

# Floating Point Numbers

# Floats are not Reals

Int's:

eg.  $40000 * 40000 \rightarrow 1600000000$

$600000 * 600000 \rightarrow ?$

**Floats:**

**Eg 2 : Is  $(x + y) + z = x + (y + z)$ ?**

eg

$(1e20 + -1e20) + 3.14 \rightarrow 3.14$

$1e20 + (-1e20 + 3.14) \rightarrow ??$

**Need to understand details of underlying implementations**

# Objective

- To understand the fundamentals of floating-point representation
- To know the IEEE-754 Floating Point Standard

# Patriot Missile

- Gulf War I
- Failed to intercept incoming Iraqi scud missile (Feb 25, 1991)
- 28 American soldiers killed

GAO Report: GAO/IMTEC-92-26 Patriot Missile Software Problem

<http://www.fas.org/spp/starwars/gao/im92026.htm>



# Patriot Design

- Intended to operate only for a few hours
  - Defend Europe from Soviet aircraft and missile
- Four 24-bit registers (1970s design!)
- Kept time with integer counter: incremented every 1/10 second
- Calculate speed of incoming missile to predict future positions:
$$\text{velocity} = \text{loc}_1 - \text{loc}_0 / (\text{count}_1 - \text{count}_0) * 0.1$$
- But, cannot represent 0.1 exactly!

# Floating Imprecision

- 24-bits:

$$\begin{aligned} 0.1 &= 1/2^4 + 1/2^5 + 1/2^8 + 1/2^9 \\ &\quad + 1/2^{12} + 1/2^{13} + 1/2^{16} + 1/2^{17} \\ &\quad + 1/2^{20} + 1/2^{21} \\ &= 209715 / 2097152 \end{aligned}$$

$$\text{Error is } 0.2/2097152 = 1/10485760$$

One hour = 3600 seconds

$$3600 * 1/10485760 * 10 = 0.0034\text{s}$$

$$20 \text{ hours} = 0.0687\text{s}$$

Miss target! (137 meters)

Two weeks before the incident, Army officials received Israeli data indicating some loss in accuracy after the system had been running for 8 consecutive hours. Consequently, Army officials modified the software to improve the system's accuracy. However, the modified software did not reach Dhahran until February 26, 1991--the day after the Scud incident.

GAO Report

# Review of Numbers

- Computers are made to deal with numbers
- What can we represent in N bits?

- Unsigned integers:

