

# Why represent numbers

Modeling uses data-types that are encoded in a computer. The details of the encoding impact the efficiency of algorithms we use to understand the systems we are modeling and the impacts of these algorithms on the people using the systems. Case study: how to encode numbers?

## Fixed width definition

**Definition** For  $b$  an integer greater than 1,  $w$  a positive integer, and  $n$  a nonnegative integer \_\_\_\_\_, the **base  $b$  fixed-width  $w$  expansion of  $n$**  is

$$(a_{w-1} \cdots a_1 a_0)_{b,w}$$

where  $a_0, a_1, \dots, a_{w-1}$  are nonnegative integers less than  $b$  and

$$n = \sum_{i=0}^{w-1} a_i b^i$$

## Fixed width example

| Decimal<br>$b = 10$ | Binary<br>$b = 2$ | Binary fixed-width 10<br>$b = 2, w = 10$ | Binary fixed-width 7<br>$b = 2, w = 7$ | Binary fixed-width 4<br>$b = 2, w = 4$ |
|---------------------|-------------------|--|--|--|
| $(20)_{10}$         |                   |  |  |  |

# Fixed width fractional definition

**Definition** For  $b$  an integer greater than 1,  $w$  a positive integer,  $w'$  a positive integer, and  $x$  a real number the **base  $b$  fixed-width expansion of  $x$  with integer part width  $w$  and fractional part width  $w'$**  is  $(a_{w-1} \cdots a_1 a_0 . c_1 \cdots c_{w'})_{b,w,w'}$  where  $a_0, a_1, \dots, a_{w-1}, c_1, \dots, c_{w'}$  are nonnegative integers less than  $b$  and

$$x \geq \sum_{i=0}^{w-1} a_i b^i + \sum_{j=1}^{w'} c_j b^{-j} \quad \text{and} \quad x < \sum_{i=0}^{w-1} a_i b^i + \sum_{j=1}^{w'} c_j b^{-j} + b^{-w'}$$

|   |  |
|---|--|
| 3.75 in fixed-width binary,<br>integer part width 2,<br>fractional part width 8 |  |
| 0.1 in fixed-width binary,<br>integer part width 2,<br>fractional part width 8  |  |

```
welcome $jshell
| Welcome to JShell -- Version 10.0.1
| For an introduction type: /help intro

[jshell> 0.1
$1 ==>

[jshell> 0.2
$2 ==>

[jshell> 0.1 + 0.2
$3 ==>

[jshell> Math.sqrt(2)
$4 ==>

[jshell> Math.sqrt(2)*Math.sqrt(2)
$5 ==>

[jshell> █
```

Note: Java uses floating point, not fixed width representation, but similar rounding errors appear in both.

# Negative int expansions

**Representing negative integers in binary:** Fix a positive integer width for the representation  $w$ ,  $w > 1$ .

|                | To represent a positive integer $n$   | To represent a negative integer $-n$  |
|----------------|---|---|
| Sign-magnitude | $[0a_{w-2} \cdots a_0]_{s,w}$ , where $n = (a_{w-2} \cdots a_0)_{2,w-1}$<br><br>Example $n = 17$ , $w = 7$ :  | $[1a_{w-2} \cdots a_0]_{s,w}$ , where $n = (a_{w-2} \cdots a_0)_{2,w-1}$<br><br>Example $-n = -17$ , $w = 7$ :            |
| 2s complement  | $[0a_{w-2} \cdots a_0]_{2c,w}$ , where $n = (a_{w-2} \cdots a_0)_{2,w-1}$<br><br>Example $n = 17$ , $w = 7$ : | $[1a_{w-2} \cdots a_0]_{2c,w}$ , where $2^{w-1} - n = (a_{w-2} \cdots a_0)_{2,w-1}$<br><br>Example $-n = -17$ , $w = 7$ : |

## Calculating 2s complement

For positive integer  $n$ , to represent  $-n$  in 2s complement with width  $w$ ,

- Calculate  $2^{w-1} - n$ , convert result to binary fixed-width  $w - 1$ , pad with leading 1, or
- Express  $-n$  as a sum of powers of 2, where the leftmost  $2^{w-1}$  is negative weight, or
- Convert  $n$  to binary fixed-width  $w$ , flip bits, add 1 (ignore overflow)

*Challenge: use definitions to explain why each of these approaches works.*

# Representing zero

## Representing 0:

So far, we have representations for positive and negative integers. What about 0?

|                | To represent a <b>non-negative</b> integer $n$  | To represent a <b>non-positive</b> integer $-n$  |
|----------------|---|--|
| Sign-magnitude | $[0a_{w-2} \cdots a_0]_{s,w}$ , where $n = (a_{w-2} \cdots a_0)_{2,w-1}$<br>Example $n = 0, w = 7$ :  | $[1a_{w-2} \cdots a_0]_{s,w}$ , where $n = (a_{w-2} \cdots a_0)_{2,w-1}$<br>Example $-n = 0, w = 7$ :            |
| 2s complement  | $[0a_{w-2} \cdots a_0]_{2c,w}$ , where $n = (a_{w-2} \cdots a_0)_{2,w-1}$<br>Example $n = 0, w = 7$ : | $[1a_{w-2} \cdots a_0]_{2c,w}$ , where $2^{w-1} - n = (a_{w-2} \cdots a_0)_{2,w-1}$<br>Example $-n = 0, w = 7$ : |