# Information Technology for Digital Humanities
## Lecture 4

Mario Verdicchio

Università degli Studi di Bergamo

Academic Year 2023-2024

# Lecture 4 (October 4 2023)

- Fundamental concepts: technology
  - Previous lecture: hardware
  - Now: how the electronic signals the hardware works with are turned into something meaningful for the human users

# Information Technology for **Digital** Humanities
## Lecture 4

Mario Verdicchio

Università degli Studi di Bergamo

Academic Year 2023-2024

# Tools for **data** processing

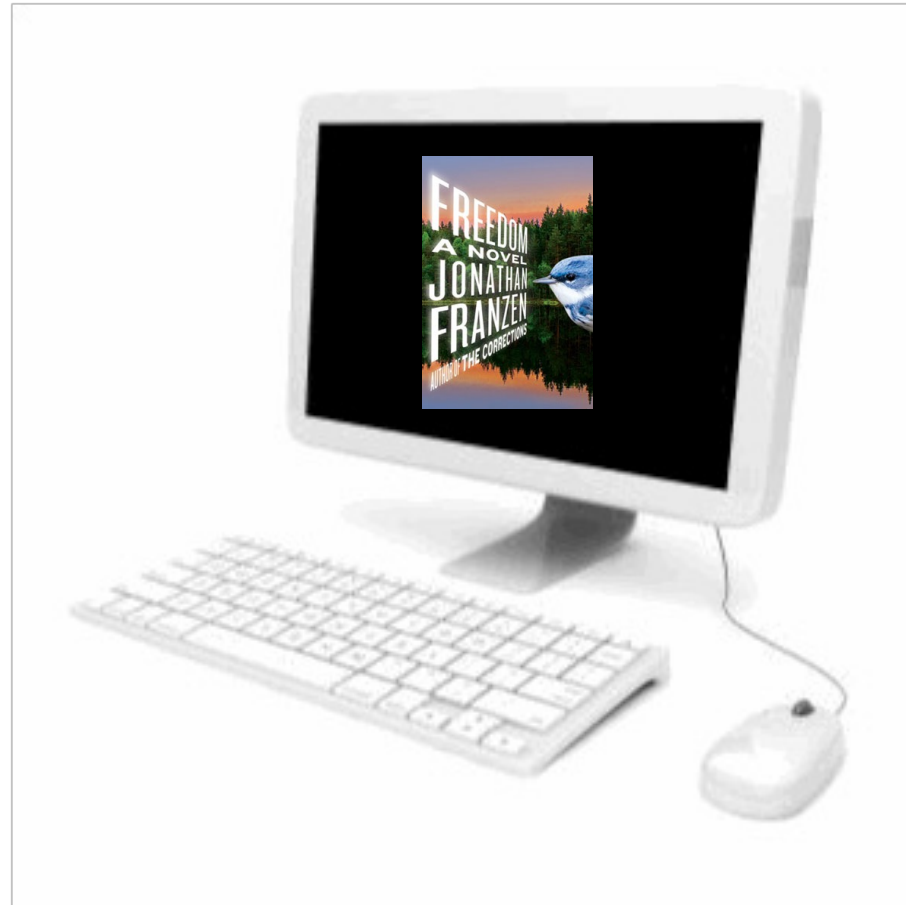This lesson focuses on what data/symbols/signs can be processed using computers.

**symbol**

# 8

**symbol**

In computer science, the term "symbol" has a much weaker meaning than in the humanities: by "symbol" we mean any sign, such as a figure, without any associated meaning, as happens instead for the lion, to which we associate by symbolism the concept of "courage".

Having simple numbers on a screen seems a bit limiting compared to what we usually see happening on the computers around us.

Authors, for example, write entire novels on their computer.

And, certainly, at least these three activities take place on the computers of the students of the course: social networks, listening to music, and watching films.

