

⑤ Bin to 8-3 format

(no denormal)

bias: 3 total = 8

EXP = 3

$\frac{0}{+} \frac{00010010}{e=1} \frac{0010}{m}$

$\rightarrow +1.0010 \times 2^{e-\text{bias}}$

$\rightarrow +1.0010 \times 2^{1-3}$

$\rightarrow +1.0010 \times 2^{-2}$

$\rightarrow 0.010010$
 $\quad \quad \frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{32}$

$\frac{1}{4} + \frac{1}{32} =$

③ Bin to 8-3 format
 (denormal)
 bias: 3 total = 8
 Exp = 3 M: 4

0 0000 0010
 S(+) denormal M

if(denormal)
 $\pm 0, M \times 2^{1 - \text{bias}}$

$\pm 0, 0010 \times 2^{1-3}$

$\rightarrow 0, 0010 \times 2^{-2}$

$\rightarrow 0, \underset{\frac{1}{2}}{0} \underset{\frac{1}{2}}{0} \underset{\frac{1}{2}}{0} 10 \rightarrow \frac{1}{32} \rightarrow 0, 03125$

④ Bin to float (Small)

total length = 12

Exp = 4

bias = 5 (denormal included)

M = 7

00000|00|000
+ α
→

$\neg(\alpha)$

→ $+0, M \times 2^{1-bias}$

→ $+0, 100|000 \times 2^{1-5}$

→ $+0, 100|000 \times 2^{-4}$

→ 0.0000|00|000

$$\rightarrow \frac{1}{32} + \frac{1}{256} =$$

⑤


HW1.5. Float Ordering (8-3 format)

Consider a floating point format with

- 8 total bits
- 1 bit for the sign
- 3 bits for the exponent
- a bias of 3
- 4 bits for the mantissa
- *no* denormal encodings

Given the following floating point numbers in binary format, order them from **most negative** (top) to **most positive** (bottom).

(Try and see if you can do it without calculating the numbers.)

Drag from here: 

10101100	E: 2 - bias = -1
11011101	E: 5 - bias = 2
00011101	
11001101	E: 4 - bias = 1

Construct your solution here:

11011101
 11001101
 10101100
 00011101

Sign

if (E₁ = E₂) cmp M₁, M₂

⑥ Float to Bin

$$F = 8 \quad \text{bias} = 3$$

$$E = 3$$

$$M = 4$$

denormal included

$$0.109375 = \frac{7}{64}$$

$$= \frac{1}{64} + \frac{2}{64} + \frac{4}{64}$$

$$= \frac{1}{64} + \frac{1}{32} + \frac{1}{16}$$

$$0. \underbrace{000}_{2^{-2}} \underbrace{111}_{2^{-2}}$$

$$2^{-2}$$

$$2^{x-\text{bias}} = 2^{-2}$$

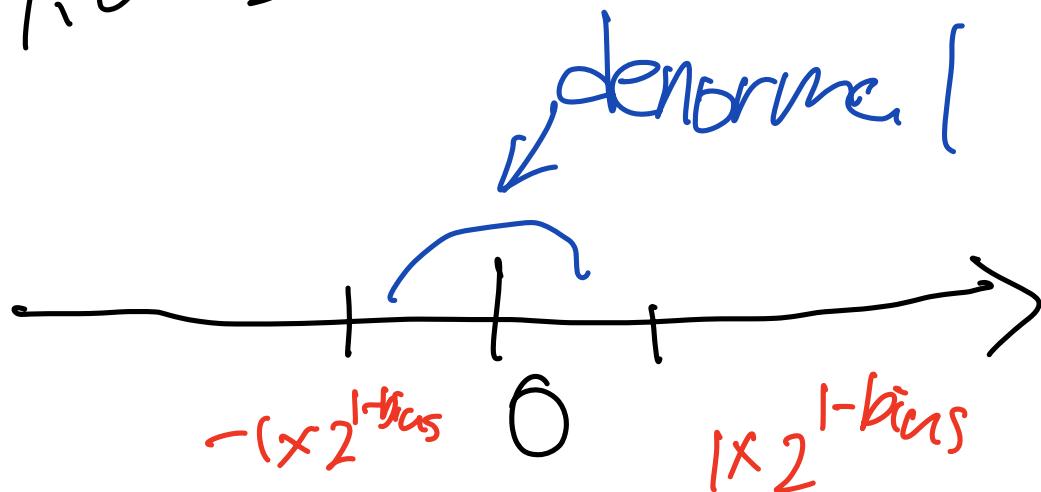
$$x = 1$$

$$0.0111 \times 2^{-2}$$

$$\begin{array}{c} \textcircled{0} \\ \hline S \end{array} \quad \begin{array}{c} 000 \\ \hline E \end{array} \quad \begin{array}{c} 0111 \\ \hline M \end{array}$$

check if denormal

$$n: \neq 0 \times 2^{1-\text{bias}}$$



⑥.2

$$-0.1875 = -\frac{3}{16}$$

$$= -\frac{1}{16} - \frac{2}{16}$$

$$= -\frac{1}{16} - \frac{1}{8}$$

denormal ✓

-0.0011

$2^{1-\text{bias}}$

$$= 2^{-2}$$



-0.11 $\times 2^{-2}$

floating format



$\overbrace{1}^{(-)}$ $\overbrace{000}^E$

$\overbrace{1100}^M$ \downarrow $\frac{1}{2} \text{ch } 0$

⑦ Floor to Bin
(no denormal)

$$T=11$$

$$E=4$$

$$\text{bias}=7$$

$$n=6$$

$$-0.34375 = -\frac{11}{32}$$

$$= -\left(\frac{5}{32} + \frac{2}{32} + \frac{8}{32}\right)$$

$$= -\left(\frac{1}{52} + \frac{1}{16} + \frac{1}{4}\right)$$

$$\rightarrow -0.01011$$

$$\rightarrow -1.011 \times 2^{-2}$$

$$E - \text{bias} = -2$$

$$\rightarrow E - 7 = -2$$

$$\rightarrow E = 5$$

$$\rightarrow 0101$$

$$\begin{array}{ccc} \frac{1}{S} & \frac{0101}{E} & \frac{011000}{M=6} \end{array}$$

⑧

Float to Bin
(include denormal)

$$T=9$$

$$E=3$$

$$\text{bias}=3$$

$$M=5$$

$$0.21875 = \frac{7}{32}$$

$$= \frac{1}{32} + \frac{2}{32} + \frac{4}{32}$$

$$= \frac{1}{32} + \frac{1}{16} + \frac{1}{8}$$

$$+ 0.000111$$

$$2^{1-bias} = 2^{-2} = \frac{1}{4}$$

$$0.2075 < \frac{1}{4}$$

it's denormal

+ 0.00011

2^{1-bias}
→ $\frac{1}{4}$

$$0.111 \times 2^{-2}$$

$$\frac{0}{S} \quad \frac{000}{E} \quad \frac{11100}{M}$$

⑨ FP convert

format 1 : $T=8$
 $E=3$
bias = 3
 $M=4$

denormal
included

→

format 2

$T=8$
 $E=3$
bias = 7
 $M=4$

denormal
included

convert 0x02

↓

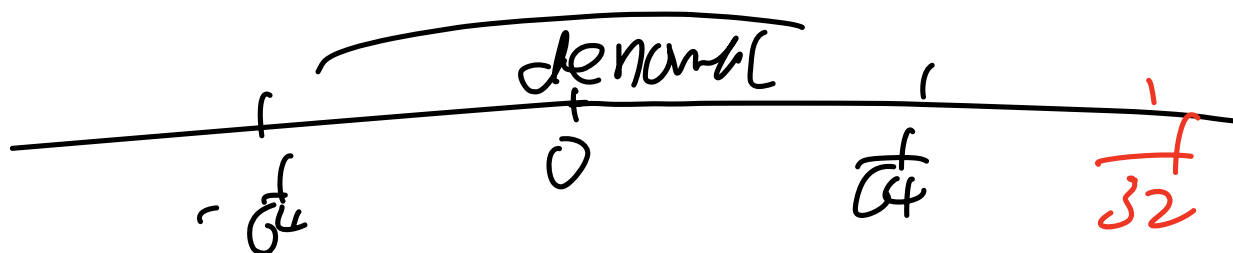
00000010

S E M

$$\begin{aligned}
 F1: & + 0.0010 \times 2^{1-3} \\
 & = 0.0010 \times 2^{-2} \\
 & = 0.00010 \\
 & = \frac{1}{32}
 \end{aligned}$$

