

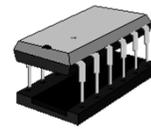
Question

Do you have your laptop here?

A yes **B** no **C** what's a laptop **D** where is here?
E none of the above

Floating Point Numbers

6.626068×10^{-34}



- Eddie Edwards
- eedwards@doc.ic.ac.uk
- <https://www.doc.ic.ac.uk/~eedwards/compsys>
- Heavily based on notes from Naranker Dulay

Eddie Edwards 2008

Floating Point Numbers

7.2

Learning Outcomes

- At the end of this lecture you should
 - Understand the representation of real numbers in other bases (e.g 2)
 - Know the mantissa/exponent representation (in base 10, 2 etc.)
 - Be able to express numbers in normalised/un-normalised form
 - Be able to convert fractions/decimals between bases
 - Know the IEEE 754 floating point format (32 and 64 bit)
 - Know the special values and when they should occur
 - Understand the issues of accuracy in floating point representation

Eddie Edwards 2008

Floating Point Numbers

7.3

Number Representation - recap

- We have seen how to represent integers
 - positive integers as binary, octal and hexadecimal
 - negative integers as one's complement, two's complement, Excess-n
 - BCD, ASCII.....
- We have also seen how to perform arithmetic
 - Addition
 - by adding the binary bits
 - overflow conditions
 - Multiplication/division
 - same "long hand" techniques as base 10
 - slightly complicated in two's compliment
 - can take the absolute values, perform calculation, then sort out the sign

Eddie Edwards 2008

Floating Point Numbers

7.4

Reals vs. Floating Point Numbers

	Mathematical Real	Floating-point Number
Range	-Infinity .. +Infinity	Finite
No. of Values	Infinite	Finite
Spacing	Constant & Infinite	Gap between numbers varies
Errors	?	Incorrect results are possible

Eddie Edwards 2008

Floating Point Numbers

7.9

Normalised Floating Point Numbers

➤ Floating Point Numbers can have multiple forms, e.g.

$$\begin{aligned}
 0.232 \times 10^4 &= 2.32 \times 10^3 \\
 &= 23.2 \times 10^2 \\
 &= 2\,320 \times 10^0 \\
 &= 232\,000 \times 10^{-2}
 \end{aligned}$$

➤ For hardware implementation its desirable for each number to have a unique representation => **Normalised Form**

➤ We'll normalise Mantissa's in the Range [1 .. R) where R is the Base, e.g.:

[1 .. 10) for DECIMAL
 [1 .. 2) for BINARY

Eddie Edwards 2008

Floating Point Numbers

7.10

Normalised Forms (Base 10)

Number	Normalised Form
23.2×10^4	2.32×10^5
-4.01×10^{-3}	-4.01×10^{-3}
$343\,000 \times 10$	3.43×10^5
$0.000\,000\,098\,9 \times 10^0$	9.89×10^{-8}

Eddie Edwards 2008

Floating Point Numbers

7.11

Binary & Decimal Fractions

Binary	Decimal
0.1	0.5
0.01	0.25
0.001	0.125
0.11	0.75
0.111	0.875
0.011	0.375
0.101	0.625

Eddie Edwards 2008

Floating Point Numbers

7.12

Floating Point Multiplication

$$\begin{aligned} N1 \times N2 &= (M1 \times 10^{E1}) \times (M2 \times 10^{E2}) \\ &= (M1 \times M2) \times (10^{E1} \times 10^{E2}) \\ &= (M1 \times M2) \times (10^{E1+E2}) \end{aligned}$$

i.e. We multiply the Mantissas and Add the Exponents

Example: $20 * 6 = (2.0 \times 10^1) \times (6.0 \times 10^0)$
 $= (2.0 \times 6.0) \times (10^{1+0})$
 $= 12.0 \times 10^1$

We must also normalise the result, so the final answer = 1.2×10^2

Eddie Edwards 2008

Floating Point Numbers

7.17

Truncation and Rounding

➤ For many computations the result of a floating point operation can be too large to store in the Mantissa.

➤ Example: with a 2-digit mantissa

$$2.3 \times 10^1 * 2.3 \times 10^1 = 5.29 \times 10^2$$

➤ TRUNCATION $\Rightarrow 5.2 \times 10^2$ (Biased Error)

➤ ROUNDING $\Rightarrow 5.3 \times 10^2$ (Unbiased Error)

Eddie Edwards 2008

Floating Point Numbers

7.18

Floating Point Addition

➤ A floating point addition such as $4.5 \times 10^3 + 6.7 \times 10^2$ is not a simple mantissa addition, unless the exponents are the same
 \Rightarrow we need to ensure that the mantissas are aligned first.

$$\begin{aligned} N1 + N2 &= (M1 \times 10^{E1}) + (M2 \times 10^{E2}) \\ &= (M1 + M2 \times 10^{E2-E1}) \times 10^{E1} \end{aligned}$$

➤ To align, choose the number with the smaller exponent & shift mantissa the corresponding number of digits to the right.

