

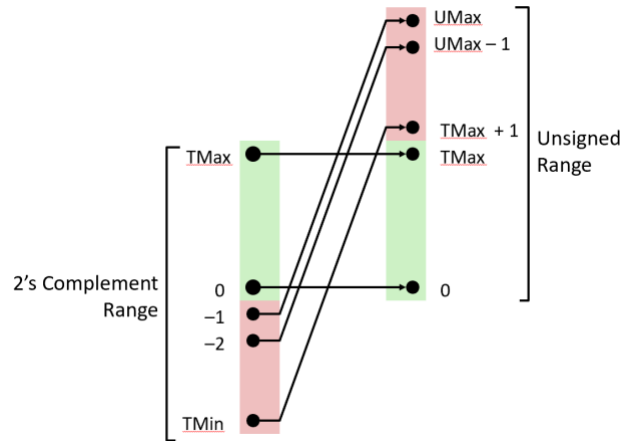
# CSE 351 Section 3 – Integers and Floating Point

Welcome back to section, we're happy that you're here ☺ . . . . .

## Integers and Arithmetic Overflow

**Arithmetic overflow** occurs when the result of a calculation can't be represented in the current encoding scheme (*i.e.*, it lies outside of the representable range of values), resulting in an incorrect value.

- **Unsigned overflow:** the result lies outside of [UMin, UMax]; an indicator of this is when you add two numbers and the result is smaller than either number.
- **Signed overflow:** the result lies outside of [TMin, TMax]; an indicator of this is when you add two numbers with the same sign and the result has the opposite sign.



### Exercises:

1) Assuming these are all signed two's complement 6-bit integers, compute the result of each of the following additions. For each, indicate if it resulted in overflow. [Spring 2016 Midterm 1C]

$$\begin{array}{r} 001001 \\ + 110110 \\ \hline 111111 \end{array}$$

No

$$\begin{array}{r} 110001 \\ + 111011 \\ \hline \text{1} 101100 \end{array}$$

No

$$\begin{array}{r} 011001 \\ + 001100 \\ \hline 100101 \end{array}$$

Yes

$$\begin{array}{r} 101111 \\ + 011111 \\ \hline \text{1} 001110 \end{array}$$

No

2) Find the largest 8-bit unsigned numeral (answer in hex) such that  $c + 0x80$  causes NEITHER signed nor unsigned overflow in 8 bits. [Autumn 2019 Midterm 1C]

Unsigned overflow will occur for  $c > 0x80$ . Signed overflow can only happen if  $c$  is negative (also  $> 0x80$ ). Largest is therefore,  $0x7F$

3) Find the smallest 8-bit numeral (answer in hex) such that  $c + 0x71$  causes signed overflow, but NOT unsigned overflow in 8 bits. [Autumn 2018 Midterm 1C]

For signed overflow, need  $(+) + (+) = (-)$ . For no unsigned overflow, need no carryout from MSB. The first  $(-)$  encoding we can reach from  $0x71$  is  $0x80$ .  $0x80 - 0x71 = 0xF$ .

## Goals of Floating Point

Representation should include: [1] a large range of values (both very small and very large numbers), [2] a high amount of precision, and [3] real arithmetic results (e.g.  $\infty$  and NaN).

## IEEE 754 Floating Point Standard

The value of a real number can be represented in scientific binary notation as:

$$\text{Value} = (-1)^{\text{sign}} \times \text{Mantissa}_2 \times 2^{\text{Exponent}} = (-1)^S \times 1.M_2 \times 2^{E-\text{bias}}$$

The binary representation for floating point values uses three fields:

- **S**: encodes the *sign* of the number (0 for positive, 1 for negative)
- **E**: encodes the *exponent* in **biased notation** with a bias of  $2^{w-1}-1$
- **M**: encodes the *mantissa* (or *significand*, or *fraction*) – stores the fractional portion, but does not include the implicit leading 1.

	S	E	M
float	1 bit	8 bits	23 bits
double	1 bit	11 bits	52 bits

How a float is interpreted depends on the values in the exponent and mantissa fields:

E	M	Meaning
0b0...0	0b0...0	+/- 0
0b0...0	non-zero	denormalized number
everything else	anything	normalized number
0b1...1	0b0...0	+/- $\infty$
0b1...1	non-zero	Not-a-Number (NaN)

### Exercises:

#### Bias Notation

- 1) Suppose that instead of 8 bits, E was only designated 4 bits. What is the bias in this case?  $2^{(4-1)} - 1 = 7$
- 2) Compare these two representations of E for the following values:

Exponent	E (4 bits)	E (8 bits)
1	1 0 0 0	1 0 0 0 0 0 0 0
0	0 1 1 1	0 1 1 1 1 1 1 1
-1	0 1 1 0	0 1 1 1 1 1 1 0

The representations are the same except the length of number of repeating bits in the middle are different.

## Floating Point / Decimal Conversions

- 3) Let's say that we want to represent the number 3145728.125 (broken down as  $2^{21} + 2^{20} + 2^{-3}$ )
- a. Convert this number to into single precision floating point representation:

0	1	0	0	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

- b. How does this number highlight a limitation of floating point representation?  
 Could only represent  $2^{21} + 2^{20}$ . Not enough bits in the mantissa to hold  $2^{-3}$ , which caused rounding.

- 4) What are the decimal values of the following floats?

0x80000000

-0

0xFF94BEEF

NaN

0x41180000

+9.5

$0x41180000 = 0b\ 0|100\ 0001\ 0|001\ 1000\ 0\dots0$ .

$S = 0, E = 128 + 2 = 130 \rightarrow \text{Exponent} = E - \text{bias} = 3, \text{Mantissa} = 1.0011_2$

$1.0011_2 \times 2^3 = 1001.1_2 = 8 + 1 + 0.5 = 9.5$

## Floating Point Mathematical Properties

- Not associative:  $(2 + 2^{50}) - 2^{50} \neq 2 + (2^{50} - 2^{50})$
- Not distributive:  $100 \times (0.1 + 0.2) \neq 100 \times 0.1 + 100 \times 0.2$
- Not cumulative:  $2^{25} + 1 + 1 + 1 + 1 \neq 2^{25} + 4$

### Exercises:

- 5) Based on floating point representation, explain why each of the three statements above occurs.

Associative: Only 23 bits of mantissa, so  $2 + 2^{50} = 2^{50}$  (2 gets rounded off). So LHS = 0, RHS = 2.

Distributive: 0.1 and 0.2 have infinite representations in binary point ( $0.2 = 0.\overline{0011}_2$ ), so the LHS and RHS suffer from different amounts of rounding (try it!).

Cumulative: 1 is 25 powers of 2 away from  $2^{25}$ , so  $2^{25} + 1 = 2^{25}$ , but 4 is 23 powers of 2 away from  $2^{25}$ , so it doesn't get rounded off.

- 6) If  $x$  and  $y$  are variable type `float`, give two *different* reasons why  $(x + 2 * y) - y == x + y$  might evaluate to false.

(1) Rounding error: like what is seen in the examples above.

(2) Overflow: if  $x$  and  $y$  are large enough, then  $x + 2 * y$  may result in infinity when  $x + y$  does not.