

# Integers II, Floating Point I

CSE 351 Summer 2020

**Instructor:**

Porter Jones

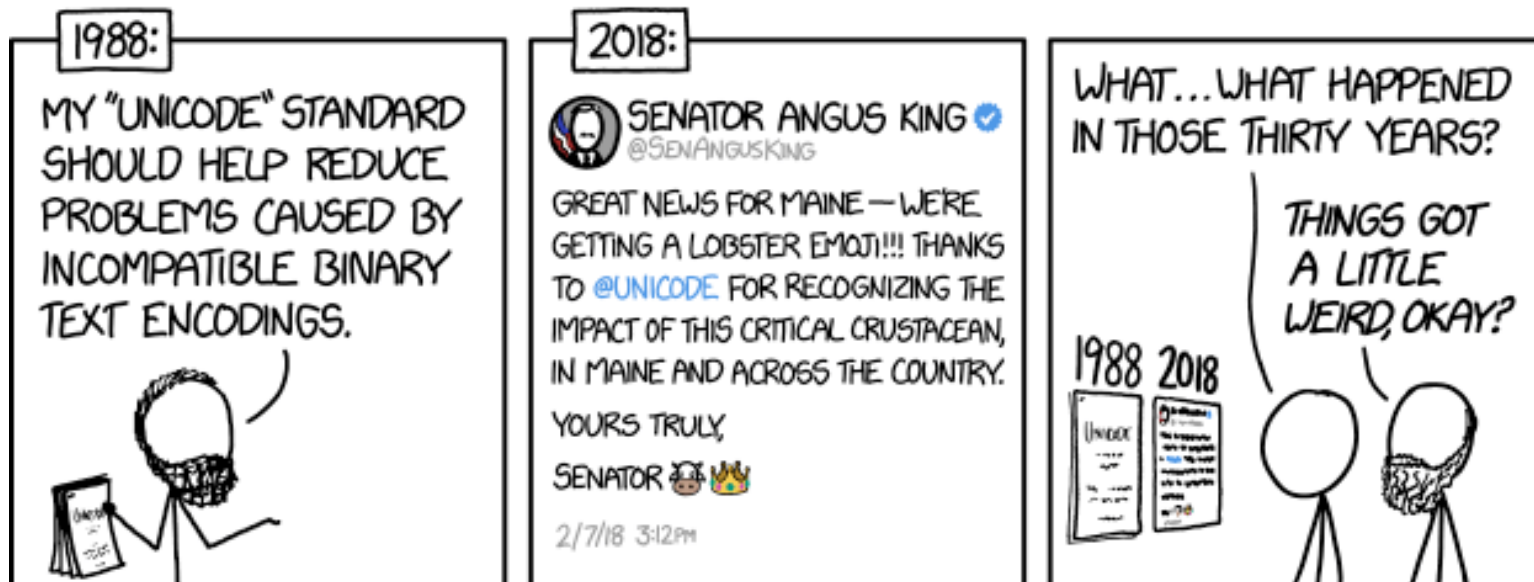
**Teaching Assistants:**

Amy Xu

Callum Walker

Sam Wolfson

Tim Mandzyuk



<http://xkcd.com/1953/>

# Administrivia

- ❖ Questions doc: <https://tinyurl.com/CSE351-7-1>
- ❖ hw4 and hw5 due Monday 7/6 – 10:30am
- ❖ hw6 and hw7 due Friday 7/10 – 10:30am
  - Will post Monday's slides later today so you can get started
- ❖ Lab 1a due Monday (7/6) (**try to finish by Friday!**)
  - Submit `pointer.c` and `lab1Areflect.txt` to Gradescope
- ❖ Lab 1b released tomorrow, due 7/10
  - Bit manipulation problems using custom data type
  - Today's bonus slides have helpful examples, tomorrow's section will have helpful examples too

# Gradescope Lab Turnin

- ❖ Make sure you pass the File and Compilation Check!
- ❖ Doesn't indicate if you passed all tests, just indicates that all the correct files were found and there were no compilation or runtime errors.
- ❖ Use the testing programs we provide to check your solution for correctness (on attu or the VM)

# Quick Aside: C Macros

- ❖ Lab1b will have you use some C macros for bit masks
- ❖ Syntax is of the form:  

```
#define NAME expression
```
- ❖ Can now use “NAME” instead of “expression” in code
- ❖ Useful to help with readability/factoring in code
  - Especially useful for defining constants such as bit masks!
- ❖ Are NOT exactly the same as a constant in Java
  - Does naïve copy and replace *before* compilation.
  - Everywhere the characters “NAME” appear in the code, the characters “expression” will now appear instead.
- ❖ See Lecture 4 (Integers I) slides for example usages

# Integers

- ❖ Binary representation of integers
  - Unsigned and signed
- ❖ Shifting and **arithmetic operations** – useful for Lab 1a
- ❖ In C: Signed, Unsigned and Casting
- ❖ Consequences of finite width representations
  - Overflow, sign extension

# Two's Complement Arithmetic

- ❖ The same addition procedure works for both unsigned and two's complement integers
  - **Simplifies hardware:** only one algorithm for addition
  - **Algorithm:** simple addition, **discard the highest carry bit**
    - Called modular addition: result is sum *modulo*  $2^w$

## ❖ 4-bit Examples:

HW	TC
0100	
+0011	
=	

HW	TC
1100	
+0011	
=	

HW	TC
0100	
+1101	
=	

# Why Does Two's Complement Work?

- ❖ For all representable positive integers  $x$ , we want:

$$\frac{\textit{bit representation of } x \\ + \textit{ bit representation of } -x}{0} \quad (\text{ignoring the carry-out bit})$$

