# Floating Point Representation

1. Fractions

10.625

$$\overbrace{1010}^{10}.10100\ 00$$
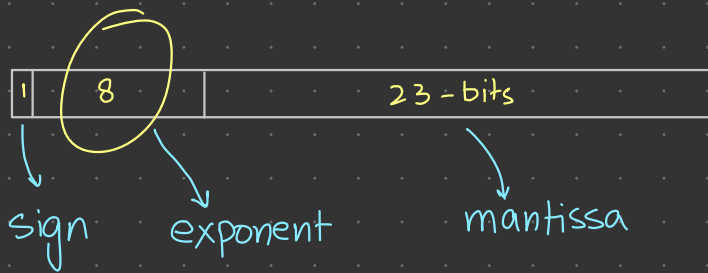
$2^1$
$2^0$

$2^{-7}$
$2^{-6}$
$2^{-5}$
$2^{-4}$
$2^{-3}$ = 0.125
$2^{-2}$
$2^{-1}$ = 0.5

= 0.5 + 0.125

= 0.625

Issue: fixed "precision" versus "size".

Solution: Use a floating point

???

32-bits — IEEE 754 Standard

| 1 | 8 | 23 - bits |
|---|---|-----------|

sign    exponent    mantissa

Double precision — $(1, 11, 52)$ = 64 - bits

First, "how".

Example:  19.59375

Step 1: Convert to binary

19 / 2 = 9    rem  1          0.59375  × 2  = 1.1875      1
9 / 2 = 4      "   1          0.1875   × 2  = 0.375       0
4 / 2 = 2      "   0          0.375    × 2  = 0.75        0
2 / 2 = 1      "   0          0.75     × 2  = 1.5         1
1 / 2 = 0      "   1          0.5      × 2  = 1.0         1

      stopping                        stopping
      point                           point.

Step 2: Find Exponent (Biased)

        left-most
            1

    1 0 0 1 1 . 1 0 0 1 1                              exp bits

Exponent = 4,        Biased exponent = 4 + 127 ~~> $(2^8 - 1)$

                                      = $(131)_{10}$

                                      = 10000011

## Step 3: Find Mantissa

$1\overset{\frown}{0011}.10011$

Mantissa is the "remainder" ⤳

| 0 | 10000011 | 0011.10011 |

Example 2: 15.0

Step 1:
$$15/2 = 7 \quad \text{rem } 1$$
$$7/2 = 3 \quad '' \quad 1$$
$$3/2 = 1 \quad '' \quad 1$$
$$1/2 = 0 \quad '' \quad 1$$

$1111.000000$

Step 2: $3 + 127 = (130)_{10}$ ⤳ $10000010$

Step 3: $0\ 10000010\ 111.0000$ ————

## Reverse Operation:

Step 1: Add leading 1 to matissa

$1.1110000000$ ————

Step 2: Find Exponent

$10000010$ ⤳ $130 - 127 = ③$

jump 3
$1\overset{\frown}{111}0000000$ ————

$1111.000000$

$15.0$

# Next Largest Possible Number:

15.0

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |

0 1 0 0 0 0 0 1 0 1 1 1.0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 |

0 1 0 0 0 0 0 1 0 1 1 1.0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1

$\frac{1}{2}$ $\frac{1}{4}$ $\frac{1}{8}$ $\frac{1}{16}$ $\frac{1}{32}$ $\frac{1}{64}$ $\frac{1}{128}$ $\frac{1}{256}$ $\frac{1}{512}$ / / /

Exponent = 3

15.0000000119209 <u>0</u>

So:  15.00000001192 <u>07</u>

$=$

15.0000000119 <u>205</u>

$=$

15.00000000000000

Let's see some code ...