

Tema 5. Aritmètica de coma flotant

Curs 2019-2020Q2 – Grup 30

Joan Manuel Parcerisa

Facultat d'Informàtica de Barcelona

Universitat Politècnica de Catalunya

Números “amb decimals”

- Com representar números no enters amb un número finit de bits? (renunciant a representar tots els racionals)
 - Coma fixa.
 - Coma flotant

Coma fixa

- **Coma fixa**

- Alguns bits per a la part entera, i alguns per a la part fraccionària
- Exemple amb 8 bits:

eeeeefff \rightarrow part entera i part fraccionària

10101,110

$$= 1 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3}$$

$$= 21,75_{10}$$

- El rang és bastant limitat
- Els números representables són equidistants

Coma flotant

- Rang molt més gran
 - però números representables no equidistants
- Notació científica normalitzada (base 10):

$$v = \pm m \times 10^e \quad \text{tal que } 1 \leq m < 10$$

$$2.34 \times 10^6 \quad \rightarrow \text{Normalitzat } (1 \leq m < 10)$$

$$0.0234 \times 10^8 \quad \rightarrow \text{No normalitzat!}$$

(la part entera de m ha de tenir 1 sol dígit, i ha de ser no-nul)

Coma flotant en base 2

$$v = \pm \underbrace{1, \overbrace{fff \dots f}^{\text{fracció (F)}}}_{\text{mantissa}} \times 2^{\underbrace{eee \dots e}_{\text{exponent (E)}}}$$

signe (S) base

- Components:

- **Signe:** 0=positiu, 1=negatiu
- **Mantissa:** conté la magnitud del número
 - Normalitzada: part entera = 1 !
 - Part entera implícita (“bit ocult”) → estalviem 1 bit
- **Exponent:** enter representat “en excés”
- **Base:** 2, implícita

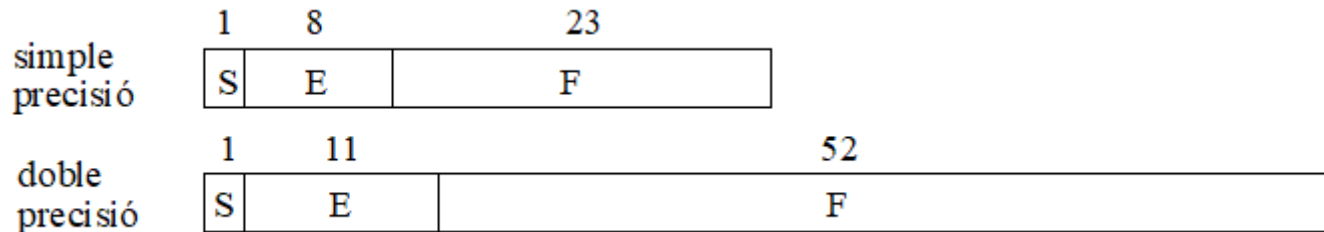
- Codificació (32 bits): signe, exponent, fracció

`v = seeeeeeeffffffffffffffffffffffffffff`

- Interpretació (S=s, E=eee...e – excés, F=0,fff...f):

$$v = (-1)^S \times (1 + F) \times 2^E$$

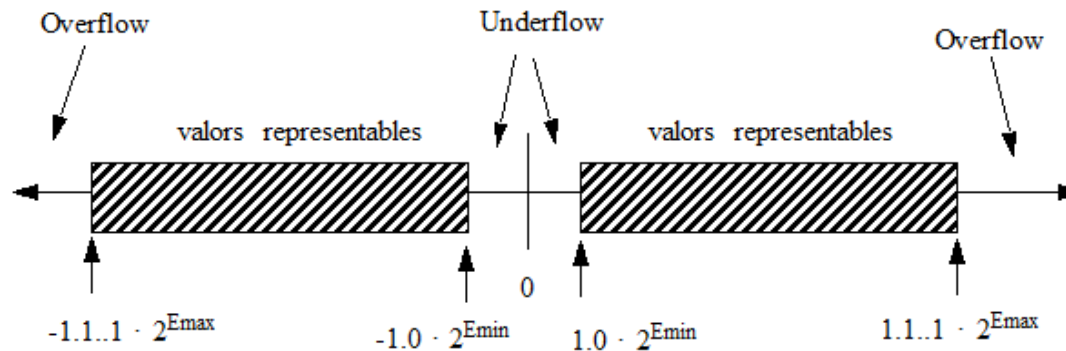
Estàndard IEEE-754 (1985-2008)