- What are the 8-bit negative encodings for the following?

$$\begin{array}{r} 00000001 \\ + \text{????????} \\ \hline 00000000 \end{array}$$

$$\begin{array}{r} 00000010 \\ + \text{????????} \\ \hline 00000000 \end{array}$$

$$\begin{array}{r} 11000011 \\ + \text{????????} \\ \hline 00000000 \end{array}$$

# Why Does Two's Complement Work?

- ❖ For all representable positive integers  $x$ , we want:

$$\frac{\text{bit representation of } x \\ + \text{bit representation of } -x}{0} \quad (\text{ignoring the carry-out bit})$$

- What are the 8-bit negative encodings for the following?

$$\begin{array}{r} 00000001 \\ + 11111111 \\ \hline 100000000 \end{array}$$

$$\begin{array}{r} 00000010 \\ + 11111110 \\ \hline 100000000 \end{array}$$

$$\begin{array}{r} 11000011 \\ + 00111101 \\ \hline 100000000 \end{array}$$

These are the bitwise complement plus 1!

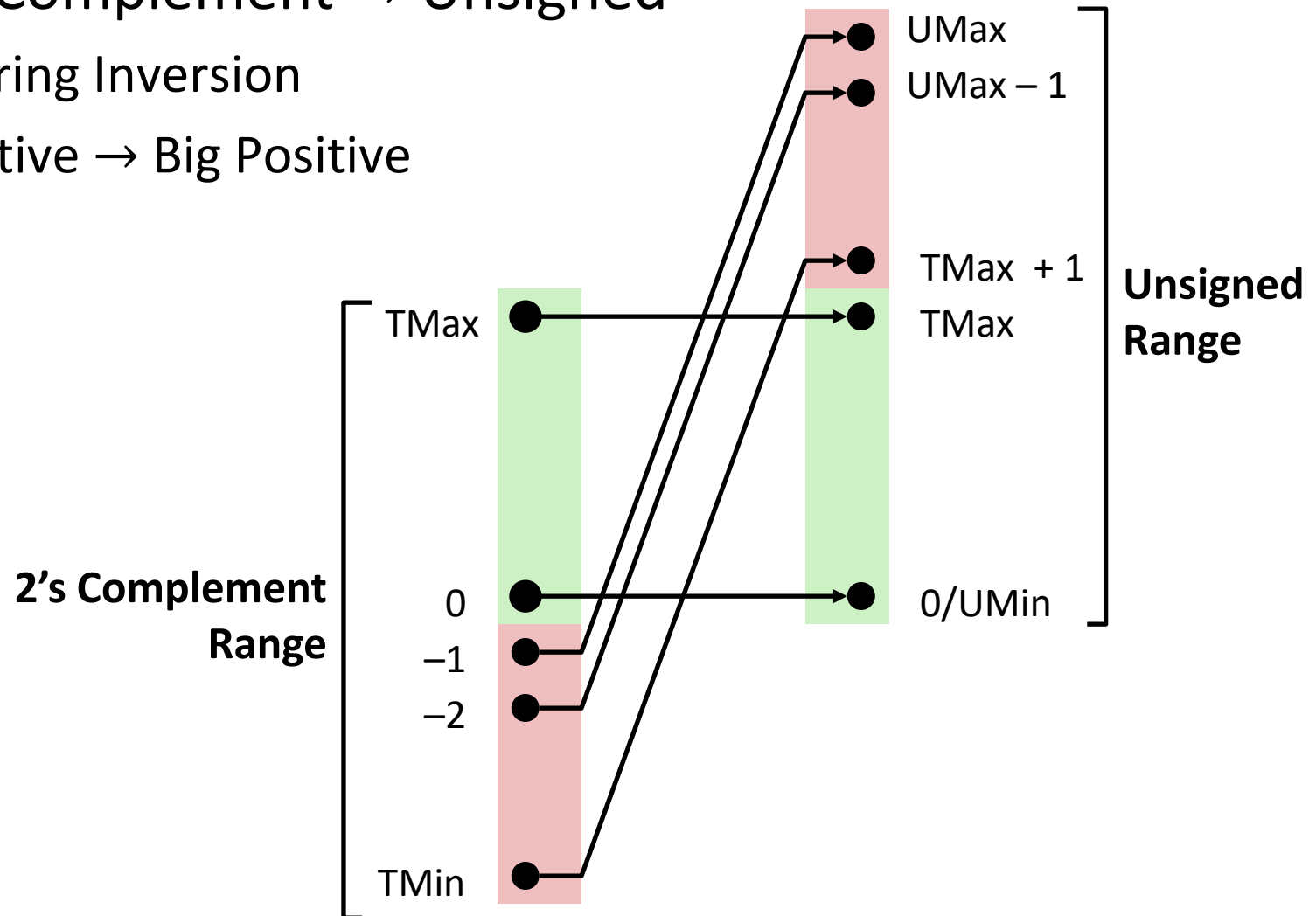
$$-x == \sim x + 1$$



# Signed/Unsigned Conversion Visualized

## ❖ Two's Complement $\rightarrow$ Unsigned

- Ordering Inversion
- Negative  $\rightarrow$  Big Positive



# Values To Remember

## ❖ Unsigned Values

- UMin = 0b00...0  
= 0
- UMax = 0b11...1  
=  $2^w - 1$

## ❖ Two's Complement Values

- TMin = 0b10...0  
=  $-2^{w-1}$
- TMax = 0b01...1  
=  $2^{w-1} - 1$
- -1 = 0b11...1

