



UNIVERSITY OF
BIRMINGHAM

Introduction to Computer Systems

More on Numbers: Representing Real Numbers



Lecture Objectives

To introduce how **real numbers** are **represented** in computer systems, and to explain some of the **limitations** of floating point representations



Lecture Outline

More on Numbers:

- ◆ Binary Arithmetic using 2's Complement
- ◆ Fixed Point Decimal to Binary
- ◆ Fixed Point Binary to Decimal
- ◆ Fixed Point Arithmetic
- ◆ Floating Point Numbers
- ◆ Numerical Precision



Binary Arithmetic – Observations

$$\begin{array}{r} +5 = 0101 \\ -5 = 1101 \\ \hline 10010 \end{array}$$

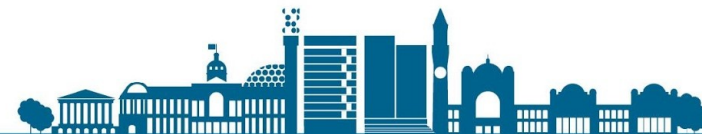
(Handwritten notes: The result 10010 is circled in pink. A bracket under the last four bits '0010' points to a green circle containing '+2'. The first '1' is also circled in pink.)

- ◆ Add a “sign” bit – assume on the left
- ◆ Then $+5 = 0101$ and $-5 = 1101$
- ◆ We know that $+5 + -5 = 0$, so this should hold in binary!
- ◆ But:

0 1 0 1
1 1 0 1
1 0 0 **1** 0

We need to Fix It!

- ◆ We need to be a bit smarter about this ...



Binary Arithmetic – 2's Complement

$$-5 = 1011$$

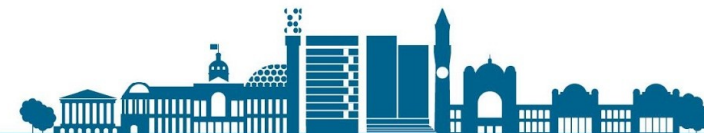
- ◆ In 2's Complement the convention is:
 - ◆ All positive numbers start with 0
 - ◆ All negative numbers start with 1
 - ◆ Negation is achieved by:
 - ◆ Flipping all the bits
 - ◆ Adding 1 to the least significant (right-most)
 - ◆ Lets see the same example again:
 - ◆ $+5 = 0101$ and $-5 = 1011$ (in 2C notation)

$$\begin{array}{r} 0101 \\ \hline 1010 \\ + 1 \\ \hline 1011 \\ = -5 \end{array}$$

$$\begin{array}{r} +5 \rightarrow 0101 \\ -5 \rightarrow 1011 \\ \hline 10000 \end{array}$$

This carry will be discarded

How about Binary Arithmetic? Subtraction etc.



How Real Numbers are Represented?

How do we represent real numbers in computer systems?

Two of the common ways to achieve this are:

7.9 ◆ **Fixed Point:** binary point is fixed e.g.

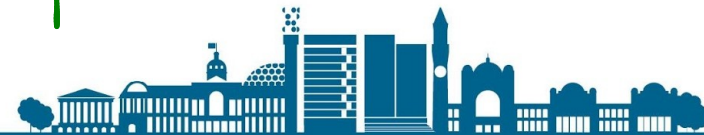
1101101.000100101

◆ **Floating Point:** binary point floats to the right of the most significant 1 and an exponent is used e.g.

1.101101000100101 $\times 2^6$

- IEEE 754 - https://en.wikipedia.org/wiki/IEEE_754

floating after the first bit



Fixed Point Decimal-to-Binary

$$(5.537)_{10}$$

- ◆ Integer part convert as before (repeated division by 2)
- ◆ Non-integer part follows the opposite process
- ◆ Repeated multiplication by 2, keeping integer part:

$$\begin{aligned} \rightarrow 0.537 \times 2 &= 1.074 \\ 0.074 \times 2 &= 0.148 \\ 0.148 \times 2 &= 0.296 \\ 0.296 \times 2 &= 0.592 \\ 0.592 \times 2 &= 1.184 \\ \rightarrow 0.184 \times 2 &= 0.368 \end{aligned}$$

So $0.537_{10} = 0.100010_2$

$$(101.100010)_2$$

3.6

10.6

◆ 0.625?

