



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)

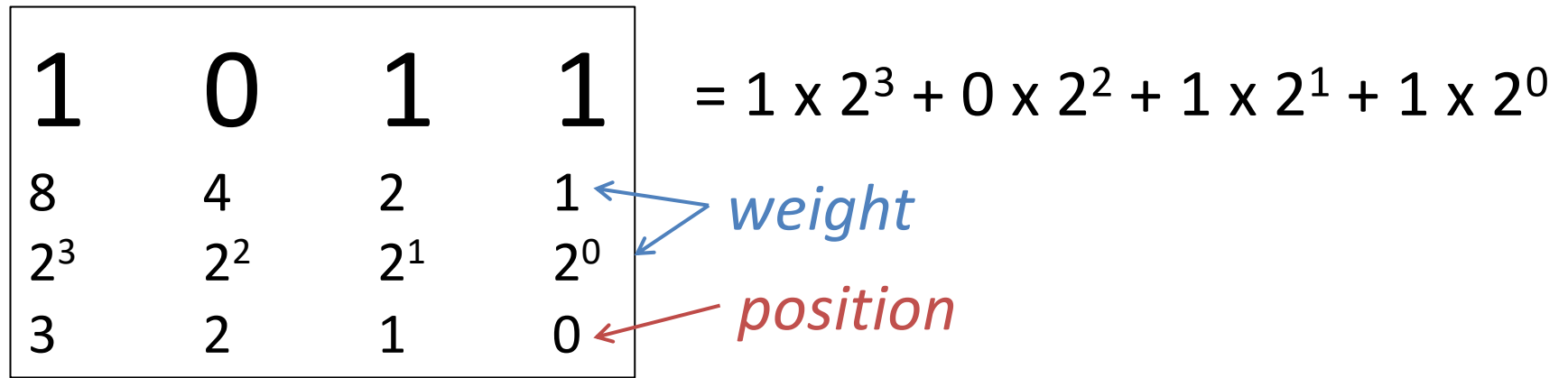
“Least-significant” bit(s)
or “low-order” bit(s)

MSB

0110010110101001

LSB

(4-bit) unsigned integer representation



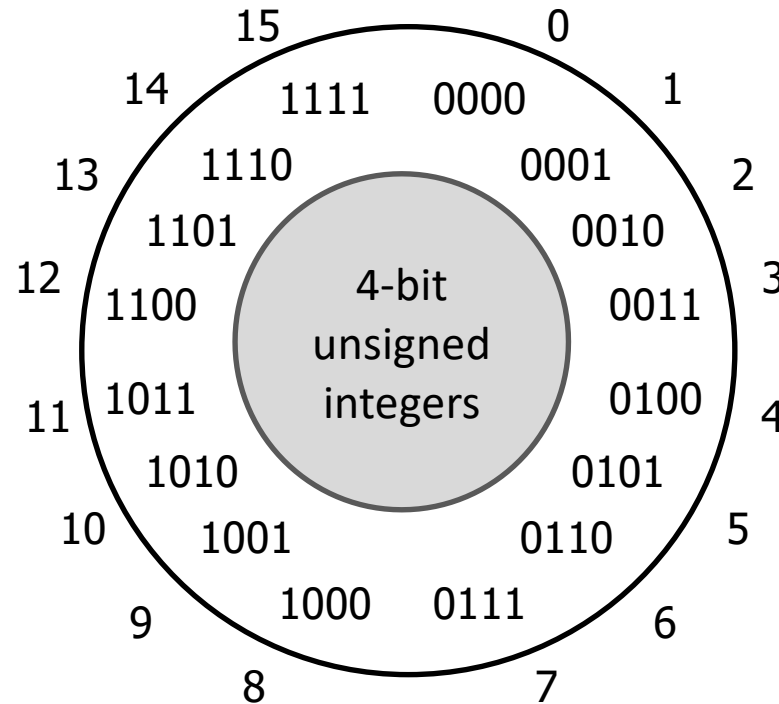
n -bit unsigned integers:

minimum =

maximum =

modular arithmetic, overflow

$$\begin{array}{r}
 11 \\
 + 2 \\
 \hline
 1011 \\
 + 0010 \\
 \hline
 1101
 \end{array}$$



$$\begin{array}{r}
 13 \\
 + 5 \\
 \hline
 1101 \\
 + 0101 \\
 \hline
 0010
 \end{array}$$

