

Lesson 6: Intro to Floating-Point Representation

Recall that for an m -bit fixed-point representation of

a type I rational $x \in \mathbb{Q}_I$, we chose two

nonnegative integer parameters $l, f \in \mathbb{Z}$ that determined

- the implicit location of the binary point when interpreting the value of the raw m -bit word
- the power in the denominator of our binary fraction (and thus the number of bits we shift the radix point left or right)

Example 6.1) Let $x = \frac{165}{16} = 10.3125$

$$\Rightarrow x = 8 + 2 + \frac{1}{4} + \frac{1}{16}$$

$$\Rightarrow x = 2^3 + 2^1 + \frac{1}{2^2} + \frac{1}{2^4}$$

Example 6.1) Let $x = \left(\frac{165}{16}\right) = \frac{165}{2^4} = 10.3125$

2⁴
 ↑
 we know $x \in \mathbb{Q}_I$

$$\Rightarrow x = \frac{128 + 32 + 4 + 1}{16}$$

$$\Rightarrow x = \left(\frac{2^7 + 2^5 + 2^2 + 2^0}{2^4} \right)_{10} = \left(2^3 + 2^1 + 2^{-2} + 2^{-4} \right)_{10}$$

$$\Rightarrow x = \left(\frac{10100101}{10000} \right)_2 = (1010.0101)_2$$

$$\Rightarrow \boxed{x} = (1010.0101)_2$$

↑
math object

Let's try to encode $x \in \mathbb{Q}_I$ exactly in fixed-point format

$$\Rightarrow \boxed{\hat{x}} = \text{vfixed}_{2(4,4)} \left(\frac{165}{16} \right) = \text{vfixed}_{2(4,4)} \underbrace{(10100101)}_{\text{raw bit string}}$$

$$= \underbrace{(1010.0101)}_{\text{interpreted value}}_2$$

the value of bit string used to store "something" related to x in a computer's memory

Example 6.1, (continued...)

$$\hat{x}_1 = \text{ufixed}_{2(s,3)}\left(\frac{165}{16}\right) = \text{ufixed}_{2(s,3)}(10100101)$$

$$= \underbrace{01010}_{l=5} . \underbrace{010}_{f=3} \boxed{1}$$

what do we do with this extra bit?

□ Chop: 01010.010

□ Round Down: 01010.010

□ Round up: 01010.011

This forces a discussion of rules for rounding...

this is the first time we're actually needed to approximate a number in this class.

≈ 01010.010 ← let's assume a naive approach and simply chop off all trailing bits that are unsupported by our data format

Prior to this lesson, all equality was exact... Welcome to numerical analysis!!

But notice if we set $\hat{x}_1 = 01010.010$ to be our approximation

then
$$\hat{x}_1 = 8 + 2 + \frac{1}{4} = \frac{32 + 8 + 1}{4} = \frac{41}{4} = \frac{164}{16}$$

⇒ we have an absolute error in our approximation of

$$|x - \hat{x}_1| = \left| \frac{165}{16} - \frac{164}{16} \right| = \left| \frac{1}{16} \right| = \frac{1}{16}$$

Example 6.1, continued...

Let's try a different method for dealing with the extra bit from our exact number that can't be stored in our $\text{ufixed}_{2(5,3)}\left(\frac{165}{16}\right)$ format.

$$\hat{X}_2 = \text{ufixed}_{2(5,3)}\left(\frac{165}{16}\right)$$

$$= \text{ufixed}_{2(5,3)}(10100101)$$

$$= \underbrace{01010}_{l=5} \cdot \underbrace{010}_{f=3} \boxed{1} \quad \text{this time, let's round up}$$

$$\approx 01010.011$$

Now let's set $\hat{X}_2 = 01010.011$ to be our second approximation

$$\Rightarrow \hat{X}_2 = 8 + 2 + \frac{1}{4} + \frac{1}{8} = \frac{64 + 16 + 2 + 1}{8} = \frac{83}{8} = \frac{166}{16} \quad \begin{array}{l} > \times \\ \text{we rounded} \\ \text{up!} \end{array}$$

\Rightarrow we have absolute error in our approximation of

$$|X - \hat{X}_2| = \left| \frac{165}{16} - \frac{166}{16} \right| = \left| \frac{-1}{16} \right| = \frac{1}{16}$$

Example 6.1, continued...)