F2: denormal 范围

$$\begin{aligned}
 2^{1-\text{bias}} &= 2^{1-7} = 2^{-6} \\
 &\pm \frac{1}{64}
 \end{aligned}$$



not denormal

F₂ (no decimal)

0.0000010

1.0 × 2⁻⁵

1.0 × 2⁻⁵

$$E - 7 = -5$$

$$E = 2 \rightarrow 010$$

$$S \rightarrow 0$$

$$\begin{array}{ccc} 0 & 010 & 00000 \\ \hline S & E & M \end{array}$$

↓ hex

0x2 0

⑨. 2

format 1:

$$T = 8$$

$$E = 3$$

$$\text{bias} = 4$$

$$n = 4$$

denormal
included

f2:

$$T = 8$$

$$E = 3$$

$$\text{bias} = 6$$

$$n = 4$$

denormal
included

$$0x05 \rightarrow \overset{f1}{\begin{array}{c|c|c} 0 & 0000 & 0101 \\ \hline \bar{s} & \bar{E} & n \end{array}}$$

Since E: 000

it's denormal $(f1)$

$$2^{1-bias} = 2^{1-4} = 2^{-3}$$

$$\rightarrow 0.M \times 2^{-3}$$

$$\rightarrow 0.0101 \times 2^{-3}$$

$$\rightarrow 0.0000101$$

1/32 1/128

$$\rightarrow 0.0390625$$

f2 denormal check

$$2^{1-f2bias} = 2^{1-6} = 2^{-5}$$

0.038...

0.03/25 0 0.03/25

not denormal

5M 1.1M x2 E-bias

1.01 x2 -5

$$E - 6 = -5$$

$$E = 1$$

$$E \rightarrow 001$$

$$S \rightarrow 0$$

$$\begin{array}{ccccccc} s & & E & & & & \text{hex} \\ \underline{0} & \underline{001} & \underline{0100} & & & & \rightarrow 14 \end{array}$$

(10) FP Convert

$$\begin{aligned} f1: \quad T &= 8 \\ E &= 3 \\ \text{bias} &= 3 \\ n &= 4 \end{aligned}$$

denormal
included

$$f2: \quad T = 10$$

$$\begin{aligned} E &= 3 \\ \text{bias} &= 3 \\ n &= 6 \end{aligned}$$

denormal
included

$$0 \times 23 \rightarrow \underbrace{0010}_{+ \quad E} \underbrace{0011}_n$$

not denormal

$$E = 2$$

E-bias

$$\rightarrow 1.00011 \times 2$$

$$\rightarrow 1.00011 \times 2^{2-3}$$

$$\rightarrow 1.00011 \times 2^{-1}$$

$$\rightarrow 0.10011$$

$\frac{1}{2} \quad \frac{1}{16} \quad \frac{1}{32}$

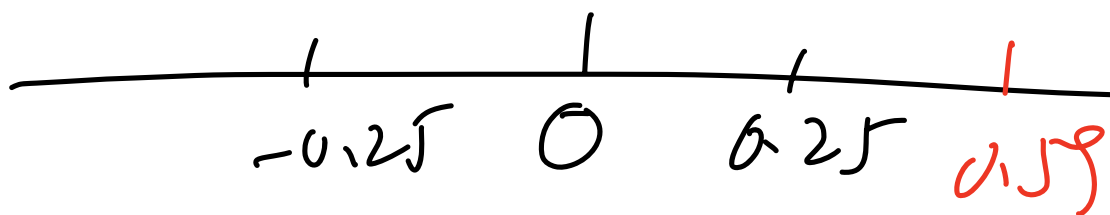
$$\rightarrow \frac{1}{2} + \frac{1}{16} + \frac{1}{32}$$

