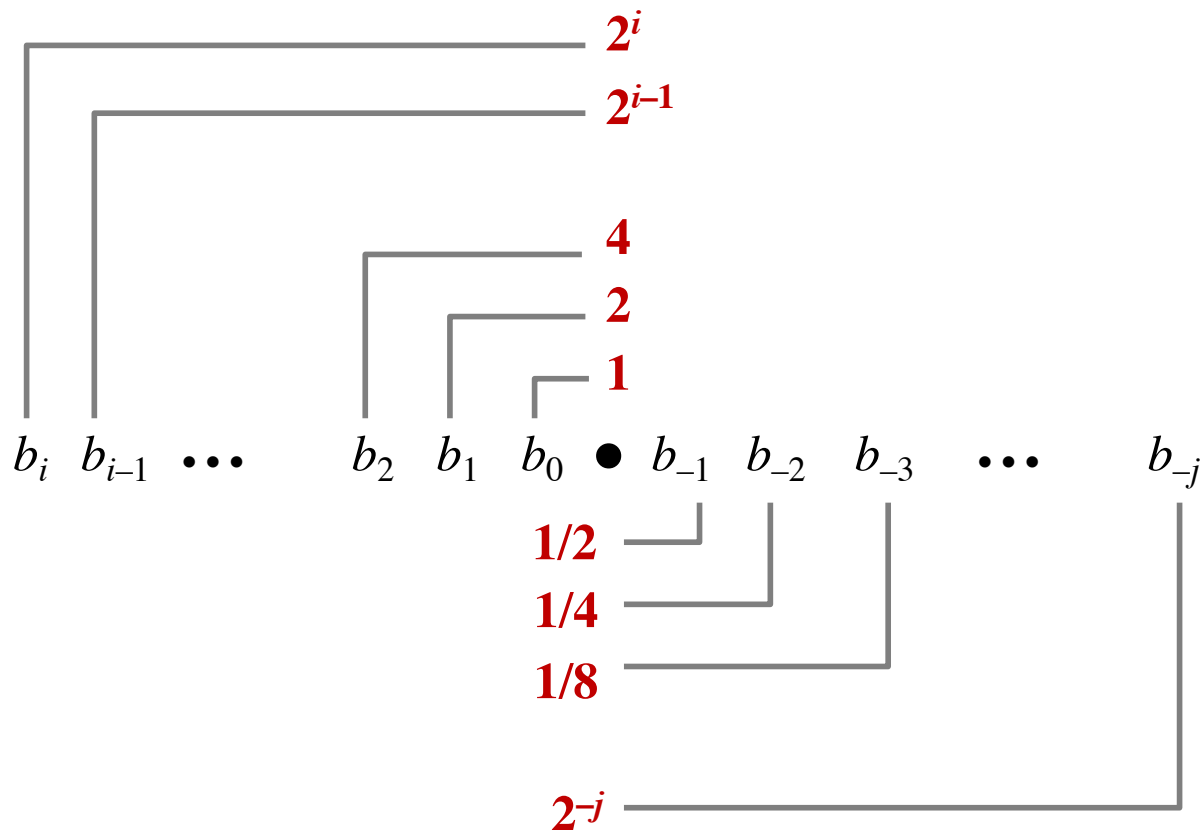# Floating Point Representation

Fractional binary numbers
IEEE floating-point standard
Floating-point operations and rounding
***Lessons for programmers***

Many more details we will skip (it's a 58-page standard...)
See CSAPP 2.4 for more detail.

# Fractional Binary Numbers

$$2^i$$
$$2^{i-1}$$
**4**
**2**
**1**

$b_i \;\; b_{i-1} \;\; \cdots \;\;\;\; b_2 \;\; b_1 \;\; b_0 \;\; \bullet \;\; b_{-1} \;\; b_{-2} \;\; b_{-3} \;\;\; \cdots \;\;\;\; b_{-j}$

**1/2**
**1/4**
**1/8**

$$2^{-j}$$

$$\sum_{k=-j}^{i} b_k \cdot 2^k$$

# Fractional Binary Numbers

Value                    Representation

    5 and 3/4

    2 and 7/8

    47/64


## Observations

    Shift left =

    Shift right =

    Numbers of the form $\texttt{0.111111}..._2$ are...?

## Limitations:

    Exact representation possible  when?