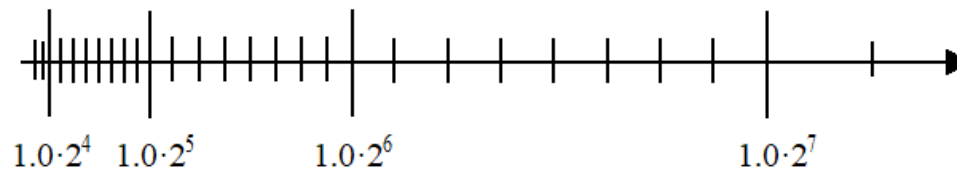
- Compromís rang vs. precisió
 - A més bits d'exponent: +rang, -precisió (i viceversa)
- Part entera = 1 (implícita, no es representa)
- Exponent
 - situat entre signe i fracció
 - en excés a 127 (simple precisió) o a 1023 (doble precisió)
 - Permet comparar magnituds amb un comparador de naturals
- En C, es representen amb els tipus `float` i `double`

IEE-754: Valors representables

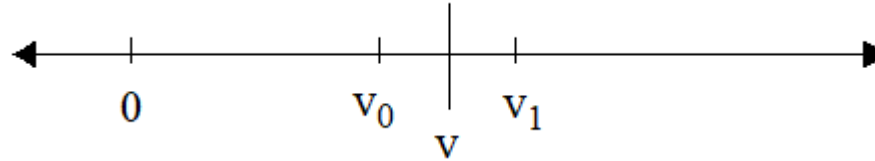
- Mantissa: $1,000\dots0 \leq \text{mantissa} \leq 1,111\dots1$
- Exponent: $E_{\min} \leq E \leq E_{\max}$
- Valors fora de rang: **Overflow**
- Valors amb magnitud $< 1.0 \cdot 2^{E_{\min}}$: **Underflow** (resultats “no fiables”)



- Amb 32 bits es poden representar 2^{32} números, igual que amb coma fixa, però no són equidistants:



Error de precisió per arrodoniment



- Resultat exacte v està entre 2 valors representables: v_0 i v_1
- Si l'arrodonim a v_0 l'error de precisió és $\varepsilon = |v - v_0|$
- Exemple: representar el número racional $\frac{1}{10}$ en simple precisió
 - És el número amb fracció periòdica:

$$v = 1, \underline{10011001100110011001100} 1100110011001100110011... \times 2^{-4}$$

23 bits de fracció

- Si l'arrodonim a v_0 (eliminant bits)

$$v_0 = 1, 10011001100110011001100 \cancel{1100110011001100110011...} \times 2^{-4}$$

- Error de precisió ($\varepsilon = |v - v_0|$)

$$\varepsilon = 0, 000000000000000000000000 \cancel{1100110011001100110011...} \times 2^{-4}$$

IEE-754: 4 modes d'arrodoniment

En l'exemple anterior...

$$v_0 = 1,10011001100110011001100 \times 2^{-4}$$

$$v = 1,10011001100110011001100110011001100110011001100110011... \times 2^{-4}$$

$$v_1 = 1,100110011001100110011001101 \times 2^{-4}$$

4.- Cap al més pròxim (o al valor “parell” en el cas equidistant)

– Mètode usat per defecte

→ En l'exemple, arrodonim a v_1 , que és més pròxim que v_0

– Màxim error: quan v és equidistant de v_0 i v_1

$$\varepsilon_{\max} = |v_1 - v_0| / 2 = 0,5 \text{ ULP}$$

– Regla pràctica: examinem el primer bit extra en 3 exemples

$$v = 1,xxxxxx...0011 \boxed{0}000101 \rightarrow v_0 = 1,xxxxxx...0011$$

$$v = 1,xxxxxx...0011 \boxed{1}000101 \rightarrow \begin{array}{r} 1,xxxxxx...0011 \\ + 0,000000...0001 \\ \hline v_1 = 1,xxxxxx...0100 \end{array}$$

$$v = 1,xxxxxx...0011 \boxed{1}000000.. \rightarrow \text{al parell} = v_0: 1,xxxxxx...0100 \boxed{0}$$

Equidistant!

IEE-754: Codificacions especials

- Exponents en valors normalitzats: $E_{\min} = 0 \dots 001$, $E_{\max} = 1 \dots 110$
- Exponents reservats: $000 \dots 0$ i $111 \dots 1$
- **\pm Zero** ($E = 000 \dots 0$, $F = 000 \dots 0$)

No es pot representar en format normalitzat...