## ❖ Example: Values for $w = 64$

	Decimal	Hex
UMax	18,446,744,073,709,551,615	FF FF FF FF FF FF FF FF
TMax	9,223,372,036,854,775,807	7F FF FF FF FF FF FF FF
TMin	-9,223,372,036,854,775,808	80 00 00 00 00 00 00 00
-1	-1	FF FF FF FF FF FF FF FF
0	0	00 00 00 00 00 00 00 00

# Integers

- ❖ Binary representation of integers
  - Unsigned and signed
- ❖ Shifting and arithmetic operations – useful for Lab 1a
- ❖ **In C: Signed, Unsigned and Casting**
- ❖ Consequences of finite width representations
  - Overflow, sign extension

# In C: Signed vs. Unsigned

## ❖ Casting

- Bits are unchanged, just interpreted differently!
  - `int tx, ty;`
  - `unsigned int ux, uy;`
- *Explicit* casting
  - `tx = (int) ux;`
  - `uy = (unsigned int) ty;`
- *Implicit* casting can occur during assignments or function calls
  - `tx = ux;`
  - `uy = ty;`



# Casting Surprises

## ❖ Integer literals (constants)

- By default, integer constants are considered *signed* integers
  - Hex constants already have an explicit binary representation
- Use “U” (or “u”) suffix to explicitly force *unsigned*
  - Examples: 0U, 4294967259u

## ❖ Expression Evaluation

- When you mixed unsigned and signed in a single expression, then **signed values are implicitly cast to unsigned** (unsigned “dominates”)
- Including comparison operators  $<$ ,  $>$ ,  $==$ ,  $<=$ ,  $>=$



# Casting Surprises

❖ 32-bit examples:

- TMin = -2,147,483,648, TMax = 2,147,483,647

Left Constant	Order	Right Constant	Interpretation
<b>0</b> 0000 0000 0000 0000 0000 0000 0000 0000		<b>0U</b> 0000 0000 0000 0000 0000 0000 0000 0000	
<b>-1</b> 1111 1111 1111 1111 1111 1111 1111 1111		<b>0</b> 0000 0000 0000 0000 0000 0000 0000 0000	
<b>-1</b> 1111 1111 1111 1111 1111 1111 1111 1111		<b>0U</b> 0000 0000 0000 0000 0000 0000 0000 0000	
<b>2147483647</b> 0111 1111 1111 1111 1111 1111 1111 1111		<b>-2147483648</b> 1000 0000 0000 0000 0000 0000 0000 0000	
<b>2147483647U</b> 0111 1111 1111 1111 1111 1111 1111 1111		<b>-2147483648</b> 1000 0000 0000 0000 0000 0000 0000 0000	
<b>-1</b> 1111 1111 1111 1111 1111 1111 1111 1111		<b>-2</b> 1111 1111 1111 1111 1111 1111 1111 1110	
<b>(unsigned) -1</b> 1111 1111 1111 1111 1111 1111 1111 1111		<b>-2</b> 1111 1111 1111 1111 1111 1111 1111 1110	
<b>2147483647</b> 0111 1111 1111 1111 1111 1111 1111 1111		<b>2147483648U</b> 1000 0000 0000 0000 0000 0000 0000 0000	
<b>2147483647</b> 0111 1111 1111 1111 1111 1111 1111 1111		<b>(int) 2147483648U</b> 1000 0000 0000 0000 0000 0000 0000 0000	

# Integers

- ❖ Binary representation of integers
  - Unsigned and signed
- ❖ Shifting and arithmetic operations – useful for Lab 1a
- ❖ In C: Signed, Unsigned and Casting
- ❖ **Consequences of finite width representations**
  - **Overflow, sign extension**

# Arithmetic Overflow

Bits	Unsigned	Signed
0000	0	0
0001	1	1
0010	2	2
0011	3	3
0100	4	4
0101	5	5
0110	6	6
0111	7	7
1000	8	-8
1001	9	-7
1010	10	-6
1011	11	-5
1100	12	-4
1101	13	-3
1110	14	-2
1111	15	-1

- ❖ When a calculation produces a result that can't be represented in the current encoding scheme
  - Integer range limited by fixed width
  - Can occur in both the positive and negative directions
- ❖ C and Java ignore overflow exceptions
  - You end up with a bad value in your program and no warning/indication... oops!



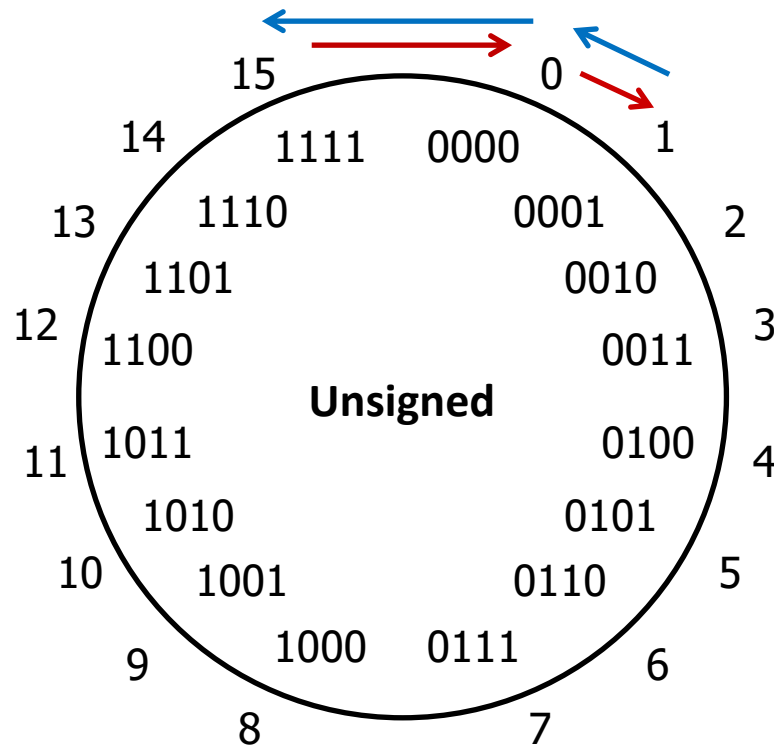
# Overflow: Unsigned

- ❖ **Addition:** drop carry bit ( $-2^N$ )

