Q1- What decimal number does the bit pattern 0×0C000000 represent if it is a floating point number? Use the IEEE 754 standard.

$(0 \times 0 C 000000)_{16}$

Mantissa

$= 0000\ 1100\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000$

- 1-bit sign (S) $= 0 \Rightarrow$ positive

- 8-bit exponent (e) $= 00011000 = 24$
- 23-bit fraction/mantissa (M) $= 0$

actual exponent $= 24 - 127 = -103$

$$= (-1)^{0} \times (1 + \overset{\text{Mantissa}}{F}) \times 2^{(e-127)}$$

$$= (-1)^{S} (1.0) \times 2^{-103}$$

$$= 1.0 \times 2^{-103}$$

Q2 Solution
0 0C000000 = 0 0 0 0 1 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 = 0 0 0 0 1 1 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0

sign is positive
exponent = 0 x18 = 24-127= -103
there is a hidden 1
mantissa = 0
answer 1.0 x $2^{-103}$

Q2- Write down the binary representation of the decimal number 63.25 assuming the IEEE 754 double precision format.

Q2 Solution

$63.25 \times 100 = 111111.01 \times 2^0$

normalize, move binary point five to the left

$1.1111101 \times 2^5$

sign = positive, exp = 1023 + 5 = 1028

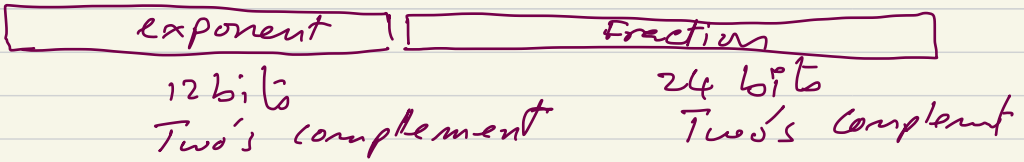Final bit pattern:

0 100 0000 0100 1111 1010 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000 0000

= 0x404FA00000000000

Write down the binary bit pattern to represent -1.5625 x 10^-1 assuming a format similar to that employed by the DEC PDP-8 (the leftmost 12 bits are the exponent stored as a two's complement number, and the rightmost 24 bits are the fraction stored as a two's complement number). No hidden 1 is used. Comment on how the range and accuracy of this 36-bit pattern compares to the single and double precision IEEE 754 standards.

**Solution**

Convert the decimal value to binary.

$$-1.5625 \times 10^{-1} = -0.15625$$

$$0.15625 = \tfrac{1}{8} + \tfrac{1}{32} = 0.00101_2$$

So $-0.15625 = -0.00101_2$.

Normalize so that the first bit after the point is 1:

$$-0.00101_2 = -0.101_2 \times 2^{-2}$$

Therefore the fraction value $= -0.101_2 = -0.625$ and the exponent $= -2$.

**Fraction (24 bits, two's complement):**

$+0.101_2$ as a 24-bit fraction: $0101\,0000\,0000\,0000\,0000\,0000$

Negate (invert and add 1 LSB):

$$1011\,0000\,0000\,0000\,0000\,0000$$

Check: $-1 + 0.011_2 = -1 + 0.375 = -0.625$ ✓

**Exponent (12 bits, two's complement):**

$$-2 = 1111\,1111\,1110$$

**Full 36-bit pattern:**

$$1111\,1111\,1110\ \ 1011\,0000\,0000\,0000\,0000\,0000$$

**Comments on range and accuracy**

- **Range:** The 12-bit two's complement exponent spans roughly $-2048$ to $+2047$, giving magnitudes from about $2^{-2048}$ to $2^{2047}$ ($\approx 10^{\pm 616}$). This far exceeds IEEE 754 single precision (8-bit exponent, $\approx 10^{\pm 38}$) and is even larger than double precision (11-bit exponent, $\approx 10^{\pm 308}$).

- **Accuracy:** The 24-bit fraction provides about 23–24 bits of significance (one bit used for sign, no hidden 1). This is comparable to single precision, which has 24 bits of significand (23 stored + 1 hidden). It is considerably less precise than double precision, which has 53 bits of significand.

In summary, this 36-bit format offers a much wider exponent range than either IEEE 754 format, but its precision is only on par with single precision and well below double precision.

$-1.5625 \times 10^{-1}$

| exponent | Fraction |
|---|---|

12 bits
Two's complement

24 bits
Two's Complement

No Hidden 1 is used

| | |
|---|---|
| $0.15625 \times 2$ | 0 |
| $0.3125 \times 2$ | 0 |
| $0.625 \times 2$ | 1 |
| $0.25 \times 2$ | 0 |
| $0.5 \times 2$ | 1 |

$-0.15625 \rightarrow -0.00101$

No leading 1
No hidden bit

$-0.101 \times 2^{-2}$

Fraction = 101

exponent = -2

Encode the fraction (24-bit two's complement)

$0.101)_2 = 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3}$

$= \frac{1}{2} + 0 + \frac{1}{8} = 0.5 + 0.125$

$= 0.625$

$- 0.101)_2 \rightarrow - 0.625)_{10}$

In Two's complement

$+ 0.625 \rightarrow 0.10100\text{-----}0$

Invert bits $\rightarrow 1.010111\text{-----}1$

add 1 $\rightarrow 1.011000\text{-----}0$

Fraction number $\rightarrow 101100\text{-----}0$

---

Encode the Exponent (12-bit two's complement)

$-2$ in 12-bit two's complement

$+2: 0\text{----}00010$

Invert: $1\text{-----}1101$

add 1: $1\text{-----}110$

Exponent $= 1\text{-----}1110$

Full 36-bit Pattern

$1111\text{----}1110\ 101100\text{-------}0$

12 bits     24 bits

**Range** : The PDP-8 format has a much larger exponent range

far beyond IEEE standard $(-2048 \rightarrow 2047)$

$\pm 127$ for single $^{VS}$ $\pm 1023$ for double

Precision : PDP-8 matches single- precision's decimal digits

$2^{-23} \rightarrow 23 \log_{10} 2 = 23 \times 0.3 \approx 6$ decimal digits

but lack the hidden bit

Double precision offers ~16 digits with 52 bits

$52 \log_{10} 2 = 52 \times 0.3 = 16$

Q3 Solution

$-1.5625 \times 10^{-1} = -0.15625 \times 10^{0}$
$= -0.00101 \times 2^{0}$
move the binary point two to the right
$= -0.101 \times 2^{-2}$
exponent $= -2$,
fraction $= -0.10100000000000000000000$
answer: 111111111110 10110000000000000000000