Now let's try storing x in the $\text{ufixed}_{2(3,5)} \left(\frac{165}{16} \right)$ format and see what happens. To this end, we set

$$\hat{X}_3 = \text{ufixed}_{2(3,5)} \left(\frac{165}{16} \right)$$
$$= \text{ufixed}_{2(3,5)} (10100101)$$

$$= \boxed{1} \underbrace{010}_{e=3} . \underbrace{01010}_{f=5}$$

what do we do
with this MSB?

- Chop: 010.01010
- Round up: 110.01010

$$\approx 010.01010$$

let's assume a naive approach again
and simply chop off the MSB

$$\Rightarrow \hat{X}_3 = 2 + \frac{1}{4} + \frac{1}{16} = \frac{32 + 4 + 1}{16} = \frac{37}{16}$$

\Rightarrow we have absolute error in our approximation of

$$|x - \hat{X}_3| = \left| \frac{165}{16} - \frac{37}{16} \right| = \left| \frac{128}{16} \right| = +8 \text{ units}$$

this makes sense
because we chopped
off $+8 = 2^3$ bit.

5

Example 6.1, continued -)

Of course, we might try to round up using the $\text{ufixed}_{2(3,5)}$ format to yield

$$\hat{X}_4 = \text{ufixed}_{2(3,5)} \left(\frac{165}{16} \right)$$

$$= \text{ufixed}_{2(3,5)} (1010\ 0101)$$

$$= \boxed{1} \underbrace{010}_{l=3} . \underbrace{01010}_{f=5}$$

$$\approx 110.01010$$

$$\Rightarrow \hat{X}_4 = 4 + 2 + \frac{1}{4} + \frac{1}{16} = \frac{64 + 32 + 4 + 1}{16} = \frac{101}{16}$$

\Rightarrow we have absolute error in our approximation of

$$|x - \hat{X}_4| = \left| \frac{165}{16} - \frac{101}{16} \right| = \left| \frac{64}{16} \right| = +4$$

this error makes sense since we rounded the exact 4-bit integer value 1010 to an approximate 3-bit integer 110 and $1010 - 110 = 100$

(6)

Example 6.1, continued...

Notice that our choice of parameters $l, p \in \mathbb{Z}$ matter. In only one case can we encode the exact value of $x = \frac{165}{16}$. In all other formats, we must use an approximation to x .

Data type and rounding rule	Binary Value	Absolute Error	Type of error
\square $\text{ufixed}_{2(4,4)}\left(\frac{165}{16}\right)$	$1010.0101 = \hat{x}$	$ x - \hat{x} = 0$	NONE
\square $\text{ufixed}_{2(5,3)}\left(\frac{165}{16}\right)$ (use chop)	$\hat{x}_1 = 01010.010$ error happens here	$ x - \hat{x}_1 = \frac{9}{16}$	LSB error
\square $\text{ufixed}_{2(5,3)}\left(\frac{165}{16}\right)$ (use round up)	$\hat{x}_2 = 01010.011$ error happens here	$ x - \hat{x}_2 = \frac{5}{16}$	LSB error
\square $\text{ufixed}_{2(3,5)}\left(\frac{165}{16}\right)$	$\hat{x}_3 = 010.0101$ error happens here	$ x - \hat{x}_3 = 8$	MSB error
\square $\text{ufixed}_{2(3,5)}\left(\frac{165}{16}\right)$	$\hat{x}_4 = 110.0101$	$ x - \hat{x}_4 = 4$	MSB error

Example 6.1 ...)

□ Given these observations, which parameter values are the most appropriate for this example?

□ what if, ^{for some reason,} we could not choose $l=4$ and $f=4$?
What would be your next favorite option? why?

□ What type of error is "preferable": LSB errors or MSB errors? why?