Example: $4.5 \times 10^3 + 6.7 \times 10^2 = 4.5 \times 10^3 + 0.67 \times 10^3$
 $= 5.17 \times 10^3$
 $= 5.2 \times 10^3$ (rounded)

Eddie Edwards 2008

Floating Point Numbers

7.19

Exponent Overflow & Underflow

➤ **EXPONENT OVERFLOW** occurs when the Result is too Large
i.e. when the Result's Exponent $>$ Maximum Exponent

Example: if Max Exponent is 99 then $10^{99} * 10^{99} = 10^{198}$ (overflow)

On Overflow \Rightarrow Proceed with incorrect value or infinity value or raise an Exception

➤ **EXPONENT UNDERFLOW** occurs when the Result is too Small
i.e. when the Result's Exponent $<$ Smallest Exponent

Example: if Min Exp. is -99 then $10^{-99} * 10^{-99} = 10^{-198}$ (underflow)

On Underflow \Rightarrow Proceed with zero value or raise an Exception

Eddie Edwards 2008

Floating Point Numbers

7.20

Comparing Floating-Point Values

- Because of the potential for producing in-exact results, comparing floating-point values should account for close results.
- If we know the likely magnitude and precision of results we can adjust for closeness (epsilon), for example, for equality we can:
 $a = b \quad a > (b - e) \text{ AND } a < (b + e)$
- $a = 1 \quad a > (1 - 0.000005) \text{ AND } a < 1 + 0.000005$
 $a > 0.999995 \text{ AND } a < 1.000005$
- Alternatively we can calculate $|a - b| < e$ e.g. $|a - 1| < 0.000005$
- A more general approach is to calculate the closeness based on the relative size of the two numbers being compared.

Eddie Edwards 2008

Floating Point Numbers

7.21

Floating point numbers - questions

Eddie Edwards 2008

Floating Point Numbers

7.22

- What is the binary notation for 3.625
- A 11.011 B 10.101 C 11.101 D 101.11 E 11.11
- What is binary 0.1101 in decimal?
- A 0.8125 B 0.8 C 0.8625 D 0.9125 E 0.7865

Eddie Edwards 2008

Floating Point Numbers

7.23

IEEE Floating-Point Standard

- IEEE: Institute of Electrical & Electronic Engineers (USA)
- Comprehensive standard for Binary Floating-Point Arithmetic
- Widely adopted => Predictable results independent of architecture
- The standard defines:
 - The format of binary floating-point numbers
 - Semantics of arithmetic operations
 - Rules for error conditions

Eddie Edwards 2008

Floating Point Numbers

7.24

Single Precision Format (32-bit)

Sign S	Exponent E	Significand F
1 bit	8 bits	23 bits

- > The mantissa is called the **SIGNIFICAND** in the IEEE standard
- > Value represented = $\pm 1.F \times 2^{E-127}$ $127 = 2^{8-1} - 1$
- > The Normal Bit (the 1.) is omitted from the Significand field => a **HIDDEN** bit
- > Single precision yields 24-bits = ~ 7 decimal digits of precision
- > Normalised Ranges in decimal are approximately:
 -10^{+38} to -10^{-38} , **0**, $+10^{-38}$ to $+10^{+38}$

Eddie Edwards 2008

Floating Point Numbers

7.25

Exponent Field

> In the IEEE Standard, exponents are stored as **Excess (Bias) Values**, not as 2's Complement Values

> **Example:** In 8-bit Excess 127
 -127 would be held as 0000 0000

...
 0 would be held as ... 0111 1111
 1 would be held as 1000 0000

...
 128 would be held as ... 1111 1111

> Excess notation allows non-negative floating point numbers to be compared using simple integer comparisons, regardless of the absolute magnitude of the exponents.

Eddie Edwards 2008

Floating Point Numbers

7.26

Double Precision Format (64-bit)

Sign S	Exponent E	Significand F
1 bit	11 bits	52 bits

- Value represented = $\pm 1.F \times 2^{E-1023}$ $1023 = 2^{11-1} - 1$
- > Yields 53 bits of precision = ~ 16 decimal digits of precision
 - > Normalised Ranges in decimal are approximately:
 -10^{+308} to -10^{-308} , **0**, $+10^{-308}$ to $+10^{+308}$
 - > Double-Precision format is preferred for its greater precision. Single-precision is useful when memory is scarce and for debugging numerical calculations since rounding errors show up more quickly.

Eddie Edwards 2008

Floating Point Numbers

7.27

Example: Conversion to IEEE format

What is 42.6875 in IEEE Single Precision Format?

