

# Integer Representation

Representation of integers: unsigned and signed

Modular arithmetic and overflow

Sign extension

Shifting and arithmetic

Multiplication

Casting

# Fixed-width integer encodings

**Unsigned**  $\subset \mathbb{N}$  non-negative integers only

**Signed**  $\subset \mathbb{Z}$  both negative and non-negative integers

$n$  bits offer only  $2^n$  distinct values.

## Terminology:

“Most-significant” bit(s)  
or “high-order” bit(s)

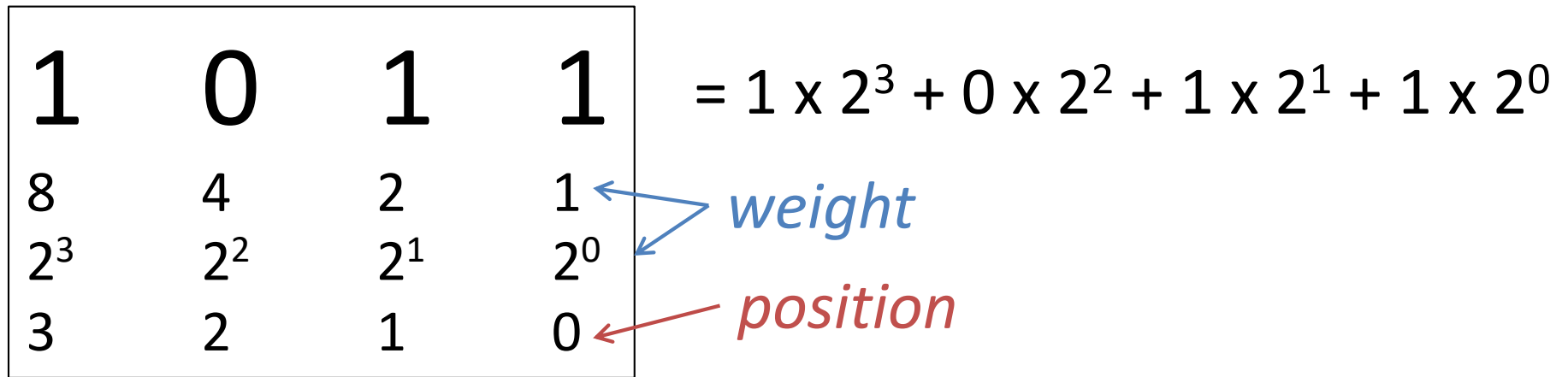
“Least-significant” bit(s)  
or “low-order” bit(s)