$x+y$ 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:

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



$$\begin{array}{r} 0000100 \\ +1000011 \\ \hline \end{array}$$

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 =

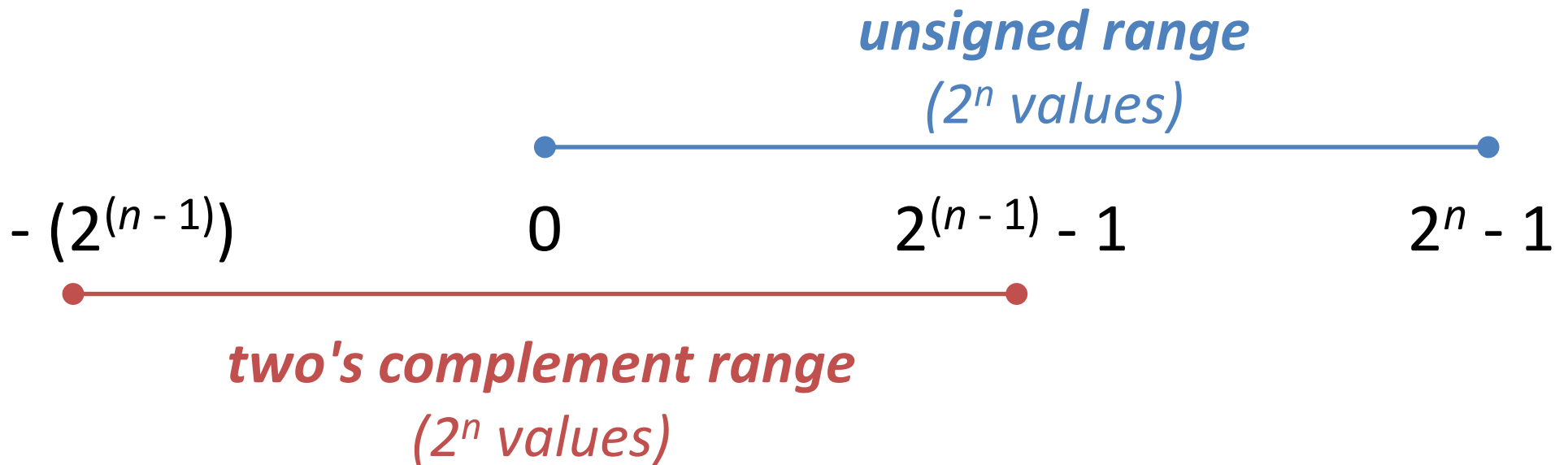
maximum =

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