$$\rightarrow 0.59375$$

check format \rightarrow denormal

$$2^{1-f_2 \text{ bias}} = 2^{1-3}$$

$$= 2^{-2}$$



not denormal

$$0.10011$$

$$\xrightarrow{1.0 \times 2^{-1.5}}$$

$$1.0011 \times 2^{-1}$$

$$E - \text{bias} = -1$$

$$E - 3 = -1$$

$$E = 2$$

(f) $\underline{00100}$ $\underline{00100}$
 $\downarrow \text{hex}$
 08C

(10). 2

$$\begin{aligned} f_1: T &= 8 \\ E &= 3 \\ \text{bias} &= 3 \\ M &= 4 \end{aligned}$$

include
denormal

$$\begin{aligned} f_2: T &= 12 \\ E &= 4 \\ \text{bias} &= 7 \\ M &= 7 \end{aligned}$$

include
denormal

$$0x8e \rightarrow \begin{array}{c|c|c} 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ \hline S & E & & & M & & & \end{array}$$

→ denormal
 f_1

$$-0.1110 \times 2^{1-\text{bias}} \quad (\text{f}_1)$$

$$\rightarrow -0.1110 \times 2^{1-3}$$

$$\rightarrow -0.1110 \times 2^{-2}$$

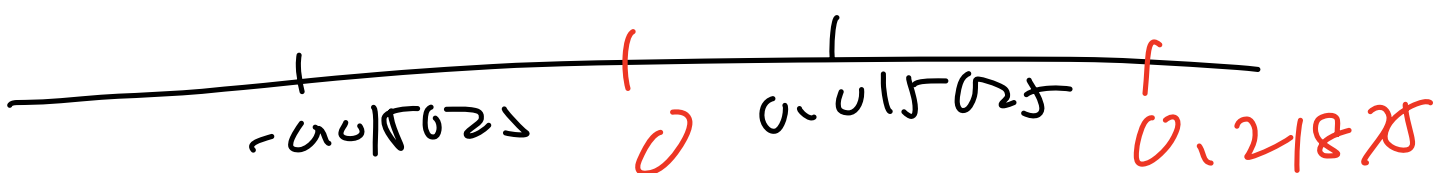
S.E.) $\rightarrow -0.001110$

$\frac{1}{8} \quad \frac{1}{16} \quad \frac{1}{32}$

$$\rightarrow 0.21875$$

check f2 denormal

$$2^{1-7} = 2^{-6} = 0.015625$$



not denormal

0.001110

$$\xrightarrow{1.M \times 2^{E-\text{bias}}} 1.\underline{110} \times 2^{-3}$$

$$E_2 - \text{bias}_2 = -3$$

$$E_2 - 7 = -3$$

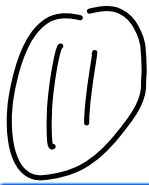
$$E_2 = 4$$

↓

(-) | 0100

11000000
↓ nex

→ A60



HW1.11. Format Comparison

Consider floating point formats Format-1 with

- 11 total bits
- 1 bit for the sign
- 4 bits for the exponent
- a bias of **6**
- 6 bits for the mantissa
- *no* denormal encodings

and Format-2 with

- 12 total bits
- 1 bit for the sign
- 4 bits for the exponent
- a bias of **4**
- 7 bits for the mantissa
- *no* denormal encodings

Which of the following formats can represent a *larger magnitude* number?

- ☐ (a) Format 1
- ☐ (b) Format 2
- ☐ (c) They can both represent the same largest magnitude number

Which of the following formats can represent a *smaller magnitude* number?

- ☐ (a) Format 1
- ☐ (b) Format 2
- ☐ (c) They can both represent the same smallest magnitude number

Suppose we changed both formats to support denormal representations. Would this change the answer to the previous (smaller magnitude) part? Pick the answer for this case.

- ☐ (a) Format 1
- ☐ (b) Format 2
- ☐ (c) They can both represent the same smallest magnitude number