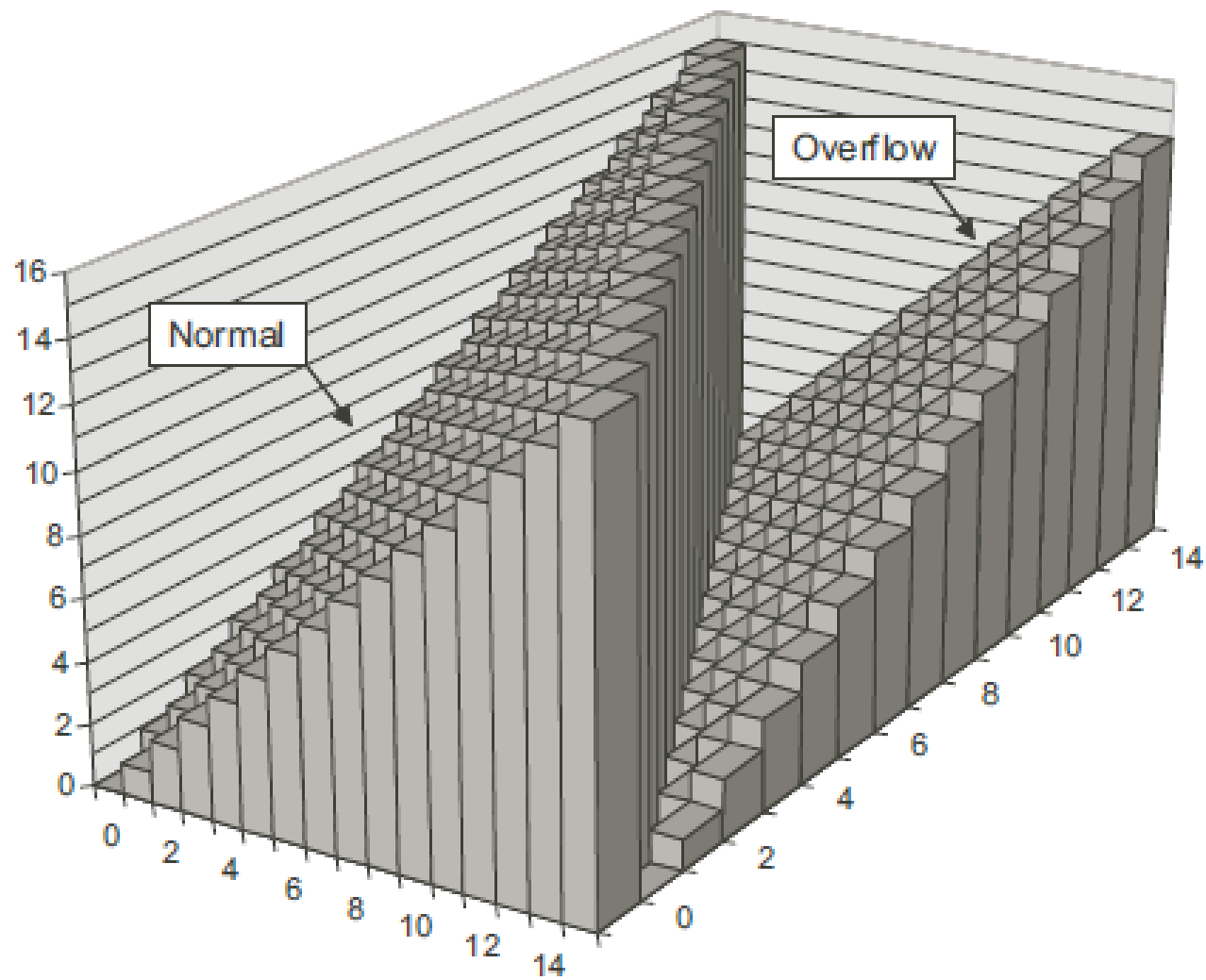
0 to  $2^N - 1$

- Signed Integers (Two's Complement)

$-2^{(N-1)}$  to  $2^{(N-1)} - 1$



### Unsigned Addition (4-bit word)



# Other Numbers

- What about other numbers?
  - Very large numbers? (seconds/century)  
 $3,155,760,000_{10}$  ( $3.15576_{10} \times 10^9$ )
  - Very small numbers? (atomic diameter)  
 $0.00000001_{10}$  ( $1.0_{10} \times 10^{-8}$ )
  - Rationals (repeating pattern)
    - $2/3$  (0.6666666666. . .)
  - Irrationals  
 $2^{1/2}$  (1.414213562373. . .)
  - Transcendentals
    - $e$  (2.718...),  $\pi$  (3.141...)
- All represented in scientific notation