Example 6.2) Let $x = \left(\frac{9965}{64}\right)_{10}$

$$\Rightarrow x = 128 + 16 + 8 + 2 + 1 + \frac{1}{2} + \frac{1}{8} + \frac{1}{16} + \frac{1}{64} = (155.703125)_{10}$$

$$\Rightarrow x = 2^7 + 2^4 + 2^3 + 2^1 + 2^0 + 2^{-1} + 2^{-3} + 2^{-4} + 2^{-6}$$

$$\Rightarrow x = \frac{2^{13} + 2^{10} + 2^9 + 2^7 + 2^6 + 2^5 + 2^3 + 2^2 + 2^0}{2^6}$$

$$\Rightarrow x = \left(\frac{10011011101101}{1000000}\right)_2$$

$$\Rightarrow x = (10011011.101101)_2$$

□ Suppose we want to store this number in $m=8$ bits.

Do the same "optimal" choice of parameters $l=4$ and $f=4$

from example 6.1 yield an optimal data encoding strategy

in this example?

□ If not, what are more strategic values of $l, f \in \mathbb{Z}$ for this case if we assume we cannot extend our range (i.e. $m=8$)?

Example 6.2, continued ...)

The first observation we might make is that ^{the} economic decision to keep $m=8$ has immediate mathematical consequences. Since we cannot store $x = \left(\frac{9965}{64}\right)_{10}$ exactly using only 8 bits and thus any ^{8-bit} fixed-point data encoding will necessarily be approximate and inaccurate.

The first thing we should do after we come to this realization is cry. Then, after we recover our sensibilities, we might ask:

□ even though any 8-bit fixed point data encoding of $x = \left(\frac{9965}{64}\right)_{10}$ will be inaccurate, what is the best we can do (i.e. which 8 bit encoding will produce an approximation with minimal error).

Example 6.2 continued...

From our work in example 6.1 and using our hard earned intuition, we might guess that $l=8, f=0$ and the round-up scheme does the trick. Let's try this on for size:

$$\hat{X}_1 = \text{ufixed}_{2(8,0)}\left(\frac{9965}{64}\right)$$

$$= \text{ufixed}_{2(8,0)}\left(\underbrace{10\ 0110\ 1110\ 1101}_{\substack{\text{raw, uninterpreted} \\ \text{14-bit word}}}\right)$$

$$= \underbrace{1001\ 1011}_{l=8} \overset{\circ}{\cdot} \boxed{101101}$$

$f=0$

we can't store 14-bits in an 8-bit word so this trailing part will be lost. But, we can round up:

$$\approx 1001\ 1100$$

$$\begin{array}{r} 1001\ 1011 \\ + \quad \quad 1 \\ \hline 1001\ 1100 \end{array}$$

Example 6.2 continued ...)

If we set

$$\hat{X}_1 = 1001\ 1100.0 \quad \leftarrow \text{implicit fractional part since } f=0$$

$$= 2^7 + 2^4 + 2^3 + 2^2$$

$$= 128 + 16 + 8 + 4$$

$$= 156 = \frac{9984}{64}$$

\Rightarrow we have absolute error in our approximation of

$$|x - \hat{X}_1| = \left| \frac{9985}{64} - \frac{9984}{64} \right| = \left| \frac{-19}{64} \right|$$