**MSB**

0110010110101001

**LSB**

# (4-bit) unsigned integer representation



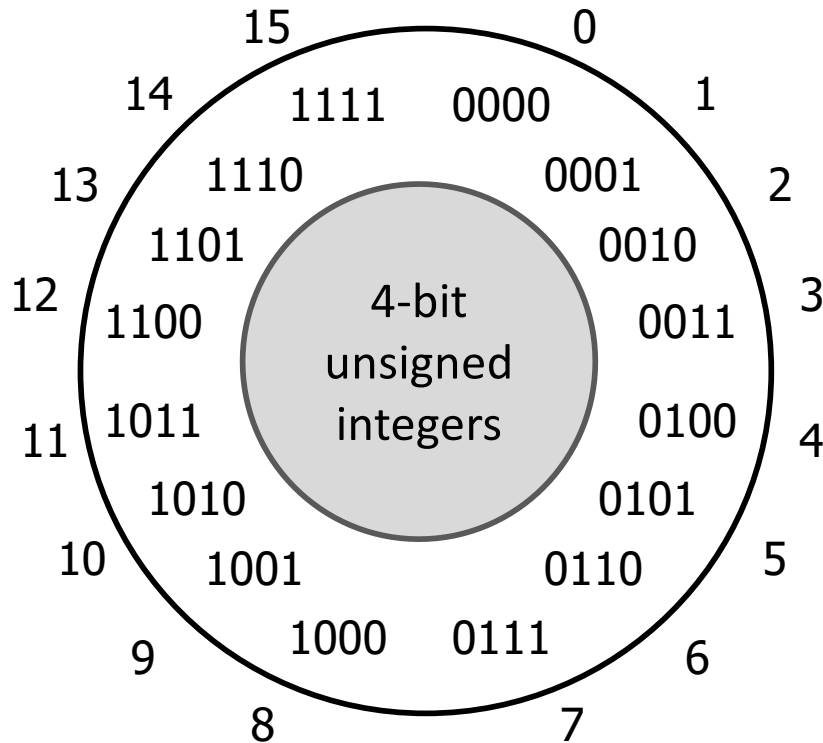
***n*-bit unsigned integers:**

**minimum =**

**maximum =**

# modular arithmetic, overflow

$$\begin{array}{r} 11 \quad 1011 \\ + 2 \quad + 0010 \\ \hline \end{array}$$



$$\begin{array}{r} 13 \quad 1101 \\ + 5 \quad + 0101 \\ \hline \end{array}$$

$x+y$  in  $n$ -bit unsigned arithmetic is

in math

*unsigned overflow* =  
=

Unsigned addition *overflows* if and only if

# sign-magnitude



Most-significant bit (MSB) is *sign bit*

0 means non-negative

1 means negative

Remaining bits are an unsigned magnitude

8-bit sign-magnitude:

00000000 represents \_\_\_\_\_

01111111 represents \_\_\_\_\_

10000101 represents \_\_\_\_\_

10000000 represents \_\_\_\_\_

Anything weird here?

## Arithmetic?

Example:

$$4 - 3 \neq 4 + (-3)$$



$$\begin{array}{r} 00000100 \\ +10000011 \\ \hline \end{array}$$

## Zero?



# (4-bit) two's complement signed integer representation



|        |       |       |       |
|--------|-------|-------|-------|
| 1      | 0     | 1     | 1     |
| $-2^3$ | $2^2$ | $2^1$ | $2^0$ |

$$= 1 \times -2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$

**4-bit two's complement integers:**

minimum =

maximum =

# two's complement vs. unsigned

|            |           |     |       |       |       |
|------------|-----------|-----|-------|-------|-------|
| —          | —         | ... | —     | —     | —     |
| $2^{n-1}$  | $2^{n-2}$ | ... | $2^2$ | $2^1$ | $2^0$ |
| $-2^{n-1}$ | $2^{n-2}$ | ... | $2^2$ | $2^1$ | $2^0$ |

*unsigned places*

*two's complement places*

What's the difference?