Pega: tenim un zero positiu i un negatiu!

- **Denormals** ($E = 000 \dots 0$, $F \neq 000 \dots 0$)

Bit ocult implícit = 0 (en lloc de 1). Exponent implícit = E_{\min}

Representen números en la “regió d’*underflow*”: $0, xxx \dots x \times 2^{E_{\min}}$

- **\pm Inf** ($E = 111 \dots 1$, $F = 000 \dots 0$)

Evitar *overflows* en càlculs intermedis: $1 / (1 + 100/x) = 0$ per a $x \rightarrow 0$

Definim algunes regles: $1 \div 0 = \text{Inf}$, $1 / \text{Inf} = 0$, $x \pm \text{Inf} = \text{Inf}$, etc.

- **NaN** (*Not-a-Number*) ($E = 111 \dots 1$, $F \neq 000 \dots 0$)

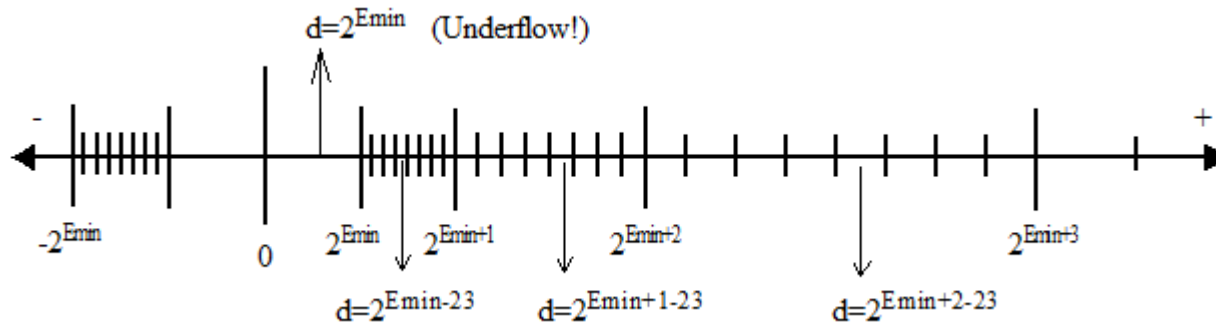
Resultat d’operacions *invàlides*: $\infty \pm \infty$, $0 \times \infty$, $\infty \div \infty$, \sqrt{x} ($x < 0$), etc.

Definim algunes regles: $x \text{ operat NaN} = \text{NaN}$

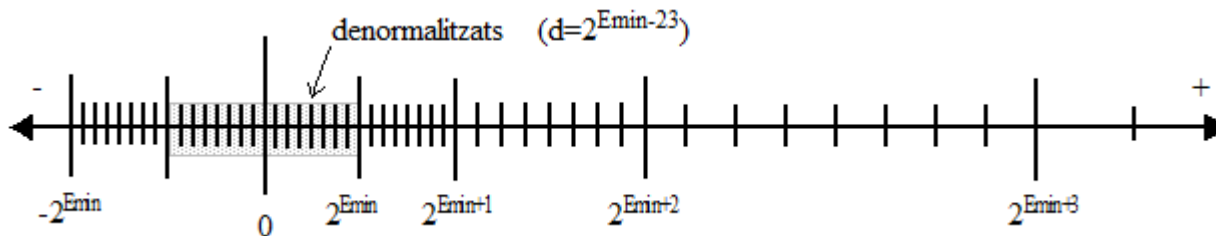
IEEE-754: resum dels formats

		Exponent (E)	
		Tot 0	Tot 1
Mantissa (F)	Tot 0	Zero	Infinit
	Altres	Denormal	NaN
		Normalitzat	

Underflow i nombres denormals



- L'error de precisió depèn de la distància (d) entre números consecutius
- No és uniforme: decreix geomètricament amb l'exponent...
- **Excepte prop del zero!**
- Prop del zero, l'error és del mateix ordre que la magnitud (l'error relatiu pot arribar a ser $= 1$). Els resultats no són fiables: Excepció d'Underflow
- Alternativa: números denormals: $0, xxxx \dots xx \times 2^{E_{min}}$



Exemple. Representar $v = -1029,68$

1. Convertir la part entera, per divisions successives

$$1029 = 10000000101 \quad (\text{té 11 bits})$$

2. Convertir la fracció, per multiplicacions successives

$$0,68 \times 2 = 1,36 \quad \rightarrow 1 \quad (\text{primer bit de fracció})$$

$$0,36 \times 2 = 0,72 \quad \rightarrow 0 \quad (\text{segon bit de fracció})$$

... etc.