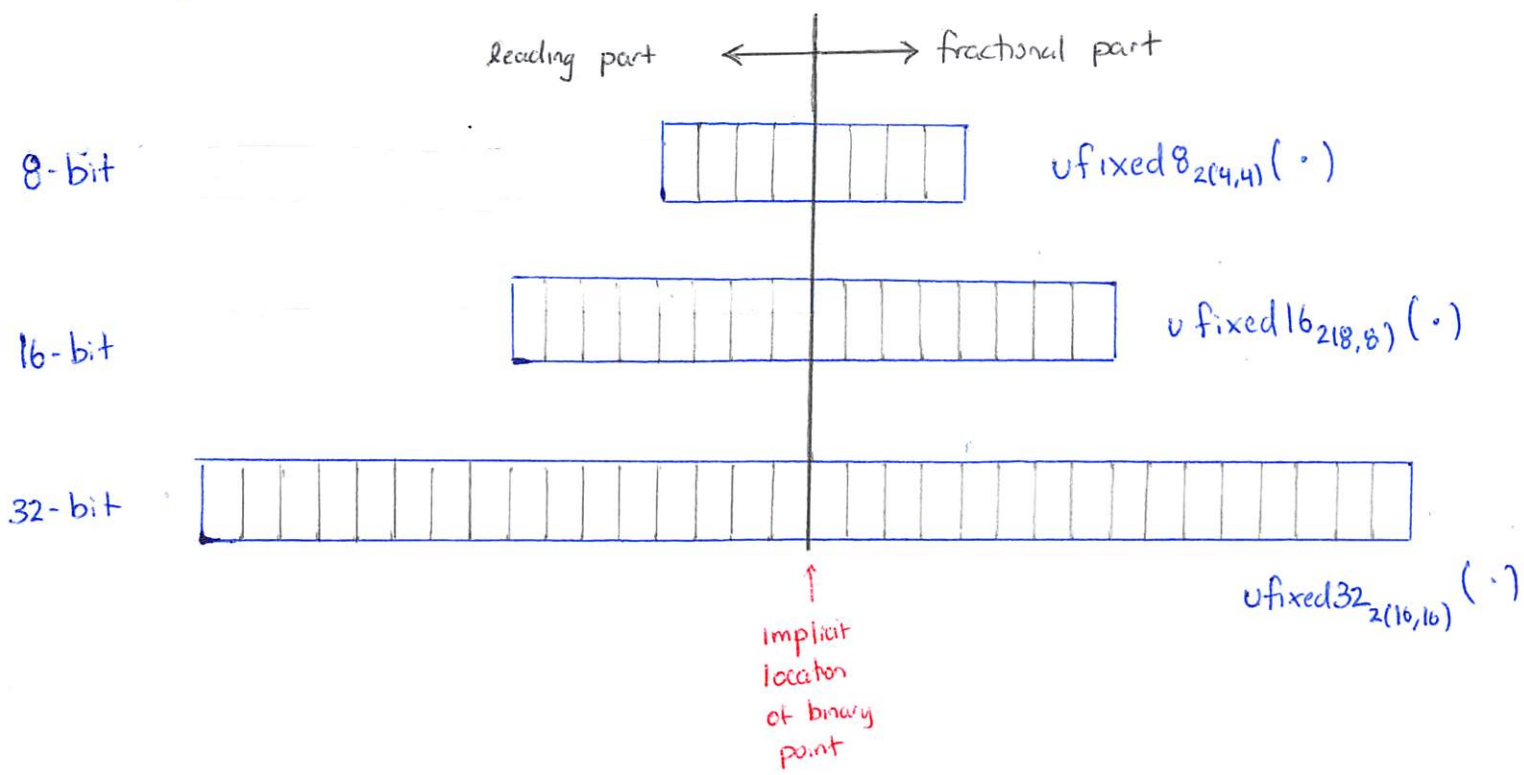
$$= \frac{+19}{64} = 16 + 2 + 1$$

$$= \frac{2^4 + 2^1 + 2^0}{2^6} = 2^{-2} + 2^{-5} + 2^{-6}$$

$$= 0.010011$$

Lesson 6, part 2: Intro to Floating-Point Continued

Based on our work in Lesson 6, part one we saw that for a chosen m -bit, unsigned fixed-point format with $m = 8, 16, 32,$ or $64,$ we might choose $l = \frac{m}{2} = f$ so that exactly half the bits are dedicated to the leading part and half the bits are dedicated to the fractional part of our number. Let's visualize this:



Let's analyze the range of each of these data formats.

$$\text{ufixed}_{2(4,4)}(\cdot) : \quad x = \text{ufixed}_{2(4,4)}(B_7 B_6 B_5 B_4 B_3 B_2 B_1 B_0) = \sum_{k=-4}^3 b_k \cdot 2^k$$

where $b_k = B_{k+4}$

Lower-bound : 0000 0000 ← raw, uninterpreted 8-bit word

interpreted value
↓

$$\Rightarrow x = \text{ufixed}_{2(4,4)}(0000 0000)$$

$$= (0000 \cdot 0000)_2$$

$$= (0)_{10}$$

Upper-bound : 1111 1111 ← raw, uninterpreted 8-bit word

$$\Rightarrow x = \text{ufixed}_{2(4,4)}(1111 1111)$$

$$= (1111 \cdot 1111)_2$$

$$= 1 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 + 1 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} + 1 \cdot 2^{-4}$$

$$= 8 + 4 + 2 + 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16}$$

$$= \frac{128 + 64 + 32 + 16 + 8 + 4 + 2 + 1}{16} = \frac{255}{16}$$

⇒ upper bound for $\text{ufixed}_8_{2(4,4)}(1111\ 1111)$ is

$$\frac{255}{16} = \frac{256-1}{16} = \frac{2^8-2^0}{2^4} = 2^4 - \frac{1}{2^4} = 2^e - \frac{1}{2^f}$$

