

lecture 2

- fixed point
- IEEE floating point standard

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Wed. January 13, 2016

Fixed point

Fixed point means we have a constant number of bits (or digits) to the left and right of the binary (or decimal) point.

Examples :

Assignment Project Exam Help

23953223.49 (base 10)
<https://powcoder.com>

Currency uses a fixed number of digits to the right.
Add WeChat powcoder

10.1101 (base 2)

Two's complement for fixed point numbers

e.g. 0110.1000 which is 6.5 in decimal

How do we represent -6.5 in fixed point ?

$$\begin{array}{r} 0110.1000 \\ 1001.0111 \quad \leftarrow \text{invert bits} \\ + \underline{0000.0001} \quad \leftarrow \text{add } 0001 \\ \hline 0000.0000 \end{array}$$

<https://powcoder.com>
Add WeChat powcoder

Thus,

$$\begin{array}{r} 1001.0111 \quad \leftarrow \text{invert bits} \\ + \underline{0000.0001} \quad \leftarrow \text{add } .0001 \\ \hline 1001.1000 \quad \leftarrow \text{answer: -6.5 in (signed) fixed point} \end{array}$$

Scientific Notation (floating point)

$$300,000,000 = 3 \times 10^8$$
$$= 3.0 E + 8$$

Assignment Project Exam Help

<https://powcoder.com>

$$0.0000456 = 4.56 E - 6$$

Add WeChat powcoder



"Normalized" : one digit to the left of the decimal point.

Scientific Notation in binary

$$(1000.01)_2 = 1.00001 \times 2^3$$

Assignment Project Exam Help

<https://powcoder.com>

$$(0.111)_2 = 1.11 \times 2^{-1}$$

Add WeChat powcoder

"Normalized" means one "1" bit to the left of the binary point. **(Note that 0 cannot be represented this way.)**

sign

"exponent"

+

|

x

2

E

"significand"

Assignment Project Exam Help

(also called

<https://powcoder.com>
"mantissa")

Add WeChat powcoder

How to represent this information ?

How to represent the number 0 ?

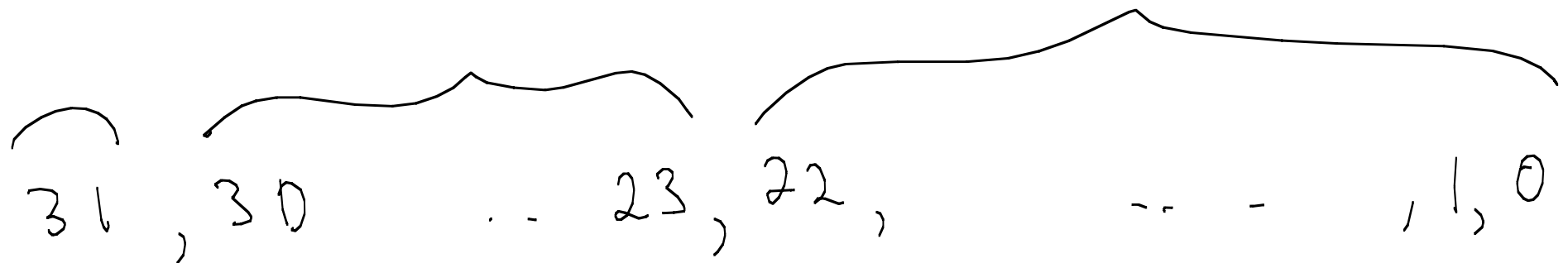
IEEE floating point standard (est. 1985)

case 1: single precision (32 bits = 4 bytes)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Let's look at these three parts, and then examples.

sign 0 for positive, 1 for negative

"significand" Assignment Project Exam Help

<https://powcoder.com>

↓ Add WeChat powcoder

You don't encode the "1" to the left of the binary point.

