# 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)

0110010110101001

**MSB**  **LSB**
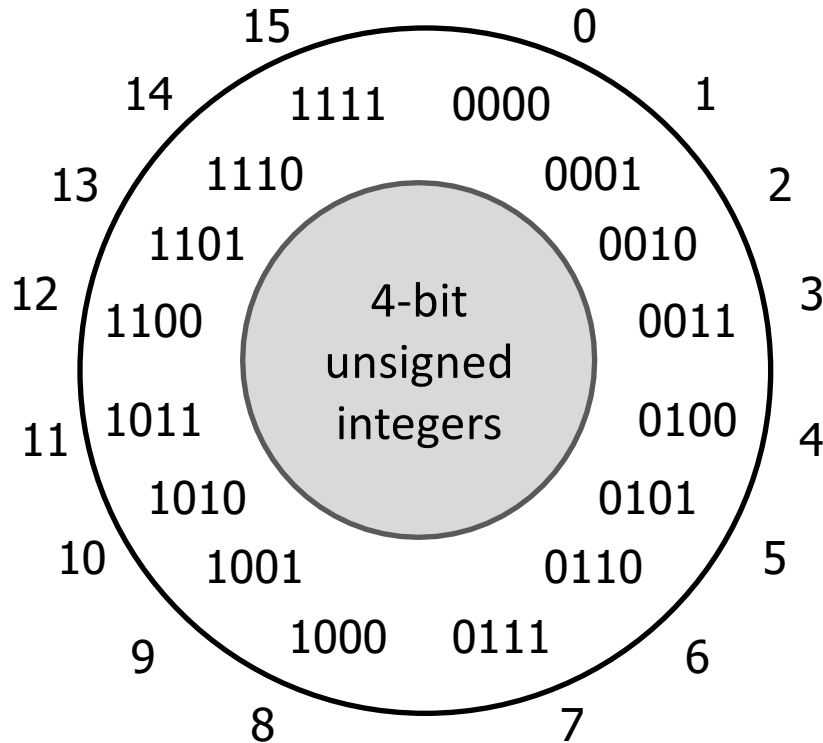
# (4-bit) **unsigned integer representation**

| 1 | 0 | 1 | 1 |
|---|---|---|---|
| 8 | 4 | 2 | 1 |
| $2^3$ | $2^2$ | $2^1$ | $2^0$ |
| 3 | 2 | 1 | 0 |

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

*weight*

*position*

*n*-bit unsigned integers:

**minimum =**

**maximum =**

# modular arithmetic, overflow

$$11 \quad 1011$$
$$\underline{+\ 2} \quad \underline{+\ 0010}$$

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

4-bit unsigned integers

$$13 \quad 1101$$
$$\underline{+\ 5} \quad \underline{+\ 0101}$$

**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:                    Anything weird here?

**0**0000000 represents _____

**Arithmetic?**

**0**1111111 represents _____                    Example:
4 - 3 != 4 + (-3)

**1**0000101 represents _____

```
  00000100
+ 10000011
```

**1**0000000 represents _____

ex

**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

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

*unsigned places*

*two's complement places*

What's the difference?

$n$-bit minimum =                    $n$-bit 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 =
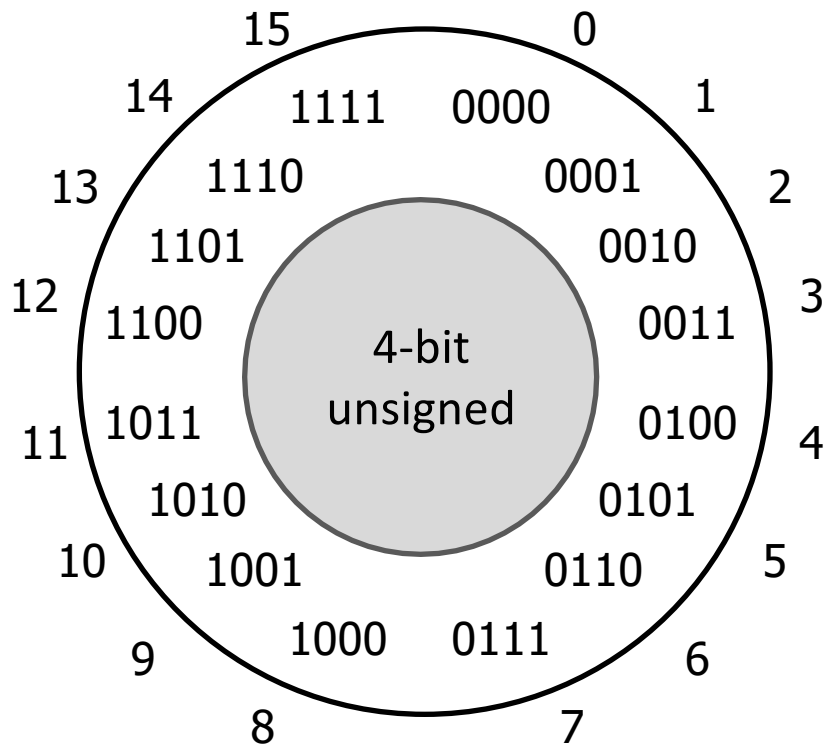
# 4-bit **unsigned** vs. 4-bit **two's complement**