# Exponential Notation

- The following are equivalent representations of **1,234**

$$123,400.0 \times 10^{-2}$$

$$12,340.0 \times 10^{-1}$$

$$1,234.0 \times 10^0$$

$$123.4 \times 10^1$$

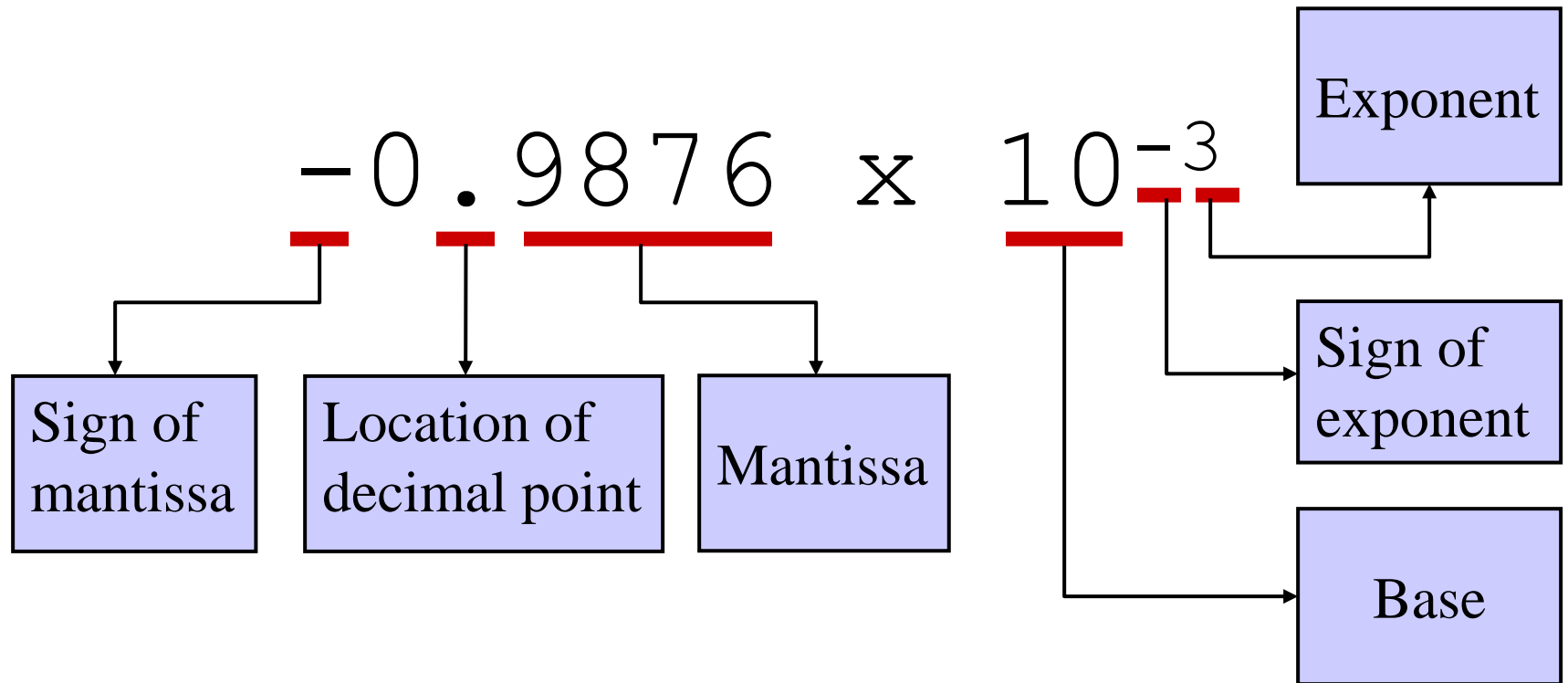
$$12.34 \times 10^2$$

$$1.234 \times 10^3$$

$$0.1234 \times 10^4$$

The representations differ in that the decimal place – the “point” -- “floats” to the left or right (with the appropriate adjustment in the exponent).

