

Exponent encoding

The single precision binary floating-point exponent is encoded using an offset binary representation, with the zero offset being 127; also known as exponent bias in the IEEE 754 standard.

- $E_{\min} = 01_H - 7F_H = -126$
- $E_{\max} = FE_H - 7F_H = 127$
- Exponent bias = $7F_H = 127$

Thus, in order to get the true exponent as defined by the offset binary representation, the offset of 127 has to be subtracted from the stored exponent.

The stored exponents 00_H and FF_H are interpreted specially.

Exponent	Significand zero	Significand non-zero	Equation
00_H	zero, -0	subnormal numbers	$(-1)^{\text{signbits}} \times 2^{-126} \times 0.\text{significandbits}$
$01_H, \dots, FE_H$	normalized value		$(-1)^{\text{signbits}} \times 2^{\text{exponentbits}-127} \times 1.\text{significandbits}$
FF_H	\pm infinity	NaN (quiet, signalling)	

The minimum positive (subnormal) value is $2^{-149} \approx 1.4 \times 10^{-45}$. The minimum positive normal value is $2^{-126} \approx 1.18 \times 10^{-38}$. The maximum representable value is $(2-2^{-23}) \times 2^{127} \approx 3.4 \times 10^{38}$.

Converting from decimal representation to binary32 format

In general refer to the IEEE 754 standard itself for the strict conversion (including the rounding behaviour) of a real number into its equivalent binary32 format.

Here we can show how to convert a base 10 real number into an IEEE 754 binary32 format using the following outline :

- consider a real number with an integer and a fraction part such as 12.375
- convert the integer part into binary as shown in Binary numeral system
- convert the fraction part using the following technique as shown here
- add the two results and adjust them to produce a proper final conversion

Conversion of the fractional part:

consider 0.375, the fractional part of 12.375. To convert it into a binary fraction, multiply the fraction by 2, take the integer part and re-multiply new fraction by 2 until a fraction of zero is found or until the precision limit is reached which is 23 fraction digits for IEEE 754 binary32 format.

$0.375 \times 2 = 0.750 = 0 + 0.750 \Rightarrow b_{-1} = 0$, the integer part represents the binary fraction digit. Re-multiply 0.750 by 2 to proceed

$0.750 \times 2 = 1.500 = 1 + 0.500 \Rightarrow b_{-2} = 1$

$0.500 \times 2 = 1.000 = 1 + 0.000 \Rightarrow b_{-3} = 1$, fraction = 0.000, terminate

We see that $(0.375)_{10}$ can be exactly represented in binary as $(0.011)_2$. Not all decimal fractions can be represented in a finite digit binary fraction. For example decimal 0.1 cannot be represented in binary exactly. So it is only approximated.

Therefore $(12.375)_{10} = (12)_{10} + (0.375)_{10} = (1100)_2 + (0.011)_2 = (1100.011)_2$

Also IEEE 754 binary32 format requires that you represent real values in $(1.x_1x_2\dots x_{23})_2 \times 2^e$ format, (see Normalized number, Denormalized number) so that 1100.011 is shifted to the left by 3 digits to become $(1.100011)_2 \times 2^3$

Finally we can see that: $(12.375)_{10} = (1.100011)_2 \times 2^3$

From which we deduce:

- The exponent is 3 (and in the biased form it is therefore $130 = 1000\ 0010$)
- The fraction is 100011 (looking to the right of the binary point)

From these we can form the resulting 32 bit IEEE 754 binary32 format representation of 12.375 as:
 $0-10000010-100011000000000000000000 = 41460000_{\text{H}}$

Note: consider converting 68.123 into IEEE 754 binary32 format: Using the above procedure you expect to get $42883\text{EF}9_{\text{H}}$ with the last 4 bits being 1001 However due to the default rounding behaviour of IEEE 754 format what you get is $42883\text{EFA}_{\text{H}}$ whose last 4 bits are 1010 .

