

IEEE 754

Floating point numbers.

fixed point.

exponent

1011.1101

mantissa

-0.072 x 10

4

i

4

f

2^{-f}

Reduction

~~-0.7~~

Normalized →

-1.72 x 10⁻²

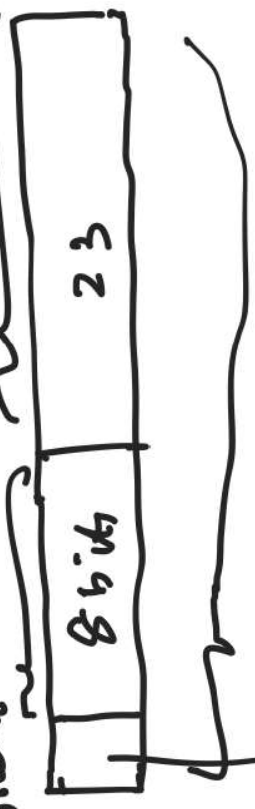
IEEE 754.

1985

Single point representation -
 Exponent
 mantissa

19.59375₁₀

float a;

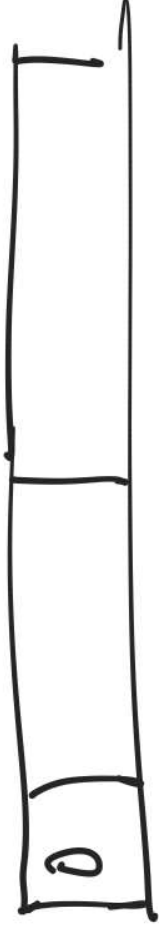


(1) Determine the Sign bit.

Sign bit $\rightarrow 0$

(for negative number $\rightarrow 1$)

Sign



(2)

$$\begin{array}{r} 2 \overline{) 19} \\ 2 \overline{) 14} \\ 2 \overline{) 4} \\ 2 \overline{) 2} \\ 2 \overline{) 0} \\ \hline 0 \end{array}$$

$$\begin{array}{r} 0.59375 \times 2 = 1.1875 \\ 1.1875 \times 2 = 2.375 \\ 2.375 \times 2 = 4.75 \\ 4.75 \times 2 = 9.5 \\ 9.5 \times 2 = 19.0 \end{array}$$

Binary representation:

$[10011, 10011]$

$i=7, f=5$

Step 3

Normalize

(place the binary point after leftmost 1)

$[1.001110011 \times 2^4]$

mantissa,

exponent.

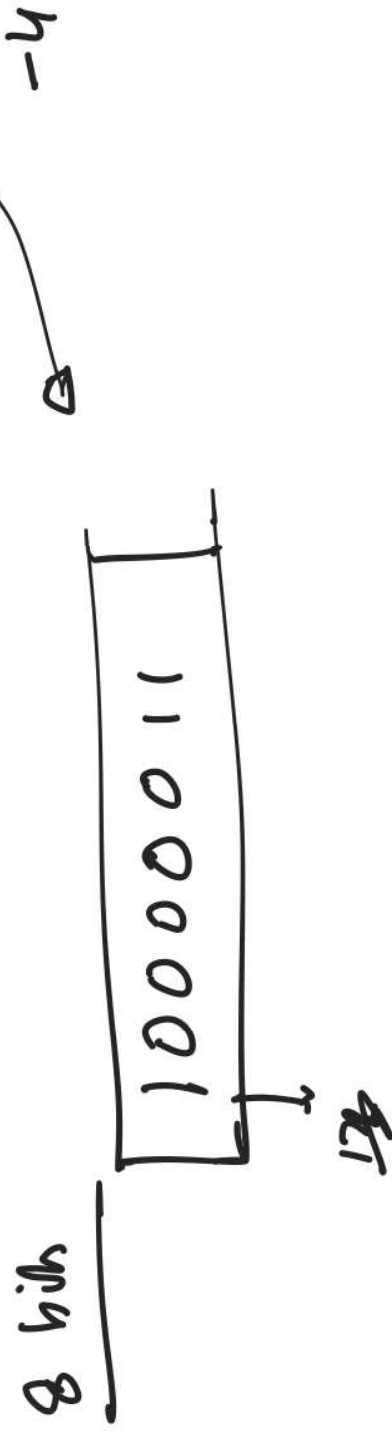
Step 4

Biased Exponent.

