

# **MATH50003**

# **Numerical Analysis**

## **II.1 Reals**

**Dr Sheehan Olver**

# **Part II: Representing Numbers**

**How do computers compute with numbers?**

**Why are there errors, eg. in divided differences?**

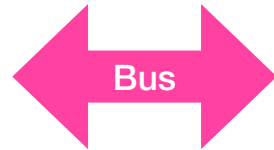
**Can we understand and bound these errors?**

# Simplified Model of a Computer

How do computers compute?

## Memory

Address (64-bits)	Data (8-bits)
000...000	11010101
000...001	00011101
000...010	10011001
⋮	⋮
111...111	11001110



## Registers

Address	Data (64-bits)
1	110...101
2	000...101
3	100...001
⋮	⋮
37	110...110

CPU

Input



Output

# Mathematical model

CPU's work on  $p$ -bits at a time

Cores take (1x or 2x)  $p$ -bits  
and return  $p$ -bits.

Operations are

$$f : \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p} \quad \text{or}$$

$$f : \mathbb{Z}_{2^p} \times \mathbb{Z}_{2^p} \rightarrow \mathbb{Z}_{2^p}$$

for  $\mathbb{Z}_m := \{0, 1, \dots, m-1\}$

$$2^p = 2^{64}$$

But how to handle  $\infty$ -cardinality sets  
integers/reals?

## Limitations

- Memory is finite
- All operations work on  $p$ -bits at a time
- No such thing as throwing an error
- Any operation that manipulates more than  $p$ -bits must be a composition of simpler functions

# Part II

## Representing Numbers

1. **Reals** via floating point
2. **Floating point arithmetic** and bounding errors
3. **Interval arithmetic** for rigorous computations

Appendix B Integers

# Ariane 5 rocket explosion

Learn floating point, or else...



## II.1.1 Real numbers in binary

We can represent any real number using binary digits

**Definition 3** (binary format). For  $B_0, \dots, B_p \in \{0, 1\}$  denote an integer in binary format by: bits base -2

$$\pm(B_p \dots B_1 B_0)_2 := \pm \sum_{k=0}^p B_k 2^k$$

Examples

$$(1)_2 = \overset{B_0}{1} \cdot 2^0 = 1$$

$$(10)_2 = 1 \cdot 2^1 + 0 \cdot 2^0 = 2$$

$$-(11)_2 = -(1 \cdot 2^1 + 1 \cdot 2^0) = -3$$

$$(10001)_2 = \underbrace{2^4}_{16} + 0 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + \underbrace{1 \cdot 2^0}_1 = 17$$

**Definition 4** (real binary format). For  $b_1, b_2, \dots \in \{0, 1\}$ , Denote a non-negative real number in *binary format* by:

$$(B_p \dots B_0 . b_1 b_2 b_3 \dots)_2 := (B_p \dots B_0)_2 + \sum_{k=1}^{\infty} \frac{b_k}{2^k}.$$

$$(101.101)_2 = 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} \\ = 5.625.$$

**Example 3** (rational in binary).

Consider

$$\frac{1}{3} = 0.333\ldots = \sum_{k=1}^{\infty} \frac{3}{10^k}$$

Claim:

$$\frac{1}{3} = (0.010101\ldots)_2 = \sum_{k=1}^{\infty} \frac{1}{2^{2k}}$$

$\underbrace{\quad}_{1/4} \quad \underbrace{\quad}_{1/16} \quad \underbrace{\quad}_{1/64}$



Use Geometric Series:

$$\sum_{k=0}^{\infty} z^k = \frac{1}{1-z} \quad \text{for } |z| < 1$$

Write  $z = \frac{1}{4}$ , so that

$$-1 + 1 + \sum_{k=1}^{\infty} \frac{1}{(z^2)^k} = \sum_{k=0}^{\infty} z^k - 1 = \frac{1}{1-\underbrace{z}_{1/4}} - 1$$

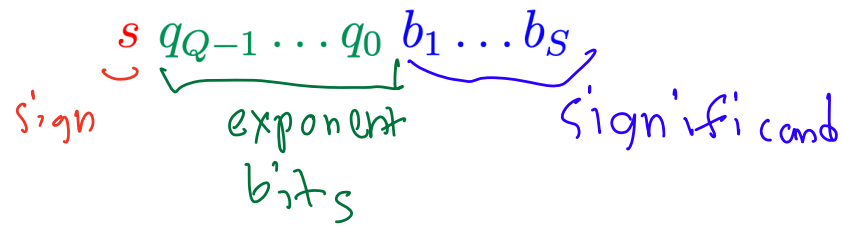
$$= \frac{1}{3}$$



## II.1.2 Floating-point numbers

How do we represent an uncountable set with only  $p$ -bits?

