

Data Representation

Number bases

Denary (or decimal) is base-10 and is the number system we are most familiar with. We have the columns of units, tens, hundreds, thousands and so on. Base-10 means that we have 10 possible values (0, 1, 2, 3, 4, 5, 6, 7, 8, 9) in each column.

Binary is base-2 and has 2 values, 0 and 1. It requires a greater number of digits in binary to represent a number than denary. This is how data and instructions are stored in a computer.

To calculate the maximum value for a given number of bits we use $2^n - 1$ where n is the number of bits. For example for 4 bits we have $2^4 - 1$ which is 15.

Bits	Max value binary	Max value denary
1	1 ₂	1 ₁₀
2	11 ₂	3 ₁₀
3	111 ₂	7 ₁₀
4	1111 ₂	15 ₁₀
5	11111 ₂	31 ₁₀
6	111111 ₂	63 ₁₀
7	1111111 ₂	127 ₁₀
8	11111111 ₂	255 ₁₀

Hexadecimal is base-16. To make up the 16 values we use the ten denary numbers in addition to 6 letters (A, B, C, D, E, F).

Denary	Hex.	Binary
0 ₁₀	0 ₁₆	0000 ₂
1 ₁₀	1 ₁₆	0001 ₂
2 ₁₀	2 ₁₆	0010 ₂
3 ₁₀	3 ₁₆	0011 ₂
4 ₁₀	4 ₁₆	0100 ₂
5 ₁₀	5 ₁₆	0101 ₂
6 ₁₀	6 ₁₆	0110 ₂
7 ₁₀	7 ₁₆	0111 ₂

Denary	Hex.	Binary
8 ₁₀	8 ₁₆	1000 ₂
9 ₁₀	9 ₁₆	1001 ₂
10 ₁₀	A ₁₆	1010 ₂
11 ₁₀	B ₁₆	1011 ₂
12 ₁₀	C ₁₆	1100 ₂
13 ₁₀	D ₁₆	1101 ₂
14 ₁₀	E ₁₆	1110 ₂
15 ₁₀	F ₁₆	1111 ₂

Hexadecimal is used a lot in computing because it much easier to read than binary. There are far fewer characters than binary. So hexadecimal is often used in place of binary as a shorthand to save space. For instance, the hexadecimal number 7BA3D456 (8 digits) is 01111011101000111101010001010110 (32 digits) in binary which is hard to read.

Hexadecimal is better than denary at representing binary because hexadecimal is based on powers of 2.

Converting between number bases

Denary to binary conversion

- Create a grid:

128	64	32	16	8	4	2	1

- Add a 1 to the corresponding cell if number contributes to target number and 0 to all the other cells

Worked example: convert 24₁₀ to binary.

128	64	32	16	8	4	2	1
0	0	0	1	1	0	0	0

$$16_{10} + 8_{10} = 24_{10}$$

The binary value is 11000₂ (we can ignore the preceding zeros)

Binary to denary conversion

Worked example: Convert 01011001₂ to denary

- Create the grid:

128	64	32	16	8	4	2	1
0	1	0	1	1	0	0	1

- Add up the cells that have a corresponding value of 1:
 $64 + 16_{10} + 8_{10} + 1 = \underline{89_{10}}$

Hexadecimal to denary conversion

- Convert the two hex values separately to denary value
- Multiply the first value by 16
- Add the second value

Worked example: Covert A3₁₆ to denary

$$A_{16} = 10_{10}$$

$$3_{16} = 3_{10}$$

$$(10_{10} \times 16_{10}) + 3_{10} = \underline{163_{10}}$$

Denary to hexadecimal conversion

- Integer divide the denary number by 16
- Take the modulus 16 of the denary number
- Convert the two numbers to the corresponding hex values.

Worked example: Convert 189₁₀ to hex

$$189_{10} / 16_{10} = 11_{10} \text{ remainder } 15_{10}$$

$$11_{10} = B_{16}$$

$$15_{10} = F_{16}$$

$$189_{10} = \underline{BF_{16}}$$

Hexadecimal to binary conversion

- Find the corresponding 4-bit binary number for the two numbers
- Concatenate the two binary values to give the final binary value

Example: convert C3₁₆ to binary

$$C_{16} = 12_{10} = 1100_2$$

$$3_{16} = 3_{10} = 0011_2$$

$$\underline{11000011_2}$$

Binary to hexadecimal conversion

- Split the binary number into groups of 4 bits: 1110₂ 1010₂
- Find the corresponding Hex value for each of the 4-bit groups

Worked example: Convert 11101010₂ to hexadecimal

$$1110_2 \mid 1010_2$$

$$1110_2 = 14_{10} = E_{16}$$

$$1010_2 = 10_{10} = A_{16}$$

$$\underline{EA_{16}}$$

Units of Information

Unit	Symbol	Number of bytes
Kilobyte	KB	10 ³ (1000)
Megabyte	MB	10 ⁶ (1 million)
Gigabyte	GB	10 ⁹ (1 billion)
Terabyte	TB	10 ¹² (1 trillion)

A bit is the fundamental unit of binary numbers. A bit is a binary digit that can be either 0 or 1.

$$1 \text{ byte} = 8 \text{ bits}$$

$$1 \text{ nibble} = 4 \text{ bits}$$

Character Encoding

Character coding schemes allows text to be represented in the computer. One such coding scheme is **ASCII**. ASCII uses 7 bits to represent each character which means that a total of 128 characters can be represented.

