

Setting a bit of a variable to 1.

- use bitwise OR with a 1

$$\text{OR} \quad \begin{array}{c} ? \\ \hline 1 \\ \hline 1 \end{array} \quad \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \quad \begin{array}{c} 1 \\ 0 \\ 1 \end{array}$$

OR'ing with a zero doesn't change anything

Example: set the rightmost bit of x to 1.

int $x = \dots$;

$x = x | 1$;

$$\begin{array}{r} x: 0 \dots 10100 \\ 1: 0 \dots 00001 \\ \hline 0 \dots 10101 \end{array}$$

changed to 1
remain the same.

Set the m^{th} bit of x to 1:

$$x = x | \underbrace{(1 \ll m)}_{\text{mask}}; \quad // \text{mask has a 1 in the } m^{\text{th}} \text{ position}$$

"Clear" a bit to 0.

- use bitwise AND with a 0.

$$\& \quad \begin{array}{c} ? \\ \hline 0 \\ \hline 0 \end{array} \quad \begin{array}{c} 0 \\ 1 \\ 0 \end{array} \quad \begin{array}{c} 1 \\ 1 \\ 1 \end{array}$$

AND'ing anything with 1 has no change.

Clearing the rightmost bit of x (to 0).

- mask would have to be $11 \dots 110$

How to create that mask?

① #define MYMASK 0xFFFFFFFF

generally better \rightarrow ② #define MYMASK ~1 // complement of 1 000...01

$x = x \& \text{MYMASK};$

Clear the m^{th} bit of x .

$x = x \& \sim(1 \ll m);$

mask: $11 \dots 101 \dots 1$
 ↑
 bit m

bit m
 $00 \dots 010 \dots 0$

Binary Arithmetic

- already seen addition.

- subtraction:

$$x - y \equiv x + (-y)$$

flip the bits of y and add 1.

- multiplication

- same way as humans, except in binary not decimal.

Decimal:

$$\begin{array}{r} 3^3 3^4 5 \\ 27 \\ \hline 2415 \\ 690 \\ \hline 9315 \end{array}$$

Binary:

11010

← 26 decimal

110

← 6 decimal

0

11010

11010

10011100

← 156 decimal

DIVISION :

- Done the same way as humans do.
- binary not decimal

Decimal: 129

$$\begin{array}{r} 56 \overline{) 7253} \\ \underline{56} \\ 165 \\ \underline{112} \\ 533 \\ \underline{504} \\ 29 \end{array}$$

129 remainder = 29

Binary:

$$\begin{array}{r} 11110 \quad \leftarrow 30 \text{ decimal} \\ 101 \overline{) 10011010} \quad \leftarrow 154 \text{ decimal} \\ \underline{101} \\ 01001 \\ \underline{101} \\ 01000 \\ \underline{101} \\ 111 \\ \underline{101} \\ 100 \quad \leftarrow \text{remainder} \\ \quad \quad \quad \leftarrow 4 \text{ decimal} \end{array}$$

5 decimal \nearrow

IEEE Floating Point Numbers

- approximation to real numbers
- 32 bits "float"
- 64 bits "double"

32 bit IEEE floating point number:

- specified using scientific notation in binary.

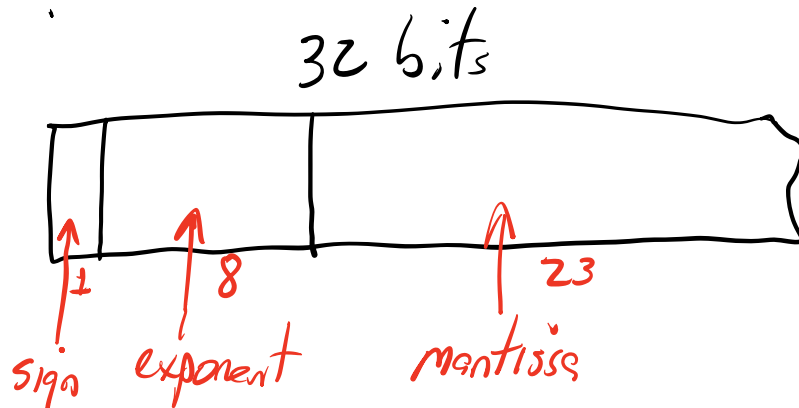
In decimal: 6.02×10^{23}
 -3.15×10^{15}
 2.01×10^{-7}

mantissa → 6.02
exponent → 23

"normalized"
- one digit before the point.

In binary: 1.0101×2^{10}
sign → -1.11×2^{-101}

The sign, exponent, and the mantissa are represented as fields in the 32-bit number.



Sign: $1 = \text{negative}$, $0 = \text{positive}$

Exponent: A bias of 127 is used.

– the number stored in the exponent field is the actual exponent plus 127.

– so actual exponent is the stored exponent $- 127$.

If stored exponent is 13, then actual exponent is -114 .

Mantissa:

- Must be normalized,
^{non-zero}
so one digit before
the point.

$$\cancel{1010.1 \times 2^5}$$

$$1.0101 \times 2^8$$

$$\cancel{0.1 \times 2^{-3}}$$

$$1.0 \times 2^{-4}$$

↑ ,

- there will always be a
1 before the decimal
point.

- except for the number
0.0

- in the stored mantissa,
the leading 1 is
implicit.

only the digits after
the point are stored

~~1.01011010...0~~

.01011010...0

Zero is represented specially
as all 0's (all 32 bits)

000...0
32 bits