Let's think.
Premise 1: computers are nothing more than automatic executors of operations on symbols.
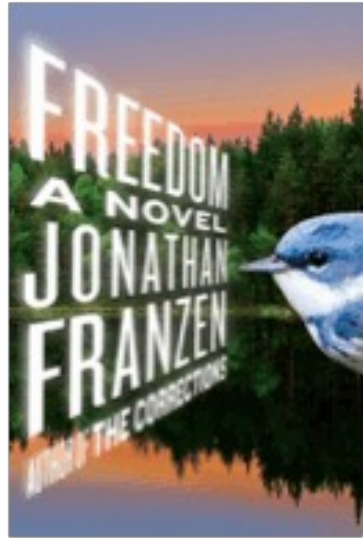Premise 2: with computers we write novels, surf social networks, listen to music, watch films.

**8**

Conclusion:
writing novels, surfing social networks, listening to music, watching films are operations on symbols.

The reasoning is flawless, but the conclusion may leave you perplexed.

To better understand the conclusion of the previous reasoning, it is necessary to clarify the link between all these activities and the symbols processed by a computer.

**8**

# encoding

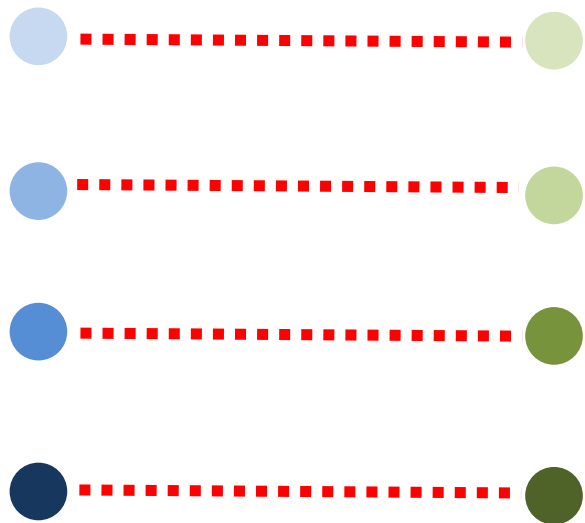This link is given by an encoding.

# ENCODING

## is

# EVERYTHING*

**\*in digital technology**

**ENCODING [ɪnˈkəʊdɪŋ]:** biunivocal correspondence between a set of entities of any kind and a set of natural numbers.
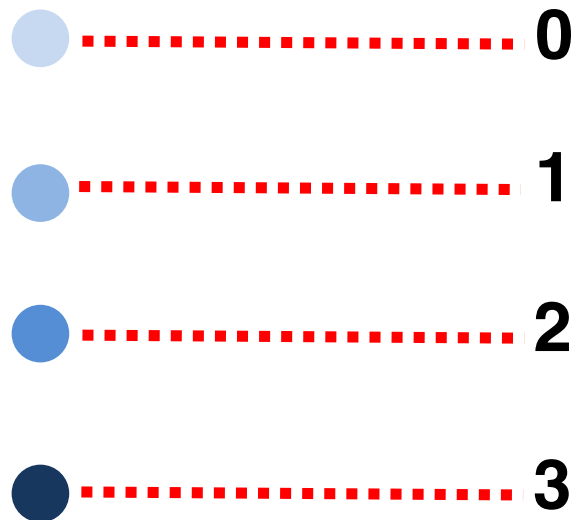
# biunivocal correspondence