15	1111
+ 2	+ 0010
<del>17</del>	<del>10001</del>
1	

- ❖ **Subtraction:** borrow ( $+2^N$ )

1	10001
- 2	- 0010
<del>-1</del>	1111
15	



$\pm 2^N$  because of modular arithmetic

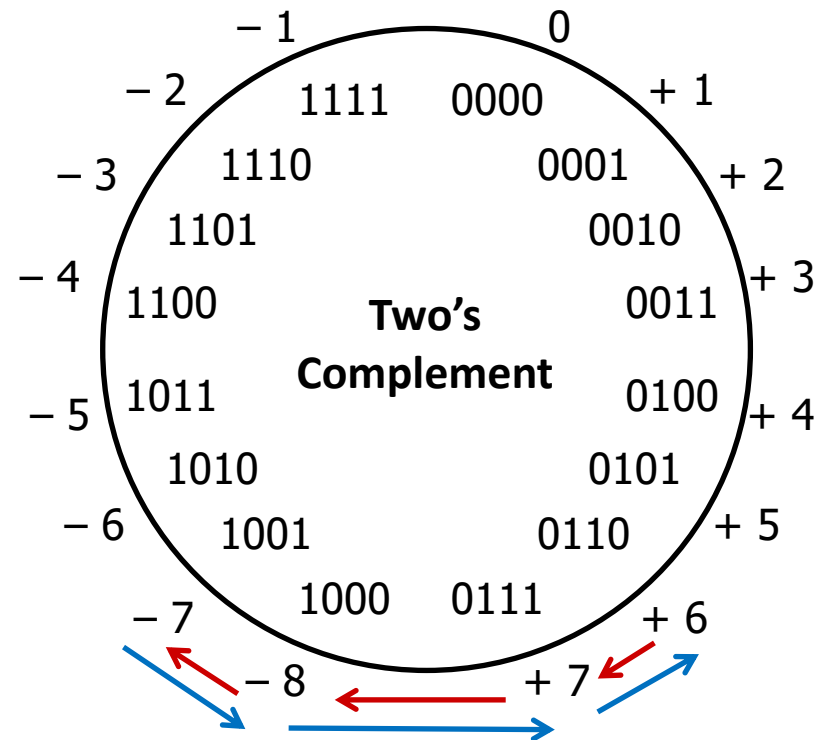
# Overflow: Two's Complement

❖ **Addition:** (+) + (+) = (-) result?

$$\begin{array}{r}
 6 \qquad 0110 \\
 + 3 \qquad + 0011 \\
 \hline
 \cancel{9} \\
 -7
 \end{array}$$

❖ **Subtraction:** (-) + (-) = (+)?

$$\begin{array}{r}
 -7 \qquad 1001 \\
 - 3 \qquad - 0011 \\
 \hline
 \cancel{-10} \\
 6
 \end{array}$$



**For signed: overflow if operands have same sign and result's sign is different**

# Sign Extension

- ❖ What happens if you convert a *signed* integral data type to a larger one?

- e.g. <sup>1 byte</sup> char → <sup>2 bytes</sup> short → <sup>4 bytes</sup> int → <sup>8 bytes</sup> long

- ❖ **4-bit → 8-bit Example:**

- Positive Case

**4-bit:**            0010    =    +2

- ✓ • Add 0's?

**8-bit:**    **00000010**    =    **+2**

- Negative Case?

# Polling Question [Int II - a]

- ❖ Which of the following 8-bit numbers has the same *signed* value as the 4-bit number **0b1100**?
  - Underlined digit = MSB
  - Vote at <http://pollev.com/pbjones>
  
- A. 0b 0000 1100
- B. 0b 1000 1100
- C. 0b 1111 1100
- D. 0b 1100 1100
- E. We're lost...

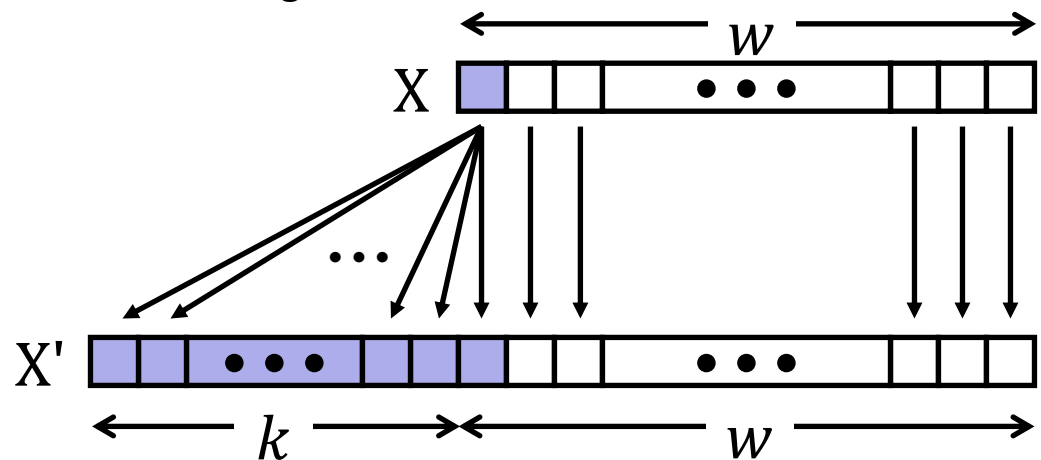
# Sign Extension

❖ **Task:** Given a  $w$ -bit signed integer  $X$ , convert it to  $w+k$ -bit signed integer  $X'$  with the same value

