



Unit 45

Fixed and floating point notation

Objectives

On completion of this unit you should understand:

- 1.** Fractional binary notation.
- 2.** Fixed point binary numbers
- 3.** Floating point binary numbers.

Fractional binary numbers

When we considered the change from denary to binary numbers in an earlier unit, we did this by writing the headings above each column.

2^8	2^7	2^6	2^5	2^4	2^3	2^2	2^1	2^0
256	128	64	32	16	8	4	2	1

Notice that each number is double the previous one as you move to the left. 1, 2, 4, 8, 16.....and so on.

To convert the denary number 35 into a binary number, we need $32 + 2 + 1$ and so we obtain, 100011_2 .

This same system applies to numbers less than one. We can continue our headings to the right of the decimal point.

2^0	.	2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}
-------	---	----------	----------	----------	----------	----------	----------------

2^{-1} is the same as $1/2^1$ or $1/2$ or 0.5.

2^{-2} is the same as $1/2^2$ or $1/4$ or 0.25.

2^{-3} is the same as $1/2^3$ or $1/8$ or 0.125.

We can therefore consider our headings as follows, starting with one and moving to the right of the decimal point.

2^0	.	2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}
1	.	$1/2$	$1/4$	$1/8$	$1/16$	$1/32$	$1/64$
or 1	.	0.5	0.25	0.125	0.0625	0.03125	0.015625....

Notice that each number is half the previous one as you move to the right.

Study these examples.

Example 1

Convert the binary number 0.101010 to a denary number.

We set up our headings as follows with the number in position. We shall use the fraction headings this time.

2^0	.	2^{-1}	2^{-2}	2^{-3}	2^{-4}	2^{-5}	2^{-6}
1	.	$1/2$	$1/4$	$1/8$	$1/16$	$1/32$	$1/64$
0	.	1	0	1	0	1	0

The number 0.101010 is therefore, $1/2 + 1/8 + 1/32 = 21/32$ or 0.65625.

Example 2

Convert the binary number 11.0110 to a denary number.

We set up our headings as follows with the number in position. We shall use the decimal headings this time.

2^1	2^0	.	2^{-1}	2^{-2}	2^{-3}	2^{-4}
2	1	.	0.5	0.25	0.125	0.0625
1	1	.	0	1	1	0

The binary number 11.00101 is therefore,
 $2 + 1 + 0.25 + 0.125 = 3.375$ as a mixed decimal or $3\frac{3}{8}$ as a mixed fraction in the denary system.

Try this exercise.

Exercise A

Convert the following binary numbers to denary numbers giving your answers as,

- a) mixed fractions,
- b) mixed decimals.

1. 11.101
2. 1.10111
3. 10.01011
4. 101.1011010
5. 1101.11
6. 100.0101010
7. 100.1001
8. 11.11
9. 0.1101
10. 101.011

Check your answers with those at the end of the unit.

Converting decimals to binary numbers

Consider these examples.

Example 3

Convert 3.75_{10} to a binary number.

We set up our headings as follows with the number in position. We use the decimal headings.

2^1	2^0	.	2^{-1}	2^{-2}	2^{-3}	2^{-4}
2	1	.	0.5	0.25	0.125	0.0625

3_{10} is written as 11_2 . We now need to convert the decimal part of the number. 0.75 is made up of $0.5 + 0.25$. so $3.75_{10} = 11.11_2$.

Example 4

Convert 8.4375_{10} to a binary number.

We set up our headings as follows with the number in position. We use the decimal headings.

2^3	2^2	2^1	2^0	.	2^{-1}	2^{-2}	2^{-3}	2^{-4}
8	4	2	1	.	0.5	0.25	0.125	0.0625

8_{10} is written as 1000_2 . We now need to convert the decimal part of the number.

0.4375 is made up of $0.25 + 0.125 + 0.0625$ so $8_{10} = 1000.0111_2$.

