

---

---

---

---

---

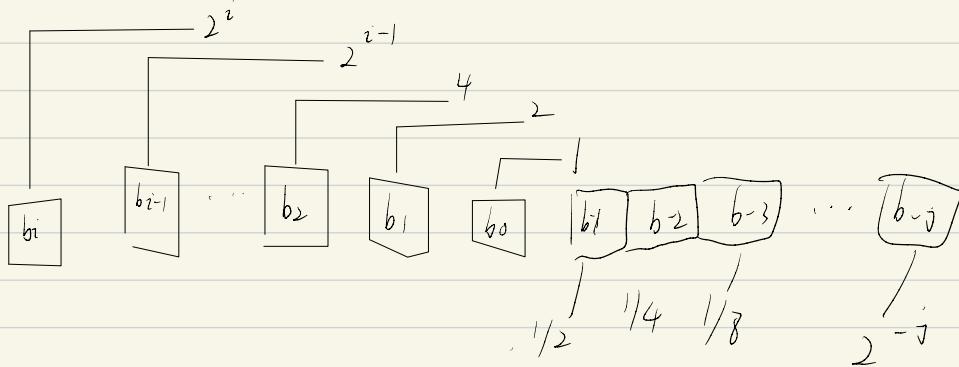


Floating point:

Fractional binary numbers

(I) What is  $1011.10_2$ ?

## Fractional Binary Numbers:



(i) Representation:

(i) Bits to right of "binary point" represent fractional powers of 2,

(ii) Represents rational numbers.

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

e.g.	Value	Representation
$5 \frac{3}{4}$		$101.11_2$
$2 \frac{7}{8}$		$10.111_2$
$1 \frac{7}{16}$		$1.011_2$

Observations:

(i) Divide by 2 by shifting right (unsigned)

(ii) Multiply by 2 by shifting left.

(iii) Numbers of form  $0.11111\ldots_2$  are just below 1.0.

$$(\star) \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots + \frac{1}{2} \rightarrow 1.0$$

(\*) Use notation  $1.0 - \epsilon$

## Limitation #1

can only exactly represent numbers of the form  $x/2^k$

Other rational numbers have repeating b.t representations

Value	Representations
-------	-----------------

1/3	0.01010101[01]... <sub>2</sub>
-----	--------------------------------

1/5	0.001100110011[0011]... <sub>2</sub>
-----	--------------------------------------

1/10	0.0001100110011[0011]... <sub>2</sub>
------	---------------------------------------

## Limitation #2.

Just one setting of binary point within the n bits

Limited range of numbers (very small values? very large?)

floating point ← trade-off

a large range

of values

much precision given

the number of bits.

# IEEE floating Point

## IEEE standard 754

Established in 1985 as uniform std. for floating point arithmetic.  
Before that, many idiosyncratic formats.  
(peculiar or individual)

### floating point Representation:

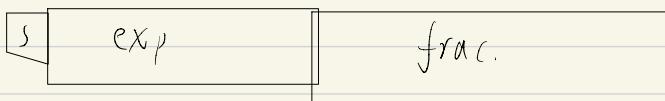
#### (I) Numerical Form:

$$(-1)^s M \cdot 2^E \quad (\text{Mantissa})$$

- (i) Sign bit's determines whether number is negative or positive.
- (ii) Significand  $M$  normally a fractional value in range [1.0, 2.0).
- (iii) Exponent  $E$  weights value by power of two.

#### (II) Encoding

- (i) MSB is sign bit  $s$ .
- (ii) exp field encodes  $E$  (but is not equal to  $E$ )
- (iii) frac field encodes  $M$  (but is not equal to  $M$ )



Single Precision    1-bit    8-bits    23-bits  
(32 bits)

Double Precision    1-bit    11-bits    52-bits,  
(64 bits)

Extended Precision    1-bit    15-bits    63 or 64-bits,  
(80 bits)

"Normalized" Values

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

(I) When :  $\text{exp} \neq 000\ldots0$  and  $\text{exp} \neq 111\ldots1$

(II) Exponent coded as a biased value:  $E = \text{Exp} - \text{Bias}$ .

(i) Exp: unsigned value of exp. field.

(ii) Bias =  $2^{k-1} - 1$ , where  $k$  is number of exponent bits.

