

# Floating-Point Numbers (the sequel)

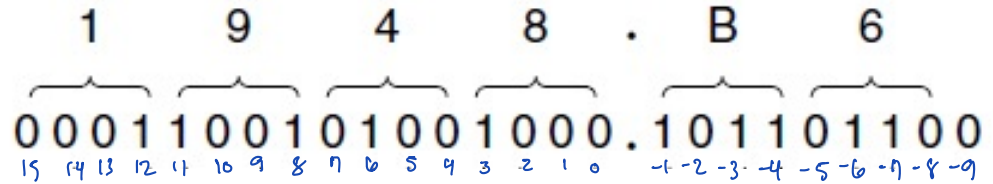


# Radix points: converting to decimal point

$$1 \cdot 16^3 + 9 \cdot 16^2 + 4 \cdot 16^1 + 8 \cdot 16^0 + 11 \cdot 16^{-1} + 6 \cdot 16^{-2}$$

**Example 1**

Hexadecimal



Binary

$$2^{12} + 2^{11} + 2^8 + 2^6 + 2^3 + 2^{-1} + 2^{-3} + 2^{-4} + 2^{-6} + 2^{-7}$$

# Converting from decimal point to binary point

Convert 45.45 to binary

$$\begin{aligned} 45 &= 22 \times 2 + 1 \\ 22 &= 11 \times 2 + 0 \\ 11 &= 5 \times 2 + 1 \\ 5 &= 2 \times 2 + 1 \\ 2 &= 1 \times 2 + 0 \\ 1 &= 0 \times 2 + 1 \end{aligned}$$

$$45_{10} = 101101_2$$

$$45.45_{10} = 101101.011100$$

$$= 1.01101011100 \times 2^5$$

$$\begin{aligned} 0.45 \times 2 &= 0.9 \\ 0.9 \times 2 &= 1.8 \\ 0.8 \times 2 &= 1.6 \\ 0.6 \times 2 &= 1.2 \\ 0.2 \times 2 &= 0.4 \\ 0.4 \times 2 &= 0.8 \end{aligned}$$

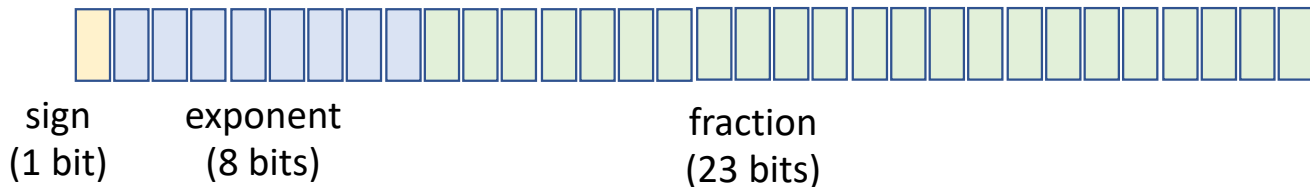
$$0.45_{10} = 0.011100_2$$

# IEEE 754: A standard for binary floating point

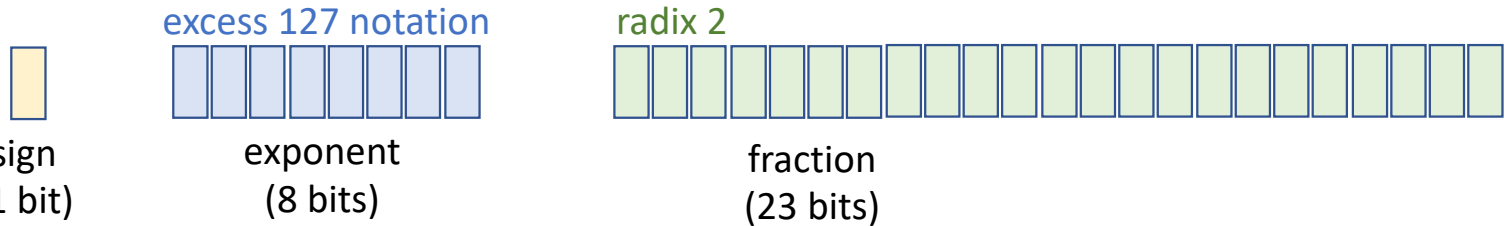
- Different formats existed until the 1980s
- Late 1970s: a committee forms to standardize floating point arithmetic!
- Mainly the work of Berkeley Math Prof. William Kahan

## Formats

- **Single precision (32 bits)**
- Double precision (64 bits)
- Extended precision (80 bits) – mainly used in ALUs



# IEEE 754: A standard for binary floating point



**significand** = 1 . fraction  
*notable implicit 1 and implicit binary point*

Representing in traditional binary point:  
 $number = significand \times 2^{exponent}$

**Example:**  $45.45_{10} = 1.01101011100 \times 2^5$   
 (The '5' in the exponent is circled in blue, and the fraction '1.01101011100' is underlined in blue.)

Convert 45.45 to IEEE 754 single precision format.

0      10000100      0110 101 1100 1100 1100 1100  
 sign      exponent (5)  
           = 132 in excess-127

= 0 100 0010 0011 0101 1100 1100 1100 1100 in IEEE 754 form.  
 OR, 4 2 3 5 C C C C in hex.



# Example 2: Converting 3F80000 (hex) from IEEE 754

excess 127 notation

radix 2

0

0 1 1 1 1 1 1 1

0 0

sign  
(1 bit)

exponent  
(8 bits)

fraction  
(23 bits)

**significand** = 1 . fraction

*notable implicit 1 and implicit binary point*

Representing in traditional binary point:

$$number = significand \times 2^{exponent}$$

Sign = positive  
 exponent = 127 in excess-127  $\Rightarrow$  0  
 fraction = 0

$$+1.\overset{\circlearrowleft}{0} \times 2^0 \leftarrow \text{exponent} = \boxed{\begin{array}{c} 1 \\ 10 \end{array}}$$

↑  
fraction



# Group work: Converting to IEEE 754

Convert 18.25 to IEEE 754 single precision format.