❖ **Rule:** Add  $k$  copies of sign bit

■ Let  $x_i$  be the  $i$ -th digit of  $X$  in binary

$$X' = \underbrace{x_{w-1}, \dots, x_{w-1}}_{k \text{ copies of MSB}}, \underbrace{x_{w-1}, x_{w-2}, \dots, x_1, x_0}_{\text{original } X}$$



# Sign Extension Example

- ❖ Convert from smaller to larger integral data types
- ❖ C automatically performs sign extension
  - Java too

```
short int x = 12345;
int      ix = (int) x;
short int y = -12345;
int      iy = (int) y;
```

Var	Decimal	Hex	Binary
x	12345	30 39	00110000 00111001
ix	12345	00 00 30 39	00000000 00000000 00110000 00111001
y	-12345	CF C7	11001111 11000111
iy	-12345	FF FF CF C7	11111111 11111111 11001111 11000111

Handwritten annotations: *0b 0011* (pointing to the last two hex digits of the Hex column) and *0b 1100* (pointing to the last two hex digits of the Hex column).

# Practice Question

For the following expressions, find a value of **signed char**  $x$ , if there exists one, that makes the expression **TRUE**. Compare with your neighbor(s)!

❖ Assume we are using 8-bit arithmetic:

■  $x == (\text{unsigned char}) x$

Example:

All solutions:

■  $x \geq 128U$

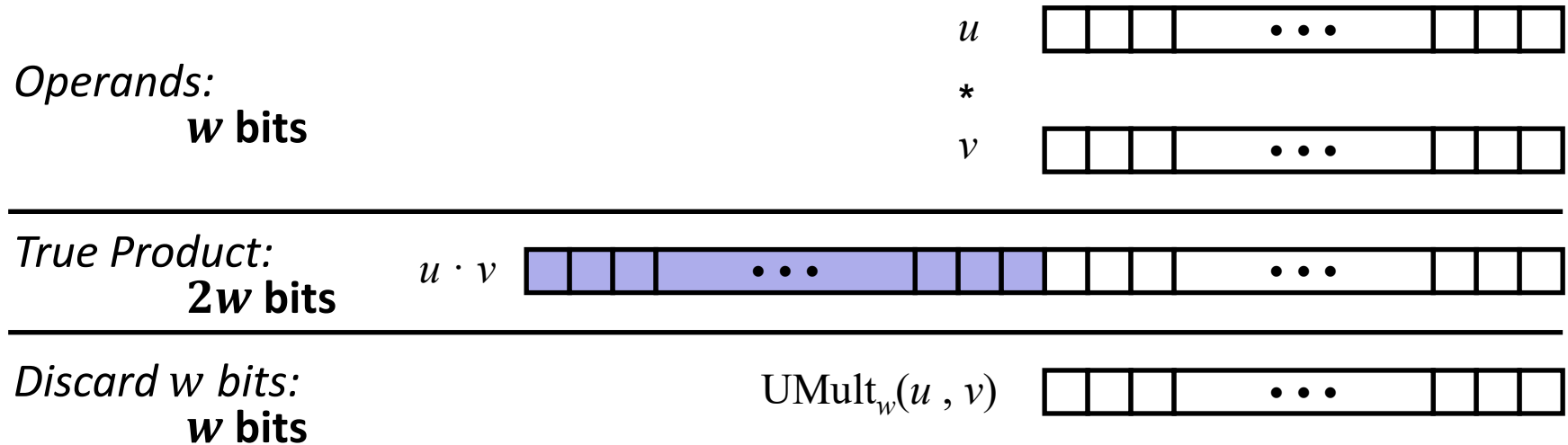
■  $x \neq (x \gg 2) \ll 2$

■  $x == -x$

- Hint: there are two solutions

■  $(x < 128U) \ \&\& \ (x > 0x3F)$

# Aside: Unsigned Multiplication in C



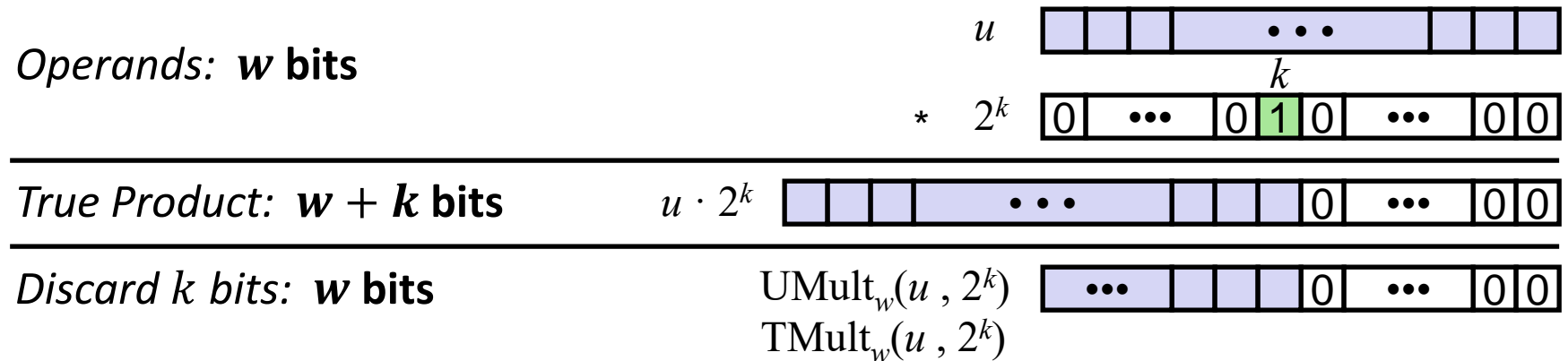
- ❖ Standard Multiplication Function
  - Ignores high order  $w$  bits
- ❖ Implements Modular Arithmetic
  - $\text{UMult}_w(u, v) = u \cdot v \pmod{2^w}$



