# Floating Point Numbers

CS 2130: Computer Systems and Organization 1

January 30, 2023

# Announcements

- TA Office Hours starting Wednesday
  - Wednesdays, Rice 011
  - Thurs-Sun, Olsson 001
- Please join our Discord server
- Lab tomorrow: hex editor
- Homework 1 due Feb 6 (Mon)

## Operations

So far, we have discussed:

- Addition: $x + y$
  - Can get multiplication
- Subtraction: $x - y$
  - Can get division, but more difficult
- Unary minus (negative): $-x$
  - Flip the bits and add 1

# Operations (on Integers)

Bit vector: fixed-length sequence of bits (ex: bits in an integer)

- Manipulated by bitwise operations

Bitwise operations: operate over the bits in a bit vector

- Bitwise not: $\sim$x - flips all bits (unary)
- Bitwise and: x & y - set bit to 1 if $x, y$ have 1 in same bit
- Bitwise or: x | y - set bit to 1 if either $x$ or $y$ have 1
- Bitwise xor: x ^ y - set bit to 1 if $x, y$ bit differs

- Logical not: !x
  - $!0 = 1$ and $!x = 0, \forall x \neq 0$
  - Useful in C, no booleans
  - Some languages name this one differently
- Left shift: x << y - move bits to the left
  - Effectively multiply by powers of 2
- Right shift: x >> y - move bits to the right
  - Effectively divide by powers of 2
  - Signed (extend sign bit) vs unsigned (extend 0)

$$1000\ 0000 \ll 1$$
$$\swarrow$$
$$00\ 0000\ 00 \gg 1$$
$$\searrow$$
$$\underline{0\ 0000\ 00}$$

What about other kinds of numbers?

Floating point numbers

- Decimal: 3.14159

  ↳ decimal point

Floating point numbers

- Decimal: 3.14159
- Binary: 11.10110

↳ binary point

Floating point numbers

- Decimal: 3.14159
- Binary: 11.10110
- With integers, the point is always fixed after all digits
- With floating point numbers, the point can move!

0110101.

Floating point numbers

- Decimal: 3.14159
- Binary: 11.10110
- With integers, the point is always fixed after all digits
- With floating point numbers, the point can move!

Challenge! only 2 symbols in binary

Convert the following decimal to scientific notation:

2130

$2.130 \times 10^3$

Convert the following binary to scientific notation:

$$101101,$$

$$1.01101 \times 2^5$$

An interesting phenomenon:

- Decimal: first digit can be any number *except* 0

$$2.13 \times 10^3$$

$0.213 \times 10^4$ (crossed out)

# Something to Notice

An interesting phenomenon:

- Decimal: first digit can be any number *except* 0

$$2.13 \times 10^3$$

- Binary: first digit can be any number *except* 0 <span style="color:red">Wait!</span>

$$1.01101 \times 2^5$$

An interesting phenomenon:

- Decimal: first digit can be any number *except* 0

$$2.13 \times 10^3$$

- Binary: first digit can be any number *except* 0 <span style="color:red">Wait!</span>

$$1.01101 \times 2^5$$

  - First digit can only be 1
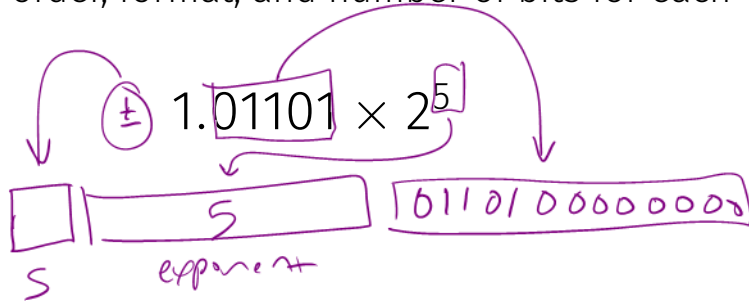
We must store 3 components

- **sign** (1-bit): 1 if negative, 0 if positive
- **fraction** or **mantissa**: (?-bits): bits after binary point
- **exponent** (?-bits): how far to move binary point

*We do not need to store the value before the binary point. Why?*

# Floating Point in Binary

How do we store them?

- Originally many different systems
- IEEE standardized system (IEEE 754 and IEEE 854)
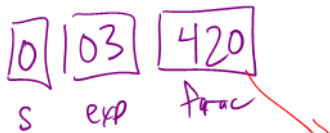- Agreed-upon order, format, and number of bits for each



$$\pm \; 1.01101 \times 2^5$$

10

# Example

A rough example in Decimal:

$$6.42 \times 10^3$$

$\frac{42}{100}$

1 sign
2 exp
3 frac

| 0 | 103 | 420 |
|---|-----|-----|
| s | exp | frac |

$\frac{420}{1000}$

How do we store the exponent?

- Exponents *can* be negative

$$2^{-3} = \frac{1}{2^3} = \frac{1}{8}$$

- Need positive and negative ints (but no minus sign)

# Exponent

How do we store the exponent?

