

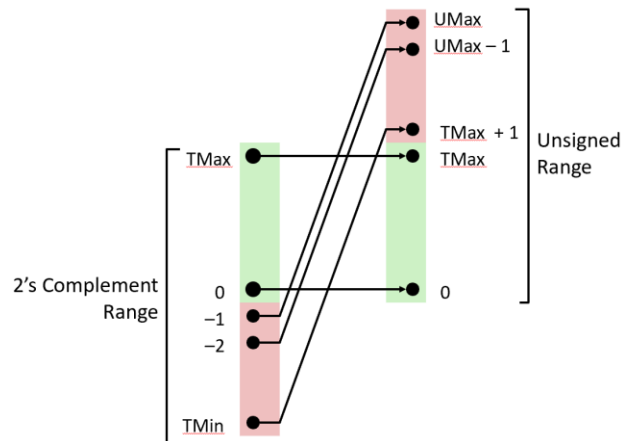
# CSE 351 Section 3 – Integers and Floating Point

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]

001001	110001	011001	101111
+ 110110	+ 111011	+ 001100	+ 011111

- 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]
- 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]

## 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
0	anything	denormalized number (denorm)
1-254	anything	normalized number
255	zero	infinity ( $\infty$ )
255	nonzero	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) Compare these two representations of E for the following values:

Exponent	E (4 bits)	E (8 bits)												
1	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td></tr></table>					<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td></tr></table>								
0	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td></tr></table>					<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td></tr></table>								
-1	<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td></tr></table>					<table border="1" style="display: inline-table; vertical-align: middle;"><tr><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td><td style="width: 20px; height: 20px;"></td></tr></table>								

Notice any patterns?



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

# IEEE 754 Float (32 bit) Flowchart