Lower case letters	26
Upper case letters	26
Numbers	10
Symbols (e.g. comma, colon)	33
Control characters	33

ASCII encoded values for some characters

A	1000001 ₂	65 ₁₀
B	1000010 ₂	66 ₁₀
a	1100001 ₂	97 ₁₀
b	1100010 ₂	98 ₁₀
“0”	0110000 ₂	48 ₁₀
“1”	0110001 ₂	49 ₁₀

- ASCII has a limited character set (7 bits, 128 characters), but **Unicode** has 16 bits and allows many more (65K) characters.
- Unicode provides a unique character for different languages and different platforms.
- It allows us to represent different alphabets for instance Greek, Mandarin, Japanese, Emojis etc.
- Unicode and ASCII are the same up to 127.

Binary addition

Binary addition rules

$$0_2 + 0_2 = 0_2$$

$$0_2 + 1_2 = 1_2$$

$$1_2 + 0_2 = 1_2$$

$$1_2 + 1_2 = 10_2 \text{ (carry 1)}$$

$$1_2 + 1_2 + 1_2 = 11_2 \text{ (carry 1)}$$

Example

$$\begin{array}{r} 1\ 0\ 1\ 0\ 1\ 0\ 0\ 1_2 \\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1_2 \\ + \quad 0\ 0\ 0\ 1\ 0\ 1\ 0\ 1_2 \\ \hline 1\ 1\ 0\ 0\ 0\ 1\ 1\ 1_2 \\ \text{carry } 1\ 1\ 1 \quad 1 \end{array}$$

Binary Shift

The binary shift operator is used to perform multiplication and division of numbers by powers of 2

<i>multiply/divide</i>	x 16	x 8	x 4	x 2	/ 2	/ 4	/ 8
<i>shift</i>	<<4	<<3	<<2	<<1	>>1	>>2	>>3

Example: Apply shift operator to 1101₂ (13₁₀)

Shift	Result	denary
<<1	11010 ₂	13 ₁₀ x 2 ₁₀ = 26 ₁₀
<<2	110100 ₂	13 ₁₀ x 4 ₁₀ = 52 ₁₀
>>1	110	13 ₁₀ // 2 ₁₀ = 6 ₁₀

Note that odd numbers are rounded down to the nearest integer when the right shift operator is applied.

Sound

Sample - Measure of the analogue signal at a given point in time

Sample rate - number of samples taken per second and is measured in Hertz.

Sample resolution - number of bits used to represent each sample

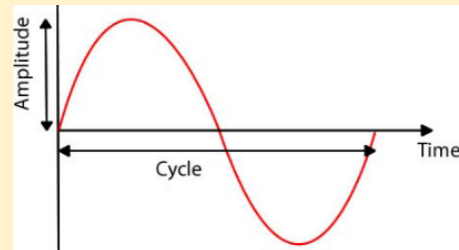
The size of sound files can be calculated using:

size of file = length (seconds) x sample rate x sampling resolution

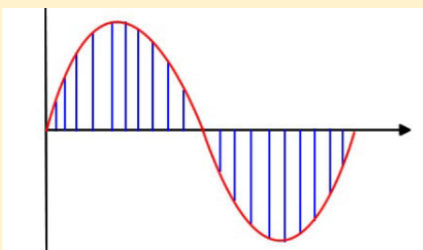
For sound to be stored digitally on a computer it needs to be converted from its continuous analogue form into a discrete binary values. The steps are:

1. Microphone detects the sound wave and converts it into an electrical (analogue) signal
2. The analogue signal is sampled at regular intervals
3. The samples are approximated to the nearest integer (quantised)
4. Each integer is encoded in binary with a fixed number of bits

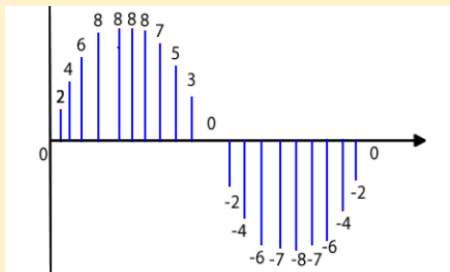
Original analogue signal



Sample signal at regular intervals



Integer values give to each sample



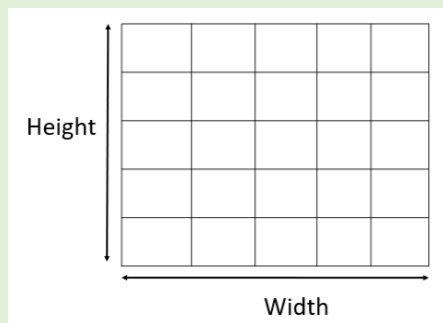
Encode as binary

0 2 4 6 8 8 8 8 7 5 3 0 ->
00000 00010 00100 01000
01000 01000 01000 00111
00101 00011 ...