**n-bit unsigned numbers:**

minimum =

maximum =

# 8-bit representations



0 0 0 0 1 0 0 1

1 0 0 0 0 0 0 1

1 1 1 1 1 1 1 1

0 0 1 0 0 1 1 1

**n-bit two's complement numbers:**

minimum =

maximum =



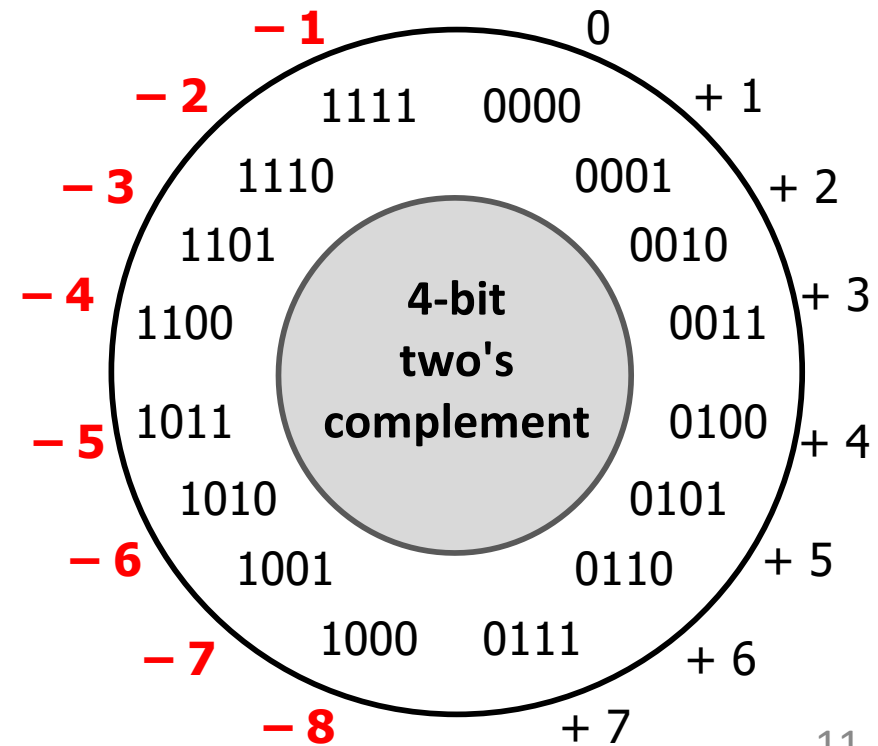
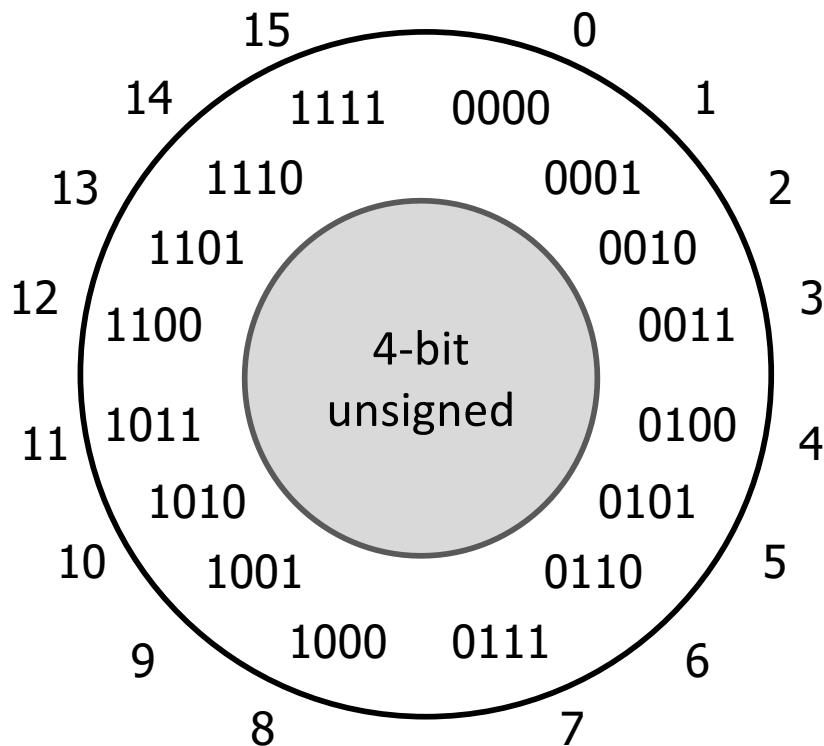
# 4-bit unsigned vs. 4-bit two's complement