- Exponents *can* be negative

$$2^{-3} = \frac{1}{2^3} = \frac{1}{8}$$

- Need positive and negative ints (but no minus sign)
- *Don't we always use Two's Complement?*

# Exponent

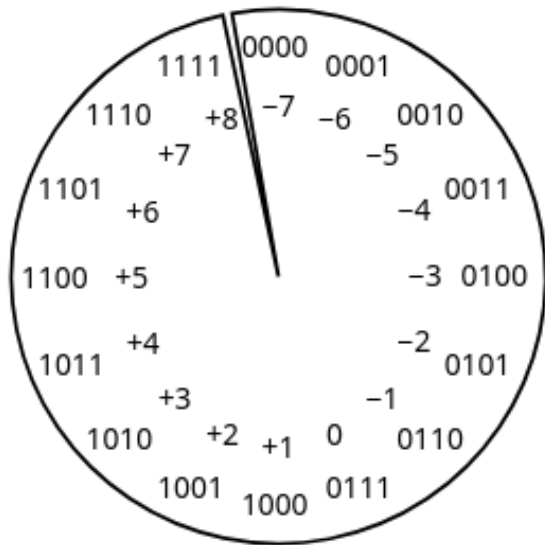How do we store the exponent?

- Exponents *can* be negative

$$2^{-3} = \frac{1}{2^3} = \frac{1}{8}$$

- Need positive and negative ints (but no minus sign)
- *Don't we always use Two's Complement?* Unfortunately Not

How do we store the exponent?

- Exponents *can* be negative

$$2^{-3} = \frac{1}{2^3} = \frac{1}{8}$$

- Need positive and negative ints (but no minus sign)
- *Don't we always use Two's Complement?* Unfortunately Not
- Biased integers
  - Make comparison operations run more smoothly
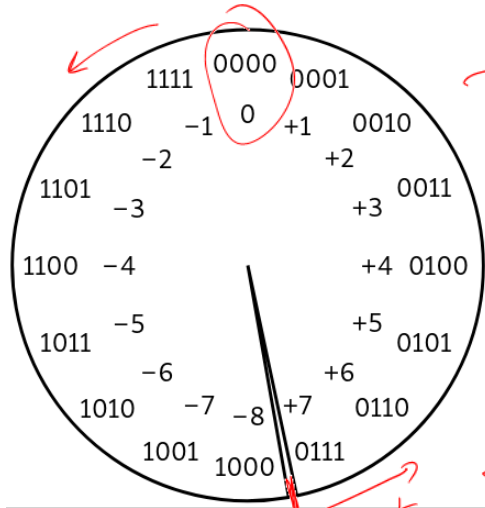  - Hardware more efficient to build
  - Other valid reasons

# Biased Integers

Similar to Two's Complement, but add **bias**

- **Two's Complement**: Define 0 as 00...0
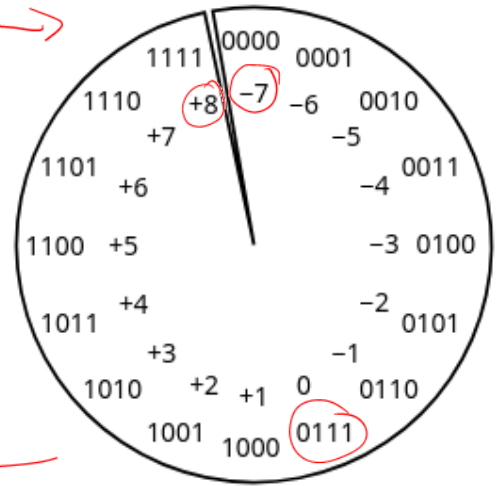- **Biased**: Define 0 as 0111...1
- Biased wraps from 000...0 to 111...1

# Biased Integers



Two's Complement

Biased

14

Calculate value of biased integers (4-bit example)

biased → $1\,\cancel{0}\,\cancel{0}\,\overset{10}{\cancel{1}}\,\overset{10}{0} = -5$

$$-\,0\,1\,1\,1$$

2's complement → $\overline{1\,0\,1\,1} = -5$

$\boxed{\sim\\+1}$
$\quad 0\,1\,0\,0$
$\quad\quad\quad 1$

binary → $\overline{0\,1\,0\,1} = 5$

$0010 = -5$

bias → $\underline{0\,0\,1\,0} = 2$

$T$

$2^s \rightarrow$
$$-\,\underline{0\,1\,1} = -7$$
$$\overline{1\,0\,1\,1}\quad -5$$

# Floating Point Example

$1.0000111$ (with $1\,000$ written above)

$+ \ 101.011_2$

$+ \ 1.01011 \times 2^2$

$8\text{-bit}$
$1 \quad \text{sign}$
$4 \quad \text{exp}$
$3 \quad \text{fraction}$

$$
\begin{array}{r}
11 \\
0010 \\
+ \ 0111 \\
\hline
1001
\end{array}
$$

$0 \ | \ 1001 \ | \ 011$

sign    exponent      frac

17

$$101.011_2$$

What does the following encode?

$$s \quad exp \quad frac$$

$$\boxed{1} \; \boxed{001110} \; \boxed{1010101}$$

$$001110 \quad \leftarrow \text{# in binary}$$
$$- \; 011111 \quad \leftarrow \text{bias}$$
$$= -17$$

$$-1.1010101 \quad \times 2^{-17}$$

$$0.00 \text{---} \quad 0110101$$

What does the following encode?

$$\boxed{1}\ \boxed{001110}\ \boxed{1010101}$$

What about 0?

# Floating Point Numbers

Four cases:

- **Normalized**: What we have seen today

$$s\ eeee\ ffff = \pm 1.ffff \times 2^{eeee - \text{bias}}$$

- **Denormalized**: Exponent bits all 0

$$s\ eeee\ ffff = \pm 0.ffff \times 2^{1 - \text{bias}}$$

- **Infinity**: Exponent bits all 1, fraction bits all 0 (i.e., $\pm\infty$)
- **Not a Number (NaN)**: Exponent bits all 1, fraction bits not all 0