**Trinity College Dublin**
**Coláiste na Tríonóide, Baile Átha Cliath**
The University of Dublin

# 4.1: Floating Point Numbers

**CSU11022 – Introduction to Computing II**

D.A.Patterson, J.L.Hennessy, "Computer Organisation and Design: ARM Edition", Morgan-Kaufmann, 2016.

(Section 3.5: Floating Point, available in the Library, doesn't have to be the ARM Edition)

Dr Jonathan Dukes / jdukes@tcd.ie
School of Computer Science and Statistics

```java
// some really small numbers and one large number
float [] vals = {
    3.7e-5f, 4.8e-5f, 1.7e-5f, 2.4e-5f,
    3.7e-5f, 4.8e-5f, 1.7e-5f, 2.4e-5f,
    3.7e-5f, 4.8e-5f, 1.7e-5f, 2.4e-5f,
    3.7e-5f, 4.8e-5f, 1.7e-5f, 2.4e-5f,
    12345.0f
};

float result;

// add the numbers first-to-last
result = 0;
for (int i = 0; i < vals.length; i++) {
  result += vals[i];
}

System.out.println("sum first-to-last: " + String.format("%.8f",result));
// output: sum first-to-last: 12345.00097656

// add the numbers last-to-first
result = 0;
for (int i = vals.length - 1; i >= 0; i--) {
  result += vals[i];
}

System.out.println("sum last-to-first: " + String.format("%.8f",result));
// output: sum first-to-last: 12345.00000000
```

32-bits … $2^{32}$ unique values that we can use to represent different things

e.g. unsigned integers

   0 … $2^{32}-1$ (or 0 … 4,294,967,295)

e.g. signed integers using 2's complement

   $-2^{31}$ … 0 … $+2^{31}-1$ (or –2,147,483,648 … 0 … +2,147,483,647)

How do we represent **real** numbers like 2½ or 3.14159265… ?

Also, how do we represent values with really large or really small magnitudes?

   e.g. 2.2 x $10^{11}$

   e.g. 1.3 x $10^{-8}$

The values $2.2 \times 10^{11}$ and $1.3 \times 10^{-8}$ are examples of (normalized) scientific notation in decimal form

$$f \times 10^e$$

Values expressed in normalised scientific notation satisfy the condition:

$$1 \le |f| < 10$$

Normalized scientific notation give us one *canonical form* in which to express a value using scientific notation and allows quick, visual comparison of magnitude

As computer scientists, we avoid expressing the same thing in different ways (a==b?)

$372.98$

$37.298 \times 10^1$

$3729.8 \times 10^{-1}$

$3.7298 \times 10^2$

Convert the following binary numbers to decimal numbers with fractions

$10010101$ **$= 1\times2^7 + 1\times2^4 + 1\times2^2 + 1\times2^0 = 149$**

$1.1$ **$= 1\times2^0 + 1\times2^{-1} = 1\tfrac{1}{2} = 1.5$**

$101000.01$ **$= 1\times2^5 + 1\times2^3 + 1\times2^{-2} = 40\tfrac{1}{4} = 40.25$**

Convert the following decimal numbers to binary floating point numbers

$5\tfrac{1}{4}$ **$= 101.01$**

$7.75$ **$= 111.11$**

$2.1$ **$= 10.000110011001100\,...$**

$9.3125$ **$= 1001.0101$**

$0.75 \times 2 \;= \mathbf{1}.5$
$0.5 \times 2 \;\;= \mathbf{1}.0$

$0.3125 \times 2 \;= \mathbf{0}.625$
$0.625 \times 2 \;\;= \mathbf{1}.25$
$0.25 \times 2 \;\;\;= \mathbf{0}.5$
$0.5 \times 2 \;\;\;\;= \mathbf{1}.0$

Like decimal values, we can express binary values using scientific notation (again, in normalized form)

e.g.