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

1 0 1 1

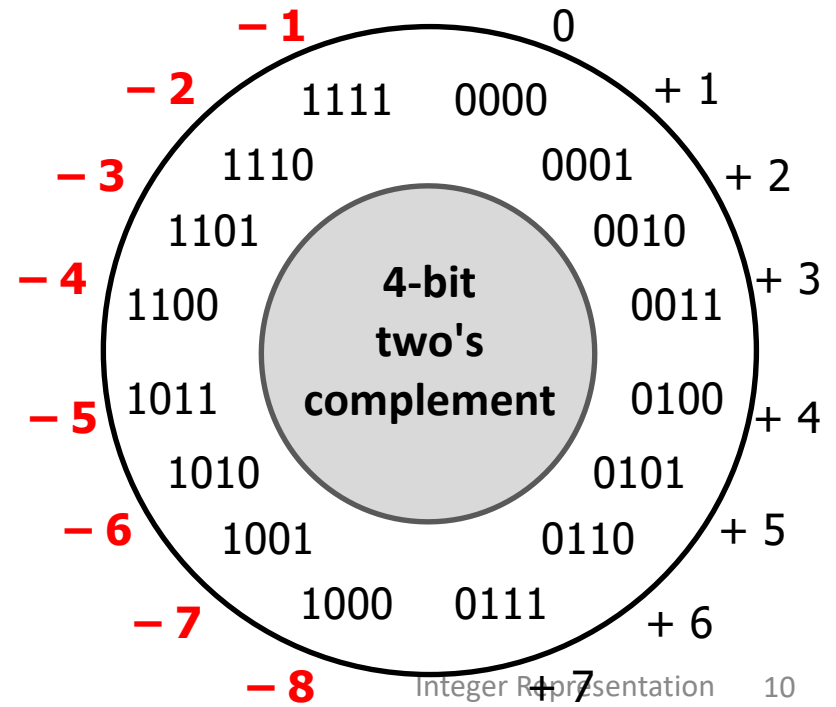
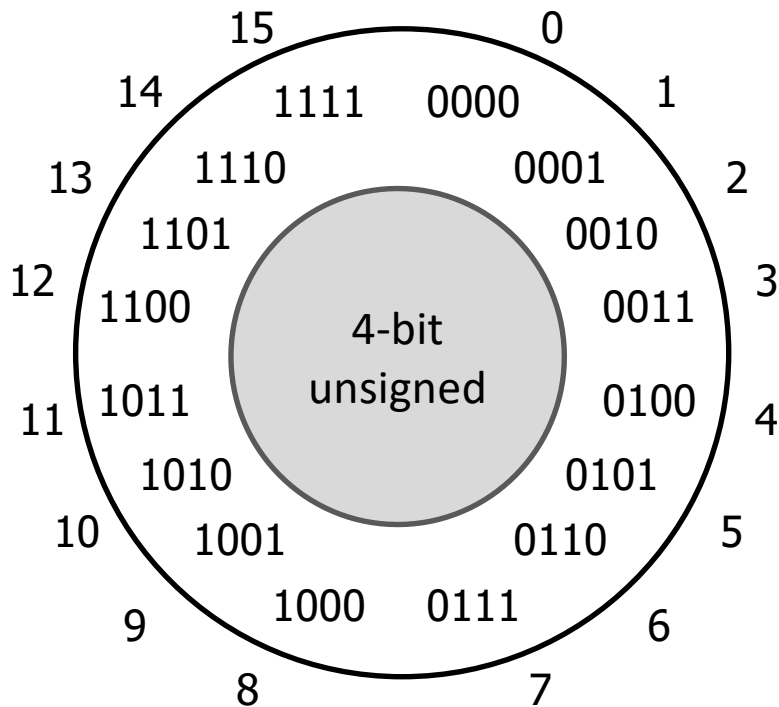
$$1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$

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

11

difference = ___ = 2__

-5



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:

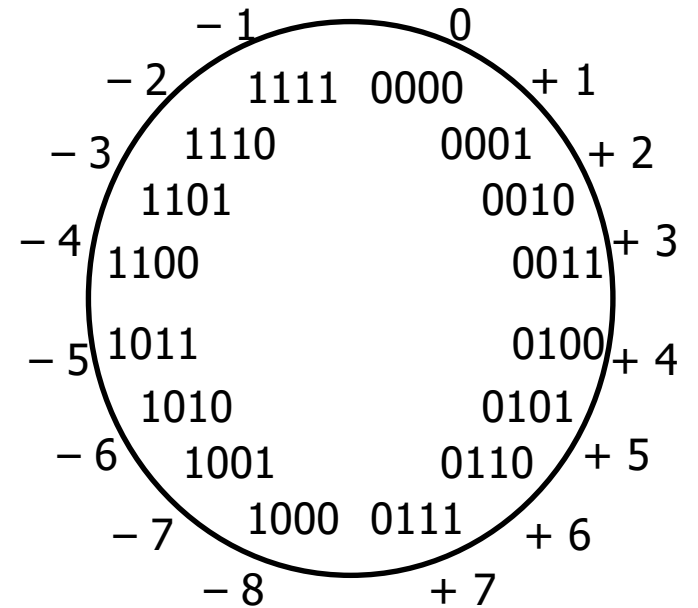
minimum =

maximum =

two's complement addition

2	0010	-2	1110
<u>+ 3</u>	<u>+ 0011</u>	<u>+ -3</u>	<u>+ 1101</u>
5		-5	

-2	1110	2	0010
<u>+ 3</u>	<u>+ 0011</u>	<u>+ -3</u>	<u>+ 1101</u>
1		-1	



Modular Arithmetic

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

Arithmetic 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.

11111111	11111111	11111111
00000001	00000010	00000011
<u>+11111111</u>	<u>+11111110</u>	<u>+11111101</u>
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

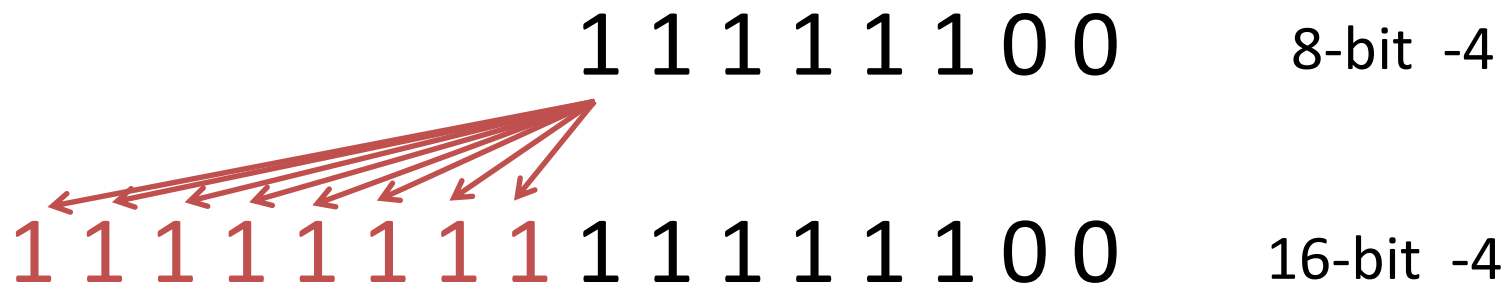
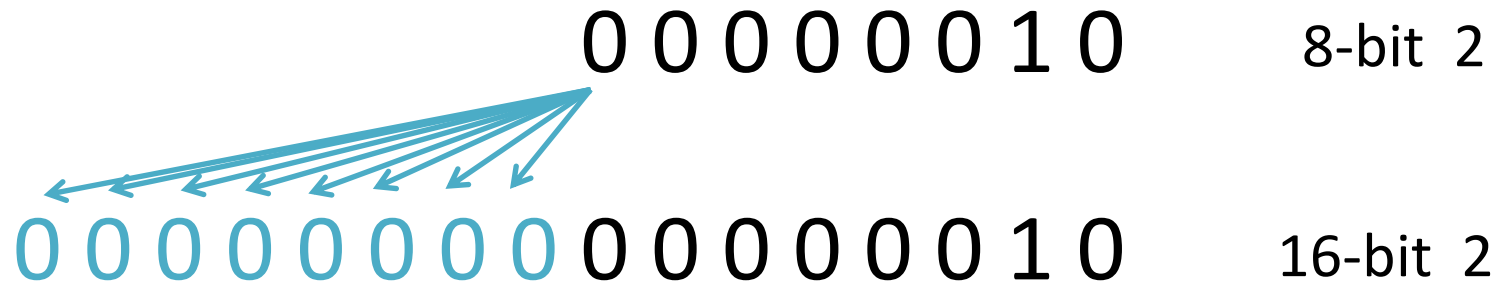
----- 0 0 0 0 0 0 1 0 16-bit 2

1 1 1 1 1 1 0 0 8-bit -4

----- 1 1 1 1 1 1 0 0 16-bit -4

Rule/name?