- Seguim fins tenir 24 bits entre part entera i fracció: $24 - 11 = 13$ bits

- Però calen alguns bits extra, suficients per discernir si arrodonir amunt o al cas equidistant (al parell). En aquest cas ens calen 5 bits extra: **10001**...

$$0,68 = 0,101011100001010001 \dots \quad (18 \text{ bits})$$

3. Ajuntar part entera i fracció ($11 + 18 = 29$ bits)

$$1029,68 = 10000000101,101011100001010001 \dots$$

Exemple (cont.)

4. Normalitzar (moure la coma 10 llocs i ajustar l'exponent)

$$1029,68 = 1,00000001011010111000010\textcolor{red}{10001}\dots \times 2^{10}$$

5. Arrodonir (al més pròxim), examinant els “bits extra” → amunt

$$1029,68 = 1,00000001011010111000010\textcolor{red}{10001}\dots \times 2^{10}$$

$$1029,68 = 1,0000000101101011100001\mathbf{1} \times 2^{10}$$

6. Codificar l'exponent en excés a 127

$$10 + 127 = 137 = 10001001$$

7. Ajuntar Signe, Exponent i Fracció (el bit implícit 1 no s'escriu!)

$$-1029,68 = 1 \ 10001001 \ 00000001011010111000011$$

8. Expressar en hexadecimal

$$-1029,68 = \mathbf{0xC480B5C3}$$

Example (cont.)

- [illegible]

Exemple: convertir $v = 0x45814140$ a base 10

1. L'escrivim en binari
 $v = 0100\ 0101\ 1000\ 0001\ 0100\ 0001\ 0100\ 0000$
2. Agrupem els bits en camps (Signe, Exponent, Fracció)
 $v = 0\ 10001011\ 00000010100000101000000$
3. Convertim l'exponent a decimal i li restem l'excés 127
 $10001011 = 139$
 $E = 139 - 127 = 12$
4. Ajuntar signe, part entera (bit ocult), fracció i potència
 $v = +1,00000010100000101000000 \times 2^{12}$
5. Moure la coma 12 posicions a la dreta i eliminar zeros finals
 $v = +1000000101000,00101000000$
6. Convertir la part entera a decimal (suma ponderada)
 $1000000101000 = 1 \cdot 2^{12} + 1 \cdot 2^5 + 1 \cdot 2^3 = 4136$
7. Convertir la fracció a decimal (movent la coma a la dreta)
 $0,00101 = 101 \times 2^{-5} = 5/32 = 0,15625$
8. Ajuntar part entera i fracció: $v = 4136,15625$

Segona part

- Operacions amb números de coma flotant
 - Bits de guarda
 - Suma (resta)
 - Multiplicació i divisió
- Coma flotant en MIPS
 - Registres
 - Instruccions
 - Subrutines i declaracions
- Associativitat

Bits de guarda

- L'estàndard IEEE-754 especifica que *el resultat de una operació aritmètica ha de ser el mateix que el que s'obtingria si es realitzés amb una precisió absoluta i al final s'arrodonís el resultat* (al més pròxim $\rightarrow \epsilon_{\max} = 0,5 \text{ ULP}$)
- El resultat intermedi d'una operació pot tenir molts (més de 250!) bits “extra”

y = X, XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX...

- Operar amb tots ells per no perdre precisió requeriria **un hardware monstruós!**
- Però es demostra que podem obtenir el mateix resultat convertint el conjunt de bits extra en tan sols tres (anomenats *bits de guarda*): **G**uard, **R**ound i **S**ticky
- Per convertir el conjunt de bits extra a 3 bits de guarda:
 - **G** i **R** són iguals als 2 primers bits extra
 - **S** (tercer bit) és la OR lògica de la resta de bits a la dreta
- Suposem el resultat intermedi

y = X,XXXXXXXXXXXXXXXXXXXXXXXXXXXXabcdefgh...