1 0 1 1

$$1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$

$$1 \times -2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$

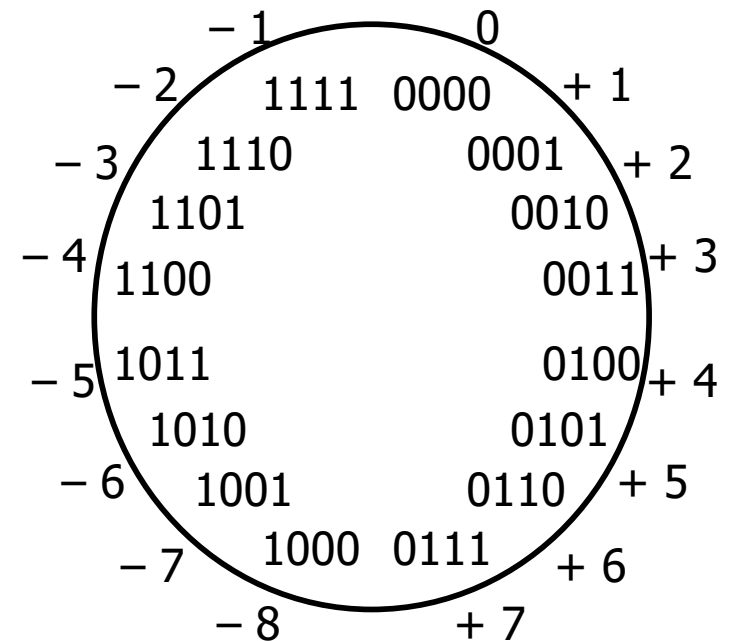
11 ← difference = \_\_\_ = 2 — → -5



# two's complement addition

|            |               |             |               |
|------------|---------------|-------------|---------------|
| 2          | 0010          | -2          | 1110          |
| <u>+ 3</u> | <u>+ 0011</u> | <u>+ -3</u> | <u>+ 1101</u> |

|            |               |             |               |
|------------|---------------|-------------|---------------|
| -2         | 1110          | 2           | 0010          |
| <u>+ 3</u> | <u>+ 0011</u> | <u>+ -3</u> | <u>+ 1101</u> |



## Modular Arithmetic

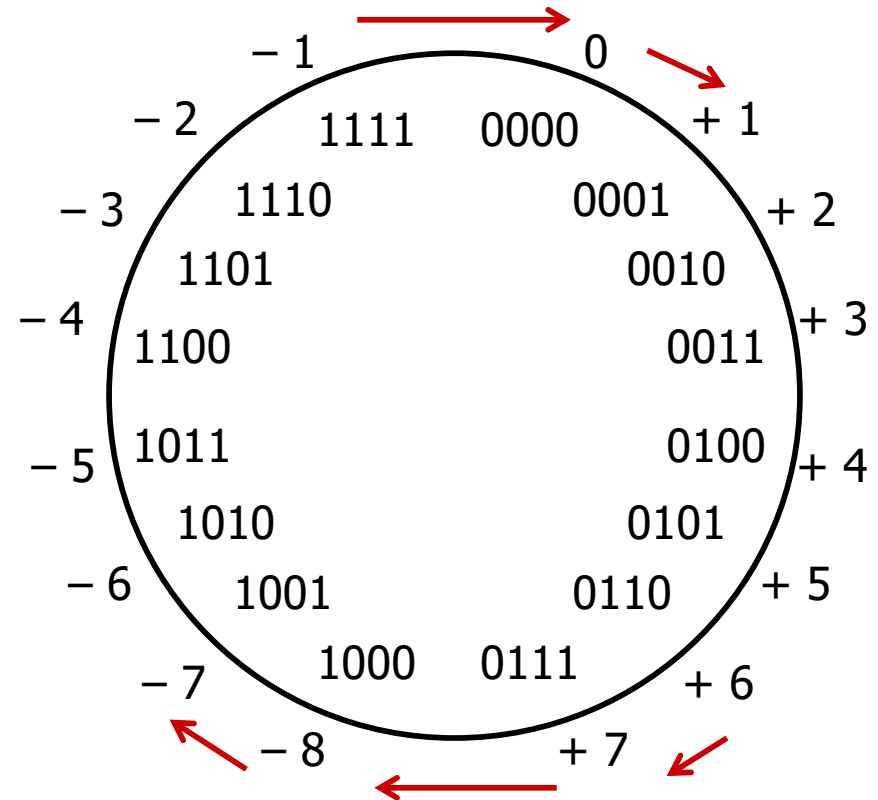
# two's complement *overflow*

## Addition *overflows*

if and only if  
if and only if

$$\begin{array}{r} -1 \quad 1111 \\ + 2 \quad + 0010 \\ \hline \end{array}$$

$$\begin{array}{r} 6 \quad 0110 \\ + 3 \quad + 0011 \\ \hline \end{array}$$



## Modular Arithmetic

Some CPUs/languages raise exceptions on overflow.  
C and Java cruise along silently... Feature? Oops? 13