It is not always possible to find an exact binary number to replace all decimals. For example if we had to find a binary number for 0.4374 , we would need to use $0.25 + 0.125$ and then 0.0624 . There is no binary digit for 0.0624 and so the nearest would be to use the binary digit for 0.0625 .

Try the following exercise.

Exercise B

Convert the following numbers to binary form.

- | | |
|-----------|-----------|
| 1. 2.0625 | 4. 3.25 |
| 2. 5.3125 | 5. 0.75 |
| 3. 0.5625 | 6. 0.8125 |

Check your answers with those at the end of the unit.

Fixed point binary numbers

Just as calculators have a fixed display of eight or ten numbers, computer registers for holding numbers are a fixed length, so compromises need to be made.

There are two major types of fixed representation. These are integer or fractional. We are going to consider an 8-bit register.

In this system, the first digit tells us whether the number is positive or negative.

A number starting with a 0 is positive.

A number starting with a 1 is negative.

The maximum positive number which can be stored using 8 bits must be 01111111. Remember that the 0 tells us that the number is positive. This number is +127.

The smallest positive number in an 8-bit register is therefore 00000001. This is +1.

Zero is considered to be neither positive nor negative.

The negative numbers in an 8-bit register range from -1 to -128. In binary form these are written as 11111111 and 10000000.

The following examples show how negative numbers are written in an 8-bit binary system, using the quick method shown in the Binary, Octal and Hexadecimal unit.

Study these examples.

Example 5

Using eight bits, and the quick method, find the binary number to represent -3_{10} .

Using eight bits, the number 3_{10} is converted to the binary number 00000011_2 .

Starting at the right hand side, rewrite the number up to and including the first 1.

00000011_2 write 1

For all the remaining digits, the 0's become 1 and the 1's become 0.

11111101 .

In this 8-bit system, -3_{10} is written 11111101_2 .

Example 6

Using the quick method find the binary number to represent -20_{10} in an 8-bit register.

The number 20_{10} , is written as 00010100_2 .

Starting at the right hand side, rewrite the number up to and including the first 1.

00010100_2 write **100**

For all the remaining digits, the 0's become 1 and the 1's become 0.

11101100.

-20_{10} is equal to the binary number 11101100_2 .

Try this exercise.

Exercise C

Using an 8-bit register, find the binary equivalent of each of the following denary numbers.

1. 27
2. 111
3. 8
4. 15
5. 39
6. -10
7. -18
8. -35
9. -101
10. -69

Check your answers with those at the end of the unit.

Floating point representation of binary numbers

Numbers can be split into two parts. The first part is called a **mantissa** and the second part is called an **exponent**.

Consider the standard form number 1.789×10^6 .

1.789 is the mantissa and 10^6 is the exponent.

A similar method can be employed for binary numbers. We shall consider a 16-bit register. This can be split into two parts. The first part carries the mantissa and the second part carries the exponent.

Consider these examples.

Example 7

A 16-bit register is split into two parts. Ten bits are to be used for the mantissa and six bits are to be used for the exponent. Convert the binary number, $0.001110011 | 000110$ into denary form.

0.001110011 is the mantissa and is using ten bits. This is a positive number because the first digit is 0.

000110 is the exponent and is using six bits. The exponent is also positive because it begins with a 0.

First consider the exponent. 000110 is equal to 6 in denary form.

This means that the point in the mantissa is to be moved six places to the right. 0001110.011

We can now convert this binary number to denary form.

$$0001110.011 = 8 + 4 + 2 + \frac{1}{4} + \frac{1}{8} = 14\frac{3}{8} \text{ or } 14.375.$$

Example 8

A 16-bit register is split into two parts. Twelve bits are to be used for the mantissa and four bits are to be used for the exponent. Convert the binary number, $0.01011100111 | 0111$ into denary form.

0.01011100111 is the mantissa and is using twelve bits. This is a positive number because the first digit is 0.

0111 is the exponent and is using four bits.

First consider the exponent. 0111 is equal to 7 in denary form.

