

a	5.00736125×10^5
b	$-2.691650390625 \times 10^{-2}$

1. Write down the binary representation of the decimal number, assuming the IEEE 754 **single** precision format.
2. Write down the binary representation of the decimal number, assuming the IEEE 754 **double** precision format.

Answer:

a)

$$5.00736125 \times 10^5 = 500736.125 = 1111010010000000000.001$$

$$= 1.111010010000000000001 \times 2^{18}$$

1. Biểu diễn theo độ chính xác đơn

$$\text{sign} = 0$$

$$\text{exponent} = 18 + 127 = 145 = 10010001_2$$

$$\text{fraction} = 111010010000000000001_2$$

0	1	0	0	1	0	0	0	1	1	1	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
s	exponent (8 bits)								fraction (23 bits)																						

3. Biểu diễn theo độ chính xác kép

$$\text{sign} = -1$$

$$\text{exponent} = 18 + 1023 = 1041 = 10000010001_2$$

$$\text{fraction} = 111010010000000000001_2$$

0	1	0	0	0	0	0	1	0	0	0	1	1110100100000000000000001000...0000															
s	Exponent (11 bits)											Fraction (52 bits) với: 21 bits cao =1110100100000000000001 31 bits thấp còn lại bằng 0															

b)

$$-2.691650390625 \times 10^{-2} = 0.02691650390625 = 441/2^{14} = 110111001 \times 2^{-14}$$

$$= 1.10111001 \times 2^{-6}$$

1. Biểu diễn theo độ chính xác đơn

$$\text{sign} = -1$$

$$\text{exponent} = -6 + 127 = 121 = 1111001_2$$

$$\text{fraction} = 10111001_2$$

Nếu đây là số có dấu dạng bù 2, số tương ứng = -1346437120
 Nếu đây là số có không dấu, số tương ứng = 2948530176

2. a)

0x24A60004 = 0010 0100 1010 0110 0000 0000 0000 0100

Nếu đây là lệnh assembly, lệnh tương ứng: *addiu \$6, \$5, 4*
 Hay *addiu \$a2, \$a1, 4*

b)

0xAFBF0000 = 1010 1111 1011 1111 0000 0000 0000 0000

Nếu đây là lệnh assembly, lệnh tương ứng: *sw \$31, 0(\$29)*
 Hay *sw \$ra, 0(\$sp)*

3. a)

0x24A60004 = 0010 0100 1010 0110 0000 0000 0000 0100

Nếu đây là số floating-point với độ chính xác đơn

0010010010100110 0000 0000 0000 0100

Bit thứ 31 = 0, là bit dấu của số floating-point → số floating-point này là số dương

Bit thứ 30 tới 23 = **01001001** = $73_{(10)}$, là phần mũ của số floating-point sau khi đã cộng thêm 127 → số floating-point này có số mũ = $73 - 127 = -54$

Bit thứ 22 tới 0 = 0100110 0000 0000 0000 0100, là phần thập phân của floating-point

⇒ số floating-point = $1.0100110\ 0000\ 0000\ 0000\ 0100 \times 2^{-54}$

b)

0xAFBF0000 = 1010 1111 1011 1111 0000 0000 0000 0000

Nếu đây là số floating-point với độ chính xác đơn

⇒ số floating-point = $-1.011\ 1111 \times 2^{-32}$

Exercise 6:

The following table shows pairs of decimal numbers

	A	B
a.	<u>-1278</u> × 10 ³	-3.90625 × 10 ⁻¹
b.	2.3109375 × 10 ¹	6.391601562 × 10 ⁻¹

Sửa lại thành
 -1.278 × 10³

1. Calculate the sum of A and B by hand, assuming that we keep 11 bits of significand and 5 bits of the exponent. (Rounding rule: add 1 if the bits to the right of the desired point is larger or equal to $100_{(2)}$). Show all the steps.
2. Calculate the sum of A and B by hand, assuming A and B are stored in the IEEE-754 single precision format. Show all the steps.