- Fem: **G=a, R=b, S=(c OR d OR e OR f OR g OR h ...)**

$y = X, \text{XXXXXXXXXXXXXXXXXXXXXXXXXXXXGRS}$

Suma (resta) en coma flotant

- Suposem un cas senzill que coneixem bé (base 10):
 - Format: mantissa normalitzada amb 4 dígits decimals: $x,xxx \times 10^{xx}$
 - Sumar: $9,999 \times 10^1 + 1,680 \times 10^{-1}$
 - Així no:
$$\begin{array}{r} 9,999 \times 10^1 \\ + 1,680 \times 10^{-1} \\ \hline = 11,679 \times 10^? \end{array}$$
 - Alinear mantisses:
$$\begin{array}{r} 9,999 \times 10^1 \\ + 0,01680 \times 10^1 \\ \hline = 10,01580 \times 10^1 \end{array}$$
 - Normalitzar:
$$1,001580 \times 10^2$$
 - Arrodonir a 4 dígits:
$$\begin{array}{r} 1,001580 \times 10^2 \\ \hline 1,002 \times 10^2 \end{array}$$

Suma (resta) en coma flotant

1. **Igualar els exponents** al major dels dos, desplaçant la coma de la mantissa del de menor exponent
2. **Sumar les magnituds** tenint en compte els signes: si són diferents cal restar-les (en aquest cas, posar al minuend la de major valor absolut!)
3. **Normalitzar el resultat** movent la coma per obtenir un 1 a la part entera
4. **Arrodonir** la mantissa
5. **Codificar** signe, exponent (en excés) i mantissa (eliminant bit ocult)

Exemple: sumar $z=x+y$

Suposem $x=0x3F40000D$, $y=0xC0800004$

1. Els escrivim en binari

$x = 0011\ 1111\ 0100\ 0000\ 0000\ 0000\ 0000\ 1101$

$y = 1100\ 0000\ 1000\ 0000\ 0000\ 0000\ 0000\ 0100$

2. Separem els camps (S, E, F)

$x = 0\ 01111110\ 1000000000000000000000001101$

$y = 1\ 10000001\ 0000000000000000000000000100$

3. Convertim exponents a decimal, afegim bit ocult, i el signe

$x = +1,1000000000000000000000001101 \times 2^{-1}$

$y = -1,0000000000000000000000000100 \times 2^2$

4. Igualem exponents al major (=2): desplacem x tres bits a la dreta

$x = +0,00110000000000000000000001\underline{101} \times 2^2$

bits extra

Exemple (cont.)

5. Sumar magnituds: signes diferents \rightarrow restar (minuend, la major = $|y|$)

$$\begin{array}{rcl}
 |y| & = & 1,000000000000000000000000100 \times 2^2 \\
 - |x| & = & 0,001100000000000000000000001 \times 2^2 \\
 \hline
 |z| & = & 0,11010000000000000000000010 \times 2^2
 \end{array}$$

GRS
101
011

6. Normalitzar mantissa (desplaçant els bits 1 posició a l'esquerra)

$$|z| = 1,10100000000000000000000010011 \times 2^1$$

7. Arrodonir mantissa (amunt)

$$\begin{array}{rcl}
 & & 1,101000000000000000000000100 \times 2^1 \\
 + & & \frac{1}{2} \\
 \hline
 |z| & = & 1,101000000000000000000000101 \times 2^1
 \end{array}$$

8. Ajuntar signe, exponent (en excés) i mantissa (sense el bit ocult)

$$\text{exponent} = 1 + 127 = 128 = 10000000$$

$$\text{signe} = 1$$

$$\begin{aligned}
 z &= 1 \ 10000000 \ 10100000000000000000000101 \\
 &= 1100 \ 0000 \ 0101 \ 0000 \ 0000 \ 0000 \ 0000 \ 0101 \\
 &= \mathbf{0xC0500005}
 \end{aligned}$$

Multiplicació (divisió) en coma flotant

- Observem que

$$m_1 \times 2^{e_1} \times m_2 \times 2^{e_2} = (m_1 \times m_2) \times 2^{(e_1 + e_2)}$$

$$m_1 \times 2^{e_1} / m_2 \times 2^{e_2} = (m_1 / m_2) \times 2^{(e_1 - e_2)}$$

- Passos

- Multiplicar (o dividir) les mantisses (igual que amb els naturals)
- Sumar exponents (o restar-los)

Multiplicació en coma flotant

- Suposem un cas senzill que coneixem bé (base 10):
 - Format: mantissa normalitzada amb 4 dígits decimals: $x,xxx \times 10^{xx}$
 - Multiplicar: $(-1,110 \times 10^{10}) \times 9,200 \times 10^{-5}$
 - Suma dels exponents: $10 + (-5) = 5$
 - Producte de les mantisses:

