

* Floating Point:

↳ Signed magnitude. X

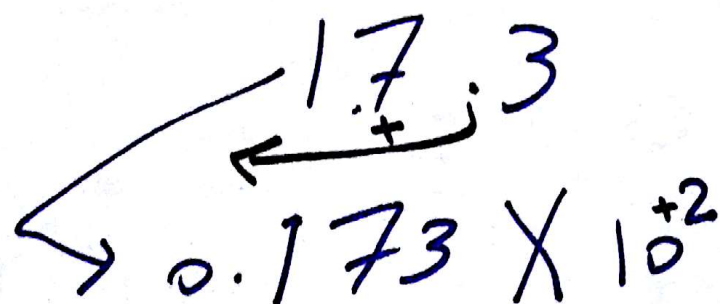
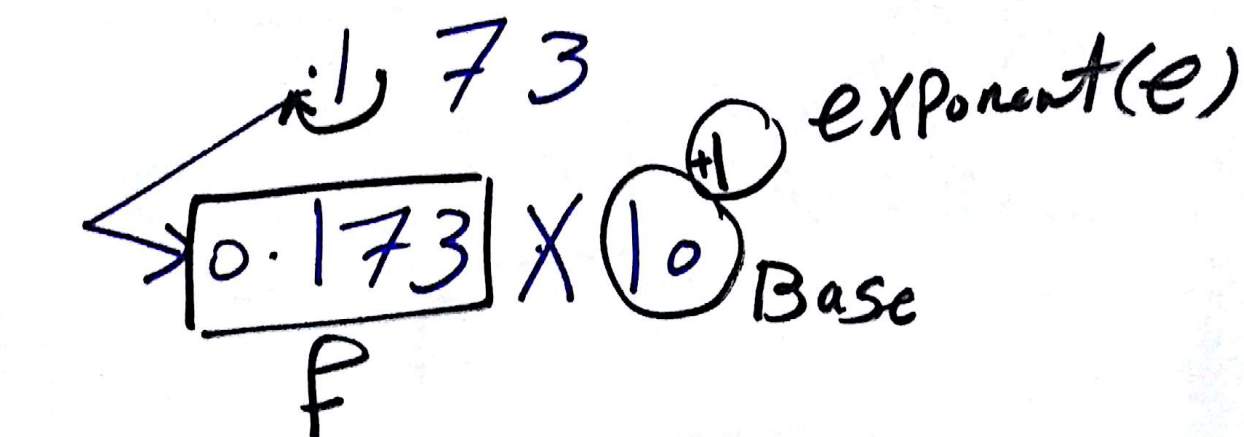
↳ 2's Complement X

↳ Excess notation. X

↳ **IEEE notation.** ✓

Code
Decode

IEEE



$$(F) \times (\text{Base})^e$$

$$0.\underline{00}173$$

$$\xrightarrow{\quad} 0.173 \times 10^{-2}$$

$$10^2$$



$$11.101$$

$$0.\underbrace{11101}_f \times 2^2 e$$

Base

$$0.\underline{0000}1101$$

$$\xrightarrow{\quad} 0.1101 \times 2^4 e$$

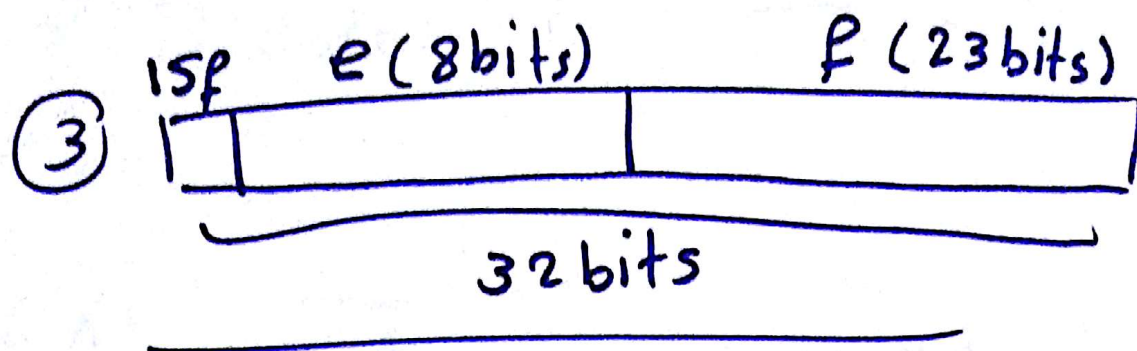
$$f \quad (f \times 2^e)$$

*IEEE

1.f x 2^e

① 1.f x 2^e

② exponent \leftrightarrow excess 127



1.1.101.