This means that the point in the mantissa is to be moved seven places to the right. 00101110.0111

We can now convert this binary number to denary form.

$$00101110.0111$$

$$= 32 + 8 + 4 + 2 + 0.25 + 0.125 + 0.0625 = 46.4375 \text{ or } 46\frac{7}{16}.$$

Try this exercise.

Exercise D

Convert the following positive 16-bit binary numbers to denary form as mixed fractions or decimals.

1. $0.010111001 \mid 000011$
2. $0.001011110 \mid 001000$
3. $0.01011110111 \mid 0110$
4. $0.00011000111 \mid 0110$
5. $0.0001110011 \mid 00111$
6. $0.1110011 \mid 00000110$
7. $0.0011100101 \mid 00100$
8. $0.01011000000 \mid 0111$
9. $0.01010100111 \mid 0101$
10. $0.00011100000 \mid 0100$

Check your answers with those at the end of the unit.

Floating point representation of negative binary numbers

We shall now consider binary numbers in a 16-bit register, where the first digit is 1. This means that the number is negative.

Consider these examples.

Example 9

Convert $1.01011110 \mid 0000111$ to denary form.

Consider the exponent, 0000111 is equal to 7.

Now consider the mantissa, 1.01011110 .

Move the decimal point seven places to the right, 10101111.0

The first digit is 1 so the number is negative. Using the two's complement quick method from the Binary Octal and Hexadecimal System unit, we can convert this number. Starting at the right hand side, rewrite the number up to and including the first 1.

10101111.0 , write, 1.0 .

For all the remaining digits the 0's become 1 and the 1's become 0 as follows, 01010001.0 . This is equal to 81

The binary number $1.01011110 \mid 0000111$ is equal to -81 in denary form.

Example 10

Convert $1.0001011100 \mid 00011$ to denary form.

Consider the exponent.

00011 is equal to 3.

Now consider the mantissa, 1.0001011100 .

Move the decimal point three places to the right, 1000.1011100 .

The first digit is 1 so the number is negative. Using the two's complement quick method from the Binary Octal and Hexadecimal System unit, we can convert this number. Starting at the right hand side, rewrite the number up to and including the first 1.

1000.1011100 , write, 100 .

For all the remaining digits the 0's become 1 and the 1's become 0 as follows, 0111.0100100 .

This is equal to $7 + 0.25 + 0.03125 = 7.28125$ or $7\frac{9}{32}$.

The binary number,

$1.0001011100 \mid 00011$ is equal to the denary number -7.28125 .

Try this exercise.

Exercise E

Convert the following negative 16-bit binary numbers to denary form. Give your answers as mixed fractions or decimals.

1. $1.010111001 \mid 000011$
2. $1.110101111 \mid 000111$
3. $1.01011110111 \mid 0110$
4. $1.00011000111 \mid 0010$
5. $1.001110011 \mid 000111$
6. $1.1110111 \mid 00000110$
7. $1.0011100101 \mid 00100$
8. $1.01011000000 \mid 0111$
9. $1.01010100000 \mid 0101$
10. $1.00011100000 \mid 0100$

Check your answers with those at the end of the unit.

Floating point representation when the exponent is negative

We know that the exponent is negative if it begins with the number 1.

$0.001011100 | 100011$ is a positive number but the exponent is negative. The mantissa is positive because it begins with a 0. This means that the number itself is positive.

The exponent is negative, so we would need to move the decimal point on the mantissa to the left

Consider this example.

Example 11

Convert $0.001011100 | 111011$ to denary form.

Consider the exponent.

111011 is negative. It begins with a 1.

Using the quick method we can write the two's complement, 000101 . This is the denary number 5, so 111011 is equal to -5 in the denary system.

Now consider the mantissa, 0.001011100 .

Move the decimal point five places to the left, 0.00000001011100 .

This is equal to $\frac{1}{256} + \frac{1}{1024} + \frac{1}{2048} + \frac{1}{4096} + \dots = \frac{23}{4096}$.

or

$0.00390625 + 0.000976562 + 0.000488281 + 0.00024414 = 0.00561523$.