# Parts of a Floating Point Number



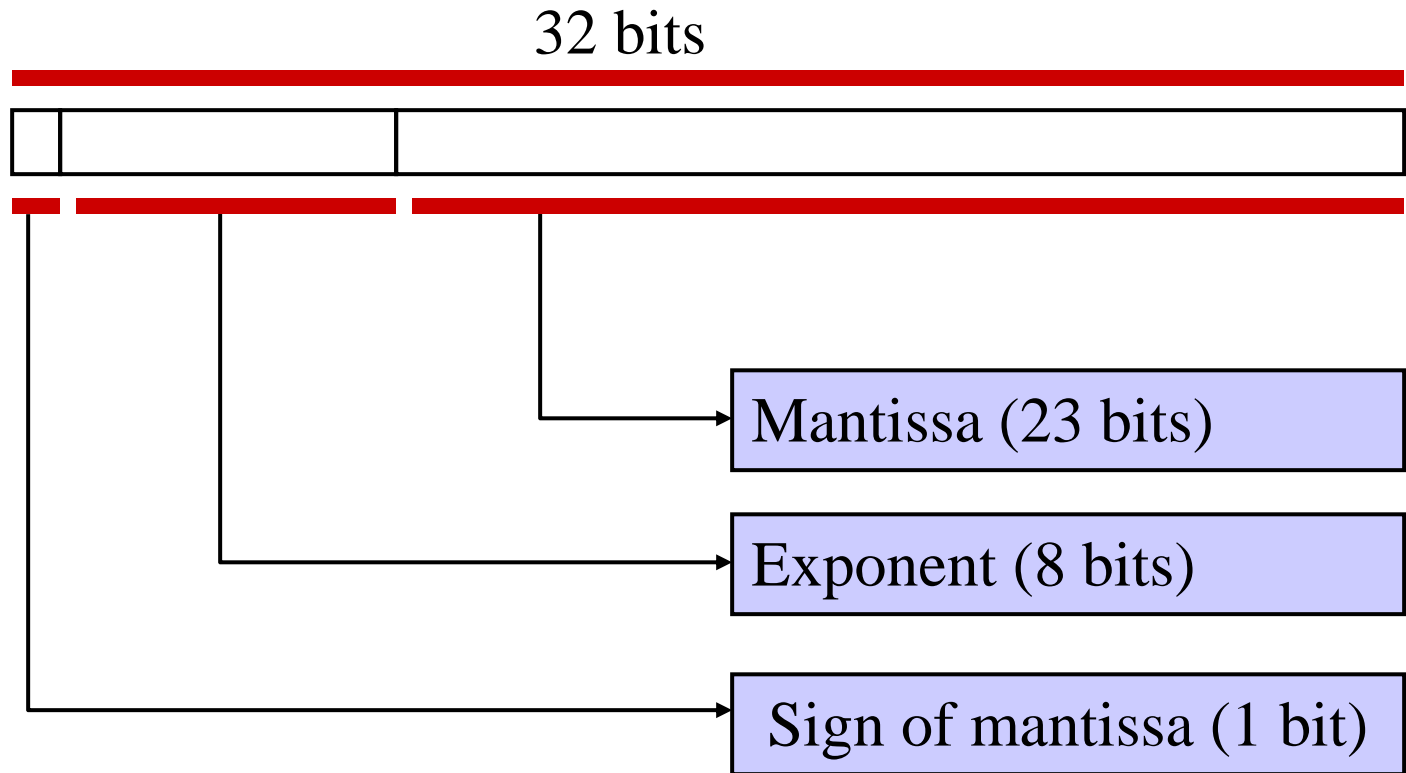
# IEEE 754 Standard

- Most common standard for representing floating point numbers
- Single precision: 32 bits, consisting of...
  - Sign bit (1 bit)
  - Exponent (8 bits)
  - Mantissa (23 bits)
- Double precision: 64 bits, consisting of...
  - Sign bit (1 bit)
  - Exponent (11 bits)
  - Mantissa (52 bits)

**Prof. Willian Kahan**



# Single Precision Format



# Normalization

- The mantissa is *normalized*
- Has an implied decimal place on left
- Has an implied “1” on left of the decimal place
- E.g.,
  - Mantissa → 101000000000000000000000000000
  - Represents...  $1.101_2 = 1.625_{10}$
- Normalized form: no leading 0s  
(exactly one digit to left of decimal point)
  - Normalized:  $1.0 \times 10^{-9}$
  - Not normalized:  $0.1 \times 10^{-8}, 10.0 \times 10^{-10}$