  1010.1 = 1.0101 x $2^3$

  0.00101 = 1.01 x $2^{-3}$

The general form is again:

$$f \times 2^{-3}$$

and in normalised form, f satisfies the following condition:

$$1_2 \leq |f| < 10_2$$

$$5.75_{10} = 101.11_2 \times 2^0$$
$$= 1.0111_2 \times 2^2$$

The normalized form of a binary number expressed using scientific notation forms the basis for its representation in a computer

# Trinity College Dublin

Coláiste na Tríonóide, Baile Átha Cliath
The University of Dublin

WeChat: cstutorcs

# 4.2: IEEE-754

Assignment Project Exam Help

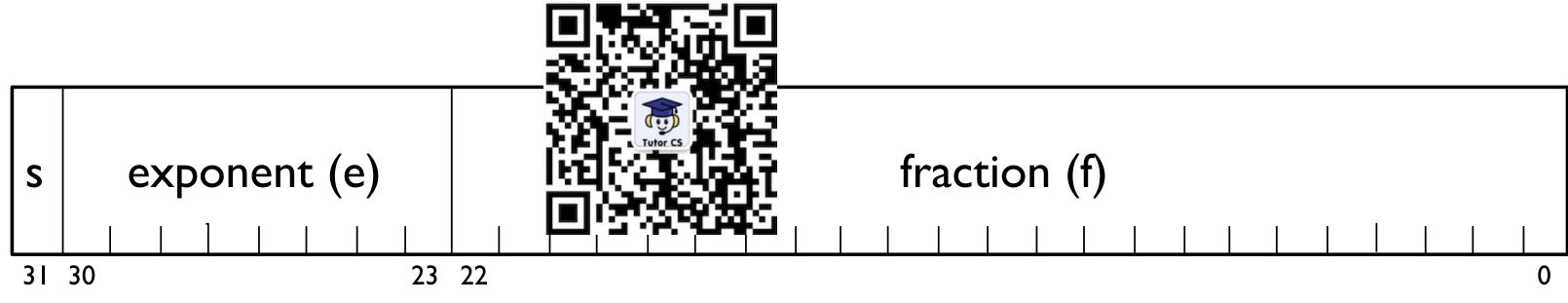CSU11022 – Introduction to Computing II

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

Dr Jonathan Dukes / jdukes@tcd.ie
School of Computer Science and Statistics

https://www.h-schmidt.net/FloatConverter/IEEE754.html

Use a different interpretation of a 32-bit value to represent floating point numbers, e.g. IEEE 754

程序代写代做 CS编程辅导

| s | exponent (e) | fraction (f) |

31 30      23 22      0

WeChat: cstutorcs

$$(-1)^s \times f \times 2^e$$

Assignment Project Exam Help

Email: tutorcs@163.com

How can we represent ...

QQ: 749389476

... positive and negative values?

https://tutorcs.com

... values with positive and negative exponents?

Where is the binary (radix) point?

Sign bit?

0 ⇒ positive floating-point number

1 ⇒ negative floating-point number

Positive and negative exponents?

Option 1: 2's Complement exponents ✗
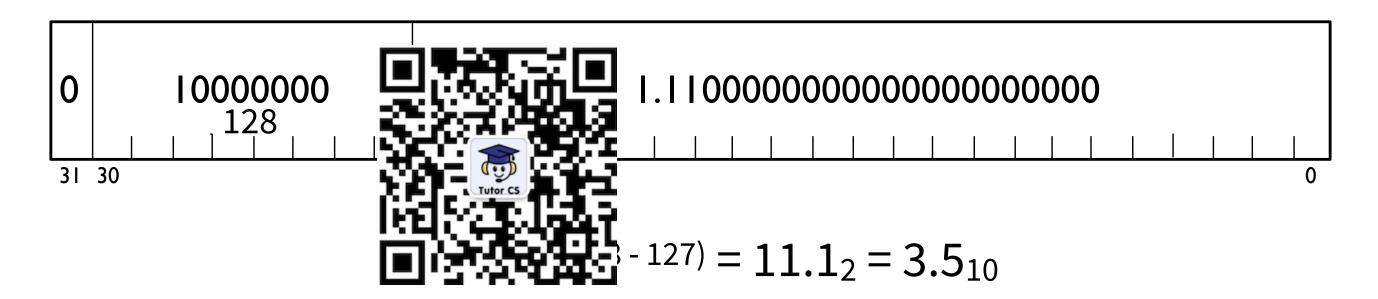
Option 2: Biased exponents ✓