      $1/3 = 0.333333..._{10} = 0.$

# Fixed-Point Representation

**Implied binary point.**

$b_7 \; b_6 \; b_5 \; b_4 \; b_3$ **[.]** $b_2 \; b_1 \; b_0$

$b_7 \; b_6 \; b_5 \; b_4 \; b_3 \; b_2 \; b_1 \; b_0$ **[.]**

**range:** difference between largest and smallest representable numbers

**precision:** smallest difference between any two representable numbers

**fixed point = fixed range, fixed precision**

# IEEE **Floating Point** Standard 754
IEEE = Institute of Electrical and Electronics Engineers

## Numerical form:

$$V_{10} = (-1)^s * M * 2^E$$

**Sign bit $s$** determines whether number is negative or positive

**Significand (mantissa) $M$** usually a fractional value in range [1.0,2.0)

**Exponent $E$** weights value by a (-/+) power of two

Analogous to scientific notation

## Representation:

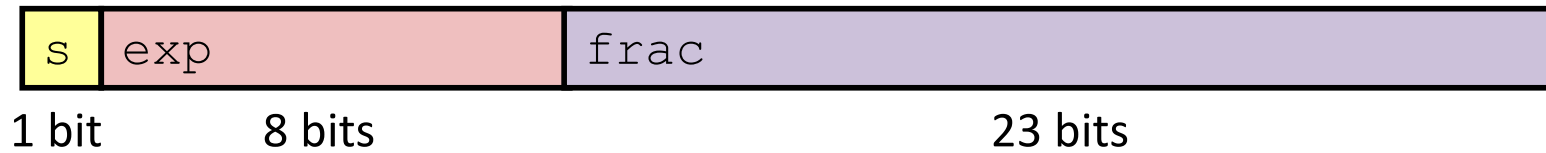MSB `s` = sign bit $s$

`exp` field encodes $E$ (but is *not equal* to E)

`frac` field encodes $M$ (but is *not equal* to M)

| s | exp | frac |
|---|-----|------|

Numerically well-behaved, but hard to make fast in hardware

# Precisions

Single precision (`float`): 32 bits

| s | exp | frac |
|---|-----|------|
| 1 bit | 8 bits | 23 bits |

Double precision (`double`): 64 bits

| s | exp | frac |
|---|-----|------|
| 1 bit | 11 bits | 52 bits |

Finite representation of infinite range…

# Three kinds of values

$$V = (-1)^s * M * 2^E$$

| s | exp | frac |
|---|-----|------|

1. **Normalized:** $M$ = 1.xxxxx…

   As in scientific notation: $0.011 \times 2^5 = 1.1 \times 2^3$

   Representation advantage?

2. **Denormalized**, near zero: $M$ = 0.xxxxx…, smallest $E$

   Evenly space near zero.

3. **Special values:**

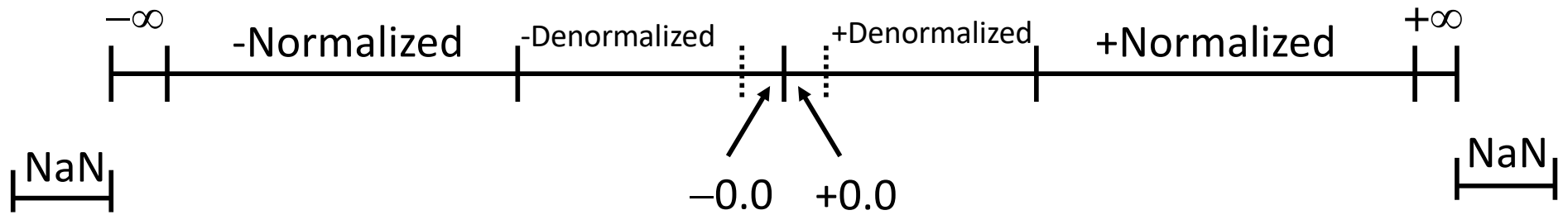   **0.0**:  **s** = 0   **exp** = 00…0   **frac** = 00…0

   **+inf, -inf:**   **exp** = 11…1   **frac** = 00…0
   division by 0.0

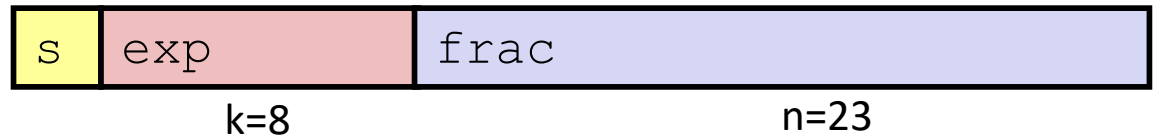   **NaN** ("Not a Number"): **exp** = 11…1   **frac** $\neq$ 00…0
   sqrt(-1), $\infty - \infty$, $\infty * 0$, etc.