**between a set of entities of any kind**

**and a set of natural numbers.**

# ENCODING

- 0
- 1
- 2
- 3

**\*from a conceptual perspective**
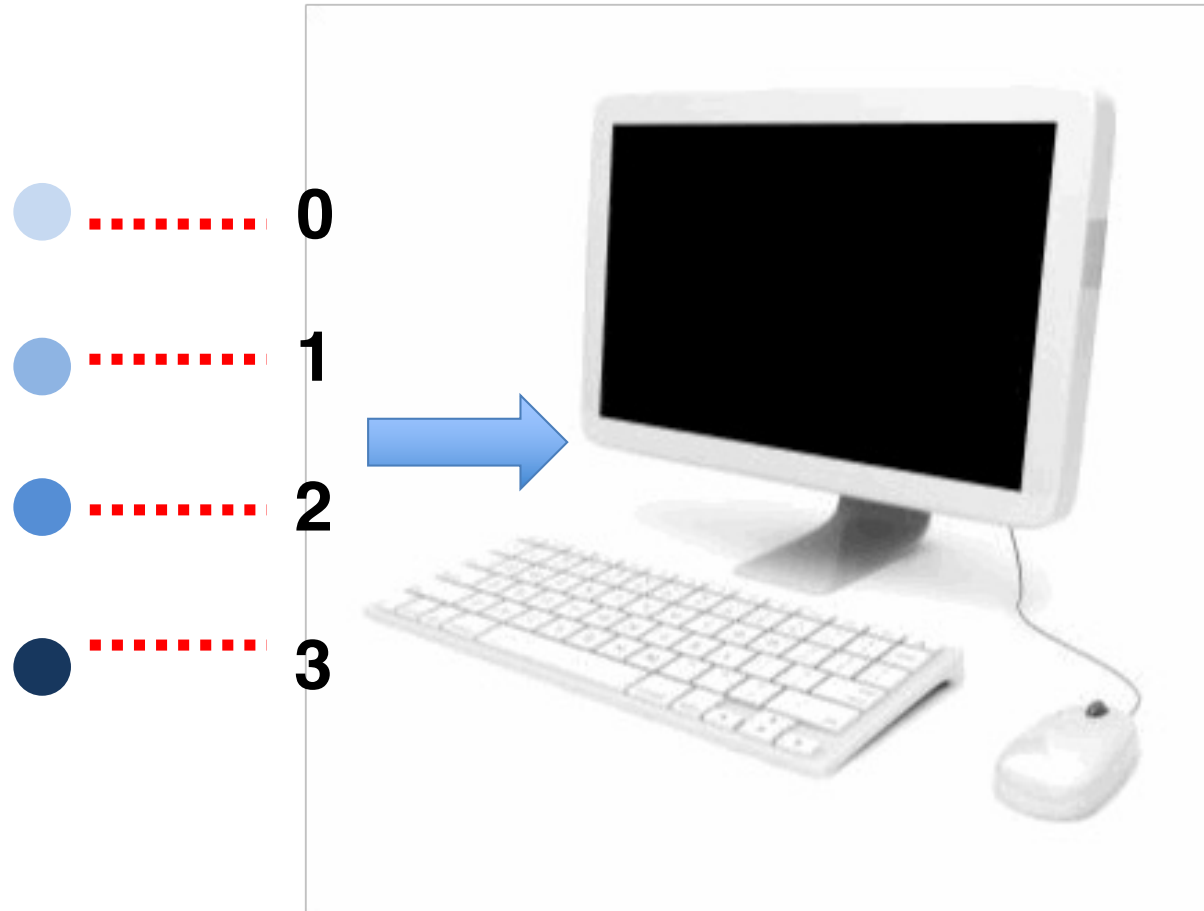
**A computer.**
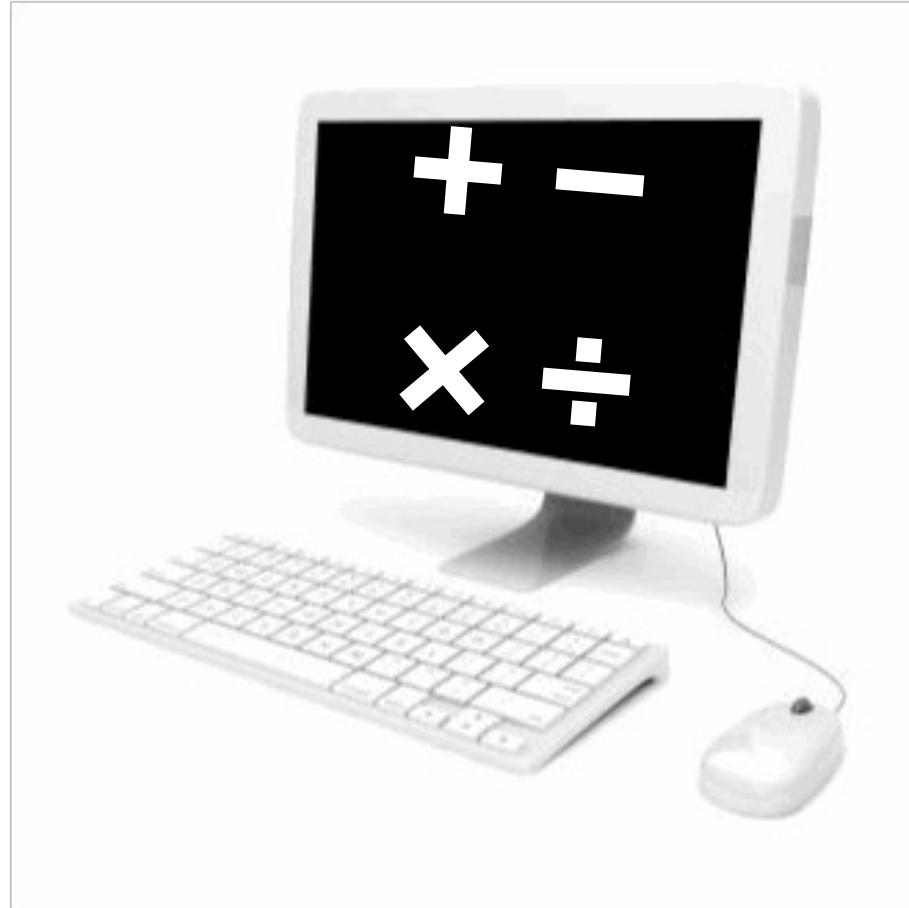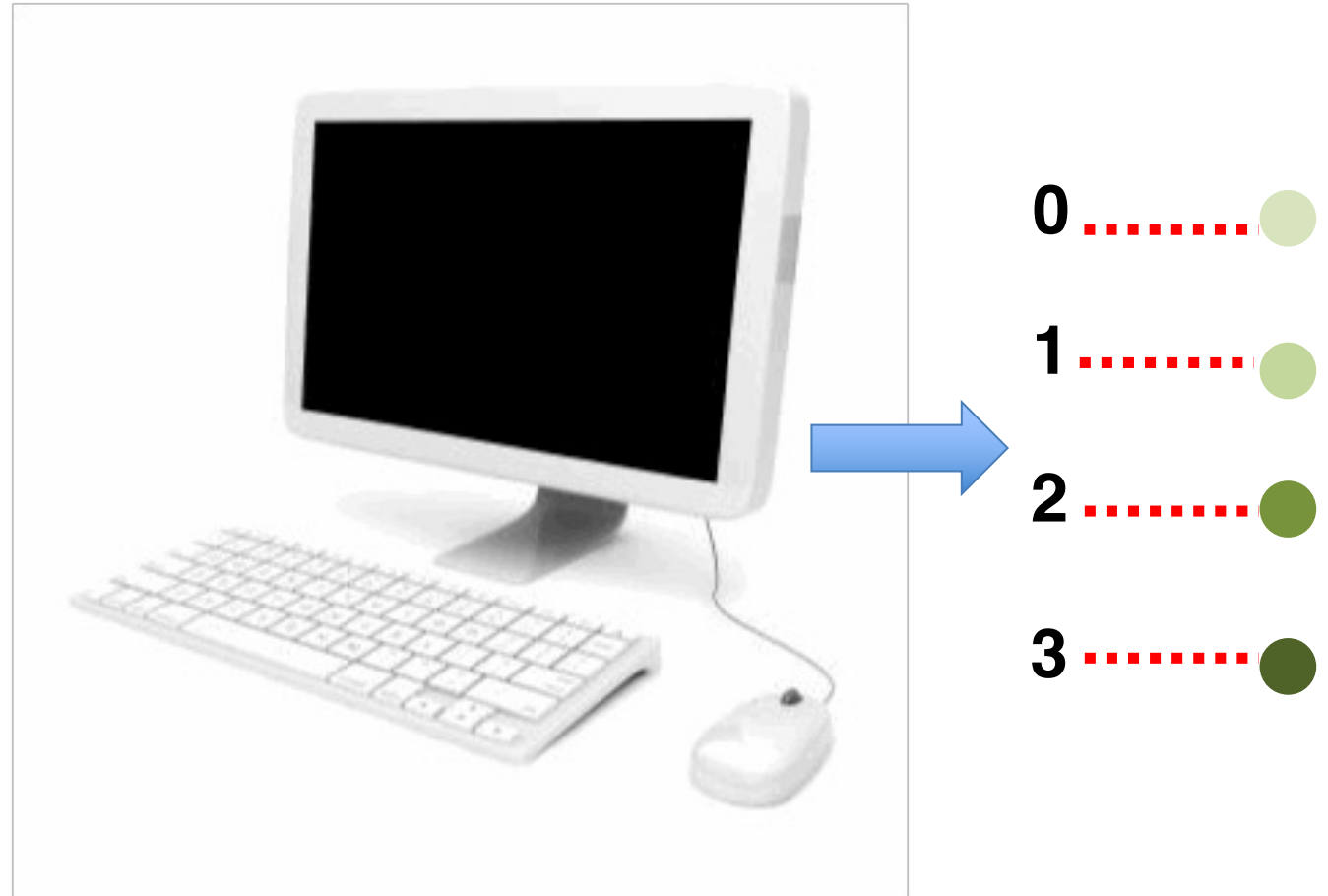
**It only works with numbers.**