Exercise: Do the same analysis for

A. $\text{ufixed}_{16}_{2(8,8)}(B)$

B. $\text{ufixed}_{32}_{2(16,16)}(B)$

C. $\text{ufixed}_{64}_{2(32,32)}(B)$

Note: there are some really beautiful & powerful features of fixed point format including

- A. fast implementation for arithmetic in binary using same circuitry as for integer arithmetic
- B. Easily defined signed data classes $\text{fixed}_8, \text{fixed}_{16}, \text{etc.}$ based on two's complement representation
- C. Relatively

We saw in class that the fixed-point format has some significant drawbacks though. Two of these include:

A. wasted bits: if we encoded the number $x = 4 + \frac{1}{4}$ in $\text{ufixed}_{2(8,8)}(\cdot)$, we would have the following bit string

0000 | 0100.0100 | 0000

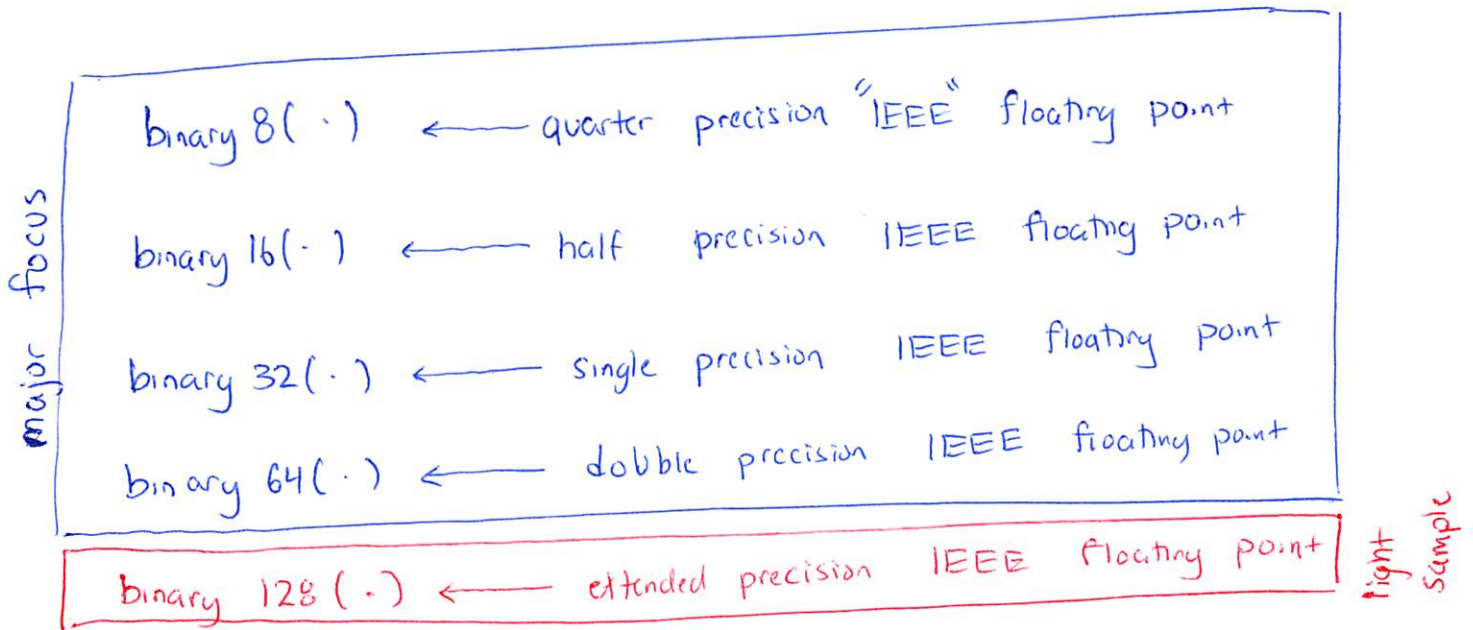
→ look at all these "wasted" bits

B. relative limited range: the range

$$0 \leq x \leq 2^l - \frac{1}{2^p}$$

is decent but falls far short of being versatile enough to encode the entire span of numbers we might want to use in scientific computing applications...

