

The IEEE standard defines a binary representation for floating point using sign, significant, and mantissa.

Sign	Exponent	Significant
1 bit	8 bits	23 bits

For normalized floats:

$$\text{Value} = (-1)^{\text{Sign}} \times 2^{(\text{Exponent} - \text{Bias})} \times 1.\text{significant}_2$$

For denormalized floats:

$$\text{Value} = (-1)^{\text{Sign}} \times 2^{(\text{Exponent} - \text{Bias} + 1)} \times 0.\text{significant}_2$$

Exponent	Significant	Meaning
0	Anything	Denorm
1-254	Anything	Normal
255	0	Infinity
255	Nonzero	NaN

HW03 2.1

$$\text{bias} = (2^{(\text{Exponent bits} - 1)} - 1)$$

$$2^{8-1} - 1 = 2^7 - 1 = 128 - 1 = 127$$

denorm stride exp=0

$$2^{-126} \times 0.0 \dots 1 = 2^{-126} \times 2^{-23} = 2^{-149} \quad \text{- Smallest Pos denorm}$$

$$2^{-126} \times 0.0 \dots 10 = 2^{-126} \times 2^{-22} = 2^{-148} \quad \text{- Second Smallest Pos de norm}$$

$$2^{-148} - 2^{-149} = 2^{-148} (1 - 2^{-1}) = 2^{-148} \cdot 2^{-1} = 2^{-149} \quad \text{denorm stride}$$

Smallest norm stride: exp=1

$$2^{-126} \times 1.0 \dots 0 = 2^{-126} \quad \text{- Smallest norm}$$

$$2^{-126} \times 1.0 \dots 1 = 2^{-126} \times (1 + 2^{-23}) = 2^{-126} + 2^{-149} \quad \text{Smallest norm stride.}$$

Same

Each time exponent increases, stride doubles so when exp=2 stride =  $2^{-149} \cdot 2 = 2^{-148}$

Stride is usually  $2^{\text{exp} - \text{bias}} \times 2^{-\text{Significant Size}}$

$$2^{x-127} \times 2^{-23} = 1$$

$$2^{x-150}$$

How many #'s between 1-2?



$$2^{127-127} \times 1.0 \dots 0 = 2^0 = 1$$

$$2^{127-127} \times 1.1 \dots 1 = 2 - (2^0 \cdot 2^{-23}) = 2 - 2^{-23} \quad \text{Stride Size}$$

$$\text{next num } 2^{128-127} \times 1.0 \dots 0 = 2^1 = 2$$

notice we have all values representable by significant between 1-2 =  $2^{23}$  things  
Same between 2-4 there are  $2^{23}$  #'s. Equally spaced ruler between Powers of 2 with  $2^{\text{sig size}}$  #'s between.

$$1.0 \overset{11 \dots 11}{=} 2 - (2^{\frac{1}{2}} \cdot 2^{-23}) = 2 - 2^{-22}$$

4-7

$$2^2 = 2^{129-127} \times 1.000 \dots 0 = 4$$

$$2^3 = 8$$

$$2^2 \times 1.0111 \dots 111$$

$$2^2 \times 1.110 \dots 0 = 7$$

$$(4-1) \cdot 2^{21} = 2^{23} - 2^{21} =$$

$$\begin{array}{r} 2^{23} - 2^{21} \\ 10 \quad 2^{21} \\ 01 \quad 2^{21} \end{array} \quad \checkmark$$

Note:

$$\sum_{i=0}^n 2^i = 2^{n+1} - 1$$

$$\sum_{i=0}^{23} 2^i = 2^{24} - 1$$

$$2-4 \quad 2^{23} - 2^{21} + 2^{22}$$

$$\begin{array}{cc} 0 & 0 \\ 1 & 1 \end{array} \left[ \begin{array}{c} 2^1 \\ 2^2 \end{array} \right] \quad 7-9$$

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

$$\begin{array}{r}
 3-7 \qquad 2^{22} \\
 2' \times \underbrace{1.10\dots0}_{22} = 3 \\
 2' \times \underbrace{1.11\dots1} \\
 2^2 \times 1.0\dots0
 \end{array}$$

decimal  $\Rightarrow$  FP

39. 5625

$$1/2 = .5$$

$$1/4 = .25$$

$$1/8 = .125$$

$$1/16 = 0.0625$$

$$1/32 = 0.03125$$

$$1/64 = 0.015625$$

① Convert to binary

$$39 = 32 + 4 + 2 + 1 \rightarrow 100111.1001$$

$0.5625 = \frac{1}{2} + \frac{1}{16}$  

(2) Convert to FP

$$2^0 \times 100111.1001$$

$$= 2^5 \times 1.001111001 \quad 132 = 128 + 4$$

$$2^5 = 2^{132-127}$$

0b 0 10000100 0011110010...

$$= 0 \times 421 \text{ E } 4000$$

8.25

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

$$2 \times 1000.01$$

$$2^3 \times 1.00001$$

Diagram illustrating the IEEE 754 floating-point format (32-bit):

