

	Signed bit	Exponent	Fraction	Decimal
i)	0	0100	1.0100	
ii)	1	1100	1101	
iii)	0	1111	1100	6.625
iv)	0	1110	0001	
v)	1	1011	1100	
vi)	0	0011	1111	

Table 7: Floating point representation

$$M = 1,2S$$

$$E = 4 - (8 - 1) = -3$$

$$-1,2S \cdot 2^{-3}$$

$$M = 1 + \frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^4} = \frac{S+8+16}{16}$$

$$E = 8 + 4 - 7 = S = \frac{29}{16}$$

$$- \frac{29}{16} \cdot 2^S$$

$$= -58$$

$$n = 1,7s$$

$$E = 1s - 7 = 8$$

$$1,7s \cdot 2^8$$

$$n = \frac{1}{16}$$

$$E = 14 - 7 = 7.$$

$$\frac{1}{16} \cdot 2^7 = \underline{\underline{136}}$$

$$n = \frac{3}{2^2} + \frac{4}{5} = \frac{7}{2^2}$$

$$E = 11 - 7 = 4.$$

$$-\frac{7}{2^2} \cdot 2^4 = -7 \cdot 4 = -28$$

$$\Omega =$$

$$E = 3 - 7 = -4$$

$$\frac{16}{16} + \frac{8}{16} + \frac{4}{16} + \frac{2}{16} \\ + \frac{1}{16}$$

$$= \frac{31}{16}$$

$$\frac{31}{16} \cdot 2^{-4} = \frac{31}{256}.$$

i

ii

$$2^n - 1$$

$\Omega + 8g\mu$

iii

$$\frac{1}{2^n}$$

iv

$$-2^n + 1$$

$$-\frac{1}{2^n}$$

(c)

2.S

$$0.S \cdot 2 = 1.0$$

$$\begin{aligned} \Pi &= (10.1)_2 = (2.S)_{10} \\ &= (.01) \cdot 2^{-1} \end{aligned}$$

$$\begin{aligned} 1 &= E - B \\ \Rightarrow E &= 1 + B = 128 \end{aligned}$$

$$x = (0 \quad 10000000 \quad 010\overbrace{0000000000000000})$$

ii

$$3.75$$

$$= (11.11)_2$$

$$(.11)_2 \cdot 2^1$$

$$E = 128$$

$$(0 \quad 100000000 \quad 1100000000 \dots)$$

iii

$$(100.1)_2$$

$$= (.001)_2 \cdot 2^2$$

$$2 = E - B$$

$$\Rightarrow E = 2 + B \\ = 129$$

$$(1 \quad 10000001 \quad 001000 \dots)$$

iv $(101)_2 = (.01) \cdot 2^1$

$$\Rightarrow (1 \ 10000000 \ 010\dots)$$

v $6.25 = (110.01)_2$

$$= (.1001)_2 \cdot 2^2$$

$$\Rightarrow (0 \ 10000001 \ 100100\dots)$$

d $0.75 = (0.11)_2$

$$(1.1)_2 \cdot 2^{-1}$$

$$E = 126$$

$$X = (1 \ 0111110 \ 10000\dots)$$

$$-1 = E - B$$

$$= E - 1023$$

$$\Leftrightarrow E = 1022$$

$$x = (0.0111111\ldots 0\ldots 100000\ldots)$$

Exercise 7

a) $8 + 2 = 10$
 $\frac{3}{4} + \frac{1}{16} = \underline{\underline{\frac{13}{16}}}$

b) $10 + \frac{13}{16} - 128$

$= -117.1875$

c) $41.0625 - 128 = -86.375$

d) $-128 + 4 + 16 = -108$ 3125
 $+ \frac{1}{2} + \frac{1}{8} + \frac{1}{16}$

$= -108 + 0,5 + 0,125 + 0,0625$

$= -107.6875$

(e)

$$1 + 4 + 16 + 64$$
$$1/2 + 1/16$$

(f)

$$2 + 8 + 16 + 64$$
$$1/2 + 1/4 + 1/8 + 1/16$$

b i $2^{n_1-1} - 1 + \left(1 - \frac{1}{2^{n_2}}\right)$

↓
req

ii $\frac{1}{2^{n_2}}$

iii -2^{n_1-1}

iv $-\frac{1}{2^{n_2}}$

c i 1.S

0001.1000

ii 0101.0100

iii -0.125

87S

500
250
-125

1111.1110

iv - 3.0

0011.0000

→ 11011111

✓ -6.75

"rajouter 1 à la
dernière pos"

0110.1100

0010

- valeur
8 oct.

(1001). (0100)⁺¹

-?