These drawbacks are significant enough that they warrant deep thought about possible alternative methods to encode $x \in \mathbb{Q}_I$. One very powerful, well-designed, and famous solution to this problem is known as IEEE Floating-Point format. In our work, we will ^{heavily} focus on four floating-point formats for signed $x \in \mathbb{Q}_I$ and lightly sample a fifth format:



Foundational Idea for Floating Point Formats

The major inspiration for floating point format is not complicated. In fact, you might be familiar with the idea already and know it as either

- scientific notation or
- exponential notation

These tools allow us to save time by representing numbers optimally

Example 6.3) Suppose we wanted to represent the number

$$X = (\underbrace{0.000}_{\text{leading zeros}}9765625)_{10}$$

$$= \left(\frac{1}{2^{10}} \right)_{10}$$

$$= (\underbrace{0.0000000001}_{\text{leading zeros}})_{2}$$

Notice how annoying it is to have to write out all the leading zeros.

Example 6.3, continued ...)

- Moreover, these leading zeros don't add any significant "value" to our representation outside of telling us where the radix point should be.
- Based on a clever use of our knowledge of arithmetic and algebra, we can capture the same information much more efficiently by allowing the radix point to float:

$$x = (0.0009765625)_{10}$$

$$= 0000\ 0000\ 0000\ 0000.0009\ 7655\ 2500\ 0000$$

↖ $ufixed32_{10(16,16)}(\cdot)$
format required
to capture exactly

$$= 9.765625 \times 10^{-4}$$

Note: to transform

$$0.0009765625 = 9.765625 \times 10^{-4}$$

↑
take radix point
and float it down the digit stream
four places to the right

Example 6.3, continued

The math goes like this

$$x = (0.0009765625)_{10}$$

$$= (0.0009765625) \cdot 1$$

$$= (0.0009765625) \cdot 10^0$$

$$= (0.0009765625) \cdot 10^{4-4}$$

$$= (0.0009765625) \cdot 10^4 \cdot 10^{-4}$$

$$= (0.0009765625 \cdot 10^4) \cdot 10^{-4}$$

float to the right
by four spaces

properly account for
shift in exponent
notation

$$= 9.765625 \times 10^{-4} \leftarrow \text{scientific notation (in decimal)}$$

This downside to fixed point lead to a very tedious & focused development of an alternate much more versatile ¹ set of data classes called floating-points

binary 8 ← quarter precision

binary 16 ← half precision

binary 32 ← single precision

binary 64 ← double precision

binary 128 ← extended precision

You already have all the genius you need to understand the main idea behind these formats

$$\begin{aligned}
 X &= \frac{1}{2^{10}} = 0.0009765625 \cdot 1 \\
 &= 0.0009765625 \cdot 10^0 \\
 &= 0.0009765625 \cdot 10^{4-4} \\
 &= (0.0009765625) \cdot 10^4 \cdot 10^{-4} \\
 &= (0.0009765625 \cdot 10^4) \times 10^{-4}
 \end{aligned}$$

Example 6.3, continued ...)

Of course, we can do the same type of analysis in binary

$$X = \frac{1}{2^{10}} = (0.0000000001)_2$$

$$= (0.0000000001) \cdot \underbrace{2^0}_{=1}$$

$$= (0.0000000001) \times \underbrace{2^{10} \cdot 2^{-10}}_{=2^0}$$

$$= (0.0000000001 \cdot 2^{10}) \times 2^{-10}$$

$$= 1.0 \times 2^{-10} \quad \leftarrow \text{scientific notation (in binary)}$$

Example 6.4) Consider the following $x \in \mathbb{Q}_I$

$$x = 2^{-4} + 2^{-7} + 2^{-8}$$

$$= (0.00010011)_2$$

If we represented this in fixed point, we would need 8 = f bits for the fractional part to capture this value exactly. However, let's transform into scientific (floating point) format as we did in Example 6.3:

$$x = (0.00010011) \cdot 2^4 \cdot 2^{-4}$$

$\overbrace{2^4}^{= 2^0 = 1}$
4

want right shift by four bits

$$= 1.0011 \times 2^{-4}$$

leading part = 1 (normalized)
new fractional part
-4 ← exponent value
binary is always base = 2