Subtract a constant bias (b = 127) from stored exponent to obtain signed exponent

| s | exponent (e) | fraction (f) |
|---|---|---|
| 31 | 30 ... 23 | 22 ... 0 |

$$(-1)^{s} \times f \times 2^{e-b}$$

The following two representations are of the same value ($3.5_{10}$)

| 0 | 10000000 | 1.110000000000000000000000 |
|---|---|---|

128

31 30 0

$\times 2^{(128 - 127)} = 11.1_2 = 3.5_{10}$

| 0 | 10000001 | 0.111000000000000000000000 |
|---|---|---|

129

31 30 23 22 0

$+0.111 \times 2^{(129 - 127)} = 11.1_2 = 3.5_{10}$ (same value!)

We don't want multiple representations of the same value!

Storing floating-point numbers in normalized form avoids this problem:

$$1_2 \leq |f| < 10_2, \text{ so f is in the form } 1 . ddddd \ldots$$

if (a == b) ...

With normalisation

**01 x 2$^{-4}$**

... becomes ...

**1.0100 x 2$^{-6}$**

adjust fraction so there is a single 1 to left of radix point

compensate by adjusting exponent accordingly

**If there is always going to be a 1 to the left of the radix point, we don't need to store it!**

Increases precision by one bit!

| s | exponent (e) | | fraction (f) |
|---|---|---|---|

31 30 23 22 0

程序代写代做 CS编程辅导

$\times 1 . f \times 2^{(e-b)}$

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

**https://www.h-schmidt.net/FloatConverter/IEEE754.html**

**Trinity College Dublin**

Coláiste na Tríonóide, Baile Átha Cliath

The University of Dublin

# 4.3: IEEE-754 Examples

**CSU11022 – Introduction to Computing II**

Dr Jonathan Dukes / jdukes@tcd.ie

School of Computer Science and Statistics

**https://www.h-schmidt.net/FloatConverter/IEEE754.html**

1.25 = $1.01_2$ (already normalised)

s = 0

e = 0 + 127 = 127 = $01111111_2$

f = $1.01_2$ or $.01_2$ after removing the hidden bit

```
s         e                    f
0    01111111    0100000000000000000000000
```

0011 1111 1010 0000 0000 0000 0000 0000

 3    F    A    0    0    0    0    0

$10.75 = 1010.11_2 \times 2^0 = 1.01011_2 \times 2^3$

$s = 0$

$e = 3 + 127 = 130 = 10000010_2$

$f = 1.01011_2$ or $.01011_2$ after removing the hidden bit

```
s      e                    f
0    10000010    01011000000000000000000

0100 0001 0010 1100 0000 0000 0000 0000
  4     1     2     C     0     0     0     0
```

$-0.125 = -0.001_2 \times 2^0 = 1.0_2 \times 2^{-3}$

s = 1

e = -3 + 127 = 124 = $01111100_2$

f = $1.0_2$ or $.0_2$ after removing the hidden bit

| s | e | f |
|---|---|---|
| 1 | 01111100 | 00000000000000000000000 |

1011 1110 0000 0000 0000 0000 0000 0000

  B     E     0     0     0     0     0     0

程序代写代做 CS编程辅导

```
s    e                 f
0   10000010    100101000000000000000
```

s = 0 (positive)

e = 130 ($2^{130-127} = 2^3$)

f = **1.**100101 (after adding the **hidden bit**)

+1.100101 x $2^3$ = +1100.101 = **+12.625**

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

程序代写代做 CS编程辅导

Special bit patterns, e.g.