0110.1100

1001.0011f

1

1001.0100

2 - 7

X

1001 0011

0,001

1, 110

1, 111

0 1 2 3

(-4) -3 -2 -1

110.00

110.00

a

i

$$(.1010010)_2 \cdot 2^3$$

$$3 = E - B$$

$$= 127$$

$$\Leftrightarrow 130 = E$$

(0 10000010 1010010...)

ii

$$(-0101101)_2 \cdot 2^3$$

$$E = 127 + 4 = 131$$

(1 10000010 0101101...)

iii

$$(-111111)_2 \cdot 2^3$$

$$E = 130$$

$$(0 \text{ } 10000010 \text{ } 11111100\dots)$$

iv

$$(.0011001)_2 \cdot 2^4$$

$$E = 131$$

$$(1 \text{ } 10000010 \text{ } 001100100\dots)$$

Exercise 8

a) I

from $(2^{24}-1) \cdot 2^{2^8} - 1 - 127$
 ~~$\rightarrow -(2^{24}-1) \cdot 2^{2^8} - 1 - 127$~~

range = $2^{24}-1$ | $2^8-1-127$
 \nearrow \nearrow
 $M \cdot 2^E \cdot (-1)^S \nearrow 1$

$$\frac{1}{2^{24}-1} \cdot 2^{8-1-127}$$

Max value for magnitude:

1. 111 1111 11111 $a_0 \cdot \left(\frac{1-q^{n+1}}{1-q}\right)$

$$\downarrow \\ 1 + \frac{1}{2} + \frac{1}{4} + \dots$$

$$= 1 \cdot \left(\frac{1-q^{24}}{1-q}\right)$$

$$= \frac{1 - \frac{1}{2^{24}}}{\frac{1}{2}}$$

-

$2 - 2^{-23}$

$$\left(2 - \frac{1}{2^{23}}\right) \cdot 2^{\epsilon}$$

$$= \left(2 - \frac{1}{2^{23}}\right) \cdot 2^{2^r}$$

$$= 2^{120} - \frac{2^{128}}{2^{23}}$$

$$= 2^{129} - 2^{105}$$

$$= 2^{105} (2^{24} - 1)$$

S and N is symmetric

$$\pm (0 \ 1 \ 2 \ 3 \dots)$$

(2)

$$\left(1 + \frac{1}{2} \frac{1 - (1/2)^{82}}{1 - (-1/2)} \right) \cdot 2^{2-1 \underbrace{-1023}_{1024}}$$

$$= \left(1 + 1 - \frac{1}{2^{82}} \right) \cdot 2^{2^{10}}$$

$$= \left(2 - \frac{1}{2^{82}} \right) \cdot 2^{1024}$$

$$= 2^{1024} \left(2 - \frac{1}{2^{82}} \right)$$

$$= 2^{1025} - 2^{972}$$

$$= 2^{972} (2^{52} - 1)$$

ii

Resolution : smallest diff.

(biased) min value

0 00000000

(1.0000000000

$$127 = 2^7 - 1$$

↓

$$E_{\min} = +1 - 2^{-127}$$

\uparrow \uparrow
bias, 127 \min
 \min is 0

$$\text{dor } E_{\min} = -2^7 + 1 = -128 + 1$$

$$E = 1 \cdot 2^{-(1 - 2^7)} = 2^{-127}$$

b) i) 2.70

500
+ 250

500
+ 125
+ 0625

→ (10.1011) .0XXX

excesso

$$= (.0110) \cdot 2^1$$

$$1 = E - 1$$

$$\Rightarrow E = 2$$

$$(0 \text{ } 10 \text{ } 0110) = x_1$$

ii) 3.9

500
+ 250] 875
+ 125] 335...
+ 0625

$$\Rightarrow (011.111)$$

$$1 = E - 1$$

$$\Rightarrow E = 2$$

$$(0 \text{ } 10 \text{ } 111)$$

iii

-0.2

on enlève 0.1875

$$\begin{aligned}
 & 0.125 \\
 & + 0.0625 \\
 & = 0.1875 \pm 0.005 \\
 & \underline{\underline{0.250 \pm 0.00}}
 \end{aligned}$$

$$\begin{aligned}
 & 0.0011 \\
 & (\cdot 1) \cdot 2^{-3} \\
 & -3 = E-1 \\
 & \Rightarrow E = -2
 \end{aligned}$$