Normalized floating-point numbers

Mathematically, we can denote the general format of a normalized floating point number as

$$X = \boxed{\pm} \left(\overset{\substack{\text{implied leading value of 1} \\ \downarrow}}{b_0} . \underbrace{b_{-1} b_{-2} \dots b_{-p}}_{\text{significand}} \right) \times 2^{\boxed{e} \leftarrow \text{exponent value}}$$

where $b_0 = 1$ when we have normalized numbers

$$\Rightarrow X = \pm \left(1 + \underbrace{0. b_{-1} b_{-2} \dots b_{-f}}_{\text{significand}} \right) \times 2^e$$

$$= \pm (1 + f) \times 2^e \quad \text{w/ } f = 0. b_{-1} b_{-2} \dots b_{-p}$$

Let's explore how we might store quarter precision

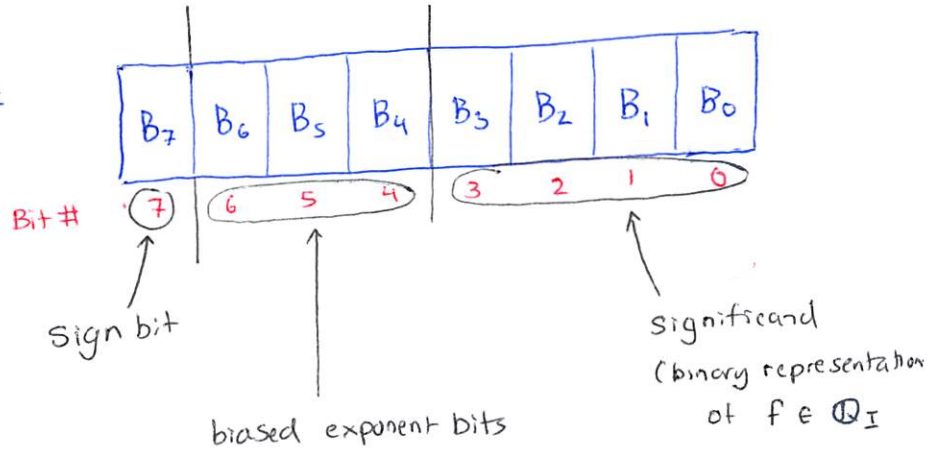
floating-point numbers using 8-bits.

Quarter-precision floats

$m = 8$ -bits & normalized

Let $x = \text{binary}_8(B_7 B_6 B_5 B_4 B_3 B_2 B_1 B_0)$

Raw, uninterpreted
8-bit word



Let x be the interpreted value resulting from 8-bit word

$$\text{sign bit: } B_7 = \begin{cases} 0 & \text{if } x \geq 0 \\ 1 & \text{if } x < 0 \end{cases}$$

biased exponent field: $B_6 B_5 B_4 \xrightarrow{\text{mapped to}}$

Example 6.5)

$$x = 1 = 2^0$$

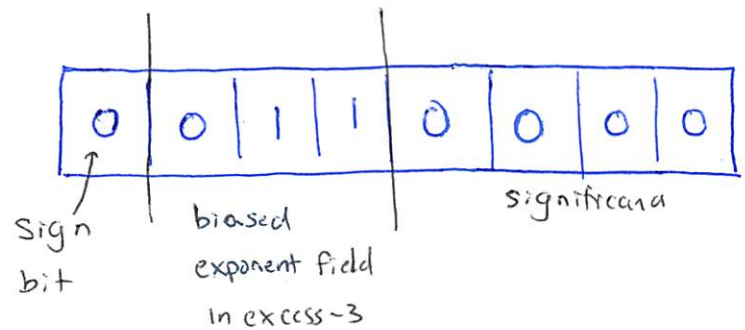
$$= 1.0 \times 2^0$$

$$= + 1. \underbrace{0000}_{\text{significant}} \times 2^0$$

$$\Rightarrow \begin{cases} B_7 = 0 & \text{since } x \geq 0 \\ f=0 \Leftrightarrow B_3 B_2 B_1 B_0 = 0000 \\ e = (0)_0 \Leftrightarrow 011 = B_6 B_5 B_4 \end{cases}$$

$$x = (1)_{10} = (1)_2$$

stored \rightarrow



$$x = (2)_{10}$$

$$= (10)_2$$

left shift

$$= (10 \cdot 2^{-1}) \cdot 2^1$$

$$= \boxed{+} \overset{\text{sign bit}}{1} . \underbrace{0000}_{\text{significant}} \times 2^{\boxed{1} \text{ exponent value}}$$