					1	1	0	
					9	2	0	0
					<hr/>			
					0	0	0	0
						0	0	0
						2	2	2
						9	9	9
						<hr/>		
					1	0	2	1
 - Ajuntar-ho tot: $-10,212000 \times 10^5$
 - Normalitzar: $-1,0212000 \times 10^6$
 - Arrodonir a 4 dígits (avall): $-1,0212000 \times 10^6$

Example: multiplicar $z=x \times y$

Suposem $x=0x3F600000$, $y=0xBED00002$

1. Els escrivim en binari

```
x = 0011 1111 0110 0000 0000 0000 0000 0000
```

$$y = 1011 \ 1110 \ 1101 \ 0000 \ 0000 \ 0000 \ 0000 \ 0010$$

2. Separem els camps (S, E, F)

```
x = 0 01111110 110000000000000000000000
```

y = 1 01111101 1010000000000000000000010

3. Convertim exponents a decimal, afegim bit ocult, i el signe

$$x = +1,11000000000000000000000000 \times 2^{-1}$$
$$y = -1,10100000000000000000000010 \times 2^{-2}$$

4. Multipliquem mantisses (sense zeros finals de x , parts enteres en negreta)

$$\begin{array}{r}
 1110\dots \\
 \times 110100000000000000000010 \\
 \hline
 1110\dots \\
 1110\dots \\
 1110\dots \\
 1110\dots \\
 \hline
 1011011000000000000000011100\dots (+20 \text{ bits } 0)
 \end{array}$$

Exemple (cont.)

5. Sumem els exponents: $-1 + (-2) = -3$

6. El producte queda $10,1101100000000000000000011100 \times 2^{-3}$

7. Normalitzem: desplacem tots els bits 1 posició a la dreta
 $|z| = 1,011011000000000000000001110 \times 2^{-2}$
el Sticky és la or lògica

8. Arrodonim la mantissa (amunt)

$$\begin{array}{r} 1,011011000000000000000001110 \times 2^{-2} \\ + 1 \\ \hline |z| = 1,0110110000000000000000010 \times 2^{-2} \end{array}$$

9. Codifiquem l'exponent: $-2 + 127 = 125 = 01111101$

10. Ajuntem signe (negatiu), exponent i mantissa (sense el bit ocult!)

$$\begin{array}{rcl} z & = & 1 \ 01111101 \ 011011000000000000000010 \\ & & 1011 \ 1110 \ 1011 \ 0110 \ 0000 \ 0000 \ 0000 \ 0010 \\ & = & \mathbf{0xBEB60002} \end{array}$$

Coma flotant en el MIPS

- Coprocesador de coma flotant (o FP) CP1
 - Procesador adjunt que extén el ISA
 - 32 registres de simple precisió: $\$f0, \$f1, \dots, \$f31$
 - O bé 16 parells per a doble precisió: $\$f0-\$f1, \$f2-\$f3, \dots$
 - Les instruccions FP operen només sobre registres FP:

- Accés a memòria

<code>lwcl ft,offset(rs)</code>	<code>ldcl ft,offset(rs)</code>
<code>swcl ft,offset(rs)</code>	<code>sdcl ft,offset(rs)</code>

- Aritmètiques

<code>add.s fd,fs,ft</code>	<code>add.d fd,fs,ft</code>
<code>sub.s fd,fs,ft</code>	<code>sub.d fd,fs,ft</code>
<code>mul.s fd,fs,ft</code>	<code>mul.d fd,fs,ft</code>
<code>div.s fd,fs,ft</code>	<code>div.d fd,fs,ft</code>

Instruccions de coma flotant

- Còpia entre registres

mfc1	rt, fs	→	copia de fs a rt
mtc1	rt, fs	→	copia de rt a fs
mov.s	fd, fs	→	copia de fs a fd

- Comparació

c.xx.s fs, ft	c.xx.d fs, ft
---------------	---------------

- on xx és un de {eq, lt, le}
- escriu el resultat al *bit de condició* (registre implícit)

- Salt

bclt	etiqueta	→	salta si el bit de condició = TRUE
bclf	etiqueta	→	salta si el bit de condició = FALSE