Sign extension for two's complement



Casting from smaller to larger signed type does sign extension.

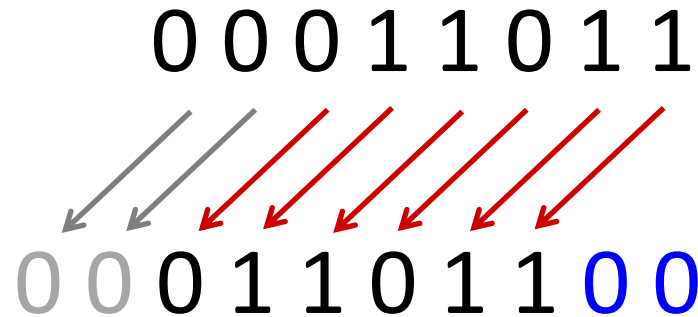
unsigned shifting and arithmetic

unsigned

$x = 27;$

$y = x \ll 2;$

$y == 108$

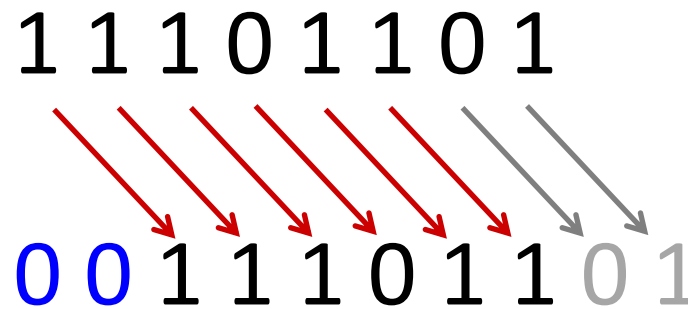


logical shift left

$n =$ shift distance in bits, $w =$ width of encoding in bits



logical shift right



unsigned

$x = 237;$

$y = x \gg 2;$

$y == 59$

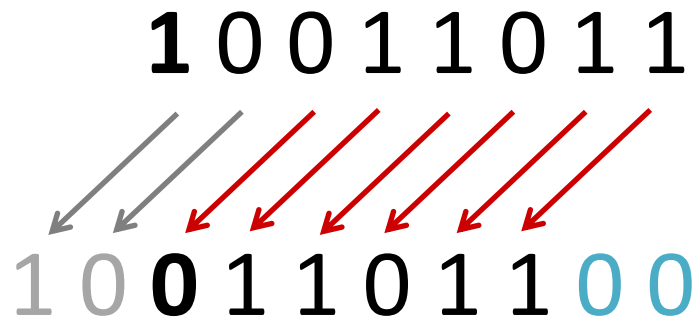
two's complement **shifting** and **arithmetic**

signed

x = -101;

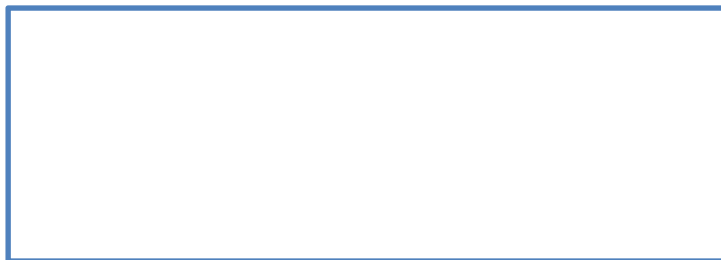
y = x << 2;

y == 108

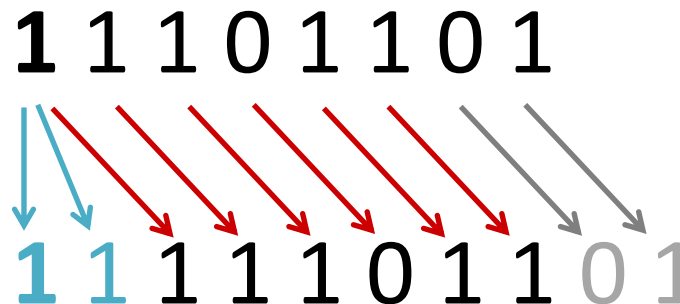


logical shift left

n = shift distance in bits, w = width of encoding in bits



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

Parenthesize shifts to be clear about precedence, which may not always be what you expect.

