
COMP1521 WK07

— number representations —

Course Updates

- Assignment 1 was due this Monday! Congrats for making it through the first half of the course 🎉
- 5 days to submit with the late penalty (before 0 is awarded)
- Assignment 2 releasing some time soon
- C Assignment

Contents

- Negative numbers - Two's Complement
 - Q1, Q2, Q3
- Floating point IEEE-754 representation
 - Q4, Q5

Two's Complement

- Negative numbers are determined by whether or not the m.s.b (most significant bit) is set to 1

Example.

-15 = 0b10001



Representing Two's Complement Numbers in Binary

1. Represent the number in binary
2. Flip all the bits (set 0s to 1s and 1s to 0s)
3. Add one to the number

Example. num = -127

1. 127 = -0b0111 1111

2. = 0b1000 0000

3. = 0b1000 0001



Tutorial Q2

Assume that the following hexadecimal values are 16-bit twos-complement. Convert each to the corresponding decimal value.

a. **0x0013** 0000 0000 0001 0011 $\rightarrow 19$

b. **0x1234** \rightarrow 0001 0010 0011 0100 \rightarrow

c. **0xffff**

d. **0x8000**

1000 0000 0000 0000 $\rightarrow -1$

1000 0000 0000 0001 $\rightarrow -(2^{15} + 1)$

Tutorial Q3

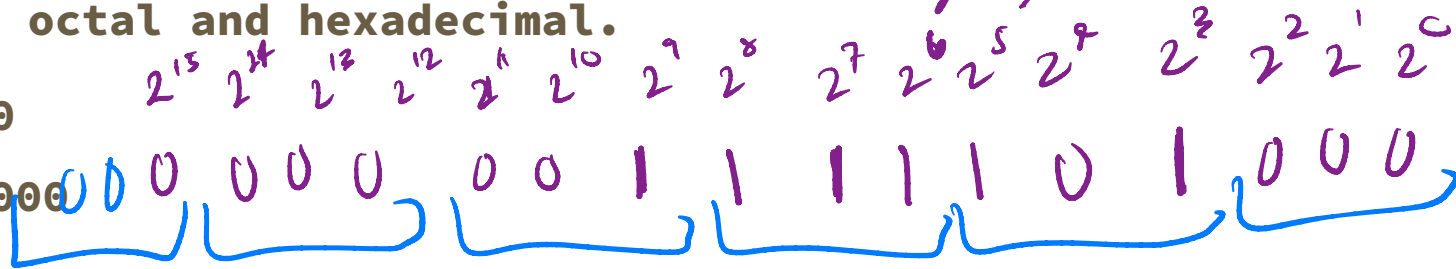
Give a representation for each of the following decimal values in 16-bit two's-complement bit-strings. Show the value in binary, octal and hexadecimal.

a. 1000

b. 100000

c. -5

d. -100



0x03C8

001750

Tutorial Q3

$$\text{max pos} \rightarrow \underline{2^{15} - 1}$$

Give a representation for each of the following decimal values in 16-bit two's-complement bit-strings. Show the value in binary, octal and hexadecimal.

a. 1000

b. 100000

c. -5

d. -100

$$2^{16} = \underline{65536}$$

Tutorial Q3

Give a representation for each of the following decimal values in 16-bit two's-complement bit-strings. Show the value in binary, octal and hexadecimal.

a. 1000 1. 0000 0000 0000 0101

b. 100000

c. -5 2. 1111 1111 1111 1011

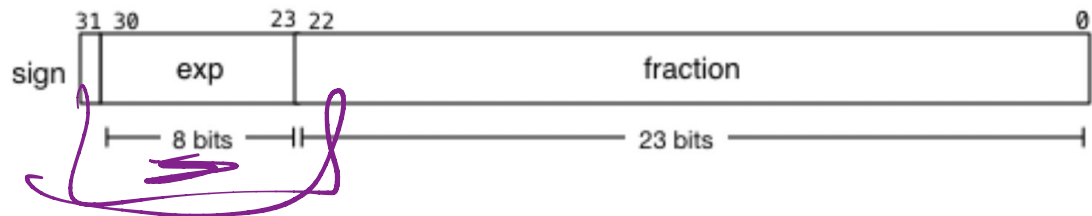
d. -100

0xFFFF 0177773

Floating Point Numbers

- We need a way to represent base 10 Real numbers in Binary
- Easy to represent Integers in base2 but Fractional numbers can be tough to get accurate (can cause floating point errors)
- Base10 has scientific notation
 - $6.022 * 10^{23}$
- Base2 also has scientific notation
- IEEE-754 is the standard for floating point arithmetic

Representing Floating Point Numbers in Binary



$$2^7 - 1$$
$$2^{(\# \text{ exp bits})} - 1$$

- **sign bit** → indicates if number is positive or negative
- **exponent bits** → unsigned 8 bit value which is interpreted from -127 to 128 by subtracting the "bias" **127**
- **fraction (mantissa) bits** → 23 bits whose values are in the range 0 to 1. We add 1 to the fractional component to convert the values to be between 1 to 2

