



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)

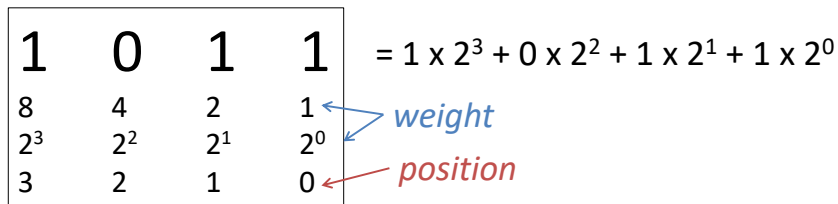
"Least-significant" bit(s)
or "low-order" bit(s)

MSB

0110010110101001

LSB

(4-bit) unsigned integer representation



n -bit unsigned integers:

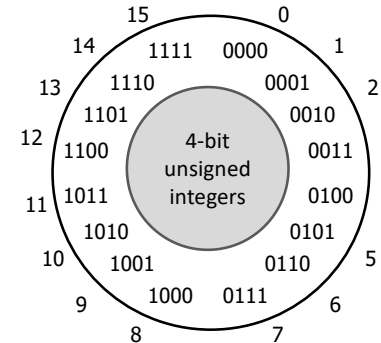
minimum =

maximum =

modular arithmetic, overflow

11 1011
+ 2 + 0010

13 1101
+ 5 + 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?

00000000 represents _____

01111111 represents _____

10000101 represents _____

10000000 represents _____

Arithmetic?

Example:
4 - 3 != 4 + (-3)



00000100
+10000011

Zero?

ex

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



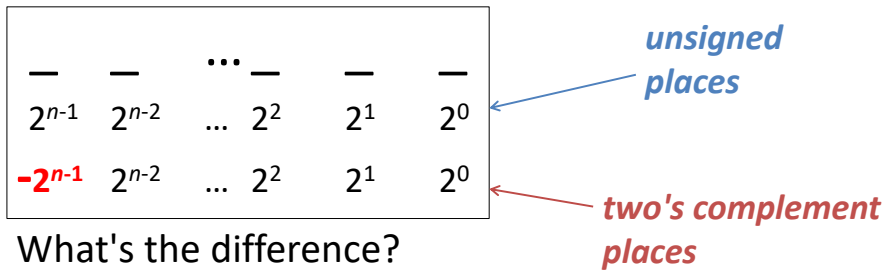
1	0	1	1	= 1 x -2 ³ + 0 x 2 ² + 1 x 2 ¹ + 1 x 2 ⁰
-2 ³	2 ²	2 ¹	2 ⁰	

4-bit two's complement integers:

minimum =

maximum =

two's complement vs. unsigned



n-bit unsigned numbers:

minimum =

maximum =

8-bit representations

ex

00001001 10000001

11111111 00100111

n-bit two's complement numbers:

minimum =

maximum =

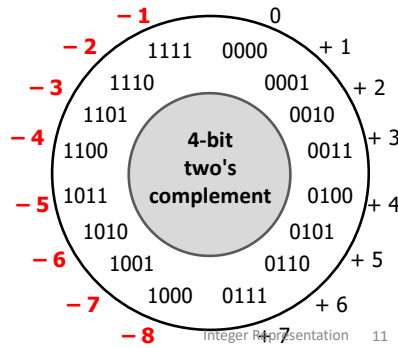
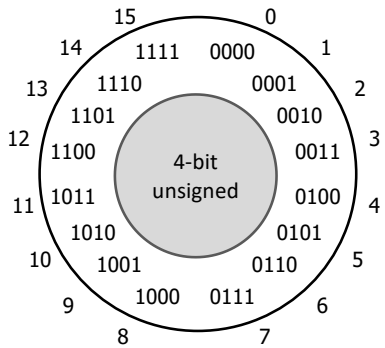
4-bit unsigned vs. 4-bit two's complement

1 0 1 1

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

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

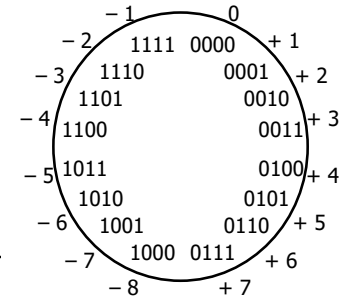
11 ← difference = = 2 — → -5



two's complement addition

$$\begin{array}{r} 2 \quad 0010 \quad -2 \quad 1110 \\ + 3 \quad + 0011 \quad + -3 \quad + 1101 \\ \hline \end{array}$$

$$\begin{array}{r} -2 \quad 1110 \quad 2 \quad 0010 \\ + 3 \quad + 0011 \quad + -3 \quad + 1101 \\ \hline \end{array}$$



Modular Arithmetic

Integer Representation 12

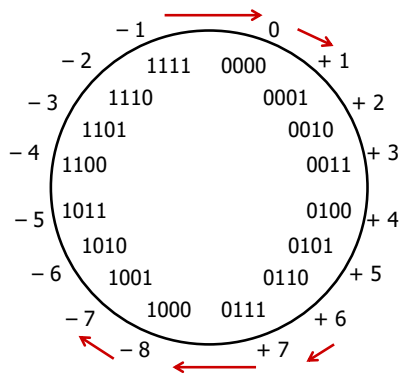
two's complement overflow

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?

