

# Floating Point



Normalized

$$(-1)^s \times M = (1.\text{frac}) \times 2^{E = \text{Exp} - \text{bias}}$$

$\overset{k}{2-1}$   
 $\uparrow$

Denormalized

$$(-1)^s \times M = (0.\text{frac}) \times 2^{E = 1 - \text{bias}}$$

$\overset{k}{2-1}$   
 $\uparrow$

★ Overflow handled as too

Special

exp  $\rightarrow$  all 1s

frac  $\rightarrow 0_s \rightarrow \infty$

Non  $0_s \rightarrow \text{NaN}$

Max fractionary  $\rightarrow 1 - \frac{1}{2^N}$

Range  $\rightarrow$

$$-(2 - 2^{-N}) \times 2^{\text{bias}}$$

$$+ (2 - 2^{-N}) \times 2^{\text{bias}}$$

---

0  $\rightarrow$  0 0  $\rightarrow$

-1  $\rightarrow$  1 0 1 1 1 1 1 0 0 ~

$$362.1875$$

$$\begin{aligned} 362 &= 161 \cdot 2 + 0 \\ 161 &= 80 \cdot 2 + 1 \\ 80 &= 40 \cdot 2 + 0 \\ 40 &= 20 \cdot 2 + 0 \\ 20 &= 10 \cdot 2 + 0 \\ 10 &= 5 \cdot 2 + 0 \\ 5 &= 2 \cdot 2 + 1 \\ 2 &= 1 \cdot 2 + 0 \\ 1 &= 0 \cdot 2 + 1 \end{aligned}$$

$$101000010$$

$$\begin{aligned} 0.1875 \cdot 2 &= 0.375 \rightarrow 0 \\ 0.375 \cdot 2 &= 0.75 \rightarrow 0 \\ 0.75 \cdot 2 &= 1.5 \rightarrow 1 \\ 0.5 \cdot 2 &= 1.0 \rightarrow 1 \end{aligned}$$

$$0011$$

$$101000010.0011$$

$$1010000100011 \cdot 2^{-8}$$

$$S = 0 \quad 119 \quad 127 + (-8)$$

64 + 9' 32 13 4 8 ← exp-bias

$$001101101$$

$$0100010 \rightarrow 0$$

$$\square 117.34$$

$$\begin{aligned} 117 &= 58 \cdot 2 + 1 \\ 58 &= 29 \cdot 2 + 0 \\ 29 &= 14 \cdot 2 + 1 \\ 14 &= 7 \cdot 2 + 0 \\ 7 &= 3 \cdot 2 + 1 \\ 3 &= 1 \cdot 2 + 1 \\ 1 &= 0 \cdot 2 + 1 \end{aligned}$$

$$1110101 \quad 1$$

$$127 + 6$$

$$133$$

$$15$$

$$\textcircled{20}$$

$$0110101.0101 \text{ all } 0000 \text{ } 10100$$

$$\textcircled{8} \quad 110001111$$

$$1101010101$$

$$\begin{aligned} (-1)^S \times M(1.\text{frac}) \times 2^{E(\text{exp-bias})} \\ (-1)^S \times M(0.\text{frac}) \times 2^{E(1-\text{bias})} \end{aligned}$$

# Addition



① match exponent  
\* manage in a way to always add (not subtract) and if necessary change the sign \* or subtraction if not 0 result =

② Add

③ Normalize

④ to FP rep if asked

-230.625

1

$$\begin{aligned} 230 &= 115 \cdot 2 + 0 & 0.625 \cdot 2 &= 1 \\ 115 &= 57 \cdot 2 + 1 & 0.25 &= 0 \\ 57 &= 28 \cdot 2 + 1 & & \\ 28 &= 14 \cdot 2 + 0 & & \\ 14 &= 7 \cdot 2 + 0 & & \\ 7 &= 3 \cdot 2 + 1 & & \\ 3 &= 1 \cdot 2 + 1 & & \\ 1 &= 0 \cdot 2 + 1 & & \end{aligned}$$

1100110101

1.1100110101 · 2<sup>7</sup>

300.250

$$\begin{aligned} 300 &= 150 \cdot 2 + 0 \\ 150 &= 75 \cdot 2 + 0 \\ 75 &= 37 \cdot 2 + 1 \\ 37 &= 18 \cdot 2 + 1 \\ 18 &= 9 \cdot 2 + 0 \\ 9 &= 4 \cdot 2 + 1 \\ 4 &= 2 \cdot 2 + 0 \\ 2 &= 1 \cdot 2 + 0 \\ 1 &= 0 \cdot 2 + 1 \end{aligned}$$

1001011000.01

-1.1100110101 · 2<sup>7</sup>

-1.1100110101 · 2<sup>7</sup>  
1.0010110000 · 2<sup>8</sup>

1.0010110000 · 2<sup>8</sup>

-0.11100110101

1.100100101 · 2<sup>7</sup>

24      7+128      134  
134      128

0 |

400.625

$$\begin{aligned} 400 &= 200 \cdot 2 + 0 & 0.625 \cdot 2 &= 1.250 \text{ (1)} \\ 200 &= 100 \cdot 2 + 0 & & \\ 100 &= 50 \cdot 2 + 0 & 0.250 \cdot 2 &= 0 \\ 50 &= 25 \cdot 2 + 0 & & \\ 25 &= 12 \cdot 2 + 1 & & \\ 12 &= 6 \cdot 2 + 0 & & \\ 6 &= 3 \cdot 2 + 0 & & \\ 3 &= 1 \cdot 2 + 1 & & \\ 1 &= 0 \cdot 2 + 1 & & \end{aligned}$$

110010000.101

128  
135

1.10010000101 · 2<sup>8</sup>

200.250

128  
2  
5      7

11001000.01

+1.100100001 · 2<sup>7</sup>

0 10000111 (1) 100100000101

0 10000111 (0) 11001000010

(10) 01011000111

128 2<sup>7</sup>

0 100001101001011000111

## Multiplication

- 1, add  $E_s$
- 2, multiply significands
- 3, convert to FP rep after!

## Division

- 1, subtract  $E_s$
- 2, divide significands
- 3, convert to FP after!

---

## Denormalized FP<sub>s</sub>

are just like normalized,  
in Arithmetic point of view, except

- 1,  $E$  can't go lower than 1-bias
- 2,  $M < 1$ , if after arithmetic  $M \geq 1$ , reinterpret result as normalized!

if this violated

set  $E = 1\text{-bias}$   
and then shift  $M$   
appropriately.

★ if  $E > -\text{bias}$   
then use denormalized.

⊗ associative

int ✓  
float X

$$1.1001\ 0000\ 101 \times 2^8$$

$$1.1001\ 00001 \times 2^7 \quad 101$$

$$\begin{array}{r} 1.1001\ 0000\ 101 \\ \times 1.1001\ 0000\ 100 \\ \hline \end{array}$$

$$\begin{array}{r} 11001000010100 \\ 110010000101000000 \\ 11001000010100000000 \\ 1100100001010000000000 \end{array}$$