

How to represent real numbers

In decimal scientific notation: 5.28×10^3 has 3 parts

- ❖ sign (i.e. +)
- ❖ fraction a/k/a mantissa (i.e. 5.28)
- ❖ base (i.e., 10) to some power (i.e. 3)

Most of the time, the usual representation has one digit to the left of decimal point

- ❖ Example: -0.1234×10^6

A number is *normalized* if the leading digit is not 0

- ❖ Example: -1.234×10^5

Real representation inside a computer

Use a representation akin to scientific notation
 $sign \times mantissa \times base^{exponent}$

where *sign* is 1 or -1

Many variations in choice of representation for

- ❖ mantissa (could be 2's complement, sign and magnitude etc.)
- ❖ base (could be 2, 8, 16 etc.; 10 is a little difficult)
- ❖ exponent (like mantissa)

Arithmetic support for real numbers is called
floating-point arithmetic

What is the 520 bridge called?

3

Floating-point representation: IEEE Std

Floating Point is necessarily an approximation

Basic choices

- ❖ A single precision number must fit into 1 word (4 bytes, 32 bits)
- ❖ A double precision number must fit into 2 words
- ❖ The base for the exponent is 2
- ❖ There should be approximately as many positive and negative exponents

Additional criteria

- ❖ The mantissa will be represented in sign and magnitude form
- ❖ Numbers will be normalized

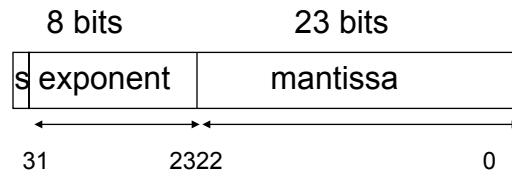
This standard was significant milestone

4

Ex: MIPS representation of IEEE Standard

A number is represented as : $(-1)^S \cdot F \cdot 2^E$

In single precision the representation is:



9.021×10^4	-9.021×10^4
5.28×10^{-3}	-5.28×10^{-3}

5

MIPS representation (continued)

Bit 31 sign bit for mantissa (0 pos, 1 neg)

Exponent 8 bits ("biased" exponent, see next slide)

mantissa 23 bits : always a *fraction* with an *implied binary point* at left of bit 22

Number is *normalized* (see implication next slides)

0 is represented by all zero's ... wouldn't have to be

Note that having the most significant bit as sign bit makes it easier to test for positive and negative

6

Biased exponent

Biased notation is yet another technique for representing signed numbers

The “middle” of the range will represent 0. For eight bits 0 is 127 (01111111)

All exponents starting with a “1” will be positive exponents

❖ Example: 10000001 is exponent 2 (10000001 - 01111111)

All exponents starting with a “0” will be negative exponents

❖ Example 01111110 is exponent -1 (01111110 - 01111111)

The largest positive exponent will be 11111111, about 10^{38}

The smallest negative exponent is about 10^{-38}

7

Normalization

Since numbers must be normalized, there is always a “one” at the left of the binary point:

1.0011011100001

No need to put it in (improves precision by 1 bit)

Often referred to as the “implied 1”

But need to reinstate it when performing operations

In summary, in MIPS a floating-point number has the value:

$$(-1)^S \cdot (1 + \text{mantissa}) \cdot 2^{(\text{exponent} - 127)}$$

8

Double precision

Takes 2 words (64 bits)

Exponent 11 bits (instead of 8) ... what's the bias?

Mantissa 52 bits (instead of 23)

Still biased exponent and normalized numbers

Still 0 is represented by all zeros

We can still have *overflow* (the exponent cannot handle super big numbers) and *underflow* (the exponent cannot handle super small numbers)

9

Floating-Point Addition

Quite “complex” (more complex than fl multiplication)

Need to know which of the addends is larger (compare exponents)

Need to shift “smaller” mantissa

Need to know if mantissas have to be added or subtracted (since it's a sign/magnitude representation)

Need to normalize the result

Correct round-off procedures are not simple (not covered in detail here)

10

One of the 4 round-off modes

Round to nearest even

❖ Example 1: in base 10. Assume 2 digit accuracy

$$3.1 * 10^0 + 4.6 * 10^{-2} = 3.146 * 10^0$$

clearly should be rounded to $3.1 * 10^0$

❖ Example 2:

$$3.1 * 10^0 + 5.0 * 10^{-2} = 3.15 * 10^0$$

By convention, round-off to nearest “even”
number $3.2 * 10^0$

Other round-off modes: towards 0, $+\infty$, $-\infty$

11

F-P add (details for round-off omitted)

1. Compare exponents. If $e_1 < e_2$, swap the 2 operands such that
 $d = e_1 - e_2 \geq 0$. Tentatively set exponent of result to e_1
2. Insert 1's at left of mantissas. If the signs of operands differ, replace 2nd mantissa by its 2's complement.
3. Shift 2nd mantissa d bits to the right (this is an arithmetic shift, i.e., insert either 1's or 0's depending on the sign of the second operand)
4. Add the (shifted) mantissas. (There is one case where the result could be negative and you have to take the 2's complement; this can happen only when $d = 0$ and the signs of the operands are different.)

12

F-P Add (continued)

5. Normalize (if there was a carry-out in step 4, shift right once; else shift left until the first “1” appears on msb)
6. Modify exponent to reflect the number of bits shifted in previous step

13

Example

Add decimal: $0.375 + 0.75$

$$3/2^3 + 3/2^2 = 0.011 + 0.11 = 1.1 \times 2^{-2} + 1.1 \times 2^{-1}$$

Now add:

$$\text{Align fractions: } 0.11 \times 2^{-1} + 1.1 \times 2^{-1}$$

$$\text{Add fractions: } 10.01 \times 2^{-1}$$

$$\text{Normalize: } 10.01 \times 2^{-1} = 1.001 \times 2^0$$

Round: Not needed

$$1.001 \times 2^0 = 1 + 1/2^3 = 1 + 1/8 = 1.125 \text{ decimal}$$

In IEEE single precision 0.75 i.e. 0.11 is

0 0111 1110 100 0000 0000 0000 0000 why?

14

Using pipelining

Stage 1

- ❖ Exponent compare

Stage 2

- ❖ Shift and Add

Stage 3

- ❖ Round-off , normalize and fix exponent

Most of the time, done in 2 stages

15

Floating-point multiplication

Conceptually easier

1. Add exponents (careful, subtract one “bias”)
2. Multiply mantissas (don't have to worry about signs)
3. Normalize and round-off and get the correct sign

16

Special Values

Allow computation to continue in face of exceptional conditions

❖ For example: divide by 0, overflow, underflow

Special value: NaN (Not a Number; e.g., $\text{sqrt}(-1)$)

❖ Operations such as $1 + \text{NaN}$ yield NaN

Special values: $+\infty$ and $-\infty$ (e.g., $1/0$ is $+\infty$)

Can also use “denormal” numbers for underflow and overflow allowing a wider range of values