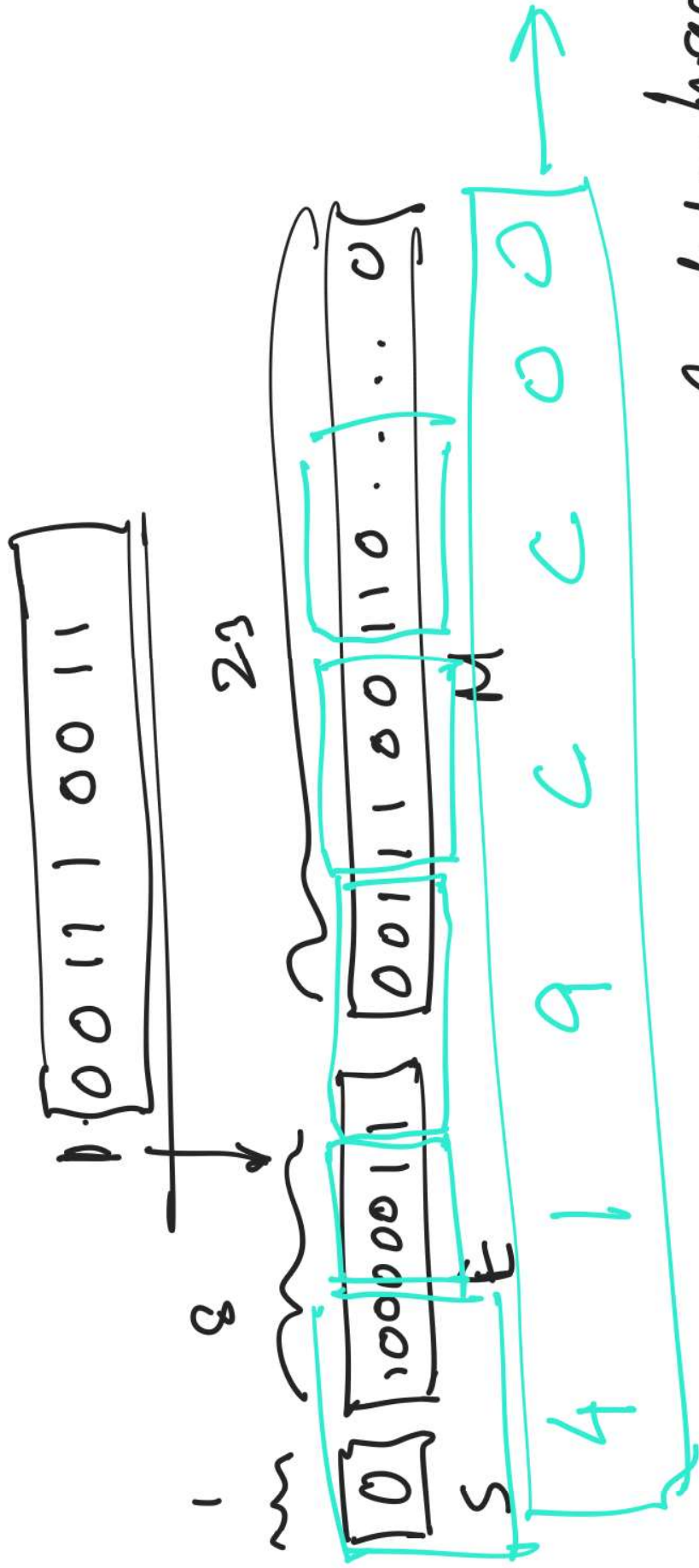
Unbiased exponent : 4.

Biased exponent : $127 + 4 = 131$

Bias/offset $\left[2^{n-1} - 1 \right]$ $\left[2^7 - 1 \right]$ $n=8$



Step 5 \rightarrow Mantissa



No. of bits in mantissa \rightarrow resolution/precision.
 " " " " exponent \rightarrow range.

0.375

0.375₁₀

0.011 × 2⁰

1.1 × 2⁻² Unbias Exponent
= -2

Sign: 0 Bias Exponent

ignore.

= 127 - 2

= 125

3 EC 000000

0	01111101	10	0 .
---	----------	------------	-----

Example

08500000

11001000010100000000

1.10100000

Sign

144

Unknown

Exponent

144-127

mantissa

-1.101×2^{17}

-1.675×2^{17}

Next 3 slides:

Ack:

<https://www.youtube.com/watch?v=RuKkePyo9zk&list=PLTd6ceoshprcpen2JvsJiuvWvqIAkzea&index=10>

2's comp

9219

629
1019

bind

[illegible]

0 1 2 3 4 5 6 7 -8 -7 -6 -5 -4 -3 -2 -1

-5	-4	-3	-2	-1	0	1	2	3	4	5	6	7	8	9	10
----	----	----	----	----	---	---	---	---	---	---	---	---	---	---	----

-9 -8 -7 -6 -5 -4 -3 -2 -1 0 1 2 3 4 5 6

8 7 6 5 4 3 2 1 0 1 2 3 4 5 6 7

7

25-5

1 → 254

$$2^{-1} = 12^7$$
$$2^{(4-1)} - 1 = 7 \text{ nodes}$$
$$\boxed{8} + 5 = 13$$

13

IEEE 754 Format	Sign	Exponent	Mantissa	Exponent Bias
32 bit single precision	1 bit	8 bits	23 bits (+ 1 not stored)	$\underline{2^{(8-1)} - 1 = 127}$
<u>64 bit double precision</u>	<u>1 bit</u>	<u>11 bits</u>	<u>52 bits (+ 1 not stored)</u>	<u>$2^{(11-1)} - 1 = 1023$</u>
128 bit quadruple precision	1 bit	15 bits	112 bits (+ 1 not stored)	$2^{(15-1)} - 1 = 16383$

double a?

Reserved Exponent Values

Exponent Value	Mantissa	Represents
✓ 11111111	All zeros	Infinity (∞)
✓ 11111111	Not all zeros	Not a number (NaN)
00000000	All zeros	Zero
00000000	Not all zeros	Subnormal (very small)

Practice problems

Convert the following decimal numbers to 32-bit IEEE 754 format by hand:

- a. 1.0
- b. -0.1
- c. 2016.0
- d. 0.00390625
- e. -3125.3125
- f. 0.33
- g. -0.67
- h. 3.14

[▼ Answer](#)

- a. 3f800000
- b. bdc00000
- c. 44fc0000
- d. 3b800000
- e. c5435500
- f. 3ea8f5c3
- g. bf2b851f
- h. 4048f5c3

Convert the following hexadecimal numbers to decimal by hand using the 32-bit IEEE 754 format:

- a. 40000000
- b. bf800000
- c. 3d800000
- d. c1804000
- e. 42c81000
- f. 3f99999a
- g. 42f6e666
- h. c25948b4

[▼ Answer](#)

- a. +2.0
- b. -1.0
- c. +0.0625
- d. -16.03125
- e. 100.03125
- f. 1.2
- g. 123.449997
- h. -54.320999