◆ 0.512?

$$(0.100010)_2 \rightarrow (?)_{10}$$



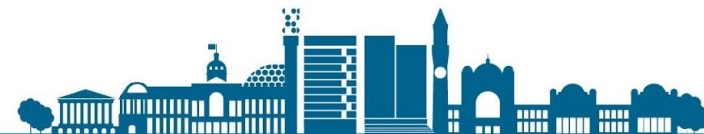
Fixed Point Binary-to-Decimal

- ◆ Allocate subset of bits to integer part, and the remainder to the non-integer part.
- ◆ For example, 4+4 bits:

$$\begin{array}{r} 1101.0101 \\ \hline 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3} + 1 \times 2^{-4} \\ = 8 + 4 + 0 + 1 + 0.0 + 0.25 + 0.0 + 0.0625 \\ = \underline{\underline{13.3125}} \end{array}$$

Handwritten annotations: Red arrows point from the fractional bits to their corresponding powers of 2. Red circles highlight the terms 0×2^{-1} and 1×2^{-2} . Red fractions $\frac{1}{2}$, $\frac{1}{4}$, and $\frac{1}{8}$ are written above the fractional part.

- ◆ 101.110?
 - ◆ 010.001?
- Handwritten annotation: A red curly brace groups the two binary numbers.*



Fixed Point Arithmetic

$$\begin{array}{r} 00010.100 \\ 11101.011 \\ + 1 \\ \hline \end{array}$$

- ◆ Everything is the same as for whole numbers
- ◆ Example: 01001.010 - 00010.100
- ◆ Take 2C and add:

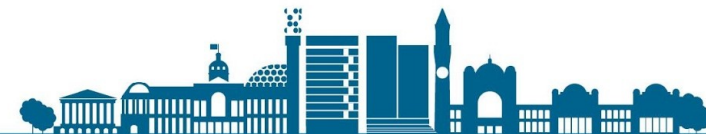
$$\begin{array}{r} 01001.010 \\ + 11101.100 \\ \hline (1)00110.110 \\ \hline \end{array}$$

(9.25)

(-2.5)

(6.75)

- ◆ 10110.101 - 00110.010 ?



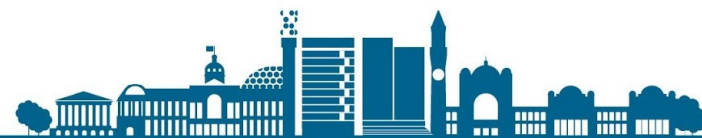
Decimal Fractions (Base 10)

e.g. $\rightarrow 42.625$

$4 \times 10 + 2 \times 1 + \frac{6}{10} + \frac{2}{100} + \frac{5}{1000}$

$= 42625 \times 10^{-3}$

$= 42 \frac{625}{1000} = 42 \frac{25}{40} = 42 \frac{5}{8}$



Infinite Decimal Fractions (Base 10)

e.g. $\frac{1}{3} = 0.3333\overline{3333} \dots$

If only fixed number of decimal places allowed, the rest is lost ✗

e.g. four places:

0.3333 $\overline{3333} \dots$

rounding error



Binary Fractions (Base 2)

Similar. e.g.

$$42\frac{5}{8}$$

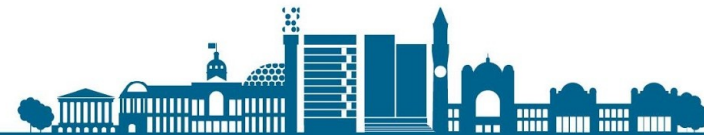
$$= \begin{array}{ccccccc} \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 0 & 1 & 0 & 1 & 0 & \cdot & 1 & 0 & 1 \\ \hline 32 & 16 & 8 & 4 & 2 & 1 & \frac{1}{2} & \frac{1}{4} & \frac{1}{8} \end{array}$$

$$= \underbrace{101010101}_{\text{binary}} \times 2^{-3}$$

$$= 0.\underbrace{000110011001100}_{\text{infinite binary fraction}}\dots$$

$$\frac{1}{10}$$

decimal



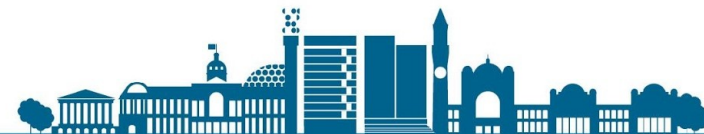
Floating Point Representation (for fractions)

General principle

- like "scientific notation", but in binary

Numbers represented as: $\pm m \cdot 2^e$

Sign (\pm)	Mantissa (m)	Exponent (e)
1 bit 0 for + 1 for –	Actual significant digits	2s complement signed binary Shows where binary point goes