# Reliability

## Ariane 5 Rocket, 1996

Exploded due to **cast** of 64-bit floating-point number to 16-bit signed number.  
**Overflow.**



## Boeing 787, 2015



"... a **Model 787 airplane** ... can lose all alternating current (AC) electrical power ... caused by a **software counter** internal to the GCUs that will **overflow** after **248 days** of continuous power. We are issuing this AD to prevent loss of all AC electrical power, which could result in **loss of control of the airplane.**"  
--FAA, April 2015

# A few reasons two's complement is awesome

Addition, subtraction, hardware

Sign

Negative one

Complement rules



# Another derivation

## How should we represent 8-bit negatives?

- For all positive integers  $x$ , we want the representations of  $x$  and  $-x$  to sum to zero.
- We want to use the standard addition algorithm.

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

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

$$\begin{array}{r} 00000011 \\ + \\ \hline 00000000 \end{array}$$

- Find a rule to represent  $-x$  where that works...

***Convert/cast signed number to larger type.***

0 0 0 0 0 0 1 0      8-bit 2

----- 0 0 0 0 0 0 1 0      16-bit 2

1 1 1 1 1 1 0 0      8-bit -4

----- 1 1 1 1 1 1 0 0      16-bit -4

Rule/name?

# unsigned shifting and arithmetic

**unsigned**

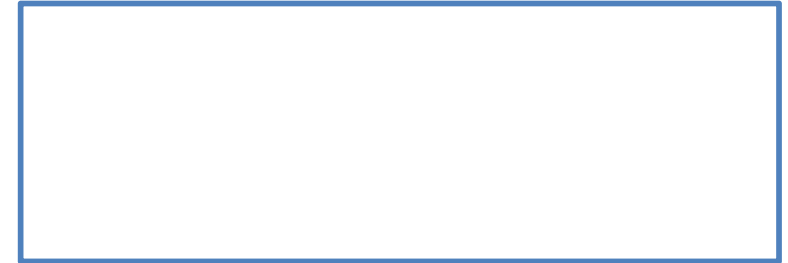
$x = 27;$

$y = x \ll 2;$

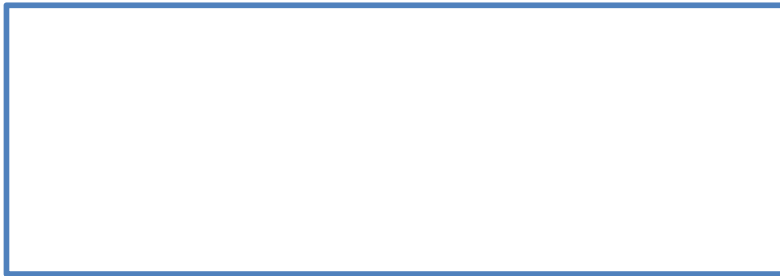
$y == 108$

0 0 0 1 1 0 1 1

0 0 0 1 1 0 1 1 0 0



logical shift left



logical shift right

1 1 1 0 1 1 0 1

0 0 1 1 1 0 1 1 0 1

**unsigned**

$x = 237;$

$y = x \gg 2;$

$y == 59$



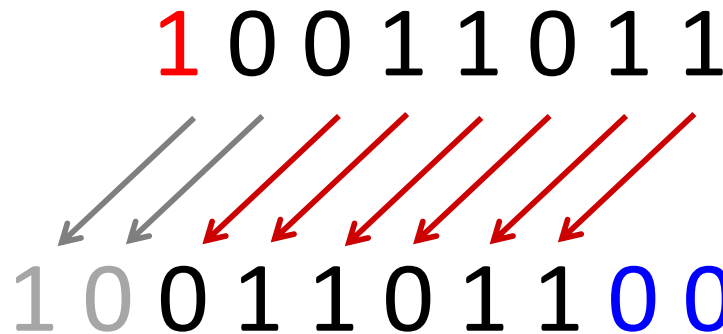
# two's complement **shifting** and **arithmetic**