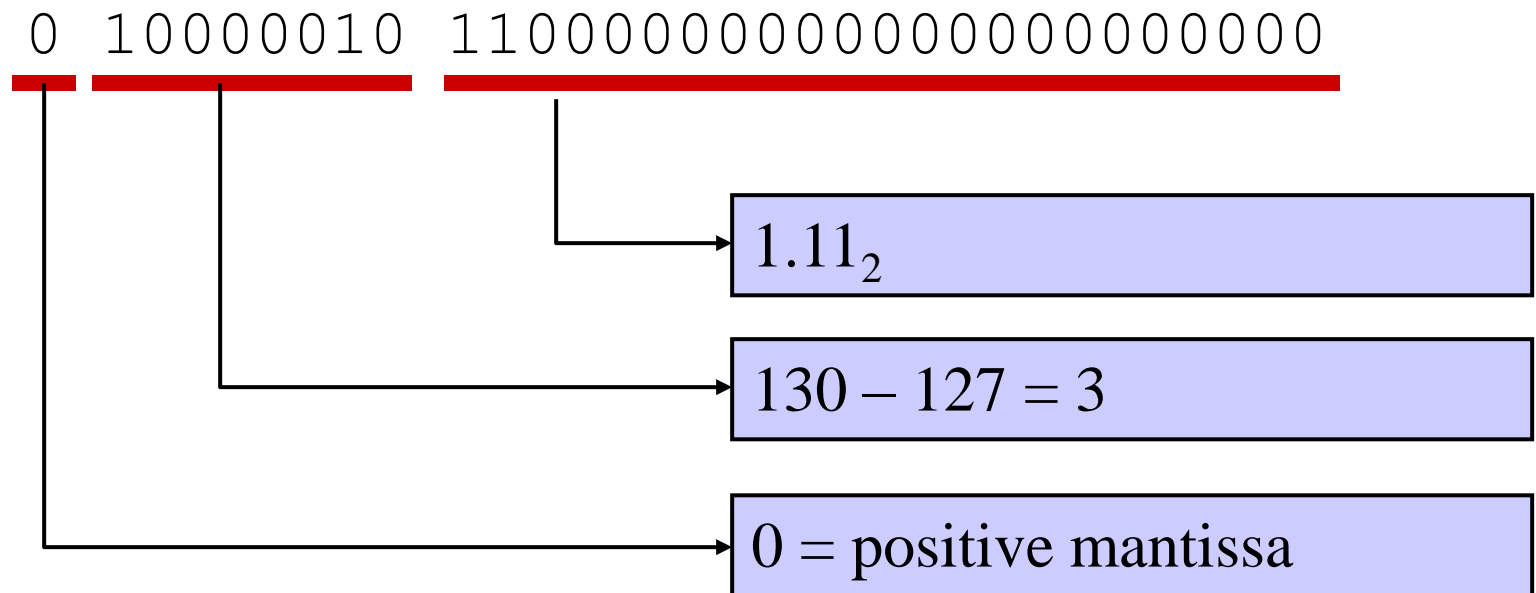
# Excess Notation

- To include +ve and –ve exponents, “excess” notation is used
- Single precision: excess 127
- Double precision: excess 1023
- The value of the exponent stored is larger than the actual exponent
- E.g., excess 127,
  - Exponent  $\rightarrow$  10000111
  - Represents...  $135 - 127 = 8$



# Example

- Single precision



$$+1.11_2 \times 2^3 = 1110.0_2 = 14.0_{10}$$

# Hexadecimal

- It is convenient and common to represent the original floating point number in hexadecimal
- The preceding example...

0	100	00001	0	110	00000	00000	00000	00000	00000
4	1	6	0	0	0	0	0	0	

# Converting from Floating Point

- E.g., What decimal value is represented by the following 32-bit floating point number?

C17B0000<sub>16</sub>

- Step 1
  - Express in binary and find S, E, and M

$$C17B0000_{16} =$$

$\underline{1} \quad \underline{10000010} \quad \underline{111101100000000000000000}_2$   
 S                      E    M

1 = negative  
 0 = positive

- Step 2
  - Find “real” exponent,  $n$
  - $n = E - 127$ 
$$= 10000010_2 - 127$$
$$= 130 - 127$$
$$= 3$$

- Step 3
  - Put S, M, and  $n$  together to form binary result
  - (Don't forget the implied "1." on the left of the mantissa.)

$$-1.1111011_2 \times 2^n =$$

$$-1.1111011_2 \times 2^3 =$$

$$-1111.1011_2$$

- Step 4
  - Express result in decimal

$$\begin{array}{r} \underline{-1111} . \underline{1011}_2 \\ -15 \end{array}$$

$2^{-1} = 0.5$   
 $2^{-3} = 0.125$   
 $2^{-4} = \underline{0.0625}$   
 $0.6875$

Answer: -15.6875

# Converting from Floating Point

- E.g., What decimal value is represented by the following 32-bit floating point number?

42808000<sub>16</sub>



# Converting to Floating Point

- E.g., Express  $36.5625_{10}$  as a 32-bit floating point number (in hexadecimal)

- Step 1
  - Express original value in binary

$$36.5625_{10} =$$

$$100100.1001_2$$

- Step 2
  - Normalize

$$100100.1001_2 =$$

$$1.001001001_2 \times 2^5$$

- Step 3
  - Determine S, E, and M

$$\underbrace{+1}_S . \underbrace{001001001}_M {}_2 \times 2^{\underbrace{5}_n}$$

$$\begin{aligned} E &= n + 127 \\ &= 5 + 127 \\ &= 132 \\ &= 10000100_2 \end{aligned}$$

$S = 0$  (because the value is positive)