→ précision trop petite

$$\begin{aligned}
 & 0.01 \quad (\cdot 0) \cdot 2^{-2} \\
 & \text{On enlève } 0.25 \\
 & \qquad \qquad \qquad \xrightarrow{\text{pas possible}}
 \end{aligned}$$

$$\begin{aligned}
 & (0.1)_2 \quad \cdot 0 \cdot 2^{-1} \quad -1 = E-1 \\
 & (\cancel{1} \quad 00 \quad 0000) \qquad \qquad \qquad \Rightarrow E=0
 \end{aligned}$$

iv

1.33

$$\begin{array}{r}
 250 \\
 + 125 \\
 \hline
 375
 \end{array}$$

$$(1.0110)_2$$

$$(\cdot 011) \cdot 2^0 \quad 0 = E-1$$

$$(\cancel{0} \quad 01 \quad 0110) \Rightarrow E=1$$

(V) -0.67

(0.1011)

$$\begin{array}{r} 0.500 \\ + 0.125 \\ \hline + 0.625 \end{array}$$

$$(-0.11) \cdot 2^{-1} - 1 = E - 1$$
$$\Rightarrow E = 0$$

(1 00 0110)

c) $(7.2)_{10}$

$$\begin{array}{r} 0.125 \\ 0.0625 \end{array}] 1875$$

$(0111.0011)_2$

ii) $8 - 6.42$
 $- 6.42 = -8 + 0.08 + 0.5 + 1$

$(1001.1001)_2$

$3125 [+ 250$
 $+ 0625$

$375 [+ 250$
 $+ 125$

iii) $8 - 3.67$

$$-3.67 = -8 + 4 + 0.33$$

$(1100.0101)_2$

iv) 5.33

$(0101.0101)_2$

$E \leftarrow 10$

A	B	SP(A)	SP(B)	A + B	SP(A) + SP(B)
0.125	0.25				
-0.375	0.5				
0.625	0.75				
-0.875	-1.0				
1.125	-1.25				

Table 12: Addition in floating point representation

0.125

(0.001)

(.0) $\cdot 2^{-3}$

$$\begin{aligned} -3 &= E - B \\ \Leftrightarrow E &= -3 + 127 \\ &= 124 \end{aligned}$$

→

$x = (0 \ 01111100 \ 0000\dots)$

0.25

$$\begin{pmatrix} 0.01 \\ -0 \end{pmatrix} \cdot 2^{-2}$$

$$+2 E-B$$

$$\hookrightarrow E = 2 + 12s$$

$$Y = (0 | 01111101 \text{ } \underbrace{\dots}_{=12s})$$

0.375

$$\begin{array}{r} 0 \quad 01111101 \quad \text{00000000} \\ + 0 \quad 01111101 \quad \text{00000000} \end{array}$$

$$- 0 \quad 01111101 \quad 10000000\dots$$

-0.375

$$\begin{array}{r} 250 \\ + 125 \\ \hline \end{array}$$

$$(0.011)_2 = (-1) \cdot 2^{-2}$$

$$\begin{aligned} -2 &= E - B \\ \Leftrightarrow E &= -2 + 127 \\ &= 125 \end{aligned}$$

$$Z = (1 \ 0111101 \ 100000\dots) \quad \checkmark$$

G.S

$$(0.1)_2 = (0)_2 \cdot 2^{-1}$$

$$E = 126$$

$$A = (0 \ 01111110 \ 00000\dots) \quad \checkmark$$

↑ + ground

hidden bit
Shift

$$\begin{array}{r}
 (0\ 01111110\ 1.0)000000\dots \\
 - (1\ 01111101\ 0.1)1000000\dots \\
 \hline
 0.01000000\dots
 \end{array}$$

*

↑ on shift de 2, l'exposant passe de -1 à -3

donc on passe de 126 à 124

$$\Rightarrow 1\ 01111100\ 000000\dots$$

* $\begin{array}{r}
 11 \\
 -100 \\
 \hline
 011
 \end{array}$ on emporte 1 $\Rightarrow 10 - 01 = 01$

0.625

$$(0.101)_2 \\ (\cdot 01)_2 \cdot 2^{-1}$$

$$-1 = E - B$$

$$\Leftrightarrow E = 126$$

$$B = (0 \ 0111110 \ 0100000\dots)$$

0.75

$$(0.11)_2 \\ (\cdot 1)_2 \cdot 2^{-1}$$

$$C = (0 \ 0111110 \ 1000000\dots)$$

(0 0111110

1. 0 1000000

0 011110

1. 1 0 000000

1 0. 1 1

$$\Rightarrow (0 0111111 \quad 011000000)$$

on passe de 126 à 127

-0.875

$$\begin{array}{r}
 500 \\
 + 250 \\
 \hline
 750
 \end{array}$$

$$\begin{aligned}
 (0.111)_2 & \\
 (-.11)_2 \cdot 2^{-1} &
 \end{aligned}$$

$D = (10111110 \quad 11000\dots)$

-1.0

$$\begin{aligned}
 (1.0)_2 & \\
 (-.0)_2 \cdot 2^0 &
 \end{aligned}$$

$$\begin{aligned}
 O &= E - 127 \\
 \Rightarrow E &= 127
 \end{aligned}$$

$E = (10111111 \quad 000000\dots)$

$$\begin{aligned}
 & (10111111 \quad 0.11000\dots) \\
 + & (10111111 \quad \text{red } 1.000000\dots)
 \end{aligned}$$