The binary number,

$0.001011100 | 111011$ is equal to the denary number $\frac{23}{4096}$ or 0.00561523.

Try this exercise.

Exercise F

Convert the following binary numbers to denary form. Give your answers as fractions or decimals.

1. $0.010101100 | 111101$
2. $0.0101010000 | 11010$
3. $0.01010110000 | 1101$
4. $0.011101100 | 11110$
5. $0.01110001100 | 1111$

Check your answers with those at the end of the unit.

Floating point representation when the exponent and the mantissa are negative

1.0001011100 | 10001 has a negative mantissa and a negative exponent. The number is negative. Because the exponent is negative we need to move the decimal point of the mantissa to the left.

Example 12

Convert 1.00101111 | 111101 to denary form.

Consider the exponent.

111101 is negative. It begins with a 1.

Using the quick method we can write the two's complement, 000011. This is the denary number 3, so 111101 is equal to -3 in the denary system.

Now consider the mantissa, 1.00101111. We can write the two's complement of this number using the quick method. 0.110100001. Move the decimal point three places to the left, 0.000110100001.

This is equal to $\frac{1}{16} + \frac{1}{32} + \frac{1}{128} + \frac{1}{4096} = \frac{417}{4096}$.

The binary number,

1.00101111 | 111101 is equal to the denary number $-\frac{417}{4096}$ or -0.1018066.

Try this exercise.

Exercise G

Convert the following binary numbers to denary form. Give your answers as fractions or decimals.

1. 1.110100000 | 111101
2. 1.0101010000 | 11011
3. 1.01010110000 | 1101
4. 1.0111101100 | 11110
5. 1.01110111100 | 1111

Check your answers with those at the end of the unit.

Answers

Exercise A

- a) $3^{5/8}$ b) 3.625
- a) $1^{23/32}$ b) 1.71875
- a) $2^{11/32}$ b) 2.34375
- a) $5^{45/64}$ b) 5.703125
- a) $13^{3/4}$ b) 13.75
- a) $4^{21/64}$ b) 4.328125
- a) $4^{9/16}$ b) 4.5625
- a) $3^{3/4}$ b) 3.75
- a) $1^{13/16}$ b) 0.8125
- a) $5^{3/8}$ b) 5.375

Exercise B

- 10.0000
- 101.0101
- 0.1001
- 11.01
- 0.11
- 0.1101

Exercise C

- 00011011
- 01101111
- 00001000
- 00001111
- 00100111
- 1110110
- 11101110
- 11011101
- 10011011
- 10111011

Exercise D

- $2^{57/64}$ or 2.890625
- 47
- $23^{23/32}$ or 23.71875
- $6^{7/32}$ or 6.21875
- $14^{3/8}$ or 14.375
- $57^{1/2}$ or 57.5
- $3^{37/64}$ or 3.578125
- 44
- $10^{39/64}$ or 10.609375
- $1^{3/4}$ or 1.75

Exercise E

- $-5^{7/64}$ or -5.109375
- $-20^{1/4}$ or -20.25
- $-40^{9/32}$ or -40.28125
- $-3^{313/512}$ or -3.6113281
- $-99^{1/4}$ or -99.25
- $-4^{1/2}$ or -4.5
- $-12^{27/64}$ or -12.421875
- 84
- $-21^{1/2}$ or -21.5
- $-14^{1/4}$ or -14.25

Exercise F

- $4^3/1024$ or 0.0419921
- $2^1/4096$ or 0.0051269
- $4^3/1024$ or 0.0419921
- $5^9/512$ or 0.1152343
- $2^{27}/1024$ or 0.2216796

Exercise G

- $-^3/128$ or -0.0234375
- $-^{43}/2048$ or -0.020996
- $-^{85}/1024$ or -0.0830078
- $-^{133}/1024$ or -0.1298828
- $-^{273}/1024$ or -0.2666015