Choice of Representations

$$42\frac{5}{8} = \underline{101010.101} = \underline{101010101} \times 2^{-3}$$

$$= 10101010\underline{10} \times 2^{-4}$$

$$= 10101010\underline{100} \times 2^{-5}$$

$$= \dots$$

or

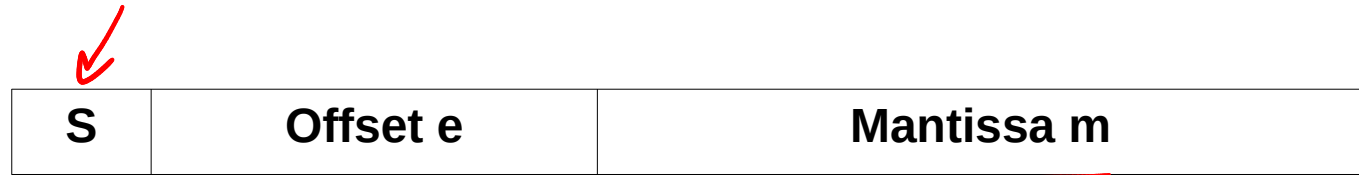
$$= \underline{101010.101} \times 2^0$$

$$= \underline{10101.0101} \times 2^1$$

$$= \dots \approx \underline{1.01010101} \times 2^5$$



Floating Point Representation in Java

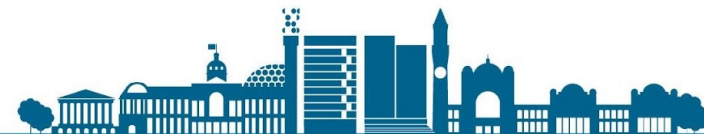


Sign
0, 1 for +, -

Exponent
2's complement
signed binary

Mantissa
Leading mantissa bit (1.) left out

- ◆ take bits of m
- ◆ put "1." at start, so normalized
- ◆ move binary point right e places
- ◆ (or left for negative e)
- ◆ note - e is stored with an "offset" added to it



Java Types for Floating Point

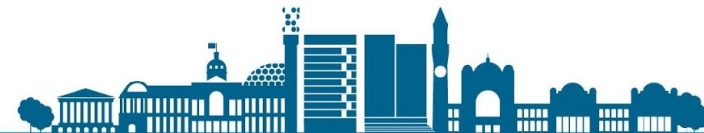
No of Bits					Bytes
Type	Sign	Mantissa	Exponent	Total	
float	1	23	8	32	4
double	1	52	11	64	8

52 bit mantissa: $2^{53} \approx 8 \times 10^{15}$ 1 extra for hidden bit

$2^{53} \Rightarrow 2^3 \times 2^{50}$

We get 15 significant decimal digits in double data type.

$$10^3 \Rightarrow 8 \times 10^{15} \approx 8 \times (2^{10})^5$$



$$5 + 127 = \underline{\underline{132}}$$

Java Types for Floating Point

$$42 \frac{5}{8} = 101010.101 = 1.01010101 \times 2^5$$

0	1000 0100	010 1010 1000 0000 0000 0000
----------	------------------	-------------------------------------

Sign
0 for +

5 + offset 127 = 132
e + 127

Mantissa
Without the leading 1.

$$\begin{aligned} \frac{1}{10} &= 0.000110011001100... \\ &= 1.10011001100... \times 2^{-4} \end{aligned}$$

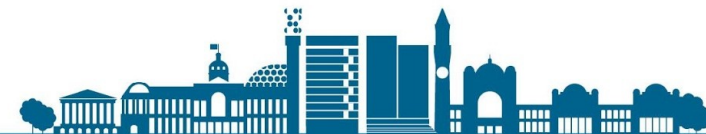
0	0111 1011	100 1100 1100 1100 1100 1101
----------	------------------	-------------------------------------

Sign
0 for +

-4 + offset 127 =
123
e + 127

Mantissa
Without the leading 1.

Rounded Up
Rounding Error!



Money as Floating Point?

Floating Point value for amount in pounds with pence as fraction?

Not a good practice!

Pence need infinite binary fractions

e.g. 10p is 0.0001100110011001100...

so we get rounding errors

Always use int or long (or BigInteger) for money.



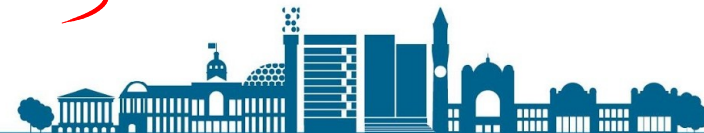
Factorial with Double

```
/**
 * Calculate factorial.
 * requires: 0 <= n
 * @param n number whose factorial is to be calculated
 * @return factorial of n
 */
public static double dfact(int n){
    double a = 1;
    for (int i = 1; i <= n; i++){
        a = a * i;
    }
    return a;
}
```

n, n!