Zero (±)

| ? | 00000000 | 00000000000000000000000 |
|---|----------|-------------------------|

31 30       23 22       0

WeChat: cstutorcs

Infinity (±)

Assignment Project Exam Help

| ? | IIIIIIII | 00000000000000000000000 |
|---|----------|-------------------------|

Email: tutorcs@163.com

31 30       23 22       0

QQ: 749389476

Not a Number (NaN)

https://tutorcs.com

| 0 | IIIIIIII | ??????????????????????? ( != 0 ) |
|---|----------|----------------------------------|

31 30       23 22       0

## 32-Bit Single Precision (bias = 127)

| s | exponent (e) | fraction (f) |
|---|---|---|

31 30 23 22 0

## 64-Bit Double Precision (bias = 1023)

63 32

| s | exponent (e) | fraction (f) |
|---|---|---|

fraction (f)

31 0

# Trinity College Dublin

**Coláiste na Tríonóide, Baile Átha Cliath**
The University of Dublin

# 4.4: Floating Point addition

**CSU11022 – Introduction to Computing II**

Dr Jonathan Dukes / jdukes@tcd.ie
School of Computer Science and Statistics

**https://www.h-schmidt.net/FloatConverter/IEEE754.html**

We can add the fractions of floating point values if their exponents are the same

If their exponents are not the same to begin, shift the fraction of the value with the smaller exponent to compensate

e.g.

$1.01101 \times 2^3 + 1.00110 \times 2^{-2}$

$= \quad 1.01101 \times 2^3 + 0.0000100110 \times 2^3$

$= \quad 1.0111000110 \times 2^3$

$$
\begin{array}{r}
1.0110100000 \\
+ \ 0.0000100110 \\
\hline
1.0111000110
\end{array}
$$

```
                          ┌─────────────────┐
                          │     Start       │
                          └────────┬────────┘
                                   │
                                   ▼
┌──────────────────────────────────────────────────────────────────────────────┐
│ Compare the exponents of the two numbers; shift the fraction of the smaller    │
│ number to the right until its exponent would match the larger exponent         │
└────────────────────────────────────┬───────────────────────────────────────────┘
                                   │
                                   ▼
                          ┌─────────────────┐
                          │  Add the fractions │
                          └────────┬────────┘
                                   │
                                   ▼
┌──────────────────────────────────────────────────────────────────────────────┐
│ Normalize the result by either shifting right and incrementing the exponent or │
│ shifting left and decrementing the exponent                                    │
└────────────────────────────────────┬───────────────────────────────────────────┘
                                   │
                                   ▼
                          ╱─────────────────╲      yes
                         ╱     Overflow/      ╲─────────────► Exception!
                         ╲     Underflow      ╱
                          ╲─────────────────╱
                                   │ no
                                   ▼
┌──────────────────────────────────────────────────────────────────────────────┐
│ Round the fraction to the appropriate number of bits                           │
└────────────────────────────────────┬───────────────────────────────────────────┘
                                   │
                                   ▼
                          ╱─────────────────╲
                   no    ╱  Still Normalised ?╲
              ◄─────────╲                     ╱
                          ╲─────────────────╱
                                   │ yes
                                   ▼
                          ┌─────────────────┐
                          │      Done       │
                          └─────────────────┘
```
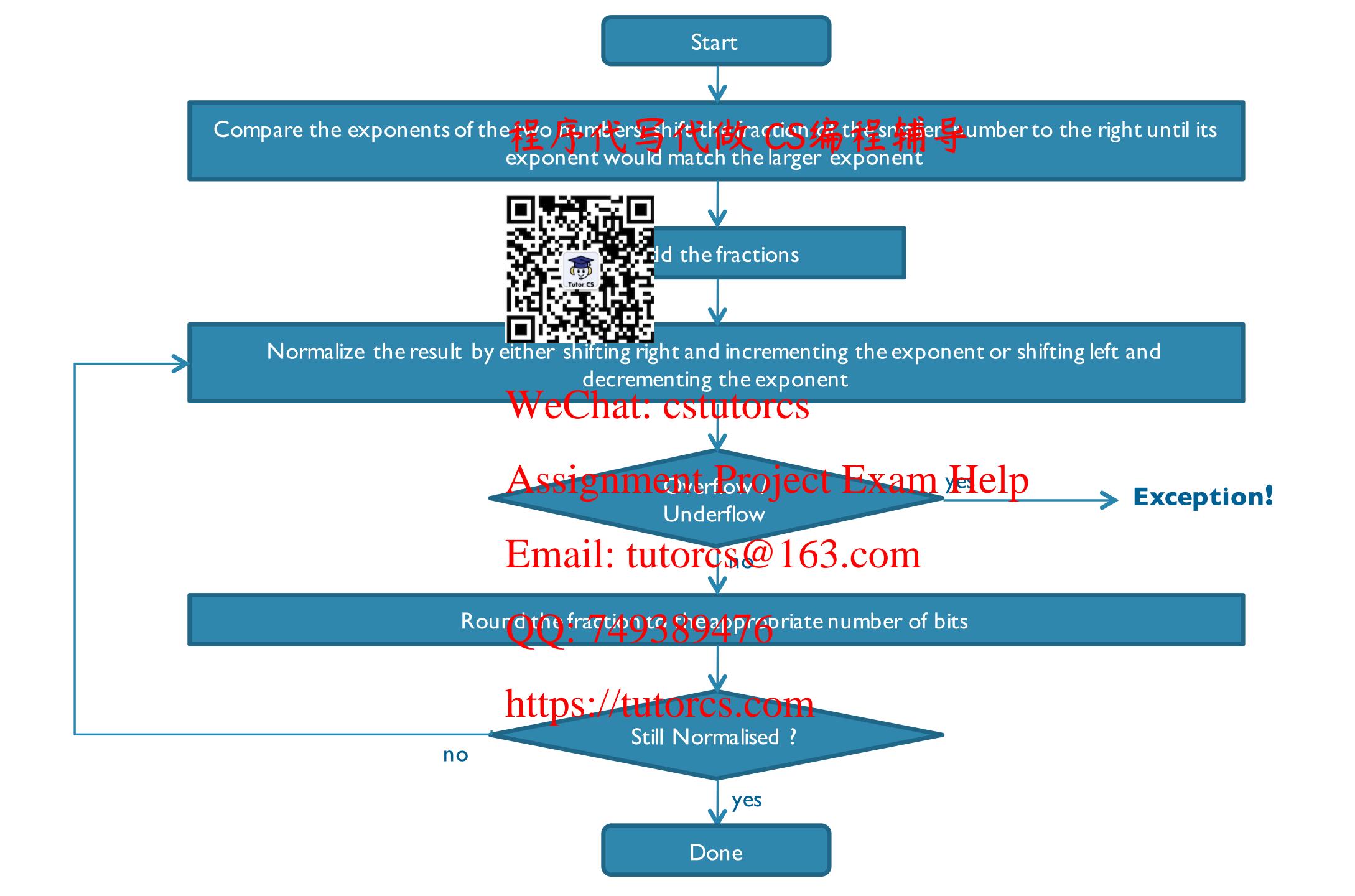
```
A = 0x3fc00000 (1.5)                    B = 0x3fe00000 (1.75)

0 01111111 10000000000000000000000     0 01111111 11000000000000000000000

s      e                  f            s      e              f



s 0 (+)                                s 0 (+)

e 01111111 (127-127=0, 2⁰)             e 01111111 (127-127=0, 2⁰)

f 1.10000 (remember hidden bit!)       f 1.110000 (remember hidden bit!)
```

