Floating-point numbers

Fractional binary numbers
IEEE floating-point standard
Floating-point operations and rounding
Lessons for programmers

Many more details we will skip (it's a 58-page standard...) See CSAPP 2.4 for a little more

Fractional Binary Numbers

Value Representation

5 and 3/4
2 and 7/8
47/64
101.11₂
10.111₂
0.101111₂

Observations

Shift left = Shift right =

Numbers of the form $0.1111111..._2$ are...?

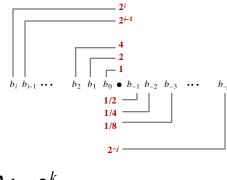
Limitations:

Exact representation possible only for numbers of the form x * $2^{\nu}\!\!\!/,$ where x and y are integers.

Other rationals have repeating bit representations

 $1/3 = 0.333333..._{10} = 0.01010101[01]..._{2}$

Fractional Binary Numbers



$$\sum_{k=-j}^{i} b_k \cdot 2^{j}$$

Fixed-Point Representation

Implied binary point. Example:

b₇ b₆ b₅ b₄ b₃ [.] b₂ b₁ b₀

Same hardware as for integer arithmetic.

 $b_7 b_6 b_5 b_4 b_3 b_2 b_1 b_0$ [.]

Fixed point = fixed range and fixed precision

range: difference between largest and smallest representable numbers precision: smallest difference between any two representable numbers

IEEE Floating Point

Analogous to scientific notation

12 000 000 1.2 x 10⁷ 1.2e7 0.000 001 2 1.2 x 10⁻⁶ 1.2e-6

IEEE Standard 754 used by all major CPUs today

IEEE = Institute of Electrical and Electronics Engineers

Driven by numerical concerns

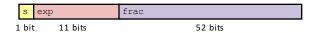
Rounding, overflow, underflow
Numerically well-behaved, but hard to make fast in hardware

Precisions

Single precision (float): 32 bits



Double precision (double): 64 bits



Finite representation of infinite range:

Not all values can be represented exactly.

Some are approximated.

Floating Point Representation

Numerical form:

$$V_{10} = (-1)^{S} * M * 2^{E}$$

Sign bit s determines whether number is negative or positive
Significand (mantissa) M normally a fractional value in range [1.0,2.0)
Exponent E weights value by a (possibly negative) power of two

Representation:

```
MSB s = sign bit s

exp field encodes E (but is not equal to E)

frac field encodes M (but is not equal to M)
```



 $V = (-1)^{S} * M * 2^{E}$ s exp

Normalization and Special Values

+inf, -inf: exp == 11...1 frac == 00...0 1.0/0.0 = -1.0/-0.0 = +inf, 1.0/-0.0 = -1.0/0.0 = -inf NaN ("Not a Number"): exp == 11...1 frac!=00...0 exp == 11...1

Denormalized/subnormal values (near 0.0) not covered here.

Floating Point Arithmetic

double
$$x = \dots$$
, $y = \dots$; double $z = x + y$;

- 1. Compute exact result.
- 2. Round, to fit:

 $V = (-1)^{S} * M * 2^{E}$

Overflow exponent if it is too wide for **exp**.

Drop LSBs of significand if it is too wide for **frac**.

Underflow if nearest representable value is 0.

frac

...

Lessons for programmers

```
float \( \neq \text{real number} \neq \text{double} \)

Rounding breaks associativity and other properties.

double \( a = \cdots \cdots \cdot b \)

if \( (abs(a - b) < epsilon) \cdots \cdots \)
```

More shortly...