# Aside: Multiplication with shift and add

❖ Operation  $u \ll k$  gives  $u * 2^k$

- Both signed and unsigned



❖ Examples:

- $u \ll 3 \quad == \quad u * 8$
- $u \ll 5 - u \ll 3 \quad == \quad u * 24$
- Most machines shift and add faster than multiply
  - Compiler generates this code automatically*

# Number Representation Revisited

- ❖ What can we represent so far?
  - Signed and Unsigned Integers
  - Characters (ASCII)
  - Addresses
  
- ❖ How do we encode the following:
  - Real numbers (*e.g.* 3.14159)
  - Very large numbers (*e.g.*  $6.02 \times 10^{23}$ )
  - Very small numbers (*e.g.*  $6.626 \times 10^{-34}$ )
  - Special numbers (*e.g.*  $\infty$ , NaN)



**Floating  
Point**

# Floating Point Topics

- ❖ **Fractional binary numbers**
- ❖ IEEE floating-point standard
- ❖ Floating-point operations and rounding
- ❖ Floating-point in C

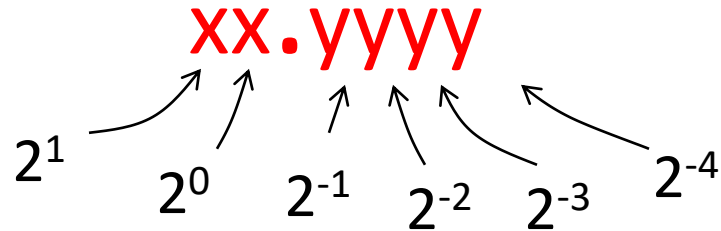


- ❖ There are many more details that we won't cover
  - It's a 58-page standard...

# Representation of Fractions

- ❖ “Binary Point,” like decimal point, signifies boundary between integer and fractional parts:

Example 6-bit  
representation:

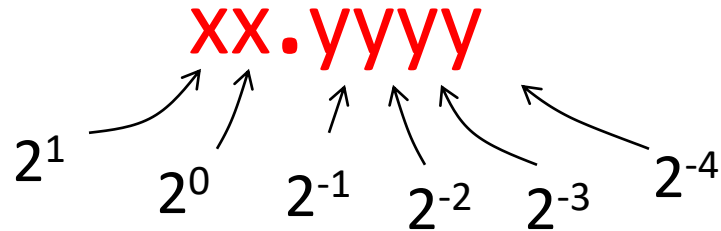


- ❖ Example:  $10.1010_2 = 1 \times 2^1 + 1 \times 2^{-1} + 1 \times 2^{-3} = 2.625_{10}$

# Representation of Fractions

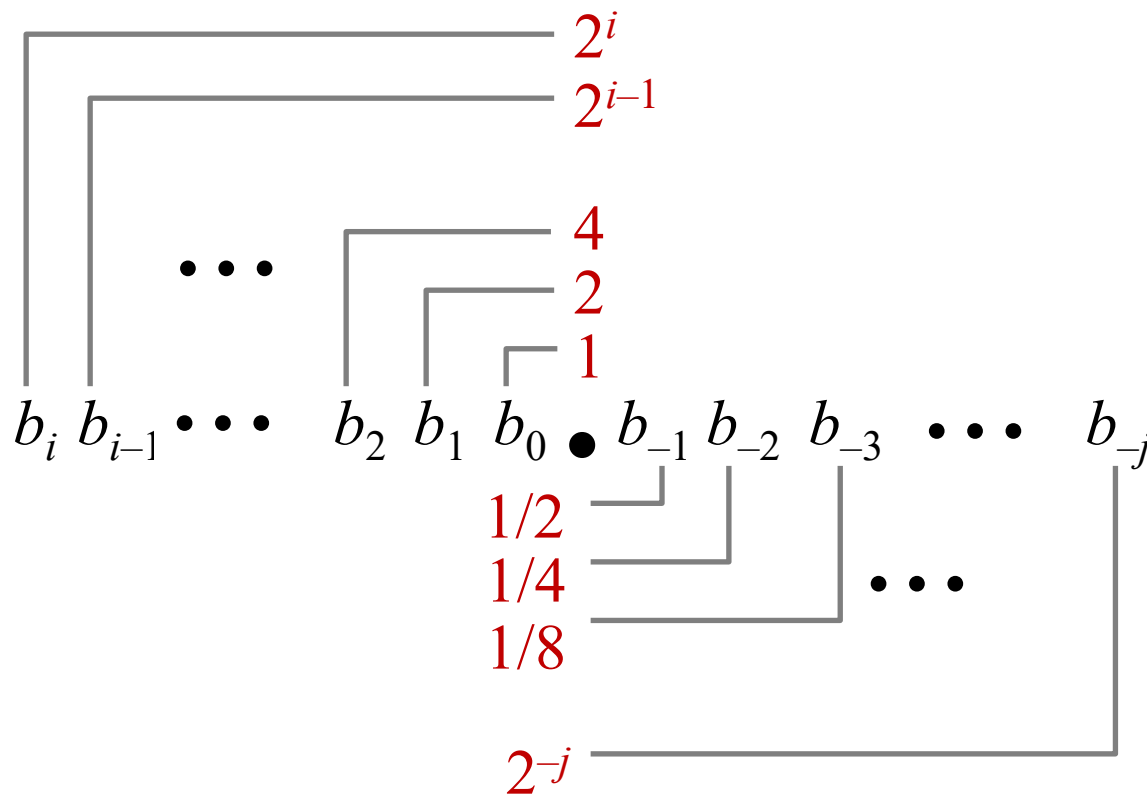
- ❖ “Binary Point,” like decimal point, signifies boundary between integer and fractional parts:

Example 6-bit  
representation:



- ❖ In this 6-bit representation:
  - What is the encoding and value of the smallest (most negative) number?
  - What is the encoding and value of the largest (most positive) number?
  - What is the smallest number greater than 2 that we can represent?