1.f x 2^e

1.1101 x 2⁺¹

* Code | ^{فرض}Represent | ^{فرض}Express

✓ $(\underline{\underline{-}})_{10} \rightarrow (\quad)_2$ in floating
Point using IEEE

✓ * Decode. $(\text{————})_2 \rightarrow (?)_{10}$

EX1
* Code $(+3.75)_{10} \rightarrow$ as binary
Pattern of floating point format
using the IEEE method.



$$\textcircled{1} \quad (+3.75)_{10}$$

↓

$$\oplus 1.1 \cdot 11$$

$$\textcircled{2} \quad 1.f \times 2^e$$

$$+ 1.\underline{111} \times 2^{\boxed{+1}}$$

f

$$\textcircled{3} \quad f = +0.111$$

$$e = (+1)_{10}$$

$$\textcircled{4} \quad eXcess\ 127$$

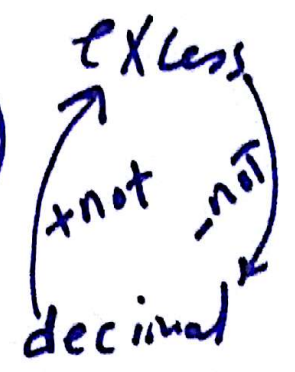
$$e = +1 + 127 = \underline{128}$$

128	64	32	16	8	4	2	1
	0	0	0	0	0	0	0

$$2 \div 133$$

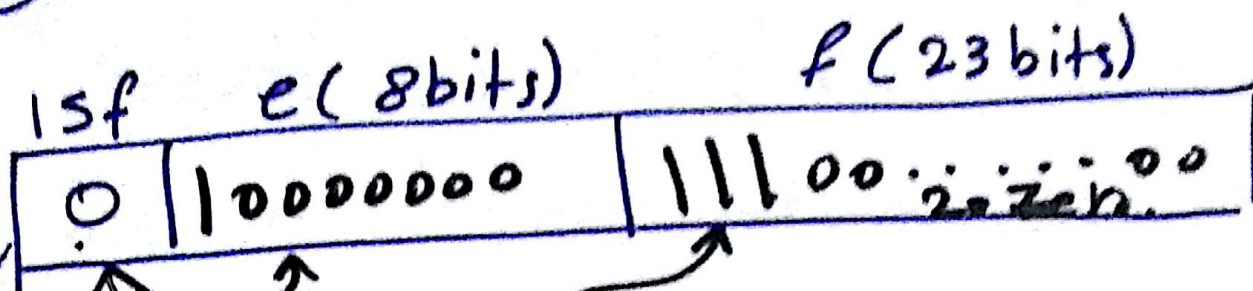
$$84 \begin{array}{r} 2+1 \\ 11 \end{array}$$

$$2^* \mid \begin{array}{r} 0.75 \\ 0.25 \boxed{1} \\ 0.5 \boxed{1} \end{array}$$



$$2 \div 128$$

⑤



$$F = (+) 0. \underline{111}$$

$$E = 10000000$$

+ \leftrightarrow 0
1 \leftrightarrow -

*The Floating Point Representation Code

$$\left(\begin{array}{c|c|c|c} 8 & 4 & 2 & 1 \\ \hline 0 & 0 & 0 & 0 \end{array} \right) \left(\begin{array}{c|c|c|c} 8 & 4 & 2 & 1 \\ \hline 0 & 0 & 0 & 0 \end{array} \right) \left(\begin{array}{c|c|c|c} 8 & 4 & 2 & 1 \\ \hline 0 & 1 & 1 & 1 \end{array} \right) \underbrace{00 \dots 00}_{2 \text{ zeros}} \right)_2$$

$$\rightarrow (4 \ 0 \ 7 \ 00000)_{16}$$



⑥

* Ex 2:

Decode $(C0A40000)_{16}$

into its equivalent decimal value.

using IEEE method.

Ans.

① $(C0A40000)_{16}$

② $1100 \ 0000 \ 1010 \ 0100 \ 0000 \dots$

A 1
B 11
C 12
D 13
E 14
F 15

②

15p	e (8b)	f (23b)
1	0000000	0100100000000000000000

12
8 4 2 1
12 → 1 1 0 0
10 → 1 0 1 0
0 1 0 0

③

f = 0.01001

e = $\begin{matrix} 128 & 64 & 32 & 16 & 8 & 4 & 2 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{matrix} - 127$

= 129 - 127

e = (+2)

exponent

(not 127)

⑦ decimal

$$④ \quad 1. \text{f} \times 2^e$$

$$-1.01001 \times 2^{+2}$$

⑤

$$- \begin{array}{c|c|c} 4 & 2 & 1 \\ \hline 1 & 0 & 1 \end{array} . \begin{array}{c|c|c} 2^{-1} & 2^{-2} & 2^{-3} \\ \hline 0 & 0 & 1 \end{array}$$