Only encode the first 23 bits to the right of the binary point.

exponent code

exponent value

00000000

reserved (explained soon)

00000001

-126

00000010

-125

00000011

-124

:

:

:

:

01111111

0

This is not two's

10000000

complement !
<https://powcoder.com>

10000001

2

Add WeChat powcoder

:

:

:

:

11111110

127

11111111

reserved (explained soon)

unsigned exponent code = exponent value + "bias"
(for 8 bits, bias is defined to be 127)

Q: What is the largest positive normalized number ?
(single precision)

1 . ~~Assignment Project Exam Help~~ $\times 2^e$

<https://powcoder.com>

Add WeChat powcoder

A:

1 . 1 1 1 1 . . . 1 1 $\times 2^{127}$

$$2^{127} \approx 10^?$$

$$2^{10} \approx 10^3$$

Assignment Project Exam Help

$$2^{127} = (2^{10})^{12} \cdot 2^7$$

<https://powcoder.com>

Add WeChat powcoder

$$= (2^{10})^{12} \cdot 2^7$$

$$\approx (10^3)^{12} \cdot 10^2$$

$$= 10^{38}$$

Q: What is the smallest positive normalized number ?
(single precision)

| ~~Assignment Project Exam Help~~ $\times 2^e$

<https://powcoder.com>

Add WeChat powcoder

A:

| 0 0 0 0 0 0 - - - 0 $\times 2^{-126}$

Exponent code 00000000 reserved for
"denormalized" numbers

$1 \pm 0 \cdot 2^{-126}$

Assignment Project Exam Help

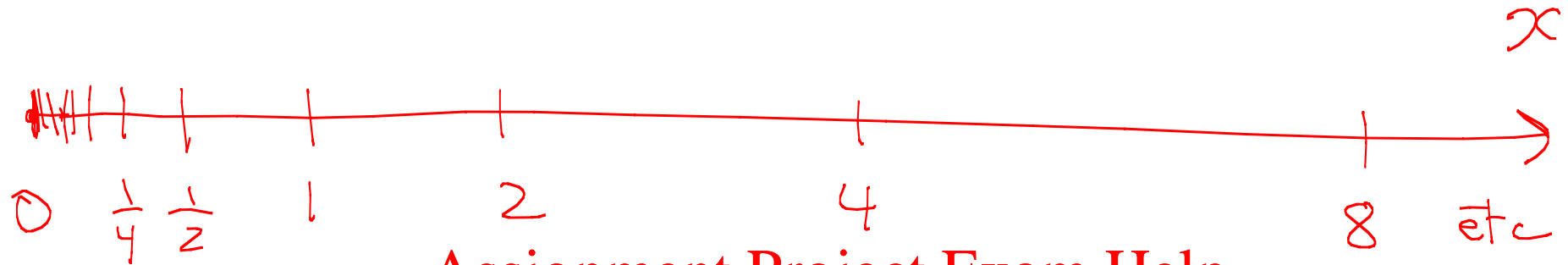
<https://powcoder.com>

belong to $(-2^{-126}, 2^{-126})$

Add WeChat powcoder

includes 0

Dividing each power of 2 interval into 2^{23} equal parts
(same for negative real numbers).