**It cannot work with anything else.**

**Input needs to be ENCODED.**

**The computer works.**

**Output needs to be DECODED.**

**COMPUTER SCIENCE**

Let's see how we can create an encoding for texts.

8

The encoding of texts is based on the encoding of the characters that compose them: if I match each letter and each punctuation mark with a number, I obtain a one-to-one correspondence between letters and numbers. For this correspondence to be useful, it must be known and shared by all those who want to use a computer to exchange the texts thus encoded.
This table shows a well-known encoding: UTF-8

| SP | ! | " | # | $ | % | & | ' | ( | ) | * | + | , | - | . | / |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 0020 | 0021 | 0022 | 0023 | 0024 | 0025 | 0026 | 0027 | 0028 | 0029 | 002A | 002B | 002C | 002D | 002E | 002F |
| 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | : | ; | < | = | > | ? |
| 0030 | 0031 | 0032 | 0033 | 0034 | 0035 | 0036 | 0037 | 0038 | 0039 | 003A | 003B | 003C | 003D | 003E | 003F |
| 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 |
| @ | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 0040 | 0041 | 0042 | 0043 | 0044 | 0045 | 0046 | 0047 | 0048 | 0049 | 004A | 004B | 004C | 004D | 004E | 004F |
| 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 |
| P | Q | R | S | T | U | V | W | X | Y | Z | [ | \ | ] | ^ | _ |
| 0050 | 0051 | 0052 | 0053 | 0054 | 0055 | 0056 | 0057 | 0058 | 0059 | 005A | 005B | 005C | 005D | 005E | 005F |
| 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 |
| ` | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o |
| 0060 | 0061 | 0062 | 0063 | 0064 | 0065 | 0066 | 0067 | 0068 | 0069 | 006A | 006B | 006C | 006D | 006E | 006F |
| 96 | 97 | 98 | 99 | 100 | 101 | 102 | 103 | 104 | 105 | 106 | 107 | 108 | 109 | 110 | 111 |
| p | q | r | s | t | u | v | w | x | y | z | { | \| | } | ~ | DEL |
| 0070 | 0071 | 0072 | 0073 | 0074 | 0075 | 0076 | 0077 | 0078 | 0079 | 007A | 007B | 007C | 007D | 007E | 007F |
| 112 | 113 | 114 | 115 | 116 | 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 | 127 |