# Value distribution

# **Normalized values,** with `float` example

$$V = (-1)^s * M * 2^E$$

| s | exp | frac |
|---|-----|------|
|   | k=8 | n=23 |

**Value:** `float f = 12345.0;`

$12345_{10}$ $\qquad = 11000000111001_2$

$\qquad\qquad\qquad = 1.1000000111001_2 \times 2^{13}$  (normalized form)

**Significand:**

$M \qquad = \qquad 1.\underline{1000000111001}_2$

`frac= ` $\qquad \underline{1000000111001}0000000000_2$

**Exponent:** $E = \mathbf{exp} - \text{Bias} \rightarrow \mathbf{exp} = E + \text{Bias}$

$E \qquad = \qquad\qquad 13$

$\text{Bias} \quad = \qquad\qquad 127 \quad = \qquad 2^7 - 1 = \mathbf{2^{k-1} - 1}$ $\qquad\qquad$ Splits exponents roughly -/+

$\mathbf{exp} = \qquad\qquad 140 \quad = \qquad 10001100_2$

**Result:**

| 0 | 10001100 | 10000001110010000000000 |
|---|----------|-------------------------|
| s | exp | frac |

# Denormalized Values: near zero

"Near zero": $\mathtt{exp} = \mathtt{000...0}$

Exponent:

$E = 1 + \mathtt{exp} - \text{Bias} = 1 - \text{Bias}$    **not:**   $\mathtt{exp} - \text{Bias}$

Significand: leading zero

$M = \mathtt{0.xxx...x}_2$

     **frac = xxx…x**

Cases:

   **exp** = 000…0, **frac** = 000…0      0.0, -0.0

   **exp** = 000…0, **frac** ≠ 000…0

# Value distribution example

6-bit IEEE-like format

| s | exp | frac |
|---|-----|------|
| 1 | 3 | 2 |

Bias = $2^{3-1} - 1 = 3$

**Full Range**

frac= 00,  01,  10,  11
$M$  = 1.00, 1.01, 1.10, 1.11

s=1, **exp**=101
$E$ = 5-3 = 2



◆ Denormalized  ▲ Normalized  ■ Infinity

s=0, **exp**=110
$E$ = 6-3 = 3

**Zoom in to 0**

exp=000
$E$ = 1-3 = -2
Denormalized
= evenly spaced

s=1, **exp**=010
$E$ = 2-3 = -1

s=0, **exp**=001
$E$ = 1-3 = -2

same spacing

# Try to represent 3.14, 6-bit example

**6-bit IEEE-like format**

Bias = $2^{3-1} - 1 = 3$

| s | exp | frac |
|---|-----|------|
| 1 | 3 | 2 |

**Value:** `3.14`;

3.14 = 11.0010 0011 1101 0111 0000 1010 000…

= 1.1001 0001 1110 1011 1000 0101 0000… $_2$ x $2^1$   (normalized form)

**Significand:**

$M$     =     1.100100011110101110111100001010000… $_2$

**frac=**     $\underline{10}_2$

**Exponent:**

$E$ = 1          Bias = 3     **exp** = 4 = $100_2$

**Result:**

**0 100 10**   =   $1.10_2 \times 2^1 = 3$     next highest?

# Floating Point Arithmetic*

$$V = (-1)^S * M * 2^E$$

| s | exp | frac |
|---|-----|------|

```
double x = ..., y = ...;
double z = x + y;
```

1. **Compute exact result.**
2. **Fix/Round**, roughly:

   Adjust $M$ to fit in [1.0, 2.0)…

   > If $M$ >= 2.0: shift $M$ right, increment $E$
   > If $M$ < 1.0: shift $M$ left by k, decrement $E$ by k

   Overflow to infinity if $E$ is too wide for **exp**

   Round* $M$ if too wide for **frac**.

   Underflow if nearest representable value is 0.

   …                                            *complicated…

# Lessons for programmers

$V = (-1)^S * M * 2^E$

| s | exp | frac |
|---|-----|------|

`float` ≠ real number ≠ `double`

Rounding breaks associativity and other properties.

```
double a = ..., b = ...;
...
if (a == b) ...

if (abs(a - b) < epsilon) ...
```