

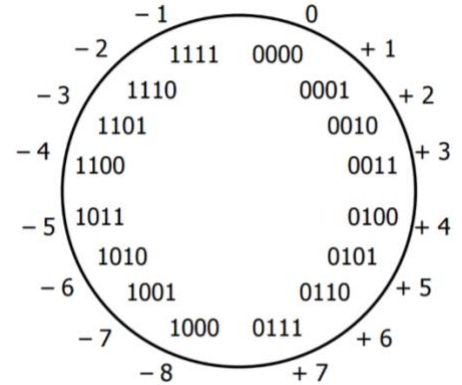
CSE 351 Section 3 – Integers and Floating Point

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

Signed Integers with Two's Complement

Two's complement is the standard for representing signed integers:

- The most significant bit (MSB) has a negative value; all others have positive values (same as unsigned)
- Binary addition is performed the same way for signed and unsigned
- The bit representation for the negative value (additive inverse) of a Two's Complement number can be found by:
flipping all the bits and adding 1 (i.e. $-x = \sim x + 1$).



The “number wheel” showing the relationship between 4-bit numerals and their Two's Complement interpretations is shown on the right:

- The largest number is 7 whereas the smallest number is -8
- There is a nice symmetry between numbers and their negative counterparts except for -8

Exercises: (assume 8-bit integers)

1) What is the **largest integer**? The **largest integer + 1**?

<u>Unsigned:</u>	<u>Two's Complement:</u>
------------------	--------------------------

2) How do you represent (if possible) the following numbers: **39, -39, 127**?

<u>Unsigned:</u>	<u>Two's Complement:</u>
39:	39:
-39:	-39:
127:	127:

3) Compute the following sums in binary using your Two's Complement answers from above. *Answer in hex.*

a. 39 -> 0b _____ + (-39) -> 0b _____ 0x __ <- 0b _____	b. 127 -> 0b _____ + (-39) -> 0b _____ 0x __ <- 0b _____
c. 39 -> 0b _____ - 127 -> 0b _____ 0x __ <- 0b _____	d. 127 -> 0b _____ + 39 -> 0b _____ 0x __ <- 0b _____

4) Interpret each of your answers above and indicate whether-or-not overflow has occurred.

a. 39+(-39) Unsigned: Two's Complement:	b. 127+(-39) Unsigned: Two's Complement:
c. 39-127 Unsigned: Two's Complement:	d. 127+39 Unsigned: Two's Complement:

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:

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

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

1EEE 754 Float (32 bit) Flowchart