(\*) Single precision: 127 ( $\text{Exp}: 1\ldots254$ ,  $E: -126\ldots127$ )  
(\*) Double precision: 1023 ( $\text{Exp}: 1\ldots2046$ ,  $E: -1022\ldots1023$ )

(III) significand coded with implied leading 1:  $M = 1.x_1x_2\ldots x_n$

(i)  $x_1x_2\ldots x_n$ : bits of frac field. ( $n=1..0$ )

(ii) Maximum when  $\text{frac} = 111\ldots1$  ( $M = 2.0 - \epsilon$ )  
get extra leading bit for "free"

Value: float  $F = 15213.0$ .

$$\begin{aligned} 15213_{10} &= 111011011010_2 \\ &= 1.11011011010_2 \times 2^{13} \end{aligned}$$

Significant:

$$M = 111011011010_2$$

$$\text{frac} = 1011010100000000000_2$$

throw this away

$$\begin{aligned} \text{Exponent: } E &= 13 & \text{Bias} &= 127 \Rightarrow \text{Exp} = 140 = 10001100_2 \end{aligned}$$

0

1 0 0 0 1 1 0 0

1 1 0 1 0 1 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0

]

Exp

frac

$$0 \leq \text{Exp} \leq 2^8 - 1$$

$$0 \leq \text{Exp} \leq 255 \quad \text{shift.}$$

$$-127 \leq E \leq (255 - 127 = 128)$$

Denormalized Value

$$v = (-1)^s M_2^E$$

$$E = 1 - \text{Bias}$$

(I) Condition:  $\exp = 000\ldots 0$

(II) Exponent value:  $E = 1 - \text{Bias}$  (instead of  $E = 0 - \text{Bias}$ )

(III) Significant coded with implied leading 0:  $M = 0.XXX\ldots X_2$

(IV) Cases:

(i)  $\exp = 000\ldots 0, \quad \text{frac} = 000\ldots 0$

\* Represents zero value

\* Note distinct values: +0 and -0 (why?)

(ii)  $\exp = 000\ldots 0, \quad \text{frac} \neq 000\ldots 0$

\* Numbers closest to 0.0

\* Equispaced.

## Special Values

(I) Condition:  $\exp = 111\dots 1$

(II) Case:  $\exp = 111\dots 1$ ,  $\text{frac} = 000\dots 0$

(i) Represents value  $\infty$  (infinity)

(ii) Operation that overflows

(iii) Both positive and negative

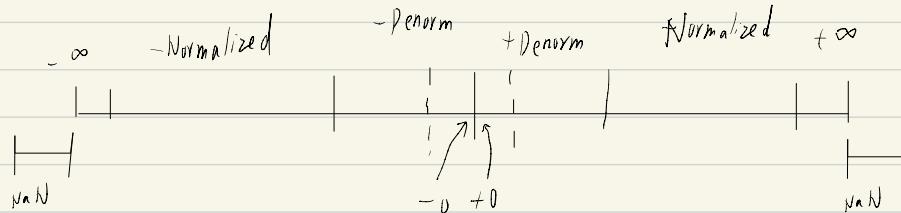
(iv) E.g.  $1.0/0.0 = +\infty$ ,  $1.0/-0.0 = -\infty$

(III) Case:  $\exp = 111\dots 1$ ,  $\text{frac} \neq 000\dots 0$

(i) Not-a-Number (NaN)

(ii) Represents  $\text{NaN}$  when no numeric value can be determined.

(iii) E.g.  $\sqrt{-1}$ ,  $\infty - \infty$ ,  $\infty \times 0$



## Today: Floating Point:

- (I) BG: Fractional binary numbers
- (II) IEEE floating point std.: Definition.
- (III) Example and properties
- (IV) Rounding, addition, multiplication
- (V) Floating point in C
- (VI) Summary

e.g.

## Tiny Floating Point Example

s	exp	frac
1	4-bits	3-bits

(I) 8 bits Floating Point Representation:

- (i) the sign bit is in the most significant bit
- (ii) the next four bits are the exponent, with a bias of 7
- (iii) the last three bits are the frac

(II) same general form as IEEE Format:

- (i) normalized, denormalized
- (ii) representation of 0, NaN, infinity

# Dynamic Range (Positive Only)

	s	exp	frac	E	Value	
Denormalized numbers	0	0000	000	-6	0	
	0	0000	001	-6	$1/8 * 1/64 = 1/512$	closest to zero.
	0	0000	010	-6	$2/8 * 1/64 = 2/512$	
	0	0000	110	-6	$6/8 * 1/64 = 6/512$	
	0	0000	111	-6	$7/8 * 1/64 = 7/512$	largest denom
Normalized numbers	0	0001	000	-6	$8/8 * 1/64 = 8/512$	smallest norm
	0	0001	001	-6	$9/8 * 1/64 = 9/512$	
	0	0110	110	-1	$14/8 * 1/2 = 14/16$	closest to 1
	0	0110	111	-1	$15/8 * 1/2 = 15/16$	below
	0	0111	000	0	$8/8 * 1 = 1$	
	0	0111	001	0	$9/8 * 1 = 9/8$	closest to 1
	0	0111	010	0	$10/8 * 1 = 10/8$	above
	0	1110	110	7	$14/8 * 128 = 224$	
	0	1110	111	7	$15/8 * 128 = 240$	largest norm
	0	1111	000	N/A	inf	

## Distribution of Values.

(I) 6-bit IEEE-like format

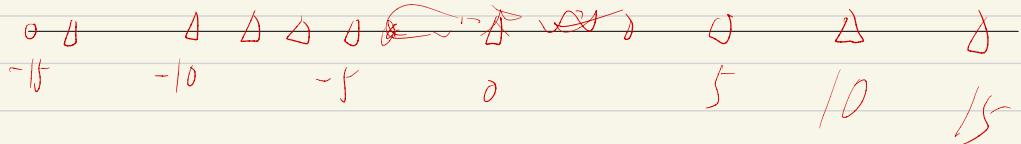
e = 3 exponent bits

f = 2 fraction bits

$$\text{Bias: } 2^{3-1} - 1 = 3.$$

s	exp	frac
	3-bits	2-bits

(II) Notice how the distribution gets denser toward zero.



## Special properties of the IEEE Encoding:

(I) FP Zero same as Integer Zero

(i) All bits = 0

(II) (an almost) Use Unsigned Integer comparison.

(i) Must first compare sign bits

(ii) Must consider  $-0 = 0$ .

(iii) NaNs problematic

Will be greater than any other values.

What should comparison yield?

(iv) otherwise Ok

Denorm vs. normalized,  
normalized vs. infinity

## Floating Point Operations: Basic Idea.

$$x + f y = \text{Round}(x+y)$$

$$x x, y = \text{Round}(x x y)$$

(I) Basic idea:

- (i) First compute exact result
- (ii) Make it fit into desired precision.

\* possibly overflow if exponent too large  
\* possibly round to fit into fraction

## Rounding

Rounding modes (illustration in \$ rounding)



• towards zero

\$1.40	\$1.60	\$1.50	\$2.50	\$ -1.50
\$1	\$1	\$1	\$2	-\$1

Round down ( $-\infty$ )

\$1	\$1	\$1	\$2	-\$2
\$1	\$1	\$1	\$2	-\$2

Round up ( $+\infty$ )

\$2	\$2	\$2	\$3	-\$1
\$2	\$2	\$2	\$3	-\$1

Nearest Even (default)

\$1	\$2	\$2	\$2	\$2
\$1	\$2	\$2	\$2	\$2

nearest even

↑  
if equ. to even number

Closer look at Round-to-Even.

(I) Default Rounding Mode.

- (i) Hard to get any other kind without dropping into assembly.
- (ii) All others are statistically biased.

A sum of set of positive numbers will consistently be over- or under-estimated.

(II) Applying to other Decimal Places/Bit Operations

When exactly halfway between two possible values.

Round so that least significant digit is even.

e.g. round to the nearest hundredth.

7.8949999

7.89

(Less than half way)

7.895000

7.90

(Greater than half way)

7.895000

7.90

(Half way - round up)

7.895000

7.88

(Half way - round down)

# Rounding Binary Numbers,

## (i) Binary Fractional Numbers

- (i) "Even" when least significant bit is 0
- (ii) "Half way" when bits to right of rounding position = 100

## (II) Examples

(i) Round to nearest  $\frac{1}{4}$  (2 bits right of binary point)

Value	Binary	Rounded	Action	Rounded Value
$2 \frac{3}{32}$	$10.00\text{0}\cancel{1}_2$	$10.00_2$	( $< 1/2$ - down)	2
$2 \frac{3}{16}$	$10.00\cancel{1}0_2$	$10.01_2$	( $> 1/2$ - up)	$2 \frac{1}{4}$
$2 \frac{7}{8}$	$10.\cancel{1}100_2$	$11.00_2$	( $1/2$ - up)	3
$2 \frac{5}{8}$	$10.10\cancel{1}0_2$	$10.10_2$	( $1/2$ - down)	$2 \frac{1}{2}$

## Fp Multiplication.

$$(2) (-1)^{s_1} M_1 2^{E_1} \times (-1)^{s_2} M_2 2^{E_2}$$

(II) Exact Result:  $(-1)^{s_1+s_2} M_1 M_2 2^{E_1+E_2}$

Sign s:  $s_1 \wedge s_2$

Significand M:  $M_1 \times M_2$

Exponent E:  $E_1 + E_2$ .