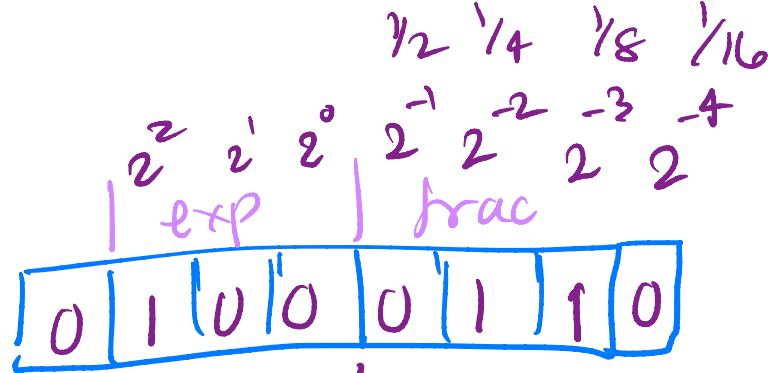
$$\text{sign} * (1 + \text{fraction}) * 2^{(\text{exponent} - \text{bias})}$$

8 bit Floating Point Example

- **sign bit** → 1 bit
- **exponent bits** → 3 bits (bias = 3)
- **fraction (mantissa) bits** → 4 bits

Example: Convert 01000110

$$\begin{aligned} \left(1 + \frac{3}{8}\right) \times 2^{4-3} &= \left(1 + \frac{3}{8}\right) \times 2 \\ &= 2.75 \leftarrow \end{aligned}$$



sign = 0 → pos

$$\text{exp} = 4$$

$$\text{frac} = \frac{1}{4} + \frac{1}{8} = \frac{3}{8}$$

Tutorial Q4

What decimal numbers do the following single-precision IEEE 754-encoded bit-strings represent?

a. 0 00000000 00000000000000000000000000000000 $\rightarrow 0$

b. 1 00000000 00000000000000000000000000000000

c. 0 01111111 10000000000000000000000000000000

d. 0 01111110 00000000000000000000000000000000 \rightarrow

e. 0 01111110 11111111111111111111111111111111

f. 0 10000000 01100000000000000000000000000000

g. 0 10010100 10000000000000000000000000000000

h. 0 01101110 101000001010000010100000

$$\text{exp} = 0$$

$$\text{frac} = 0$$

$$(1 + 0) \times 2^{0-127}$$

$$\frac{1}{2^{127}} \approx 0$$

$$b = 0$$

Each of the above is a single 32-bit bit-string, but partitioned to show the sign, exponent and fraction parts.

Tutorial Q4

What decimal numbers do the following single-precision IEEE 754-encoded bit-strings represent?

- a. 0 00000000 000000000000000000000000
- b. 1 00000000 000000000000000000000000
- c. 0 01111111 100000000000000000000000
- d. 0 01111110 000000000000000000000000
- e. 0 01111110 111111111111111111111111
- f. 0 10000000 011000000000000000000000
- g. 0 10010100 100000000000000000000000
- h. 0 01101110 101000001010000010100000

sign = 0 → pos

exp = 127

frac = $\frac{1}{2}$

$$\left(1 + \frac{1}{2}\right) \times 2^{127-127} = 1.5 \times 2^0 = 1.5$$

Each of the above is a single 32-bit bit-string, but partitioned to show the sign, exponent and fraction parts.

Tutorial Q4

What decimal numbers do the following single-precision IEEE 754-encoded bit-strings represent?

- a. 0 00000000 000000000000000000000000
- b. 1 00000000 000000000000000000000000
- c. 0 01111111 100000000000000000000000
- d. 0 01111110 000000000000000000000000
- e. 0 01111110 111111111111111111111111
- f. 0 10000000 011000000000000000000000
- g. 0 10010100 100000000000000000000000
- h. 0 01101110 101000001010000010100000

$$\text{sign} = 0 \rightarrow \text{pos}$$

$$\text{exp} = 126$$

$$\text{frac} = \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} + \dots + \frac{1}{2^{23}}$$
$$\approx 0.9999\dots$$

$$(1 + 0.999\dots) \times 2^{126-127}$$

Each of the above is a single 32-bit bit-string, but partitioned to show the sign, exponent and fraction parts.

$$= 1.999 \times 2^{-1} = \frac{1.999\dots}{2} \approx 1$$

Tutorial Q5

$$(1 + \text{frac}) \times 2^{\text{exp - bias.}}$$

Convert the following decimal numbers into IEEE 754-encoded bit-strings:

a. 2.5 =

$$(1 + 0.25) \times 2^1$$

$$\begin{aligned} \text{exp} - 127 &= 1 \\ \text{exp} &= \underline{\underline{128}} \end{aligned}$$

b. 0.375

c. 27.0

d. 100.0

$$\text{frac} = \frac{1}{4} \rightarrow \underline{0100000000000000\dots}$$

$$010000000001000000000000\dots$$

Tutorial Q5

$$(1 + \text{frac}) \times 2^{\text{exp} - 127}$$

Convert the following decimal numbers into IEEE 754-encoded bit-strings:

a. 2.5

b. 0.375

c. 27.0

d. 100.0

$$(1 + 0.5) \times 2^{-2}$$

-2

$$\text{exp} = -2 + 127 = 125$$

$$\text{frac} = 1000000 \dots$$

$$\text{exp} = 01111101$$

$$0 \ 01111101 \ 1000000 \dots$$

Special Floating Point Numbers

0000000000 11 00 110.

- **0** → exponent bits are 0

- **Infinity** → exponent bits are all 1's, fraction is 0

00000000

- **NaN** → exponent bits are all 1's, fraction is not 0