165,	5.423910666131586E295
166,	9.003691705778433E297
167,	1.5036165148649983E300
168,	2.526075744973197E302
169,	4.2690680090047027E304
170,	7.257415615307994E306

171,	Infinity
172,	Infinity
173,	Infinity
174,	Infinity



Accuracy Issues – Even in Double Floating Point

$$n = 170$$

$$n! = 7.257415615307994E306$$

$$= 725741561530799\mathbf{4}E291$$

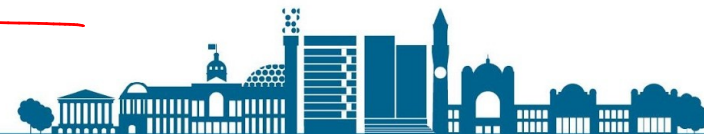
This digit is wrong!

291 more digits after it

Windows calculator  says:

7.25741561530799 $\mathbf{8}$ 9673967282111293e+306

- ◆ Floating point arithmetic loses accuracy in least significant digits.
- ◆ Most significant, and overall size, are OK.



Why is 171! too big?

2^{10}

$$\frac{170!}{171!} \approx 7.2 \times 10^{306}$$

$$\frac{170!}{171!} \approx \boxed{171 \times 7.2} \times 10^{306}$$

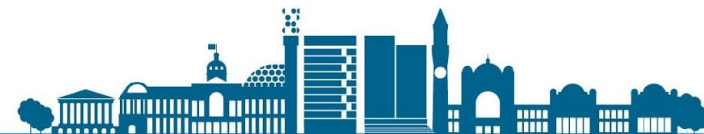
$$\approx 1231 \times 10^{306}$$

$\approx 10^3$

$$\approx 10^{309} = (10^3)^{103} \approx 2^{1030}$$

Needs binary exponent $\approx 1030 > 1024 = 2^{10}$

Exponent too big to fit in 11 bits (signed) for Java double



Floating Point Overflow

In Java:

If result is too big for datatype,

it's a special value **POSITIVE_INFINITY**

More precisely:

`Float.POSITIVE_INFINITY`

or

`Double.POSITIVE_INFINITY`

Other special values:

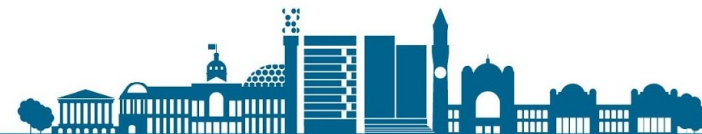
If too big but negative: **NEGATIVE_INFINITY**

Indistinguishable from 0 but known to be negative: -0.0

Impossible number (e.g. `Sqrt(-1)`) **NaN** "Not a Number"

These special values allow you to check for overflow in a program

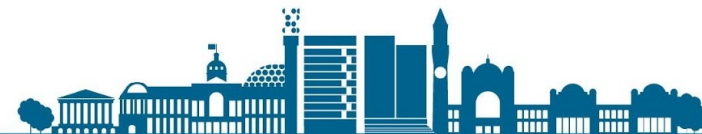
- unlike the case for integer arithmetic



Summary – Java Floating Point

1222754

- ✓◆ It normalizes where possible
- ✓◆ Unnormalised for smallest numbers (offset $e = 0$)
- ✓◆ Special representations for NaN etc.
- ◆ Details in API for:
 - ◆ `java.lang.Float.intBitsToFloat`
 - ◆ `java.lang.Double.longBitsToDouble`



Numerical Precision

- ◆ Fixed point is convenient and intuitive but has two problems

1) Numerical precision

- Only values that are an integer multiple of the smallest power of two can be represented exactly

2) Numerical range

- Increased precision of non-integer part is at the expense of numerical range
- ◆ Floating point representation effectively addresses these issues.



Summary – Numbers

- ◆ Representing numbers in the computer
 - Whole numbers in binary and hexadecimal notation
 - Positive real numbers in fixed-point binary
- ◆ Binary arithmetic is like decimal arithmetic
 - Our lack of binary practice makes it hard!
- ◆ Negative numbers are tricky things
 - But we can use a few of our own tricks – 2C
- ◆ Floating point is an alternative, but is very unnatural for us!