Implied normalized leading bit

$$\Rightarrow B_7 = 0 \quad \text{is} \quad x \geq 0$$

$$f = 0 \Rightarrow B_3 B_2 B_1 B_0 = 0000$$

$$e = 1 \Rightarrow B_6 B_5 B_4 = 100$$

Note: unsigned math

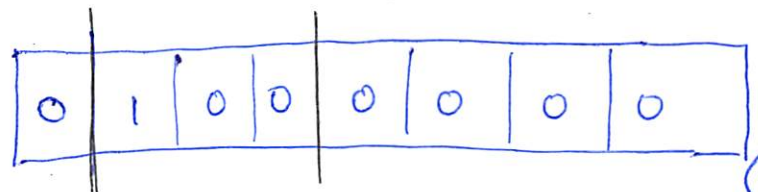
$$e = 100 - 011 = 4 - 3 = 1$$

$$\textcircled{U} - K = e$$

↑
K-base

Unsigned integer value of exponent field

$$x = (2)_{10} = (10)_2 \xrightarrow{\text{stored}}$$



$$\begin{aligned}
 X &= (11)_{10} \\
 &= 8 + 2 + 1 \\
 &= 2^3 + 2^1 + 2^0 \\
 &= (1011)_2
 \end{aligned}$$

$$= 1.011 \times 2^3$$

$$= \boxed{+} \overset{\text{implied leading value}}{\downarrow} 1.\underbrace{0110}_{\text{significant}} \times 2^{\boxed{3}} \leftarrow \text{exponent}$$

↑ sign

$$\Rightarrow B_7 = 0 \quad \text{since } X \geq 0$$

$$f = \quad \Rightarrow B_3 B_2 B_1 B_0 = 0110$$

$$e = 3$$

Note: unsigned math

$$e = 110 - \underbrace{011}_K = 6 - 3 = 3 \quad \checkmark$$

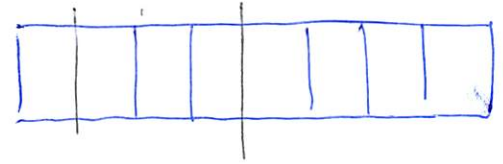
$$e = u - K \quad \text{where } u = \text{uint3}(B_6 B_5 B_4)$$

$$\Rightarrow u = e + K = (e + 3)_{10}$$

$$= (e + 011)_2$$

$$x = (15)_{10}$$

→ stored as



$$= 8 + 4 + 2 + 1$$

$$= (1111)_2$$

$$= 1.111 \times 2^3$$

Implied leading bit

$$= \boxed{+} \overset{\downarrow}{1} . \underbrace{1110}_{\text{significant}} \times 2^{\boxed{3} \text{ exponent value}}$$

sign bit

$$e = 3 = \text{uint3}(B_6 B_5 B_4) - 3$$

$$\Rightarrow e = 011 = u \quad - \underbrace{011}$$

$$= u - k$$

$$\Rightarrow u = e + k$$

$$= 011 + 011$$

$$= (110)_2$$

$$= 6$$

$$\begin{array}{r} 011 \\ + 011 \\ \hline 110 \end{array}$$

unsigned dec	exponent code = u	exponent value
	raw, uninterpreted value for 3-bit exponent field $B_6 B_5 B_4$	3-bit, excess-3 interpreted, <u>biased</u> value of e (in decimal)
0	000	-3
1	001	-2
2	010	-1
3	011	0
4	100	1
5	101	2
6	110	3
7	111	4

Sometimes called
offset binary

$$e = u - k$$

$$= \text{uint3}(B_6 B_5 B_4) - 3$$

since $k = 3$
(binary 8 has an excess-3
offset)