Answer:

1. Đề bài yêu cầu tính tổng A và B bằng tay (tức chạy từng bước) với giả sử số floating-point chỉ cho phép dùng 11 bits cho phần significand và 5 bits cho phần exponent
a)

$$A = -1.278 \times 10^3 = -1278 = -10011111110 \\ = -1.0011111110 \times 2^{10}$$

(kiểm tra số floating-point A đã đúng chuẩn chưa:
10 nằm trong phạm vi số 5 bits của phần mũ
và phần significand '1.0011111110' đúng 11 bits cho phép
→ Đúng chuẩn)

$$B = -3.90625 \times 10^{-1} = -0.390625 = -25/2^6 = -11001 \times 2^{-6} \\ = -1.1001 \times 2^{-2}$$

(kiểm tra số floating-point B này đã đúng chuẩn chưa:
-2 nằm trong phạm vi số 5 bits của phần mũ
và phần significand '1.1001' cũng không vượt quá 11 bits
→ Đúng chuẩn)

$$A + B = - (1.0011111110 \times 2^{10} + 1.1001 \times 2^{-2}) \\ = - (1.0011111110 \times 2^{10} + 0.000000000011001 \times 2^{10}) \\ = -1.0011111110\mathbf{011001} \times 2^{10}$$

Do phần significand chỉ được phép chứa 11 bits, nên A + B phải được làm tròn, phần cắt bỏ là $011001_{(2)} > 100_{(2)}$ nên $1.0011111110\mathbf{011001} \approx 1.0011111111$

$$\text{Vậy } A + B = -1.0011111111 \times 2^{10} \\ = -1001111111_{(2)} = 1279$$

b)

$$A = 2.3109375 \times 10^1 = 23.109375 \\ = 23 + 7/2^6 \\ = 10111.000111_{(2)} \\ = 1.0111000111 \times 2^4$$

(kiểm tra số floating-point A đã đúng chuẩn chưa:
4 nằm trong phạm vi số 5 bits của phần mũ
và phần significand '1.0111000111' cũng không vượt quá 11 bits
→ Đúng chuẩn)

$$B = 6.391601562 \times 10^{-1} = 0.6391601562 \\ = 1309/2^{11} = 10100011101 \times 2^{-11} \\ = 1.0100011101 \times 2^{-1}$$

(kiểm tra số floating-point B này đã đúng chuẩn chưa:
-1 nằm trong phạm vi số 5 bits của phần mũ

và phần significand '1.0100011101' cũng không vượt quá 11 bits
 ➔ Đúng chuẩn)

$$\begin{aligned} A + B &= 1.0111000111 \times 2^4 + 1.0100011101 \times 2^{-1} \\ &= 1.0111000111 \times 2^4 + 0.000010100011101 \times 2^4 \\ &= 1.01111011111101 \times 2^4 \end{aligned}$$

Do phần significand chỉ được phép chứa 11 bits, nên A + B phải được làm tròn, phần cắt bỏ là 11101₍₂₎ > 100₍₂₎ nên 1.01111011111101 ≈ 1.0111110000
 Vậy A + B = 1.0111110000 × 2⁴
 = 23.75₍₁₀₎

2. Đề bài yêu cầu tính tổng A và B bằng tay (tức chạy từng bước) với giả sử số floating-point dùng format IEEE độ chính xác đơn

a)

$$\begin{aligned} A &= -1.278 \times 10^3 = -1278 = -10011111110 \\ &= -1.0011111110 \times 2^{10} \end{aligned}$$

(kiểm tra số floating-point A đã đúng chuẩn chưa:

(10 + 127) nằm trong phạm vi số 8 bits của phần mũ
 và phần fraction '0011111110' cũng không vượt quá 23 bits
 ➔ Đúng chuẩn)

$$\begin{aligned} B &= -3.90625 \times 10^{-1} = -0.390625 = -25/2^6 = -11001 \times 2^{-6} \\ &= -1.1001 \times 2^{-2} \end{aligned}$$

(kiểm tra số floating-point B này đã đúng chuẩn chưa:

(-2 + 127) nằm trong phạm vi số 8 bits của phần mũ
 và phần fraction '1001' cũng không vượt quá 23 bits
 ➔ Đúng chuẩn)

$$\begin{aligned} A + B &= - (1.0011111110 \times 2^{10} + 1.1001 \times 2^{-2}) \\ &= - (1.0011111110 \times 2^{10} + 0.0000000000011001 \times 2^{10}) \\ &= -1.0011111110011001 \times 2^{10} \end{aligned}$$

Phần fraction này chứa 16 bits, không vượt quá 23 bits của IEEE độ chính xác đơn, nên:

$$\text{Vậy } A + B = -1.0011111110011001 \times 2^{10}$$

b)

$$\begin{aligned} A &= 2.3109375 \times 10^1 = 23.109375 \\ &= 23 + 7/2^6 \\ &= 10111.000111_{(2)} \\ &= 1.0111000111 \times 2^4 \end{aligned}$$

(kiểm tra số floating-point A đã đúng chuẩn chưa:

(4 + 127) nằm trong phạm vi số 8 bits của phần mũ
 và phần fraction '0111000111' cũng không vượt quá 23 bits
 ➔ Đúng chuẩn)

$$B = 6.391601562 \times 10^{-1} = 0.6391601562$$

$$= 1309/2^{11} = 10100011101 \times 2^{-11}$$

$$= 1.0100011101 \times 2^{-1}$$

(kiểm tra số floating-point B này đã đúng chuẩn chưa:
 (-1+127) nằm trong phạm vi số 8 bits của phần mũ
 và phần fraction '0100011101' cũng không vượt quá 23 bits
 → Đúng chuẩn)

$$A + B = 1.0111000111 \times 2^4 + 1.0100011101 \times 2^{-1}$$

$$= 1.0111000111 \times 2^4 + 0.000010100011101 \times 2^4$$

$$= 1.011110111111101 \times 2^4$$

Phần fraction này chứa 15 bits, chưa vượt quá 23 bits nên
 $A + B = 1.011110111111101 \times 2^4$

Exercise 7:

The following table shows pairs of decimal numbers

	A	B
a.	5.66015625×10^0	8.59375×10^0
b.	6.18×10^2	5.796875×10^1

1. Calculate $A \times B$ by hand, assuming that we keep 11 bits of significand and 5 bits of the exponent. (Rounding rule: add 1 if the bits to the right of the desired point is larger or equal to $100_{(2)}$). Show all the steps.
2. Calculate $A \times B$ by hand, assuming A and B are stored in the IEEE-754 single precision format. Show all the steps.

Answer:

1. Đề bài yêu cầu tính $A \times B$ bằng tay (tức chạy từng bước) với giả sử số floating-point chỉ cho phép dùng 11 bits cho phần significand và 5 bits cho phần exponent a)

$$A = 5.66015625 \times 10^0 = 1.0110101001 \times 2^2$$

(kiểm tra số floating-point A đã đúng chuẩn cho phép chưa:
 → Đúng chuẩn)

$$B = 8.59375 \times 10^0 = 1.0001001100 \times 2^3$$

(kiểm tra số floating-point B này đã đúng chuẩn cho phép chưa:
 → Đúng chuẩn)

$$A \times B = (1.0110101001 \times 2^2) \times (1.0001001100 \times 2^3)$$

Exponent của $A \times B = 2 + 3 = 5$

Significand

```

      1.0110101001
    × 1.0001001100
    -----
      00000000000
      00000000000
     10110101001
     10110101001
    00000000000
    00000000000
   10110101001
   00000000000
  00000000000
  00000000000
 10110101001
1.10000101001000101100
  
```

Do phần significand chỉ được phép chứa 11 bits, nên significand của $A \times B$ phải được làm tròn, phần cắt bỏ là $1000101100_{(2)} > 100_{(2)}$ nên $1.10000101001000101100_{(2)} \approx 1.1000010101$

Vậy $A \times B = 1.1000010101 \times 2^5$

b)

	A	B
a.	5.66015625×10^0	8.59375×10^0
b.	6.18×10^2	5.796875×10^1

$$A = 6.18 \times 10^2 = 618 = 1001101010_{(2)} \\ = 1.001101010 \times 2^9$$

(kiểm tra số floating-point A đã đúng chuẩn cho phép chưa:

→ Đúng chuẩn)

$$B = 5.796875 \times 10^1 = 57.96875 = 1.1100111111 \times 2^5$$

(kiểm tra số floating-point B này đã đúng chuẩn cho phép chưa:

→ Đúng chuẩn)

$$A \times B = (1.001101010 \times 2^9) \times (1.1100111111 \times 2^5)$$

Exponent của $A \times B = 9 + 5 = 14$

Significand

```

      1.0011010100
    × 1.1100111111
    -----
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      00000000000
      00000000000
      10011010100
      10011010100
      10011010100
      100010111110000101100

```

Significand của $A \times B = 10.00101111110000101100 \rightarrow$ phải chuẩn hóa lại

$$\begin{aligned}
 A \times B &= 10.00101111110000101100 \times 2^{14} \\
 &= 1.0001011111 \mathbf{10000101100} \times 2^{15}
 \end{aligned}$$

Do phần significand chỉ được phép chứa 11 bits, nên significand của $A \times B$ phải được làm tròn, phần cắt bỏ là $\mathbf{10000101100}_{(2)} > 100_{(2)}$ nên $1.0001011111 \mathbf{10000101100}_{(2)} \approx 1.0001100000$
 Vậy $A \times B = 1.0001100000 \times 2^{15}$

2. Đề bài yêu cầu tính $A \times B$ bằng tay (tức chạy từng bước) với giả sử số floating-point dùng format IEEE độ chính xác đơn

$$\begin{aligned}
 A &= 5.66015625 \times 10^0 = 1.0110101001 \times 2^2 \\
 &\text{(kiểm tra số floating-point A đã đúng chuẩn cho phép chưa:} \\
 &\quad \rightarrow \text{Đúng chuẩn)}
 \end{aligned}$$

$$B = 8.59375 \times 10^0 = 1.0001001100 \times 2^3$$

(kiểm tra số floating-point B này đã đúng chuẩn cho phép chưa:
 \rightarrow Đúng chuẩn)

$$A \times B = (1.0110101001 \times 2^2) \times (1.0001001100 \times 2^3)$$

$$\text{Exponent của } A \times B = 2 + 3 = 5$$

Significand

```

      1.0110101001
    × 1.0001001100
    -----
      000000000000
      000000000000
      10110101001
      10110101001
      000000000000
      000000000000
      10110101001
      000000000000
      000000000000
      000000000000
      10110101001
    1.10000101001000101100

```

Do phần fraction được phép chứa 23 bits, nên fraction của A x B từ kết quả trên thỏa mãn, không cần làm tròn

Vậy $A \times B = 1.10000101001000101100 \times 2^5$

b)

	A	B
a.	5.66015625×10^0	8.59375×10^0
b.	6.18×10^2	5.796875×10^1

$$A = 6.18 \times 10^2 = 618 = 1001101010_{(2)} \\ = 1.001101010 \times 2^9$$

(kiểm tra số floating-point A đã đúng chuẩn cho phép chưa:

→ Đúng chuẩn)

$$B = 5.796875 \times 10^1 = 57.96875 = 1.1100111111 \times 2^5$$

(kiểm tra số floating-point B này đã đúng chuẩn cho phép chưa:

→ Đúng chuẩn)

$$A \times B = (1.001101010 \times 2^9) \times (1.1100111111 \times 2^5)$$

Exponent của A x B = $9 + 5 = 14$

Significand

```

      1.0011010100
×   1.1100111111
-----
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      10011010100
      00000000000
      00000000000
      10011010100
      10011010100
      10011010100
      100010111110000101100

```

Significand của $A \times B = 10.0010111110000101100 \rightarrow$ phải chuẩn hóa lại

$$\begin{aligned}
 A \times B &= 10.0010111110000101100 \times 2^{14} \\
 &= 1.00010111110000101100 \times 2^{15}
 \end{aligned}$$

Do phần fraction được phép chứa 23 bits, nên fraction của $A \times B$ trên hợp lệ

$$\text{Vậy } A \times B = 1.00010111110000101100 \times 2^{15}$$