- Sign (S): 1 bit
- Exponent (Exp): 8 bits
- Significand (Sig): 23 bits

The diagram shows a 32-bit word structure with the following bit patterns:

- Sign (S): 1
- Exponent (Exp): 1000 0001
- Significand (Sig): 111 0...

$$(-1) \times 2^{129-127} \times 1.111$$

$$= -2^2 \times 1.111 = -11.11 = -7.5$$

# Example: Representing 1/3

## 1/3

$$= 0.33333..._{10}$$

$$= 0.25 + 0.0625 + 0.015625 + 0.00390625 + \dots$$

$$= 1/4 + 1/16 + 1/64 + 1/256 + \dots$$

$$= 2^{-2} + 2^{-4} + 2^{-6} + 2^{-8} + \dots$$

$$= 0.0101010101..._2 * 2^0$$

$$= 1.0101010101..._2 * 2^{-2}$$

▫ **Sign:** 0

▫ **Exponent** =  $-2 + 127 = 125 = 01111101$

▫ **Significand** = 0101010101...

0 0111 1101 0101 0101 0101 0101 0101 0101 0101

## Understanding the Significand (1/2)

### Method 1 (Fractions):

- In decimal:  $0.340_{10} \Rightarrow 340_{10}/1000_{10}$   
 $\Rightarrow 34_{10}/100_{10}$
- In binary:  $0.110_2 \Rightarrow 110_2/1000_2 = 6_{10}/8_{10}$   
 $\Rightarrow 11_2/100_2 = 3_{10}/4_{10}$
- Advantage: less purely numerical, more thought oriented; this method usually helps people understand the meaning of the significand better

## Understanding the Significand (2/2)

### Method 2 (Place Values):

- Convert from scientific notation
- In decimal:  
 $1.6732 = (1 \times 10^0) + (6 \times 10^{-1}) + (7 \times 10^{-2}) + (3 \times 10^{-3}) + (2 \times 10^{-4})$
- In binary:  
 $1.1001 = (1 \times 2^0) + (1 \times 2^{-1}) + (0 \times 2^{-2}) + (0 \times 2^{-3}) + (1 \times 2^{-4})$
- Interpretation of value in each position extends beyond the decimal/binary point
- Advantage: good for quickly calculating significand value; use this method for translating FP numbers

# Floating Point Fallacy

Not associated

- FP add associative?
  - $x = -1.5 \times 10^{38}$ ,  $y = 1.5 \times 10^{38}$ , and  $z = 1.0$
  - $x + (y + z) = -1.5 \times 10^{38} + (1.5 \times 10^{38} + 1.0)$   
 $= -1.5 \times 10^{38} + (1.5 \times 10^{38}) = 0.0$
  - $(x + y) + z = (-1.5 \times 10^{38} + 1.5 \times 10^{38}) + 1.0$   
 $= (0.0) + 1.0 = 1.0$
- Therefore, Floating Point add is not associative!
  - Why? FP result approximates real result!
  - This example:  $1.5 \times 10^{38}$  is so much larger than 1.0 that  $1.5 \times 10^{38} + 1.0$  in floating point representation is still  $1.5 \times 10^{38}$

4.75

$$1/2 = .5$$

$$1/4 = .25$$

$$1/8 = .125$$

$$1/16 = 0.0625$$

$$1/32 = 0.03125$$

$$1/64 = 0.015625$$

100.11

$$= 2^2 \times 1.0011 \quad 2^{1025-1023}$$

s 0

q

Exp: 1000000001

$$2^{1025-1023}$$

$$2^{16}$$

neg

pos

n	2 <sup>n</sup>
-1	0.5
-2	0.25
-3	0.125
-4	0.0625
-5	0.03125
-6	0.015625
-7	0.0078125
-8	0.00390625
-9	0.001953125
-10	0.0009765625
-11	0.00048828125
-12	0.000244140625
-13	0.0001220703125
-14	0.00006103515625
-15	0.000030517578125
-16	0.0000152587890625
-17	0.00000762939453125
-18	0.000003814697265625
-19	0.0000019073486328125
-20	0.00000095367431640625
-21	0.000000476837158203125
-22	0.0000002384185791015625
-23	0.00000011920928955078125
-24	0.000000059604644775390625
-25	0.0000000298023223876953125

2 <sup>0</sup> = 1
2 <sup>1</sup> = 2
2 <sup>2</sup> = 4
2 <sup>3</sup> = 8
2 <sup>4</sup> = 16
2 <sup>5</sup> = 32
2 <sup>6</sup> = 64
2 <sup>7</sup> = 128
2 <sup>8</sup> = 256
2 <sup>9</sup> = 512
2 <sup>10</sup> = 1024
2 <sup>11</sup> = 2048
2 <sup>12</sup> = 4096
2 <sup>13</sup> = 8192
2 <sup>14</sup> = 16384
2 <sup>15</sup> = 32768
2 <sup>16</sup> = 65536
2 <sup>17</sup> = 131072
2 <sup>18</sup> = 262144
2 <sup>19</sup> = 524288
2 <sup>20</sup> = 1048576
2 <sup>21</sup> = 2097152
2 <sup>22</sup> = 4194304
2 <sup>23</sup> = 8388608
2 <sup>24</sup> = 16777216
2 <sup>25</sup> = 33554432