Subrutines i declaracions

- Pas de paràmetres i resultats a subrutines
 - Reglamentació molt complexa per al pas de *floats*, *doubles* i mescles de tipus enters amb flotants
 - Sols veurem el cas amb 1 o 2 paràmetres *float*: s'usen `$f12` i `$f14`
 - El resultat es retorna en `$f0`
- Registres “segurs”: del `$f20` al `$f31`
- Declaració de variables de coma flotant en MIPS

```
.data
var1: .float 3.1416, 3.5E2, ...
var2: .double 3E350, 0.0038, ...
```

 - Alineen a adreces múltiples de 4 i de 8, respectivament

Exemple: traduir la funció func ()

<pre>float func (float x) { if (x<1.0) return x*x; else return 2.0-x; }</pre>	<pre>.data const1: .float 1.0 .text ... func: la \$t0,const1 lwcl \$f16,0(\$t0) # \$f16 = 1.0 c.lt.s \$f12,\$f16 # cond = (x<1.0) bclf else # branch if cond=false mul.s \$f0,\$f12,\$f12 # x*x b fisi else: add.s \$f16,\$f16,\$f16 # 1.0+1.0 sub.s \$f0,\$f16,\$f12 # 2.0-x fisi: jr \$ra</pre>
--	---

Associativitat

- Verifica la propietat associativa la suma de números en coma flotant?

$$x + (y + z) = (x + y) + z$$

- Suposem $x = -1,5 \times 10^{38}$, $y = 1,5 \times 10^{38}$, $z = 1,0$

$$\begin{aligned}x + (y + z) &= -1,5 \times 10^{38} + (1,5 \times 10^{38} + 1,0) \\&= -1,5 \times 10^{38} + 1,5 \times 10^{38} \\&= 0,0\end{aligned}$$

$$\begin{aligned}(x + y) + z &= (-1,5 \times 10^{38} + 1,5 \times 10^{38}) + 1,0 \\&= 0,0 + 1,0 \\&= 1,0\end{aligned}$$

NO!

Exercici: sumar $z=x+y$

Suposem $x=3DC00046$, $y=0xC0800004$

Exercici: sumar $z=x+y$

Suposem $x=3DC00046$, $y=0xC0800004$

1. Els escrivim en binari

$x = 0011\ 1101\ 1100\ 0000\ 0000\ 0000\ 0100\ 0110$

$y = 1100\ 0000\ 1000\ 0000\ 0000\ 0000\ 0000\ 0100$

Exercici: sumar $z=x+y$

Suposem $x=3DC00046$, $y=0xC0800004$

1. Els escrivim en binari

$x = 0011\ 1101\ 1100\ 0000\ 0000\ 0000\ 0100\ 0110$

$y = 1100\ 0000\ 1000\ 0000\ 0000\ 0000\ 0000\ 0100$

2. Separem els camps (S, E, F)

$x = 0\ 01111011\ 100000000000000001000110$

$y = 1\ 10000001\ 00000000000000000000000100$

Exercici: sumar $z=x+y$

Suposem $x=3DC00046$, $y=0xC0800004$

1. Els escrivim en binari

$x = 0011\ 1101\ 1100\ 0000\ 0000\ 0000\ 0100\ 0110$

$y = 1100\ 0000\ 1000\ 0000\ 0000\ 0000\ 0000\ 0100$

2. Separem els camps (S, E, F)

$x = 0\ 01111011\ 100000000000000001000110$

$y = 1\ 10000001\ 00000000000000000000000100$

3. Convertim exponents a decimal, afegim bit ocult, i el signe

$x = +1,10000000000000000001000110 \times 2^{-4}$

$y = -1,00000000000000000000000100 \times 2^2$

Exercici: sumar $z=x+y$

Suposem $x=3DC00046$, $y=0xC0800004$

1. Els escrivim en binari

$x = 0011\ 1101\ 1100\ 0000\ 0000\ 0000\ 0100\ 0110$

$y = 1100\ 0000\ 1000\ 0000\ 0000\ 0000\ 0000\ 0100$

2. Separem els camps (S, E, F)

$x = 0\ 01111011\ 1000000000000000001000110$

$y = 1\ 10000001\ 00000000000000000000000100$

3. Convertim exponents a decimal, afegim bit ocult, i el signe

$x = +1,1000000000000000001000110 \times 2^{-4}$

$y = -1,00000000000000000000000100 \times 2^2$

4. Igualem exponents al major (=2): desplacem x sis bits a la dreta

$x = +0,000001100000000000000001$ 000110 $\times 2^2$

bits extra

Exercici: sumar $z=x+y$

Suposem $x=3DC00046$, $y=0xC0800004$

1. Els escrivim en binari

$x = 0011\ 1101\ 1100\ 0000\ 0000\ 0000\ 0100\ 0110$
 $y = 1100\ 0000\ 1000\ 0000\ 0000\ 0000\ 0000\ 0100$