Bit Format:



# bits:  $p = 1 + Q + S$

**Definition 5** (floating-point numbers). Given integers  $\sigma$  (the *exponential shift*),  $Q$  (the number of *exponent bits*) and  $S$  (the *precision*), define the set of *Floating-point numbers* by dividing into normal, sub-normal, and *special number* subsets:

$$F_{\sigma,Q,S} := \underbrace{F_{\sigma,Q,S}^{\text{normal}}}_{\subset \mathbb{R} \setminus \{0\}} \cup \underbrace{F_{\sigma,Q,S}^{\text{sub}}}_{\subset \mathbb{R} \cup \{-0\}} \cup \underbrace{F^{\text{special}}}_{\{-\infty, \infty, \text{NaN}\}}.$$

How do bits dictate whether its normal/sub/special?

Look at **exponent**. 3 examples:

Normal  
0 10000 1010000000

3.125

Sub  
1 00000 1100000000  
 $\underbrace{\hspace{1cm}}_{=0}$   
 $-3 \times 2^{-16}$


Special  
1 11111 0000000000  
 $\underbrace{\hspace{1cm}}_{=1}$   
 $-\infty$

for  $F_{15,5,10}$ .

## II.1.3 IEEE float-point numbers


What exponent shift/number of bits/precision is used in practice?

Half  
(ML)

$$F_{16} := F_{15,5,10}$$



DON'T

Single  
(GPUs)

$$F_{32} := F_{127,8,23}$$


MEMORISE

Double  
(Standard)

$$F_{64} := F_{1023,11,52}$$


Half-precision  
 $F_{16} := F_{15,5,10}$

$$F_{\sigma,Q,S}^{\text{normal}} := \{\pm 2^{q-\sigma} \times (1.b_1b_2b_3\dots b_S)_2 : 1 \leq q < 2^Q - 1\}.$$

$\nwarrow$  is a 1 by default.

**Example 4** (interpreting 16-bits as a float). Consider the number with bits

$\overset{q_4 - q_0}{\text{0} \text{ 10000} \text{ 1010000000}}$   
 $\text{0 implies} \quad \text{encode} \quad b_1 b_2 \dots b_S$   
 $+$   
 $q = (10000)_2 = 2^4 = 16$

$\Rightarrow$  encodes

$$\begin{aligned}
 &+ 2^{16-15} \times (1.101000000000)_2 \\
 &= 2 \times \left(1 + \frac{1}{2} + \frac{1}{8}\right) = 3.25
 \end{aligned}$$

**Example 5** (rational to 16-bits). How is the number  $1/3$  stored in  $F_{16}$ ?

$$\frac{1}{3} = (0.0101\_\_\_\_)_2$$

$$= 2^{-1} \times (1.0101\_\_\_\_)_2$$

?   
 we want = 1

$$= 2^{13-15} \times (1.0101\_\_\_\_)_2 \notin F_{16}$$

$$\approx + 2^{13-15} \times (1.0101010101)_2$$

round  
(more detail)  
later

Bits: 0 01101 0101010101

$$(01101)_2 = 13$$

## II.1.4 Sub-normal and special numbers

Sub-normal have exponent bits all 0, special have all 1

If  $q = (00000)_2 = 0$ , this becomes 1, not  $q$

$$F_{\sigma, Q, S}^{\text{sub}} := \{ \pm 2^{1-\sigma} \times (0.b_1 b_2 b_3 \dots b_S)_2 \}.$$

$\uparrow$   
 0, not 1 anymore

**Example 6** (subnormal in 16-bits). Consider the number with bits

$\overbrace{1\ 00000}^{\text{all 0}}$   $\underbrace{1100000000}_{b_1 b_2 \dots b_{10}}$   
 means

encodes

$$\begin{aligned}
 - 2^{1-15} \times (0.11)_2 &= - 2^{-14} \left( \frac{1}{2} + \frac{1}{4} \right) \\
 &= - 3 \times 2^{-16}
 \end{aligned}$$





$$F^{\text{special}} := \{\infty, -\infty, \text{NaN}\}$$

When  $q = (11 \dots 1)_2$  number is special

If all  $b_k = 0$ ,  $\pm \infty$

If any  $b_k \neq 0$  NaN (LOTS OF NaNs!)

**Example 7** (special in 16-bits). The number with bits

1 11111 0000000000  
 - 11111 0000000000

$-\infty$

On the other hand, the number with bits

1 11111 0000000001  
 11111 0000000001  
 $b_{10} \neq 0$

NaN

**Time for code.**