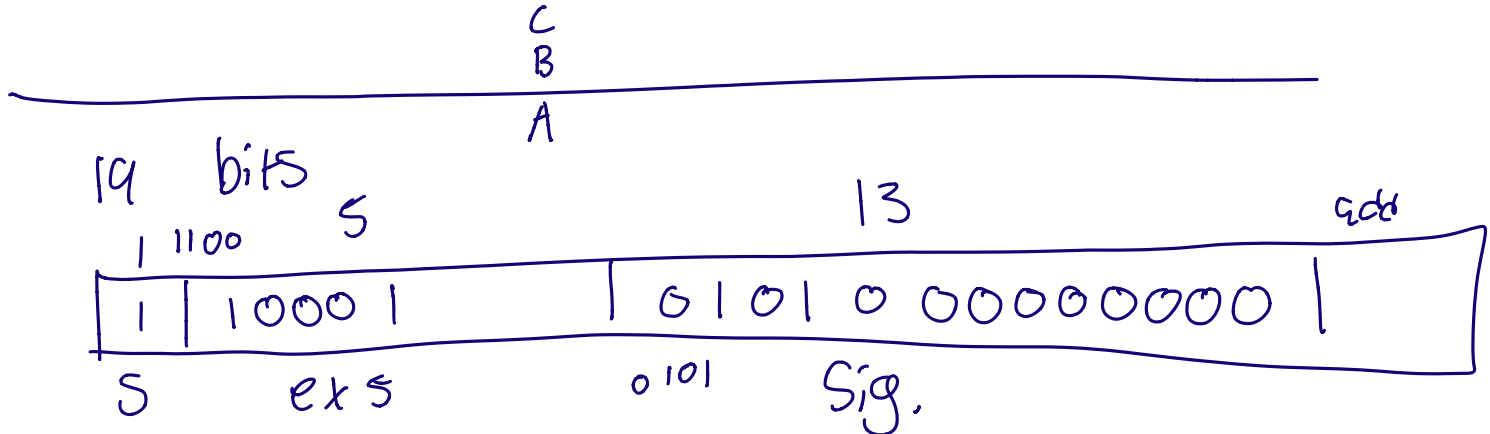
4096-15 4080

3 exp bias = -3

4 mantissa.  $2^{3-3} \times 1.0...0 = 2^0 = 1$

$2^0 \times 1.0...1 = 1(1 + 2^{-4})$

111



bias = -15

OX

-5.25

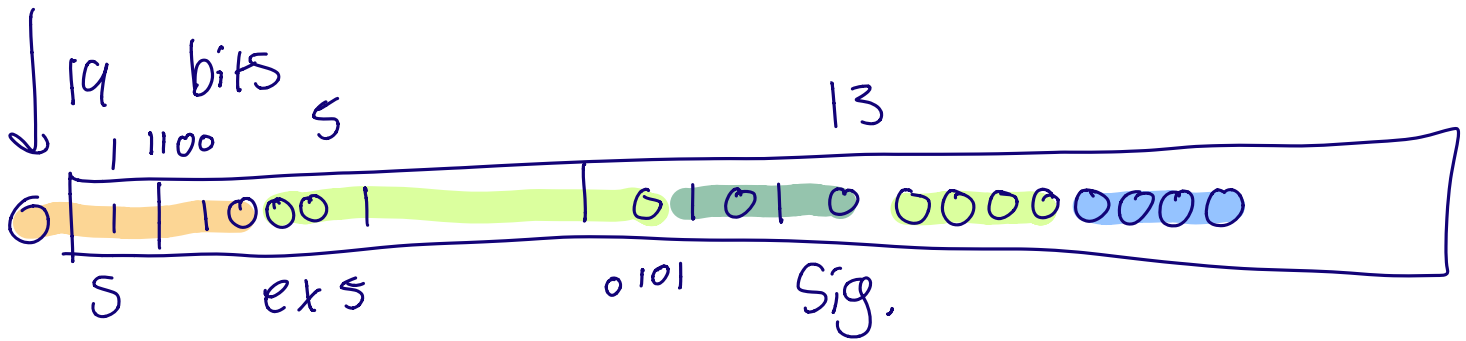
$2^0 \times 101.01$

$2^2 \times 1.0101$

0b 10001

$2^{ex-15}$   $ex = 17$

Left pad to convert



0110 0010 1010 0000 0000  
 0x 6 2 A 0 0

Part B

$$6 = 4 \cdot 1.5$$

$$2^2 \times 1.0 \dots 0 = 4$$

$$2^2 \times 1.10 \dots 0 = 4(1 + \frac{1}{2}) = 6$$

1.0 ... 1

1010

decimal  $\Rightarrow$  FP

39.5625