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;
}
```



What does this function compute?

Downsize to fake unsigned nybble type (4 bits) to make this easier to write...

```
nybble puzzle(nybble x, nybble y) {  
    nybble result = 0;  
    for (nybble i = 0; i < 4; i++){  
        if (y & (1 << i)) {  
            result = result + (x << i);  
        }  
    }  
    return result;  
}
```

	y_2	x_2

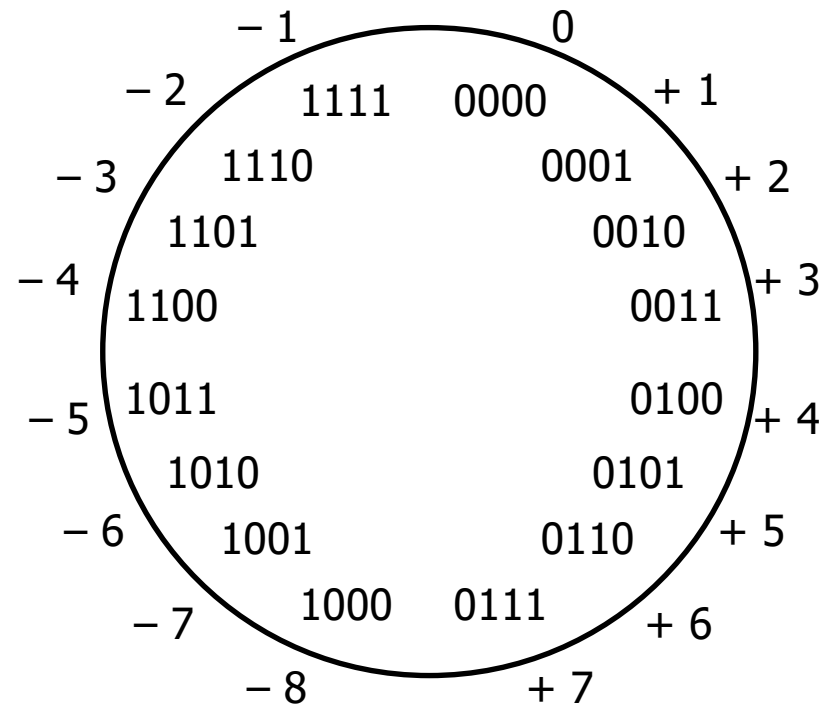
i_{10}	$y \& (1 \ll i)_2$	$result_2$
0		0 0 0 0
1		
2		
3		
4		

See Bits assignment prep exercise.

multiplication

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

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

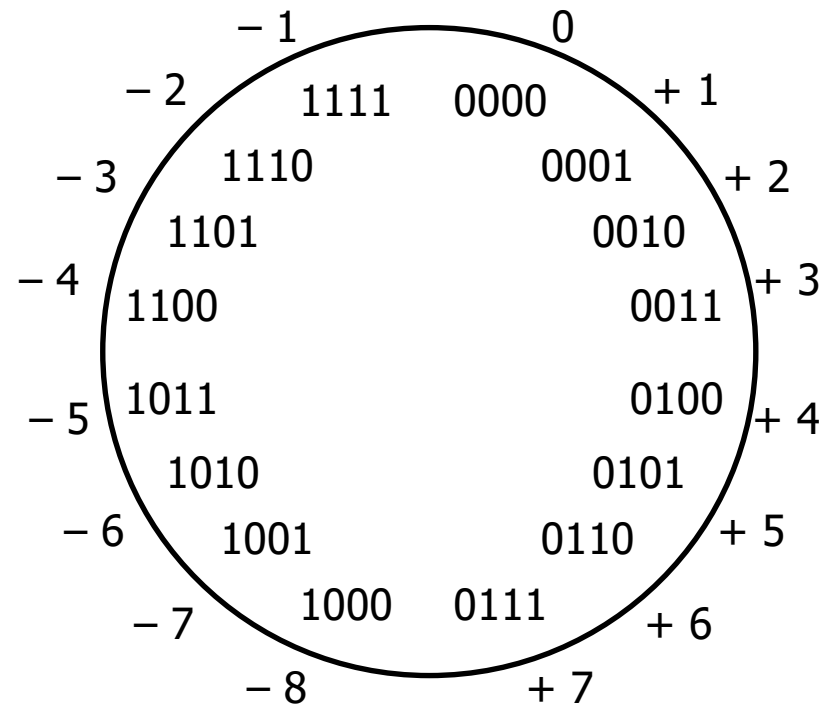


Modular Arithmetic

multiplication