# Fractional Binary Numbers



## ❖ Representation

- Bits to right of “binary point” represent fractional powers of 2
- Represents rational number:

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

# Fractional Binary Numbers

- ❖ Value                      Representation
  - 5 and 3/4               $101.11_2$
  - 2 and 7/8               $10.111_2$
  - 47/64                     $0.101111_2$
  
- ❖ Observations
  - Shift left = multiply by power of 2
  - Shift right = divide by power of 2
  - Numbers of the form  $0.111111\dots_2$  are just below 1.0
    - $1/2 + 1/4 + 1/8 + \dots + 1/2^i + \dots \rightarrow 1.0$
    - Use notation  $1.0 - \epsilon$

# Limits of Representation

## ❖ Limitations:

- Even given an arbitrary number of bits, can only **exactly** represent numbers of the form  $x * 2^y$  ( $y$  can be negative)
- Other rational numbers have repeating bit representations

### Value:

### Binary Representation:

- $1/3 = 0.333333..._{10} = 0.01010101[01]..._2$
- $1/5 = 0.001100110011[0011]..._2$
- $1/10 = 0.0001100110011[0011]..._2$



# Fixed Point Representation

- ❖ Implied binary point. Two example schemes:

#1: the binary point is between bits 2 and 3

$b_7 b_6 b_5 b_4 b_3 \text{ [.] } b_2 b_1 b_0$

#2: the binary point is between bits 4 and 5

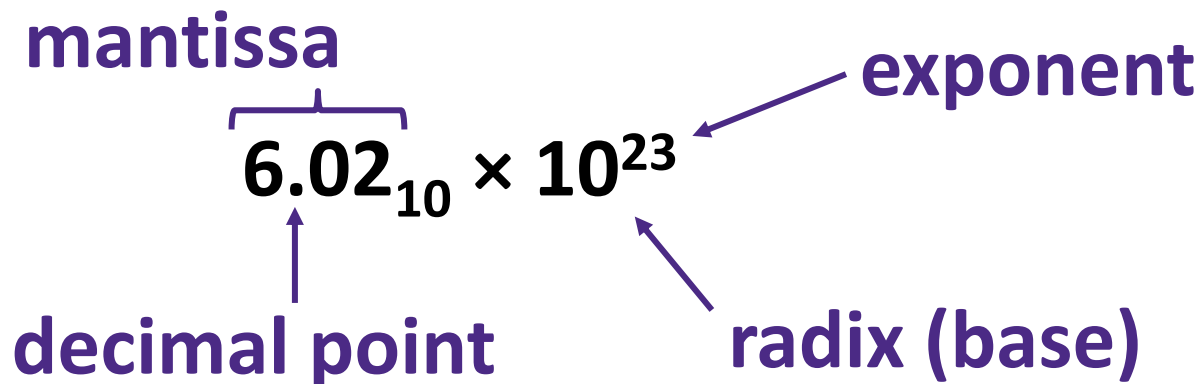
$b_7 b_6 b_5 \text{ [.] } b_4 b_3 b_2 b_1 b_0$

- ❖ Wherever we put the binary point, with fixed point representations there is a trade off between the amount of range and precision we have
- ❖ Fixed point = fixed *range* and fixed *precision*
  - range: difference between largest and smallest numbers possible
  - precision: smallest possible difference between any two numbers
- ❖ Hard to pick how much you need of each!

# Floating Point Representation

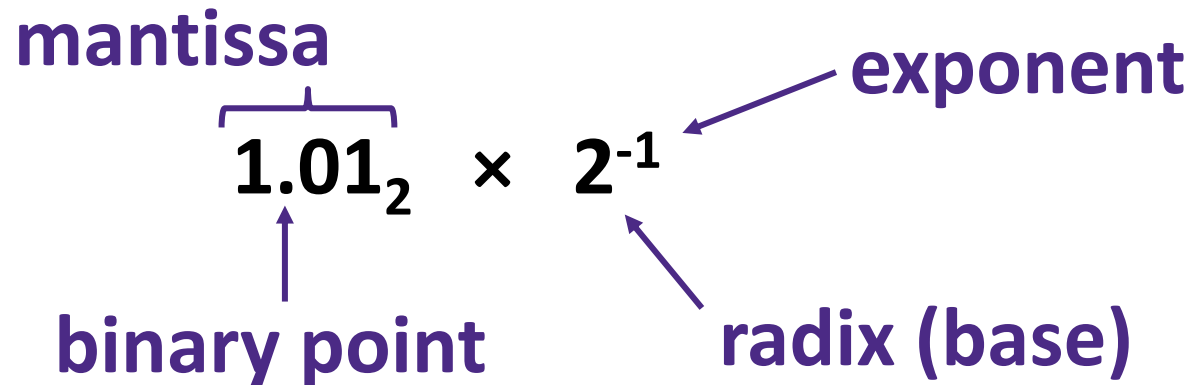
- ❖ Analogous to scientific notation
  - In Decimal:
    - Not 12000000, but  $1.2 \times 10^7$  In C: 1.2e7
    - Not 0.0000012, but  $1.2 \times 10^{-6}$  In C: 1.2e-6
  - In Binary:
    - Not 11000.000, but  $1.1 \times 2^4$
    - Not 0.000101, but  $1.01 \times 2^{-4}$
- ❖ We have to divvy up the bits we have (e.g., 32) among:
  - the sign (1 bit)
  - the mantissa (significand)
  - the exponent

