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.
◆ 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

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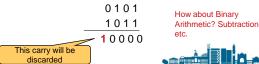◆ We know that +5 + -5 = 0, so this should hold in binary!
◆ But:

```
  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

◆ 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)

```
  0 1 0 1
  1 0 1 1
 ---------
1 0 0 0 0
```

How about Binary Arithmetic? Subtraction etc.

This carry will be discarded

---

How Real Numbers are Represented?

How do we represent real numbers in computer systems?
Two of the common ways to achieve this are:
◆ **Fixed Point**: binary point is fixed e.g.
  1101101.0001001
◆ **Floating Point**: binary point floats to the right of the most significant 1 and an exponent is used e.g.
  $1.1011010001001 \times 2^6$
  ● IEEE 754 - https://en.wikipedia.org/wiki/IEEE_754

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

$$0.537 \times 2 = \mathbf{1}.074$$
$$0.074 \times 2 = \mathbf{0}.148$$
$$0.148 \times 2 = \mathbf{0}.296$$
$$0.296 \times 2 = \mathbf{0}.592$$
$$0.592 \times 2 = \mathbf{1}.184$$
$$0.184 \times 2 = \mathbf{0}.368$$

So $0.537_{10} = 0.100010_2$

- 0.625?
- 0.512?

---

## 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:

1101.0101

$$1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3} + 1 \times 2^{-4}$$
$$= 8 + 4 + 0 + 1 + 0.0 + 0.25 + 0.0 + 0.0625$$
$$= 13.3125$$

- 101.110?
- 010.001?

---

## Fixed Point Arithmetic

2's Complemen

11101.011

+1

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

```
    0 1 0 0 1 0 1 0
  + 1 1 1 0 1 1 0 0
  -----------------
  (1) 0 0 1 1 0 1 1 0
```

- 00110.101 – 10110.010?

---

## Decimal Fractions (Base 10)



e.g.

$$42.625$$
$$10 \quad 1 \quad \tfrac{1}{10} \quad \tfrac{1}{100} \quad \tfrac{1}{1000}$$
$$4 \times 10 + 2 \times 1 + \tfrac{6}{10} + \tfrac{2}{100} + \tfrac{5}{1000}$$
$$= 42625 \times 10^{-3}$$
$$= 42\tfrac{625}{1000} = 42\tfrac{25}{40} = 42\tfrac{5}{8}$$

---

## Infinite Decimal Fractions (Base 10)



e.g. $\tfrac{1}{3} = 0.3333 | 33333 \cdots$ ✗

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

e.g. four places:

0.3333 | 33333 ...

rounding error

---

## Binary Fractions (Base 2)



Similar. e.g.

$$42\tfrac{5}{8} = 101010.101 = 101010101 \times 2^{-3}$$
$$32 \ 16 \ 8 \ 4 \ 2 \ 1 \ \tfrac{1}{2} \ \tfrac{1}{4} \ \tfrac{1}{8} \qquad \text{binary}$$

$$\tfrac{1}{10} = 0.000110011001100\cdots$$
decimal

infinite binary fraction

## Floating Point Representation (for fractions)

General principle
- like "scientific notation", but in binary

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

| Sign (±) | Mantissa (m) | Exponent (e) |
|---|---|---|
| 1 bit<br>0 for +<br>1 for – | Actual significant digits | 'Signed' binary<br>'Shows' where binary point goes |

## Choice of Representations

$$42\tfrac{5}{8} = 101010.101 = 101010101 \times 2^{-3}$$
$$= 1010101010 \times 2^{-4}$$
$$= 10101010100 \times 2^{-5}$$
$$\approx \ldots$$
$$\text{or} \quad = 101010.101 \times 2^0$$
$$= 10101.0101 \times 2^1$$
$$= \ldots \approx 1.01010101 \times 2^5$$

## Floating Point Representation in Java

| S | Offset e | Mantissa m |
|---|---|---|

| Sign<br>0, 1 for +, - | Exponent:<br>Biased binary | Mantissa:<br>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 | | | |
|---|---|---|---|---|---|
| Type | Sign | Mantissa | Exponent | Total | Bytes |
| float | 1 | 23 | 8 | 32 | 4 |
| double | 1 | 52 | 11 | 64 | 8 |

1 extra for hidden bit

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

We get 15 significant decimal digits in double data type.

## Java Types for Floating Point

Offset: $2^{(n-1)} -1$

$2^{(8-1)} - 1 = 127$

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

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

| Sign<br>0 for + | 5 + offset 127 = 132<br>e + 127 | Mantissa<br>Without the leading 1. |
|---|---|---|

$1/_{10} = 0.000110011001100...$
$= 1.10011001100... \times 2^{-4}$

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

| Sign<br>0 for + | -4 + offset 127 = 123<br>e + 127 | Mantissa<br>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

$n = 170$

$n! = 7.257415615307994E306$

$= 725741561530799\textbf{4}E291$

This digit is wrong!    291 more digits after it

Windows calculator says:
$7.25741561530799\textbf{8}9673967282111293e+306$

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

## Why is 171! too big?

$$170! \approx 7.2 \times 10^{306}$$
$$171! \approx 171 \times 7.2 \times 10^{306}$$
$$\approx 1231 \times 10^{306}$$
$$\approx 10^{309} = (10^3)^{103} \approx 2^{1030}$$

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

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

## Floating Point Overflow

More precisely:
Float.POSITIVE_INFINITY
or
Double.POSITIVE_INFINITY

In Java:
If result is too big for datatype,
it's a special value **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

◆ It normalizes where possible
◆ Unnormalised for smallest numbers (offset e = o)
◆ 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!