

①

CONVERT THE DECIMAL NUMBER

75.125 TO 32 BIT IEEE FLOATING POINT FORMAT. EXPRESSING THE

ANSWER IN HEXADECIMAL FOR CONVENIENCE

① CONVERT TO BINARY

2	75		
2	37	R	1
2	18	R	1
2	9	R	0
2	4	R	1
2	2	R	0
2	1	R	0
	0	R	1
	?		

$$.125 \times 2 = \boxed{0}.25$$

$$.25 \times 2 = \boxed{0}.5$$

$$.5 \times 2 = \boxed{1}.0$$

0

1001011.00₂

(2)

② NORMALISE (EXPRESS IT IN SCIENTIFIC NOTATION)

1001011.001

MOVE . TO AFTER FIRST 1

1.001011001

HOW MANY PLACES DOES THE POINT HAVE TO BE MOVED

1.001011001 $\times 2^6$ FORWARD TO GET TO

WHERE IT WAS

ORIGINALLY? (6) PLACES

③ IDENTIFY SIGN EXPONENT AND MANTISSA

Sign 0 FOR + NO, 1 FOR NEG NO,

EXPONENT

MANTISSA

1.001011001

(BUT DROP LEADING 1

AND BINARY POINT TO

SAVE SPACE)

001011001

MANTISSA

0
6

③

TABLE — SOMETHING YOU NEED
TO KNOW

~~BINARY~~

HEXADECIMAL

BINARY

0

0000

1

0001

2

0010

3

0011

4

0100

5

0101

6

0110

7

0111

8

1000

9

1001

A

1010

B

1011

C

1100

D

1101

E

1110

F

1111

0

④

④ ADD 127 TO EXPONENT AND CONVERT TO BINARY

$$6 + 127 = 133$$

2	133		
2	66	R	1
2	33	R	0
2	16	R	1
2	8	R	0
2	4	R	0
2	2	R	0
2	1	R	0
	0	R	1

EXPONENT



10000101

(6)

EXERCISE

Convert the decimal number

123.0625_{10} to IEEE 32 bit

floating point format.

① CONVERT TO BINARY

$$2 \overline{) 123}$$

$$2 \overline{) 61} \text{ R } 1$$

$$2 \overline{) 30} \text{ R } 1$$

$$2 \overline{) 15} \text{ R } 0$$

$$2 \overline{) 7} \text{ R } 1$$

$$2 \overline{) 3} \text{ R } 1$$

$$2 \overline{) 1} \text{ R } 1$$

$$0 \overline{) 0} \text{ R } 1$$

$$.0625 \times 2 = \boxed{0}.125$$

$$.125 \times 2 = \boxed{0}.25$$

$$.25 \times 2 = \boxed{0}.5$$

$$.5 \times 2 = \boxed{1}.0$$

$$0 \overline{) 0}$$

1111011.0000_2

⑦

② NORMALIZE (EXPRESS IN SCIENTIFIC NOTATION)

1111011.0001

MOVE . TO AFTER THE FIRST
OCCURENCE OF A 1

1.1110110001 $\times 2^6$

(THE POINT HAS TO BE MOVED FORWARD 6 PLACES
TO GET TO WHERE IT WAS ORIGINALLY)

③ IDENTIFY:

SIGN 0 (POSITIVE)

EXPONENT 6

MANTISSA 1.1110110001

(DROP LEADING 1
AND BINARY POINT) \rightarrow 1110110001

⑧

④ ADD 127 TO EXPONENT AND
CONVERT TO BINARY

$$6 + 127 = 133$$

$$2 \overline{) 133}$$

$$2 \overline{) 66} \text{ R } 1$$

$$2 \overline{) 33} \text{ R } 0$$

$$2 \overline{) 16} \text{ R } 1$$

$$2 \overline{) 8} \text{ R } 0$$

$$2 \overline{) 4} \text{ R } 0$$

$$2 \overline{) 2} \text{ R } 0$$

$$2 \overline{) 1} \text{ R } 0$$

$$0 \text{ R } 1$$

EXPONENT

100000101

9

- ⑤ ASSEMBLE THE NO. AND CONVERT TO
HEXADECIMAL (HEX.) FOR CONVENIENCE

1 BIT	8 BITS	23 BITS
SIGN	EXPONENT	MANTISSA
0	10000101	11101100010000000000000

Express THE ABOVE NO. IN HEXADECIMAL

0 1 0 0 0 0 1 0 1 1 1 1 0 1 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
4 2 F 6 2 0 0 0

$42F62000_{16}$

EXERCISE

(10)

CONVERT -12345.625_{10} TO 32 BIT

IEEE FLOATING POINT FORMAT.