**signed**

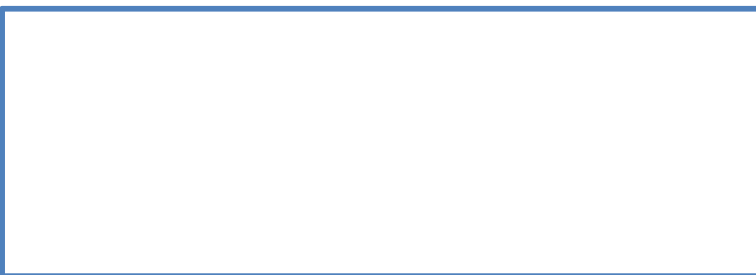
x = -101;

y = x << 2;

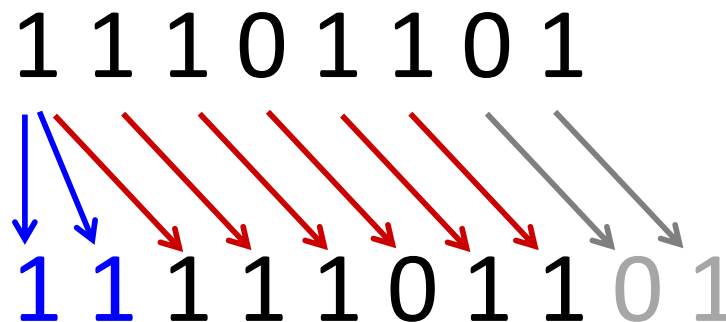
y == 108



logical shift left



arithmetic shift right



**signed**

x = -19;

y = x >> 2;

y == -5

# *shift-and-add*



## Available operations

$x \ll k$

implements  $x * 2^k$

$x + y$

Implement  $y = x * 24$  using only  $\ll$ ,  $+$ , and integer literals

# What does this function compute?

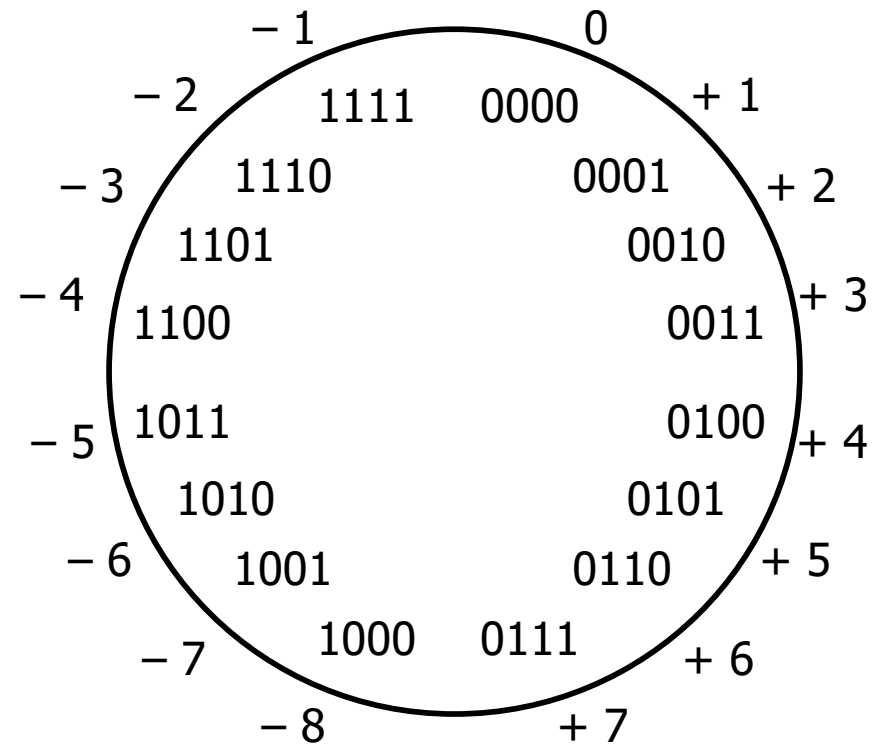


```
unsigned puzzle(unsigned x, unsigned y) {
    unsigned result = 0;
    for (unsigned i = 0; i < 32; i++) {
        if (y & (1 << i)) {
            result = result + (x << i);
        }
    }
    return result;
}
```