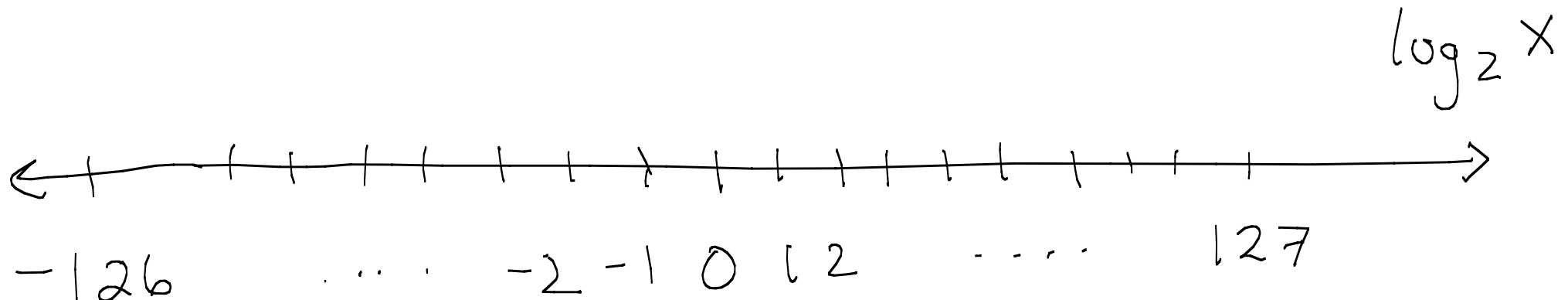


Assignment Project Exam Help

<https://powcoder.com>

Note the power of 2 intervals themselves are equally spaced on a log scale.

Add WeChat powcoder



Exponent code 11111111 also reserved.

if significand is all 0's

then value is \pm infinity (depending on sign bit)

else value is NaN ("not a number")

e.g. variable is declared but hasn't been
assigned a value

<https://powcoder.com>
Add WeChat powcoder

This is the stuff you put on an exam crib sheet.
(Yes, you can bring a crib sheet for the quizzes.)

Example: write 8.75 a single precision float (IEEE).

First convert to binary.

8.75

$$= (1000)_2 \cdot (75)_{10}$$

$$= (10001)_2 \cdot (5)_{10} \times 2^{-1}$$

$$= 100011.0 \times 2^{-2}$$

$$= 1.00011 \times 2^3$$

$$(8.75)_{10} = (1.00011)_2 \times 2^3$$

23 bit significand: 000110000000000000000000000

exponent value: $e = 3$

exponent code = exponent value (e) + bias

Thus, exponent code is unsigned $3 + 127$.

$$(130)_{10} = (10000010)_2$$

So, the 32 bit representation is :

0 10000010 000110000000000000000000

0 x 4 1 0 c 0 0 0 0

Recall last lecture: 0.05 cannot be represented exactly.

```
float x = 0;  
for (int ct = 0; ct < 20; ct ++ ) {  
    x += 1.0 / 20;  
    System.out.println( x );  
}
```

Assignment Project Exam Help

0.05

0.1

0.15

0.2

0.25

0.3

0.35000002

0.40000004

0.45000005

0.50000006

etc

<https://powcoder.com>

Add WeChat powcoder

Floating Point Addition

$$x = 1.00100100010000010100001 * 2^2$$

$$y = 1.10101000000000000101010 * 2^{-3}$$

$$x + y = ?$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Floating Point Addition

$$x = 1.00100100010000010100001 * 2^2$$

$$y = 1.10101000000000000101010 * 2^{-3}$$

Assignment Project Exam Help

$$x + y = ?$$

<https://powcoder.com>

Add WeChat powcoder

$$x = 1.0010010001000001010000100000 * 2^2$$

$$y = .00001101010000000000000101010 * 2^2$$

but the result $x+y$ has more than 23 bits of significand

How many *digits* (base 10) of precision can we represent with 23 *bits* (base 2) ?

$$2^{23}$$

Assignment Project Exam Help
 $\approx (2^{10})^2 \cdot 2^3 = (1024)^2 \cdot 2^3$
<https://powcoder.com>