Ex: Consider decimal 1 We can see that : $(1)_{10} = (1.0)_2 \times 2^0$

From which we deduce :

- The exponent is 0 (and in the biased form it is therefore $127 = 0111\ 1111$)
- The fraction is 0 (looking to the right of the binary point in 1.0 is all 0 = 000...0)

From these we can form the resulting 32 bit IEEE 754 binary32 format representation of real number 1 as:
 $0-01111111-000000000000000000000000 = 3\text{f}800000_{\text{H}}$

Ex: Consider a value 0.25 . We can see that : $(0.25)_{10} = (1.0)_2 \times 2^{-2}$

From which we deduce :

- The exponent is -2 (and in the biased form it is $127+(-2)= 125 = 0111\ 1101$)
- The fraction is 0 (looking to the right of binary point in 1.0 is all zeros)

From these we can form the resulting 32 bit IEEE 754 binary32 format representation of real number 0.25 as:
 $0-01111101-000000000000000000000000 = 3\text{e}800000_{\text{H}}$

Ex: Consider a value of 0.375 . We saw that $0.375 = (1.1)_2 \times 2^{-2}$

Hence after determining a representation of 0.375 as $(1.1)_2 \times 2^{-2}$ we can proceed as above :

- The exponent is -2 (and in the biased form it is $127+(-2)= 125 = 0111\ 1101$)
- The fraction is 1 (looking to the right of binary point in 1.1 is a single $1 = x_1$)

From these we can form the resulting 32 bit IEEE 754 binary32 format representation of real number 0.375 as:
 $0-01111101-100000000000000000000000 = 3\text{e}\text{c}00000_{\text{H}}$

Single precision examples

These examples are given in bit *representation*, in hexadecimal, of the floating point value. This includes the sign, (biased) exponent, and significand.

```

3f80 0000    = 1
c000 0000    = -2

7f7f ffff    ≈ 3.4028234 × 1038 (max single precision)

0000 0000    = 0
8000 0000    = -0

7f80 0000    = infinity
ff80 0000    = -infinity

3eaa aaab    ≈ 1/3

```

By default, 1/3 rounds up instead of down like double precision, because of the even number of bits in the significand. So the bits beyond the rounding point are 1010... which is more than 1/2 of a unit in the last place.

Converting from single precision binary to decimal

We start with the hexadecimal representation of the value, 41c80000, in this example, and convert it to binary

$$41c8\ 0000_{16} = 0100\ 0001\ 1100\ 1000\ 0000\ 0000\ 0000\ 0000_2$$

then we break it down into three parts; sign bit, exponent and significand.

Sign bit: 0

Exponent: $1000\ 0011_2 = 83_{16} = 131$

Significand: $100\ 1000\ 0000\ 0000\ 0000\ 0000_2 = 480000_{16}$

We then add the implicit 24th bit to the significand

$$\text{Significand: } \mathbf{1}100\ 1000\ 0000\ 0000\ 0000\ 0000_2 = C80000_{16}$$

and decode the exponent value by subtracting 127

Raw exponent: $83_{16} = 131$

Decoded exponent: $131 - 127 = 4$

Each of the 24 bits of the significand, bit 23 to bit 0, represents a value, starting at 1 and halves for each bit, as follows

bit 23 = 1

bit 22 = 0.5

bit 21 = 0.25

bit 20 = 0.125

bit 19 = 0.0625

.

.

bit 0 = 0.00000011920928955078125

The significand in this example has three bits set, bit 23, bit 22 and bit 19. We can now decode the significand by adding the values represented by these bits.

$$\text{Decoded significand: } 1 + 0.5 + 0.0625 = 1.5625 = C80000/2^{23}$$

Then we need to multiply with the base, 2, to the power of the exponent to get the final result

$$1.5625 \times 2^4 = \mathbf{25}$$

Thus

$$41c8\ 0000 = 25$$

This is equivalent to:

where s is the sign bit, x is the exponent, and m is the significand in base 10.

External links

- Online calculator ^[2]
- Online converter for IEEE 754 numbers with single precision ^[3]
- C source code to convert between IEEE double, single, and half precision can be found here ^[4]

References

- [1] <http://java.sun.com/docs/books/tutorial/java/nutsandbolts/datatypes.html>
- [2] <http://www.h-schmidt.net/FloatApplet/IEEE754.html>
- [3] http://www.binaryconvert.com/convert_float.html
- [4] <http://www.mathworks.com/matlabcentral/fileexchange/23173>
-

Article Sources and Contributors

Single precision floating-point format *Source:* <http://en.wikipedia.org/w/index.php?oldid=448992900> *Contributors:* -

Image Sources, Licenses and Contributors

Image:Float example.svg *Source:* http://en.wikipedia.org/w/index.php?title=File:Float_example.svg *License:* GNU Free Documentation License *Contributors:* en:User:Fresheneesz, traced by User:Stannered

License

Creative Commons Attribution-Share Alike 3.0 Unported
<http://creativecommons.org/licenses/by-sa/3.0/>