# Scientific Notation (Decimal)



- ❖ *Normalized form*: exactly one digit (non-zero) to left of decimal point
- ❖ Alternatives to representing  $1/1,000,000,000$ 
  - **Normalized:**  $1.0 \times 10^{-9}$
  - Not normalized:  $0.1 \times 10^{-8}, 10.0 \times 10^{-10}$

# Scientific Notation (Binary)



- ❖ Computer arithmetic that supports this called **floating point** due to the “floating” of the binary point
  - Declare such variable in C as `float` (or `double`)

# Scientific Notation Translation

$$\begin{aligned}2^{-1} &= 0.5 \\2^{-2} &= 0.25 \\2^{-3} &= 0.125 \\2^{-4} &= 0.0625\end{aligned}$$

- ❖ Convert from scientific notation to binary point
  - Perform the multiplication by shifting the decimal until the exponent disappears
    - Example:  $1.011_2 \times 2^4 = 10110_2 = 22_{10}$
    - Example:  $1.011_2 \times 2^{-2} = 0.01011_2 = 0.34375_{10}$
  
- ❖ Convert from binary point to *normalized* scientific notation
  - Distribute out exponents until binary point is to the right of a single digit
    - Example:  $1101.001_2 = 1.101001_2 \times 2^3$
  
- ❖ **Practice:** Convert  $11.375_{10}$  to normalized binary scientific notation

# Summary

- ❖ Sign and unsigned variables in C
  - Bit pattern remains the same, just *interpreted* differently
  - Strange things can happen with our arithmetic when we convert/cast between sign and unsigned numbers
    - Type of variables affects behavior of operators (shifting, comparison)
- ❖ We can only represent so many numbers in  $w$  bits
  - When we exceed the limits, *arithmetic overflow* occurs
  - *Sign extension* tries to preserve value when expanding
- ❖ Floating point approximates real numbers
  - We will discuss more details on Monday!

# BONUS SLIDES

Some examples of using shift operators in combination with bitmasks, which you may find helpful for Lab 1.

- ❖ Extract the 2<sup>nd</sup> most significant byte of an `int`
- ❖ Extract the sign bit of a signed `int`
- ❖ Conditionals as Boolean expressions

# Using Shifts and Masks

- ❖ Extract the 2<sup>nd</sup> most significant *byte* of an `int`:
  - First shift, then mask:  $(x \gg 16) \ \& \ 0xFF$

<b>x</b>	00000001	00000010	00000011	00000100
<b>x &gt;&gt; 16</b>	00000000	00000000	00000001	00000010
<b>0xFF</b>	00000000	00000000	00000000	11111111
<b>(x &gt;&gt; 16) &amp; 0xFF</b>	00000000	00000000	00000000	00000010

- Or first mask, then shift:  $(x \ \& \ 0xFF0000) \gg 16$

<b>x</b>	00000001	00000010	00000011	00000100
<b>0xFF0000</b>	00000000	11111111	00000000	00000000
<b>x &amp; 0xFF0000</b>	00000000	00000010	00000000	00000000
<b>(x &amp; 0xFF0000) &gt;&gt; 16</b>	00000000	00000000	00000000	00000010



# Using Shifts and Masks

- ❖ Extract the *sign bit* of a signed `int`:
  - First shift, then mask:  $(x \gg 31) \ \& \ 0x1$ 
    - Assuming arithmetic shift here, but this works in either case
    - Need mask to clear 1s possibly shifted in

<b>x</b>	<b>0</b> 0000001 00000010 00000011 00000100
<b>x&gt;&gt;31</b>	00000000 00000000 00000000 0000000 <b>0</b>
<b>0x1</b>	00000000 00000000 00000000 00000001
<b>(x&gt;&gt;31) &amp; 0x1</b>	00000000 00000000 00000000 00000000

<b>x</b>	<b>1</b> 0000001 00000010 00000011 00000100
<b>x&gt;&gt;31</b>	11111111 11111111 11111111 1111111 <b>1</b>
<b>0x1</b>	00000000 00000000 00000000 00000001
<b>(x&gt;&gt;31) &amp; 0x1</b>	00000000 00000000 00000000 00000001

# Using Shifts and Masks

## ❖ Conditionals as Boolean expressions

- For `int x`, what does `(x<<31)>>31` do?

<code>x=!!123</code>	00000000 00000000 00000000 00000000 <b>1</b>
<code>x&lt;&lt;31</code>	<b>1</b> 00000000 00000000 00000000 00000000
<code>(x&lt;&lt;31)&gt;&gt;31</code>	<b>11111111 11111111 11111111 11111111</b>
<code>!x</code>	00000000 00000000 00000000 00000000 <b>0</b>
<code>!x&lt;&lt;31</code>	<b>0</b> 00000000 00000000 00000000 00000000
<code>(!x&lt;&lt;31)&gt;&gt;31</code>	<b>00000000 00000000 00000000 00000000</b>

- Can use in place of conditional:

- In C: `if (x) {a=y;} else {a=z;} equivalent to a=x?y:z;`
- `a = ((x<<31)>>31) & y | ((!x<<31)>>31) & z;`