## (III) Fixing:

- (i) If  $M \geq 2$ , shift M right, increment E
- (ii) If E out of range, overflow
- (iii) Round M to fit frac precision.

## (IV) Implementation.

- (i) Biggest chore is multiplying significands.

## Floating Point Addition:

$$(-1)^{s_1} M_1 2^{E_1} + (-1)^{s_2} M_2 2^{E_2}$$

Assume  $E_1 > E_2$

(I) Exact Result:  $(-1)^s M \cdot 2^E$

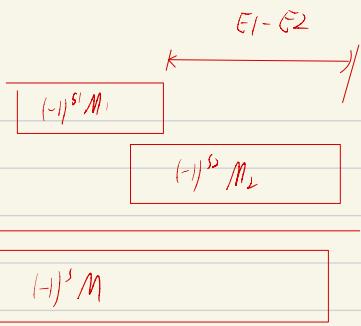
(i) sign  $s$ , significand  $M$ :

Result of signed align & add

Exponent  $E: E_1$

## (II) Fixing:

- (i) If  $M > 1$ , shift  $M$  right, increment  $E$ .
- (ii) If  $M < 1$ , shift  $M$  left  $k$  positions, decrement  $E$  by  $k$ .
- (iii) Overflow if  $E$  out of range
- (iv) Round  $M$  to fit frac precision.



## Mathematic Properties of FP Add

(I) Compare to those of Abelian group.

(i) Closed under addition Yes

But may generate infinity or NaN

(ii) Commutative? Yes

(iv) Associative? No

• Overflow and inexactness of rounding.

$$\bullet (3.14 + 1e10) - 1e10 = 0, \quad 3.14 + (1e10 - 1e10) = 3.14.$$

(vi) 0 is additive identity? Yes

(vii) Every element has additive inverse

Yes, but except for infinities & NaNs

## (II) Monotonicity

(i)  $a \geq b \Rightarrow a+c \geq b+c$  Almost

Except for infinities & NaNs

# Mathematical Properties of FP Mult?

extra properties?

Multiplication distributes over addition? No.

Possibility of overflow, inexactness of rounding,

$$\text{A } 1e20 * (1e20 - 1e20) = 0.0, \quad 1e20 * 1e20 - 1e20 \neq 1e20 = 0$$

Monotonicity,

$$a \geq b \oplus c > 0 \quad a * c \geq b * c \quad \text{Almost}$$

Except for NaNs & infinities

## Floating Point in C:

(i) C guarantees two levels;

float              single precision  
double            double precision

(II) Conversion / casting:

(i.) casting between int, float, and double changes bit representation,

(ii.) double / float → int

truncate fractional part

tie rounding toward zero.

Not defined when out of range or NaN; generally sets to TM.

(iii.) int → double

Exact conversion, as long as int has ≤ 53 bit word size.

(iv.) int ≠ float

will round according to rounding mode.

int x = ... float f = ... double d = ...

Assume neither d nor f is NaN.

x == (int)(float)x      TRUE

x == (int)(double)x      TRUE

f == (float)(double)f      TRUE

d == (double)(float)d      FALSE

f == -(-f)      TRUE

2/3 == 2/3.0      FALSE

d < 0.0       $\Rightarrow (d \times 2) < 0.0$       TRUE

d > f       $\Rightarrow -f > -d$       TRUE (monotonicity)

d \* d > 0.0 (TRUE)

(d + f) - d == f (FALSE)

summary:

- (i) IEEE Floating Point has clear mathematical properties.
- (ii) Represents numbers of form  $M \times 2^E$
- (iii) One can reason about operations independent of implementation
  - \* As if computed with perfect precision and then rounded.
- (iv). Not the same as real arithmetic.
  - | (i) violates associativity
  - | (ii) makes life difficult for compilers & serious numerical applications programmers