# multiplication

$$\begin{array}{r} 2 \\ \times 3 \\ \hline 6 \end{array} \quad \begin{array}{r} 0010 \\ \times 0011 \\ \hline 00000110 \end{array}$$

$$\begin{array}{r} -2 \\ \times 2 \\ \hline -4 \end{array} \quad \begin{array}{r} 1110 \\ \times 0010 \\ \hline 11111100 \end{array}$$



Modular Arithmetic

# multiplication

$$\begin{array}{r} 5 \\ \times 4 \\ \hline \end{array}$$

~~20~~

4

-3

$$\begin{array}{r} -3 \\ \times 7 \\ \hline \end{array}$$

~~-21~~

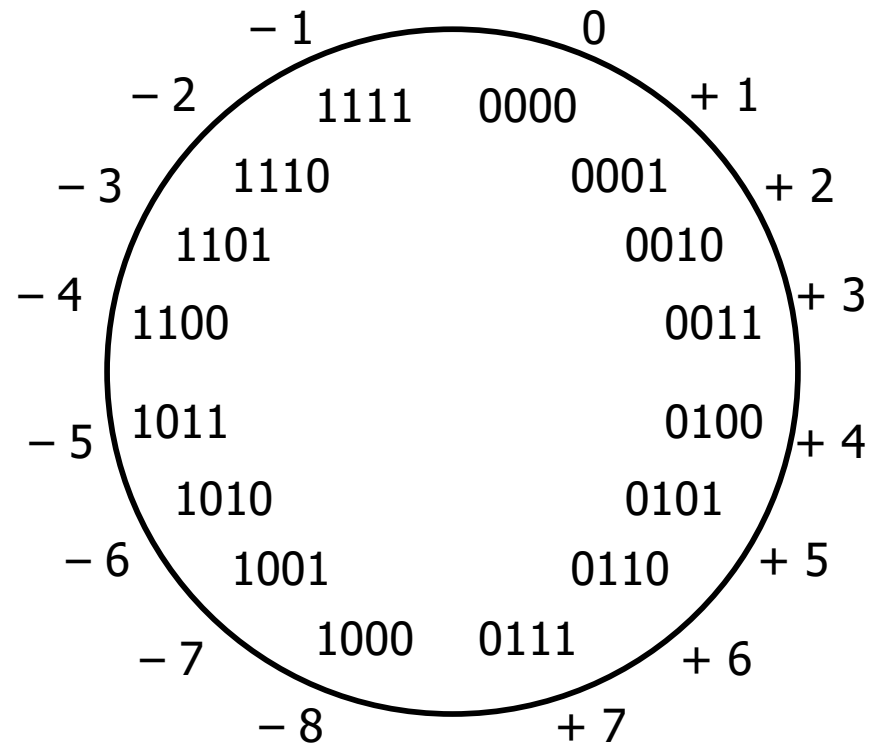
-5

$$\begin{array}{r} 0101 \\ \times 0100 \\ \hline \end{array}$$

00010100

$$\begin{array}{r} 1101 \\ \times 0111 \\ \hline \end{array}$$

11101011



Modular Arithmetic

# multiplication

$$\begin{array}{r} 5 \\ \times 5 \\ \hline \end{array}$$
$$\begin{array}{r} 0101 \\ \times 0101 \\ \hline \end{array}$$