$$1.100000 \times 2^0 \text{ A}$$
$$1.110000 \times 2^0 \text{ B}$$

```
11.010000 x 2⁰ Result (not normalised)

1.1010000 x 2¹ Result (normalised)

0 10000000 10100000000000000000000 (encoding s e f)

0x40500000 (3.25)
```

程序代写代做 CS编程辅导

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

A = 0x3fc00000 (1.5)

0 01111111 10000000000000000000000

s    e                  f

B = 0x41280000 (10.5)

0 10000010 01010000000000000000000

s    e                  f

s 0 (+)

e 01111111 (127-127=0, $2^0$)

f **1.**10000 (remember hidden bit!)

s 0 (+)

e 10000010 (130-127=3, $2^3$)

f **1.0**101000 (remember hidden bit!)

0.0011000 x $2^3$ A (adjust fraction so exponents are equal)

<u>1.0101000 x $2^3$ B</u>

1.1000000 x $2^3$ Result (already normalised)

0 10000010 10000000000000000000000 (encoding s e f)

0x41400000 (12.0)

程序代写代做 CS编程辅导

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com

程序代写代做 CS编程辅导

What about adding negative values (S==1)?

Proceed as before but before adding, the fractions of values with S==1 should be converted to their 2's Compliment

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

https://tutorcs.com