$$1 \quad 0 \quad 1 \quad 1$$

$$1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \qquad\qquad 1 \times -2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$
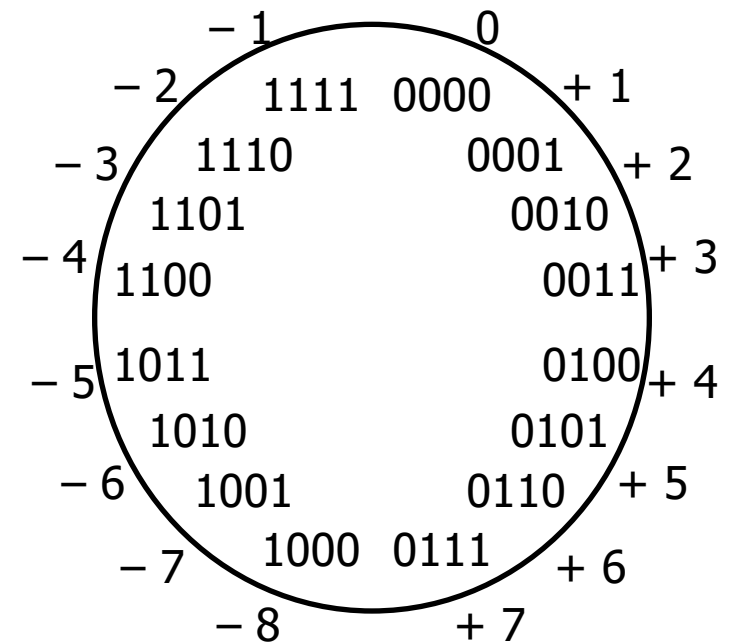
11 $\leftarrow$ **difference = ___ = 2** $\dashrightarrow$ -5



11

# two's complement **addition**

|    |          |     |         |
|----|----------|-----|---------|
| 2  | 0010     | -2  | 1110    |
| + 3| + 0011   | + -3| + 1101  |

|     |          |     |         |
|-----|----------|-----|---------|
| -2  | 1110     | 2   | 0010    |
| + 3 | + 0011   | + -3| + 1101  |



**Modular Arithmetic**

12

# two's complement *overflow*

**Addition *overflows***

if and only if
if and only if

-1        1111

+ 2      + 0010

6         0110

+ 3      + 0011

- 1      0
- 2   1111   0000   + 1
- 3    1110      0001   + 2
   1101         0010
- 4   1100        0011   + 3
- 5   1011       0100   + 4
   1010       0101
- 6   1001    0110   + 5
    1000   0111   + 6
- 7    - 8      + 7

# Modular Arithmetic

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

# 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

**ex**

## 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.

```
  00000001          00000010          00000011
+ _____        + _____        + _____
  00000000          00000000          00000000
```

- 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?