Add WeChat powcoder

$$\approx 10^6 \cdot 10^3$$

$$= 10^7$$

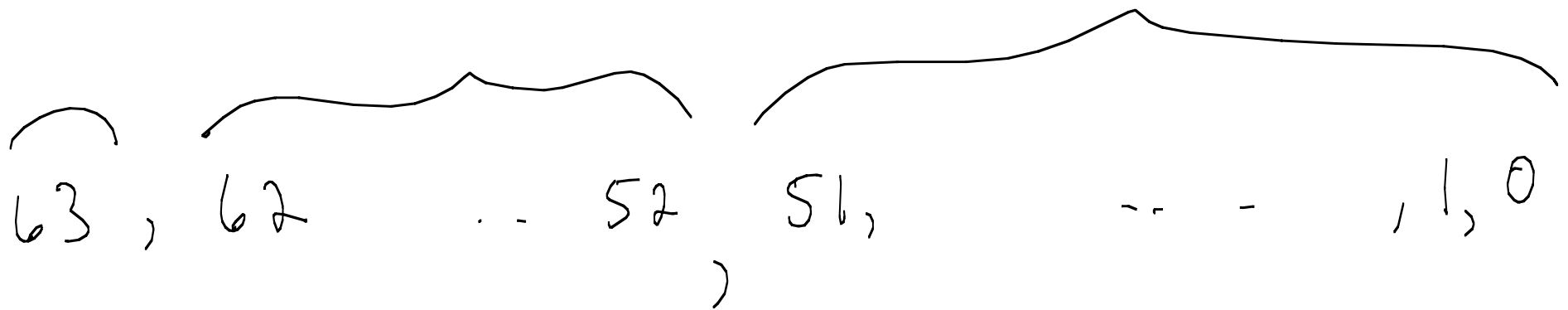
case 2: double precision (64 bits = 8 bytes)

Assignment Project Exam Help

<https://powcoder.com>

52

Add WeChat powcoder



exponent code

exponent value

unsigned exponent code = exponent value + bias

For 11 bits, bias is defined to be $2^{10} - 1 = 1023$.

000000000000

reserved

000000000001

-1022

000000000010

-1021

000000000011

-1020

:

<https://powcoder.com>

:

:

:

011111111111

Add WeChat powcoder

0

100000000000

1

100000000001

2

:

:

:

:

111111111110

1023

111111111111

reserved

Example

$$(8.75)_{10} = (1.00011)_2 \times 2^3$$

significand (52 bits)

[illegible]

Assignment Project Exam Help

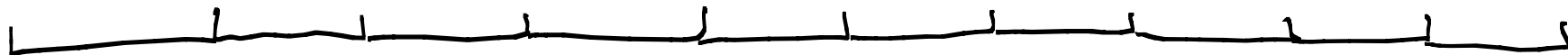
exponent = 3, code using 11 bits:

<https://powcoder.com>

3 + 1023 = 1026 = (100000000010)₂
Add WeCfrat powcoder

double precision float (64 bits)

0 100000000010 0001100000000000000000000000000000...



0 x 4 0 2 1 8 0 0 0 0 0 000000

Q: What is the largest positive normalized number ?
(double precision)

| ~~Assignment Project Exam Help~~ $\times 2^e$

<https://powcoder.com>

Add WeChat powcoder

A:

1023
2

11
(2¹⁰)¹⁰² 2³

Assignment Project Exam Help

<https://powcoder.com>

22
(10³)¹⁰² 10

Add WeChat powcoder

11
10³⁰⁷

Approximation Errors (Java/C/...)

```
double x = 0;
for (int ct=0; ct < 10; ct++) {
    x += 1.0 / 10;
    System.out.println( x );
}
```

Assignment Project Exam Help

0.1

<https://powcoder.com>

0.2

Add WeChat

Add WeChat powcoder

0.4

0.5

0.6

0.7

0.7999999999999999

0.8999999999999999

0.9999999999999999

How many *digits* of precision can we represent with 52 *bits* ?

$$\begin{aligned} & 2^{52} \\ & \approx (2^{10})^5 \cdot 2^2 \\ & \approx (10^3)^5 \cdot 10^0 \\ & \approx 10^{16} \end{aligned}$$

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

52 bits covers about the same "range" as 16 digits.

That is why the print out on the previous slide had up to (about) 16 digits to the right of the decimal point.