~~25~~  
-7

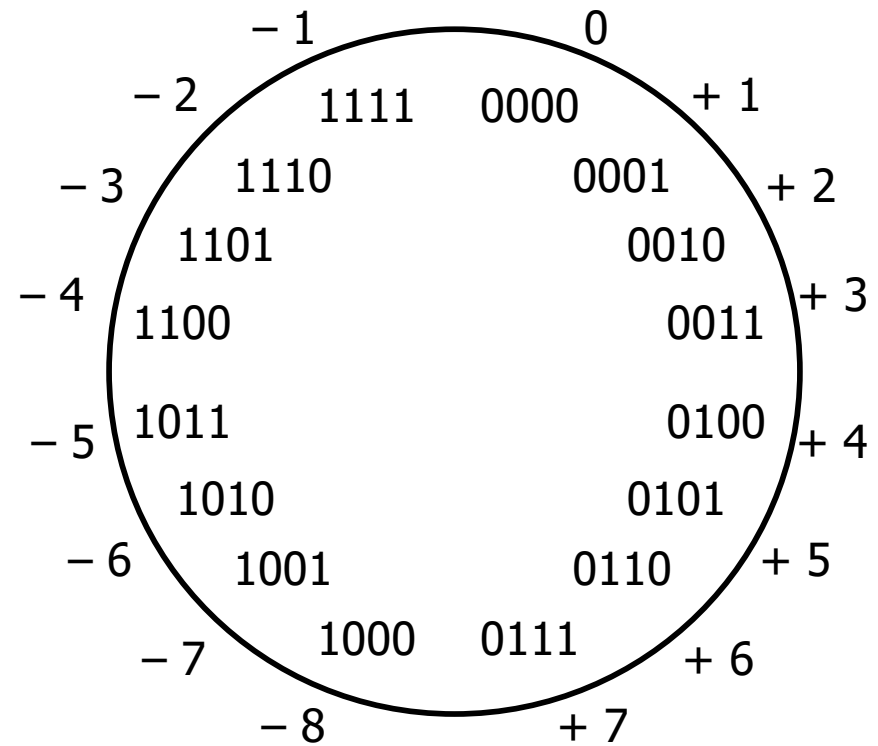
00011001

$$\begin{array}{r} -2 \\ \times 6 \\ \hline \end{array}$$
$$\begin{array}{r} 1110 \\ \times 0110 \\ \hline \end{array}$$

~~-12~~

11110100

4



## Modular Arithmetic

# Casting Integers in C



Number literals: **37** is signed, **37U** is unsigned

**Integer Casting:** *bits unchanged, just reinterpreted.*

Explicit casting:

```
int tx = (int) 73U;           // still 73
unsigned uy = (unsigned) -4; // big positive #
```

Implicit casting:

Actually does

```
tx = ux;           // tx = (int)ux;
uy = ty;           // uy = (unsigned)ty;
void foo(int z) { ... }
foo(ux);           // foo((int)ux);
if (tx < ux) ...  // if ((unsigned)tx < ux) ...
```



# More Implicit Casting in C

If you mix **unsigned** and **signed** in a single expression, then *signed values are implicitly cast to unsigned*.

How are the argument bits interpreted?

| Argument <sub>1</sub> | Op | Argument <sub>2</sub> | Type     | Result |
|-----------------------|----|-----------------------|----------|--------|
| 0                     | == | 0U                    | unsigned | 1      |
| -1                    | <  | 0                     | signed   | 1      |
| -1                    | <  | 0U                    | unsigned | 0      |
| 2147483647            | <  | -2147483647-1         |          |        |
| 2147483647U           | <  | -2147483647-1         |          |        |
| -1                    | <  | -2                    |          |        |
| (unsigned) -1         | <  | -2                    |          |        |
| 2147483647            | <  | 2147483648U           |          |        |
| 2147483647            | <  | (int) 2147483648U     |          |        |

**Note:**  $T_{min} = -2,147,483,648$       $T_{max} = 2,147,483,647$

$T_{min}$  must be written as  $-2147483647-1$  (see pg. 77 of CSAPP for details)