$$= \begin{pmatrix} 1 & 0111111 \\ & 1.1110000\ldots \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0111111 \\ 1 & 1110000\ldots \end{pmatrix}$$

1.125

$$(1.001)_2$$
$$(.001) \cdot 2^0$$

$$E = 127$$

$$F = (0\ 0111111\ 001000\dots)$$

-1.25

$$\begin{array}{r} 1 \\ 0 \\ -1 \\ \hline \end{array}$$

$$(1.01)_2$$

$$(.01) \cdot 2^0$$

$$G = (-1\ 0111111\ 010000\dots)$$

$$\begin{array}{r}
 (1 \ 0111111 \quad 1,0100000 \dots) \\
 - (0 \ 0111111 \quad 1,0010000 \dots) \\
 \hline
 (1 \ 01111100 \quad 001.00000 \dots)
 \end{array}$$

3 shift



Note, si on avait fait l'inverse, on aurait eu une retenue infinie et on aurait dû shifté vers la gauche jusqu'à avoir le premier 1 devant.

6

$$\begin{array}{r}
 625 \\
 -8+6 \\
 \hline
 -8+2+1
 \end{array}$$

0

$$\begin{array}{r}
 -1 \\
 \hline
 1
 \end{array}$$

A	B	FiP(A)	FiP(B)	A + B	FiP(A) + FiP(B)
1.125	3.25	0001.0010	0111.0100	4.375	
-4.375	6.5	1011.1010	0110.1000	2.125	
2.625	7.75	0010.1010	0111.1100	10.375	
-0.875	-1.125	1111.0010	1110.1110		
6.25	-3.875	0110.0100	1100.0010	2.375	

Table 13: Addition in fixed point representation

$$\begin{array}{r}
 11 \\
 0001.0010 \\
 + 0011.0100 \\
 \hline
 0100.0110
 \end{array}$$

$$\begin{array}{r}
 2.5 - 0.375 \\
 - 2.125
 \end{array}$$

$$\begin{array}{r}
 1111 \\
 1011.1010 \\
 + 0110.1000 \\
 \hline
 0010.0010
 \end{array}$$

$$\begin{array}{r}
 111111 \\
 1111.0010 \\
 + 1110.1110 \\
 \hline
 1110,0000
 \end{array}$$

$$\begin{array}{r}
 1111 \\
 0010.1010 \\
 + 0111.1100 \\
 \hline
 1010.0110
 \end{array}$$

$$\begin{array}{r}
 1 \\
 0110.0100 \\
 + 1100.0010 \\
 \hline
 0010.0110
 \end{array}$$

Exercise 3

16 64
↓ ↓
32 ↓

1110.1000001

(.1101000001) · 2³

B = 1S

$$3 = E - 1S$$

$$\hookrightarrow E = 18$$

10010

x = (0 10010 101000001)

$$-21.75 = -32 + 8 + 2 + 0.25$$

¹⁶
 (10101.1100)

$$(0.010111) \cdot 2^4$$

$(1 \text{ } 10011 \text{ } 01011100\dots)$

$$U = E - NS$$

7.45

$$\begin{array}{r} 250 \\ + 125 \\ \hline 375 \end{array} = 375$$
$$\begin{array}{r} 0.625 \\ + 0.03125 \\ \hline 0.65625 \end{array} = 0.65625$$

$$(111.011101)$$

0.5 0.25 0.125
0.0625

=

$$(1.11011100110) \cdot 2^2$$

$$l = E - B$$
$$\Rightarrow 17 = E$$

$$(0 \ 10001 \ 110110011)$$

-2.725

$$1 = E - 25$$

$$10.10111 \cdot 2^1$$

$$A = (-1 \ 10000 \ 0101110011)$$

Max: 2^4 (19)

14.51

$$2^4 \rightarrow 2^3$$

on shift right de 3

$$(0 \ 10011 \ 101000001) \\ 0111010000)$$

$$2^4 \rightarrow 2^2$$

on shift right de
19 - 17 = 2

7.4S

(0 10011 00110110)

-2.72S

(1 10011 000101011)