18

# unsigned **shifting** and **arithmetic**

**unsigned**
x = 27;

y = x << 2;

0 0 0 1 1 0 1 1

logical shift left

y == 108   0 0 0 1 1 0 1 1 0 0

logical shift right

1 1 1 0 1 1 0 1

**unsigned**

x = 237;

y = x >> 2;

0 0 1 1 1 0 1 1 0 1   y == 59

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

**signed**

x = -101;

y = x << 2;

y == 108

1 0 0 1 1 0 1 1

1 0 0 1 1 0 1 1 0 0

logical shift left

arithmetic shift right

1 1 1 0 1 1 0 1

1 1 1 1 1 0 1 1 0 1

**signed**

x = -19;

y = x >> 2;

y == -5

# *shift*-and-*add*

**Available operations**

    `x << k`           implements   $x * 2^k$

    `x + y`

**Implement** `y = x * 24` using only $<<$, $+$, and integer literals
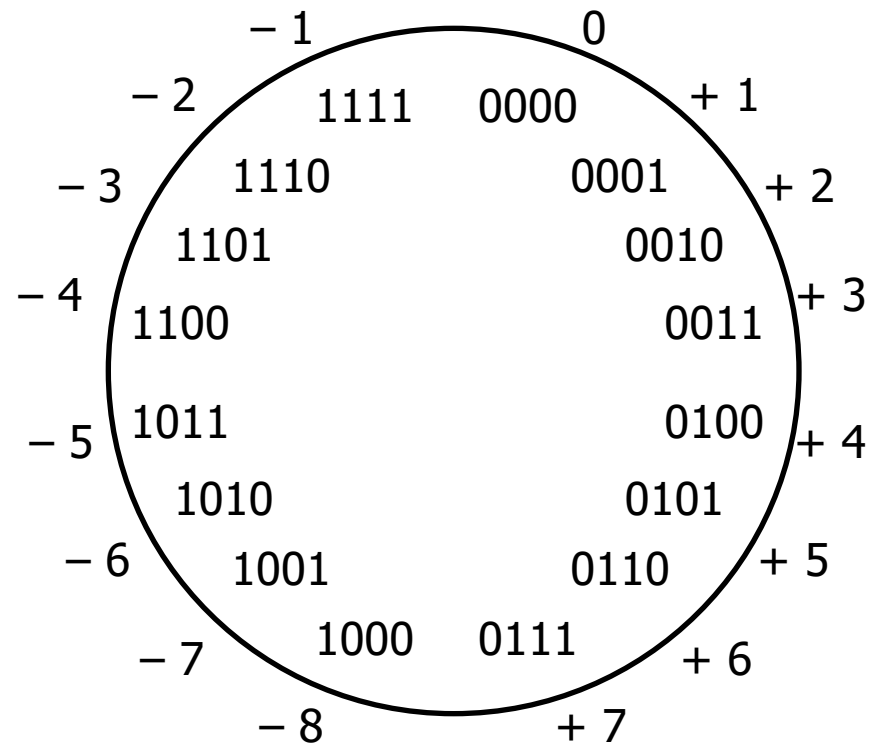
# What does this function compute?

```
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

| | | | |
|---|---|---|---|
| 2 | | 0010 | |
| x 3 | | x 0011 | |
| 6 | | 0000**0100** | |

| | | | |
|---|---|---|---|
| -2 | | 1110 | |
| x 2 | | x 0010 | |
| -4 | | **1111**1100 | |



## Modular Arithmetic

# multiplication

5         0101

x 4      x 0100

~~20~~      00010100

4

-3       1101

x 7     x 0111

~~-21~~   11101011

-2

Modular Arithmetic

25

# multiplication

| | | |
|---|---|---|
| 5 | 0101 | |
| x 5 | x 0101 | |
| ~~25~~ | <span style="color:red">0001</span>1001 | |
| <span style="color:red">-7</span> | | |

| | | |
|---|---|---|
| -2 | 1110 | |
| x 6 | x 0110 | |
| ~~-12~~ | <span style="color:red">1111</span>0100 | |
| <span style="color:red">4</span> | | |



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** underline{unsigned}*.

> How are the argument bits interpreted?

| Argument₁ | Op | Argument₂ | Type | Result |
|---|---|---|---|---|
| `0` | `==` | `0U` | unsigned | 1 |
| `-1` | `<` | `0` | signed | 1 |
| `-1` | `<` | `0U` | unsigned | **0** |
| `2147483647` | `<` | `-2147483648` | | |
| `2147483647U` | `<` | `-2147483648` | | |
| `-1` | `<` | `-2` | | |
| `(unsigned)-1` | `<` | `-2` | | |
| `2147483647` | `<` | `2147483648U` | | |
| `2147483647` | `<` | `(int)2147483648U` | | |

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