Integer Representation 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

Integer Representation 14

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.

$$\begin{array}{r} 00000001 \\ + \\ \hline 00000000 \end{array} \quad \begin{array}{r} 00000010 \\ + \\ \hline 00000000 \end{array} \quad \begin{array}{r} 00000011 \\ + \\ \hline 00000000 \end{array}$$

- 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?

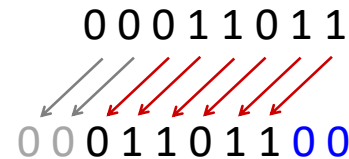
unsigned shifting and arithmetic

unsigned

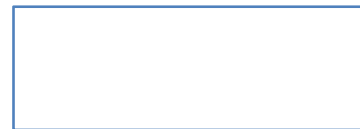
$x = 27;$

$y = x \ll 2;$

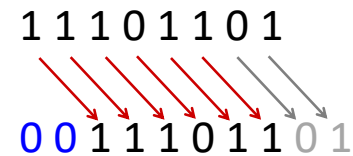
$y == 108$



logical shift left



logical shift right



unsigned

$x = 237;$

$y = x \gg 2;$

$y == 59$

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

signed

`x = -101;`

1 0 0 1 1 0 1 1

`y = x << 2;`



`y == 108`

1 0 0 1 1 0 1 1 0 0

logical shift left



arithmetic shift right

1 1 1 0 1 1 0 1

signed

`x = -19;`

`y = x >> 2;`

`y == -5`

1 1 1 1 1 0 1 1 0 1

shift-and-add

ex

Available operations

`x << k`

implements `x * 2k`

`x + y`

Implement `y = x * 24` using only `<<`, `+`, and integer literals

What does this function compute?

ex

```
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

```

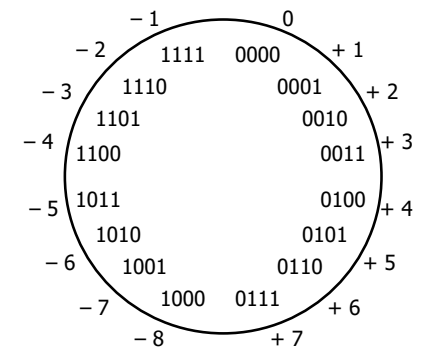
  2      0010
x 3      x 0011
-----
  6      0000110

```

```

 -2      1110
x 2      x 0010
-----
 -4      11111100

```

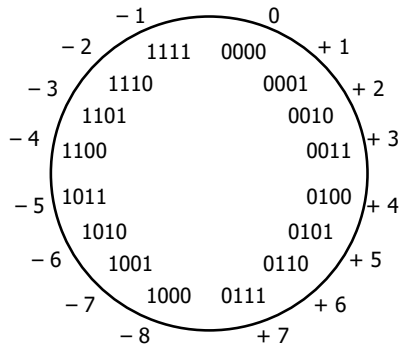


Modular Arithmetic

multiplication

```

5      0101
x 4    x 0100
-----
20  00010100
 4
-----
-3     1101
x 7    x 0111
-----
-21  11101011
 -5
  
```



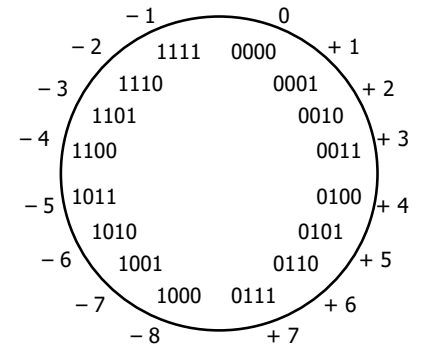
Modular Arithmetic

Integer Representation 25

multiplication

```

5      0101
x 5    x 0101
-----
25  00011001
 -7
-----
-2     1110
x 6    x 0110
-----
-12  11110100
 4
  
```



Modular Arithmetic

Integer Representation 26

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) ...
  
```

Integer Representation 27

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 ₁	Op	Argument ₂	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)

Integer Representation 28