2. Separem els camps (S, E, F)

$x = 0\ 01111011\ 100000000000000001000110$
 $y = 1\ 10000001\ 00000000000000000000000100$

3. Convertim exponents a decimal, afegim bit ocult, i el signe

$x = +1,10000000000000000001000110 \times 2^{-4}$
 $y = -1,00000000000000000000000100 \times 2^2$

4. Igualem exponents al major (=2): desplacem x sis bits a la dreta

$x = +0,0000011000000000000000001\underline{000110} \times 2^2$
bits extra

5. Determinem els 3 bits de guarda GRS

$x = +0,0000011000000000000000001001 \times 2^2$

Exemple (cont.)

6. Sumar magnituds: signes diferents \rightarrow restar (minuend, la major = $|y|$)

$ y $	=	1,000000000000000000000000100	GRS	$\times 2^2$
- $ x $	=	<u>0,0000011000000000000000001</u>	001	$\times 2^2$
$ z $	=	0,11111010000000000000000010	111	$\times 2^2$

Exemple (cont.)

6. Sumar magnituds: signes diferents \rightarrow restar (minuend, la major = $|y|$)

$$\begin{array}{rcl}
 |y| & = & 1,000000000000000000000000100 \boxed{\text{GRS}} \times 2^2 \\
 - |x| & = & 0,0000011000000000000000001 \boxed{001} \times 2^2 \\
 \hline
 |z| & = & 0,11111010000000000000000010 \boxed{111} \times 2^2
 \end{array}$$

7. Normalitzar mantissa (desplaçant els bits 1 posició a l'esquerra)

$$|z| = 1,1111010000000000000000101 \boxed{11} \times 2^1$$

Exemple (cont.)

6. Sumar magnituds: signes diferents \rightarrow restar (minuend, la major = $|y|$)

$$\begin{array}{rcl}
 |y| & = & 1,00000000000000000000000100 \boxed{\text{GRS}} \times 2^2 \\
 - |x| & = & 0,0000011000000000000000001 \boxed{001} \times 2^2 \\
 \hline
 |z| & = & 0,1111101000000000000000010 \boxed{111} \times 2^2
 \end{array}$$

7. Normalitzar mantissa (desplaçant els bits 1 posició a l'esquerra)

$$|z| = 1,111101000000000000000101 \boxed{11} \times 2^1$$

8. Arrodonir mantissa (amunt)

$$\begin{array}{rcl}
 & & 1,111101000000000000000101 \times 2^1 \\
 + & & \frac{1}{2} \\
 \hline
 |z| & = & 1,1111010000000000000001 \boxed{10} \times 2^1
 \end{array}$$

Exemple (cont.)

6. Sumar magnituds: signes diferents \rightarrow restar (minuend, la major = $|y|$)

$$\begin{array}{rcl}
 |y| & = & 1,00000000000000000000000100 \text{ GRS} \times 2^2 \\
 - |x| & = & 0,0000011000000000000000001 \text{ 001} \times 2^2 \\
 \hline
 |z| & = & 0,1111101000000000000000010 \text{ 111} \times 2^2
 \end{array}$$

7. Normalitzar mantissa (desplaçant els bits 1 posició a l'esquerra)

$$|z| = 1,111101000000000000000101 \text{ 11} \times 2^1$$

8. Arrodonir mantissa (amunt)

$$\begin{array}{rcl}
 & & 1,111101000000000000000101 \times 2^1 \\
 + & & 1 \\
 \hline
 |z| & = & 1,111101000000000000000110 \times 2^1
 \end{array}$$

9. Ajuntar signe, exponent (en excés) i mantissa (sense el bit ocult)

$$\text{exponent} = 1 + 127 = 128 = 10000000$$

$$\text{signe} = 1$$

$$\begin{aligned}
 z &= 1 \ 10000000 \ 111101000000000000000110 \\
 &= 1100 \ 0000 \ 0111 \ 1010 \ 0000 \ 0000 \ 0000 \ 0110 \\
 &= \mathbf{0xC07A0006}
 \end{aligned}$$