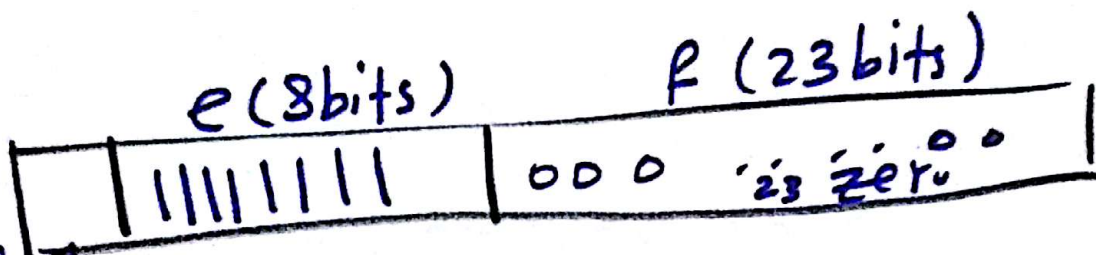


$$⑥ \quad (-5.125)$$

* Special cases in IEEE method

Decode

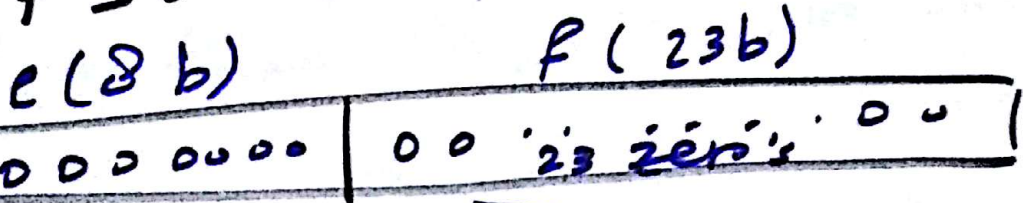
① $f=0$ and $e=1$
 23 bits \swarrow \nwarrow 8 bits



1 $\rightarrow -\infty$

0 $\rightarrow +\infty$

② $f=0$ and $e=0$



1 $\rightarrow -0$

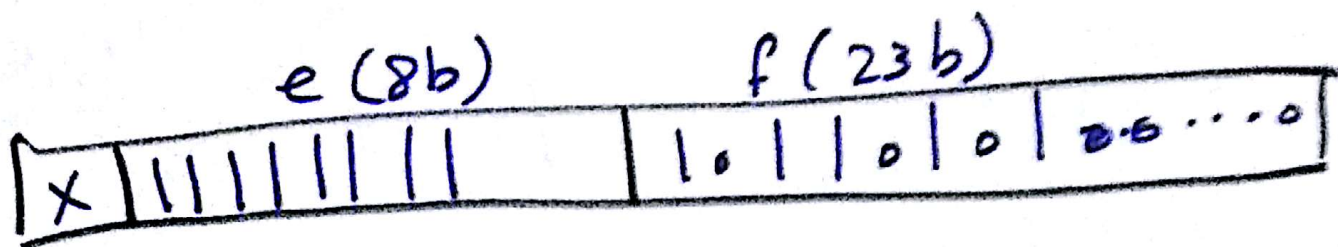
0 $\rightarrow +0$

$(80000000)_{16}$
 -0

$(00000000)_{16}$
 $+0$

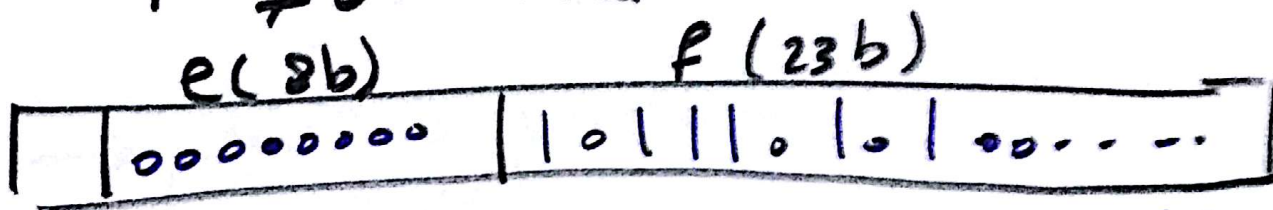
⑨

③ $f \neq 0$ and $e = 1$



→ Not a Number
(NaN)

④ $f \neq 0$ and $e = 0$



⇒ The number has a value that is less than the ^{allowed} minimum value.

$e = 1$ $\left\{ \begin{array}{l} \rightarrow f = 0 \rightarrow -\infty \text{ or } +\infty \\ \rightarrow f \neq 0 \rightarrow \text{Not a Number} \end{array} \right.$

$e = 0$ $\left\{ \begin{array}{l} \rightarrow f \neq 0 \rightarrow \text{The number} \dots \dots \dots \\ \rightarrow f = 0 \rightarrow -0 \text{ or } +0 \end{array} \right.$