EXPRESS YOUR ANSWER IN HEXADECIMAL.

ANSWER: C640E680

EXERCISE

CONVERT 8193.125_{10} TO 32 BIT

IEEE FLOATING POINT FORMAT

① CONVERT TO BINARY

$100000000000000000001.001_2$

② NORMALIZE

$1.0000000000000000001001 \times 2^{13}$

(11)

③ IDENTIFY:

SIGN

0

EXPONENT

13

MANTISSA

1.0000000000000001001

DROP LEADING 1 AND .

0000000000000001001

↓
MANTISSA

④ ADD 127 TO EXPONENT AND CONVERT TO BINARY

$$13 + 127 = 140$$

$$\begin{array}{r} 2 \overline{) 140} \end{array}$$

$$\begin{array}{r} 2 \overline{) 70} \text{ R } 0 \end{array}$$

$$\begin{array}{r} 2 \overline{) 35} \text{ R } 0 \end{array}$$

$$\begin{array}{r} 2 \overline{) 17} \text{ R } 1 \end{array}$$

$$\begin{array}{r} 2 \overline{) 8} \text{ R } 1 \end{array}$$

$$\begin{array}{r} 2 \overline{) 4} \text{ R } 0 \end{array}$$

$$\begin{array}{r} 2 \overline{) 2} \text{ R } 0 \end{array}$$

$$\begin{array}{r} 2 \overline{) 1} \text{ R } 0 \end{array}$$

$$0 \text{ R } 1$$



EXPONENT



10001100

(12)

(5)

ASSEMBLE THE NO.

(AND CONVERT TO HEX.)

1 BIT Sign	8 BITS EXPONENT	23 BITS MANTISSA
0	10001100	00000000000000010010000000

01000110000000000000010010000000
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
4 6 0 0 0 4 8 0

46000480

(13)

CONVERT -8193.125 TO 32 BIT
IEEE FLOATING POINT FORMAT

SIGN	EXPONENT	MANTISSA
1	10001100	000000000000100100000000

↑

THIS HAS CHANGED FROM PREVIOUS EXERCISE

1100

↓

C6000480

EXERCISE

CONVERT 0.03125_{10} TO 32 BIT
IEEE FLOATING POINT FORMAT

(14)

① CONVERT TO BINARY

0.03125

↓
0

$$.03125 \times 2 = \boxed{0}.0625$$

$$.0625 \times 2 = \boxed{0}.125$$

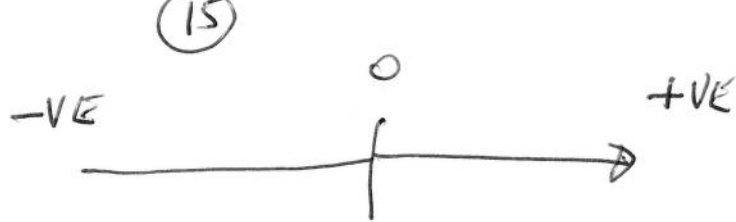
$$.125 \times 2 = \boxed{0}.25$$

$$.25 \times 2 = \boxed{0}.5$$

$$.5 \times 2 = \boxed{1}.0$$

0

0.00001₂



② NORMALIZE

0.00001

$$1.0 \times 2^{-5} \quad 5 \text{ PLACES TO THE RIGHT}$$

RIGHT \Rightarrow

NEGATIVE

EXPONENT

③ IDENTIFY SIGN EXPONENT AND MANTISSA

SIGN 0

EXPONENT -5

MANTISSA 1.0

DROP LEADING
1 AND .

0
A

MANTISSA

(16)

④ ADD 127 TO EXPONENT AND CONVERT TO BINARY

$$-5 + 127 = 122$$

$$2 \overline{) 122}$$

$$2 \overline{) 61} \text{ R } 0$$

$$2 \overline{) 30} \text{ R } 1$$

$$2 \overline{) 15} \text{ R } 0$$

$$2 \overline{) 7} \text{ R } 1$$

$$2 \overline{) 3} \text{ R } 1$$

$$2 \overline{) 1} \text{ R } 1$$

$$0 \text{ R } 1$$



7 BITS

1111010

01111010

↑ 8 BIT EXPONENT

0 ADDED HERE TO MAKE IT 8 BITS

⑦

⑤ ASSEMBLE NO. AND CONVERT TO HEX

S	E	M
0	01111010	00000000000000000000000000000000

00111101000000000000000000000000
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
3 D 0 0 0 0 0 0

3D000000