- Step 4
  - Put S, E, and M together to form 32-bit binary result

0	10000100	00100100100000000000000000000000	<sub>2</sub>
S	E	M	

$$+1.001001001_2 \times 2^5$$

- Step 5
  - Express in hexadecimal

$$\begin{array}{cccccccc}
 0 & 10000100 & 00100100 & 10000000 & 00000000 & 00000000 & 00000000 & 00000000_2 = \\
 0100 & 0010 & 0001 & 0010 & 0100 & 0000 & 0000 & 0000_2 = \\
 4 & 2 & 1 & 2 & 4 & 0 & 0 & 0_{16}
 \end{array}$$

**Answer: 42124000<sub>16</sub>**

# Converting to Floating Point

- E.g., Express  $6.5_{10}$  as a 32-bit floating point number (in hexadecimal)

# Converting to Floating Point

- E.g., Express 0.1 as a 32-bit floating point number (in hexadecimal)



Item	Single precision	Double precision
Bits in sign	1	1
Bits in exponent	8	11
Bits in fraction	23	52
Bits, total	32	64
Exponent system	Excess 127	Excess 1023
Exponent range	-126 to +127	-1022 to +1023
Smallest normalized number	$2^{-126}$	$2^{-1022}$
Largest normalized number	approx. $2^{128}$	approx. $2^{1024}$
Decimal range	approx. $10^{-38}$ to $10^{38}$	approx. $10^{-308}$ to $10^{308}$
Smallest denormalized number	approx. $10^{-45}$	approx. $10^{-324}$

**Figure B-5.** Characteristics of IEEE floating-point numbers.

# Summary: IEEE Floating Point Single Precision (32 bits)

1                  8 bits                                  23 bits



31    30                          23    22    0

Exponent values:	0	zeroes
	1-254	$\text{exp} + 127$
	255	infinities, NaN

$$\text{Value} = (1 - 2 * \text{Sign}) (1 + \text{Fraction})^{\text{Exponent} - 127}$$

# Denormalized Values

- Condition
  - $\text{exp} = 000\dots 0$
- Value
  - Exponent value  $E = -\text{Bias} + 1$
  - Significand value  $M = 0.\text{xxx}\dots\text{x}_2$ 
    - $\text{xxx}\dots\text{x}$ : bits of `frac`
- Cases
  - $\text{exp} = 000\dots 0, \text{frac} = 000\dots 0$ 
    - Represents value 0
    - Note that have distinct values  $+0$  and  $-0$
  - $\text{exp} = 000\dots 0, \text{frac} \neq 000\dots 0$ 
    - Numbers very close to 0.0
    - Lose precision as get smaller

# Special Values

- Condition
  - $\text{exp} = 111\dots 1$
- Cases
  - $\text{exp} = 111\dots 1, \text{frac} = 000\dots 0$ 
    - Represents value  $\infty$  (infinity)
    - Operation that overflows
    - Both positive and negative
    - E.g.,  $1.0/0.0 = -1.0/-0.0 = +\infty$ ,  $1.0/-0.0 = -\infty$
  - $\text{exp} = 111\dots 1, \text{frac} \neq 000\dots 0$ 
    - Not-a-Number (NaN)
    - Represents case when no numeric value can be determined
    - E.g.,  $\text{sqrt}(-1)$ ,  $\infty - \infty$

# Interesting Numbers

• Description	exp	frac	Numeric Value
• Zero	00...00	00...00	0.0
• Smallest Pos. Denorm.	00...00	00...01	$2^{-\{23,52\}} \times 2^{-\{126,1022\}}$
• Single	$\approx 1.4 \times 10^{-45}$		
• Double	$\approx 4.9 \times 10^{-324}$		
• Largest Denormalized	00...00	11...11	$(1.0 - \epsilon) \times 2^{-\{126,1022\}}$
• Single	$\approx 1.18 \times 10^{-38}$		
• Double	$\approx 2.2 \times 10^{-308}$		
• Smallest Pos. Normalized	00...01	00...00	$1.0 \times 2^{-\{126,1022\}}$
• Just larger than largest denormalized			
• One	01...11	00...00	1.0
• Largest Normalized	11...10	11...11	$(2.0 - \epsilon) \times 2^{\{127,1023\}}$
• Single	$\approx 3.4 \times 10^{38}$		
• Double	$\approx 1.8 \times 10^{308}$		