First convert to a binary number: **42.6875 = 10_1010 . 1011**

Next normalise: **1 . 0101_0101_1 x 2⁵**

Significand field is therefore: **0101_0101_1000_0000_0000_000**

Exponent field is (5+127=132): **1000_0100**

Value in IEEE Single Precision is:

Sign	Exponent	Significand
0	1000_0100	0101_0101_1000_0000_0000_000
0100__0010__0 010__1010__1100__0000__0000__0000		

In hexadecimal this value is **422A_C000**

Eddie Edwards 2008

Floating Point Numbers

7.28

Example: Conversion from IEEE format

Convert the IEEE Single Precision Value given by BECO_0000 to Decimal

BECO_0000 = 1011_1110_1 100_0000_0000_0000_0000_0000

Sign	Exponent	Significand
1	0111_1101	1000_0000_0000_0000_0000_0000

Exponent Field = 0111_1101 = 125
 True Binary Exponent = 125 - 127 = -2

Significand Field = 1000_0000_0000_0000_0000_0000
 Adding Hidden Bit = 1.1000_0000_0000_0000_0000_0000

Therefore unsigned value = $1.1 \times 2^{-2} = 0.11$ (binary)
 = 0.25 + 0.125 = 0.375 (decimal)

Sign bit = 1 therefore number is -0.375

Example: Addition

> Carry out the addition 42.6875 + 0.375 in IEEE single precision arithmetic

Number	Sign	Exponent	Significand
42.6875	0	1000_0100	0101_0101_1000_0000_0000_0000
0.375	0	0111_1101	1000_0000_0000_0000_0000_0000

> To add these numbers the exponents of the numbers must be the same => Make the smaller exponent equal to the larger exponent, shifting the mantissa accordingly.

> Note: We must restore the Hidden bit when carrying out floating point operations.

Example: Addition Contd.

> Significand of Larger No = 1.0101_0101_1000_0000_0000_0000
 Significand of Smaller No = 1.1000_0000_0000_0000_0000_0000

> Exponents differ by +7 (1000_0100 - 0111_1101). Therefore shift binary point of smaller number 7 places to the left:

> Significand of Smaller No = 0.0000_0011_0000_0000_0000_0000
 Significand of Larger No = 1.0101_0101_1000_0000_0000_0000
 Significand of SUM = 1.0101_1000_1000_0000_0000_0000

> Therefore SUM = 1.0101_1000_1 $\times 2^5 = 10_1011.0001 = 43.0625$
 Sign Exponent Significand
 0 1000_0100 0101_1000_1 000_0000_0000_0000 = 422C 4000H

Special Values

- > The IEEE format can represent five kinds of values: Zero, Normalised Numbers, Denormalised Numbers, Infinity and Not-A-Numbers (NaNs).
- > For single precision format we have the following representations:

IEEE Value Field	Sign Field	Exponent Field	Significand Exponent	True
± Zero	0 or 1	0	0 (All zeroes)	
± Denormalised No	0 or 1	0	Any non-zero bit pat.	-126
± Normalised No	0 or 1	1 .. 254	Any bit pattern	-126 .. +127
± Infinity	0 or 1	255	0 (All zeroes)	
Not-A-Number	0 or 1	255	Any non-zero bit pat.	

Denormalised Numbers

- An Exponent of All 0's is used to represent Zero and Denormalised numbers, while All 1's is used to represent Infinities and Not-A-Numbers (NaNs)
- This means that the maximum range for normalised numbers is reduced, i.e. for Single Precision the range is -126 .. +127 rather than -127 .. +128 as one might expect for Excess 127.

- Denormalised Numbers represent values between the Underflow limits and zero, i.e. for single precision we have:

$$\pm 0.F \times 2^{-126}$$

Traditionally a "flush-to-zero" is done when an underflow occurs

- Denormalised numbers allow a more gradual shift to zero, and are useful in a few numerical applications

Eddie Edwards 2008

Floating Point Numbers

7.33

IEEE 754 floating point numbers - questions

Eddie Edwards 2008

Floating Point Numbers

7.34

- What decimal is represented by the hex word C0CA0000

Answer - -6.3125

- What hex word is -0.75 in IEEE-754?

Answer - BFE8000000000000

Eddie Edwards 2008

Floating Point Numbers

7.35

Infinities and NaN's

- Infinities (both positive & negative) are used to represent values that exceed the overflow limits, and for operations like Divide by Zero
- Infinities behave as in Mathematics, e.g.

$$\text{Infinity} + 5 = \text{Infinity}, -\text{Infinity} + -\text{Infinity} = -\text{Infinity}$$

- Not-A-Numbers (NaNs) are used to represent the results of operations which have no mathematical interpretation, e.g.

$$0 / 0, +\text{Infinity} + -\text{Infinity}, 0 \times \text{Infinity}, \text{Square root of a -ve number},$$

- Operations with a NaN operand yield either a NaN result (quiet NaN operand) or an exception (signalling NaN operand)

Eddie Edwards 2008

Floating Point Numbers

7.36

This lecture - feedback

▪ The pace of the lecture was:

A. much too fast B. too fast C. about right D. too slow E. much too slow

▪ The learning objectives were met:

A. Fully B. Mostly C. Partially D. Slightly E. Not at all