Images

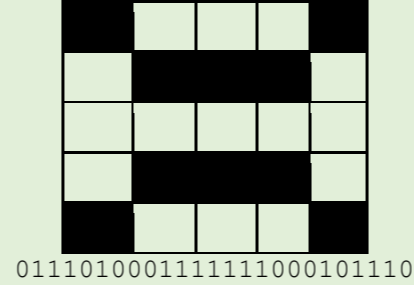
Bitmap images are made up from tiny dots called **pixels**. Each pixel will have a colour associated with it. An image can then be constructed from many of pixels which will have different colours arranged in rows and columns.

Total number of pixels in image = width in pixels x height in pixels



Colour depth is the number of bits used to represent each pixel in an image. If we have a black and white image it has two colours. Each pixel can be represented by a single pixel because a bit value of 0 is black and 1 is white.

Image and corresponding binary encoding



To represent more colours we can use more bits. For instance if we have 2-bits per pixel we can represent 4 colours because we know have 4 binary code combinations (00, 01, 10 11) where each code represents a different colour

Pixilation occurs when the image is overstretched. In these situations, the image loses quality and has a blocky and blurred appearance. This arises when the image is presented at too large a size and there are not enough pixels to reproduce the details in the image at this larger size.

Calculating the size of a bitmap image

File size in bits = width in pixels x height in pixels x colour depth

File size in bytes = width in pixels x height in pixels x colour depth / 8

Data Compression

The purpose of data compression is to make the files smaller which means that:

- Less time / less bandwidth to transfer data
- Take up less space on the disk

Given that there are 7 bits per ASCII character, the uncompressed size of an ASCII phrase is:

size = number of characters (including spaces) x 7

Run Length Encoding (RLE) is a compression method where sequences of the same values are stored in pairs of the value and the number of those values. For instance, the sequence:

0 0 0 1 1 0 1 1 1 1 0 1 1 1 1
would be represented as:
3 0 2 1 1 0 4 1 1 0 4 1

Huffman coding is a form of compression that allows us to use fewer bits for higher frequency data. More common letters are represented using fewer bits than less common letters. For instance, “a” and “e”, which occur in many words would be represented with fewer bit than “z” which occurs rarely. This allows for much more effective compression than RLE.

The steps involved in Huffman encoding as are follows:

1. Do frequency table
2. Order table
3. Create the tree
4. Add 1, 0 to the branches
5. Encode letters
6. Encode message

Worked Example: How much smaller is the phrase henry horse encoded using Huffman encoding compared with its uncompressed size.

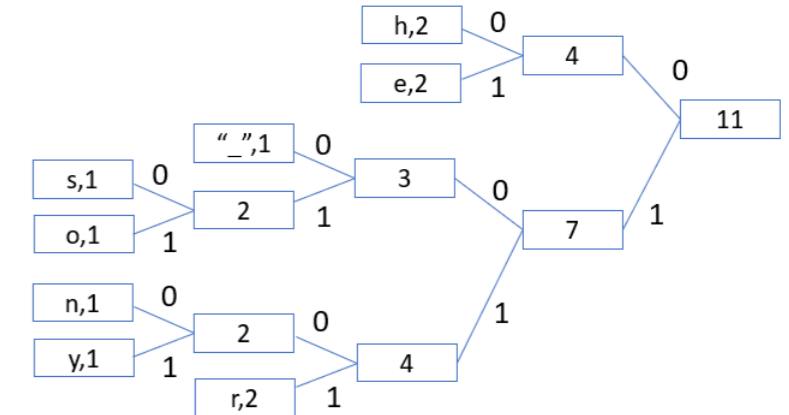
Calculate the uncompressed size

In the phrase *henry horse* there are 11 characters (including the space). Therefore the uncompressed size is $11 \times 7 = 77$ bits

Generate ordered frequency table (steps 1 and 2)

letter	frequency
e	2
h	2
r	2
<space>	1
o	1
s	1
y	1
n	1

Create the tree and add 1 and 0 to branches (steps 3 and 4)



Encode letters

Letter	encoding
e	01
h	00
r	111
<space>	100
o	1011
s	1000
n	1100
y	1101

Encode message

00 01 1100 111 1101 100 00 1011 111 1000 01 = 33 bits

Therefore by using compression we have reduced the size from 77 bits to 33 bits a saving of 44 bits.