$$\begin{array}{r}
 5 \quad \quad 0101 \\
 \times 4 \quad \quad \times 0100 \\
 \hline
 \del{20} \quad \quad \color{red}{0001}0100 \\
 4
 \end{array}$$

$$\begin{array}{r}
 -3 \quad \quad 1101 \\
 \times 7 \quad \quad \times 0111 \\
 \hline
 \del{-21} \quad \quad \color{red}{1110}1011 \\
 -5
 \end{array}$$

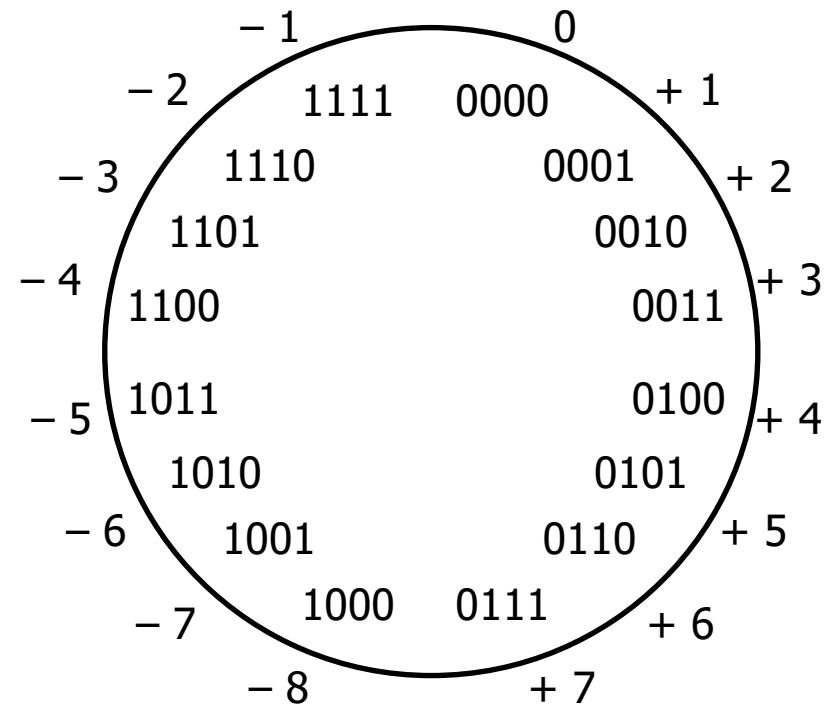


Modular Arithmetic

multiplication

$$\begin{array}{r}
 5 \\
 \times 5 \\
 \hline
 \cancel{25} \\
 -7 \\
 \hline
 -2 \\
 \times 6 \\
 \hline
 \cancel{-12} \\
 4
 \end{array}$$

$$\begin{array}{r}
 0101 \\
 \times 0101 \\
 \hline
 00011001 \\
 \\
 1110 \\
 \times 0110 \\
 \hline
 11110100
 \end{array}$$



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 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)