# UTF-8

# Universal Character Set Transformation Format – 8bit

The news about Walter Berglund wasn't picked up locally—he and Patty had moved away to Washington two years earlier and meant nothing to St. Paul now—but the urban gentry of Ramsey Hill were not so loyal to their city as not to read the *New York Times*. According to a long and very unflattering story in the *Times*, Walter had made quite a mess of his professional life out there in the nation's capital. His old neighbors had some difficulty reconciling the quotes about him in the *Times* ("arrogant," "high-handed," "ethically compromised") with the generous, smiling, red-faced 3M employee they remembered pedaling his commuter bicycle up Summit Avenue in February snow; it seemed strange that Walter, who was greener than Greenpeace and whose own roots were rural, should be in trouble now for conniving with the coal industry and mistreating country people. Then again, there had always been something not quite right about the Berglunds.

In order to be processed by a computer, the words of Jonathan Franzen's novel "Freedom" must be encoded in the form of a sequence of numbers. Each number corresponds uniquely to a character in the text.

10892735045634534765384509875309485703984562876503941753894573458934905834095844782629620423786238643782617816237276213048762763478012364023746283746763764376420870182873634756580565658276027863508271650821736582713650278135608217365657380278356238746021560609846246574568038047569345830948563094563074560384756038476501837465087314650138746507834560384756037486578346574658734506183745601837465018376456758403876573480187364571088573465783104587134653178451103874650

Let's see how social networks like Facebook can also be encoded.

8

What we see on our computer of a social network is nothing more than a set of texts and images.

We've already dealt with texts, so now let's focus our attention on images.

Looking at a very detailed photo of sweets from Bergamo, it is difficult to think that there could be a correspondence with the natural numbers.

By looking at a low-resolution image of Mario, you can instead understand how a photo can be encoded.

**pixel:**
**picture**
**element**

In the image of Mario, you can recognize basic square elements, characterized by the positions they occupy, and by their color: they are pixels.

**pixel (1/12000000)**

Image encoding is based on the pixels that compose it. In the same way as text is considered as a sequence of characters to be encoded…



…so we consider a photo a sequence of pixels. If we can encode one pixel, we can encode an entire image.

0  1  2

**encoding of the pixel's position**

We said that the first characteristic of a pixel is its position. If we imagine inserting the image in a Cartesian plane, we are able to identify the position of the pixel via its numerical coordinates

what about its color?

0   1   2

If we imagine organizing colors in a table, in the same way as pixels we are able to associate a numerical value with each color, corresponding to the coordinates of the color in the table.

The pixels that make up Mario, therefore, can be coded based on their position and their color, allowing the image, a sequence of pixels, to be expressed in the form of a sequence of numbers.

56787287298309394874897498749874 8467101
18762786287678354635413767612912 3621352
41365124387361498350238574675477 9090939
42384729384728347823749238472389 4728347
29837428364716515515625465463546 354376
47365847598475982475767151101982 4928493
84934028391809174037547478347473 6478364
73463784394010913481340946307456 0384756
03847650183746508731465013874650 7834560
38475603748657834657465873450618 3745601
83746501837645675840387657348018 7364571
08857346578310458713465317845110 3874650

We proceed in a similar way for the 12-million-pixel photo. It is no coincidence that more detailed photos take up more space in a computer memory: there are simply more encoded pixels.

10892735045634534765384509875309485703984562876503941753894573458934905834095844782629656734637463746334934898981010102949384837487348739487348176176473948739847394873847189374918374837483743814718397491384791834781347139847183947318478134718347387483478347384738748738473483748734873847384738478347834738478347834791837498137483483138746507834560384756037486578346574658734506183745601837465018376456758403876573480187364571088573465783104587134653178451103874650

0  1  2

As you will have noticed, in the table that allows color coding, not all shades are present. Just take two adjacent squares and imagine taking the intermediate color: it is not clearly present in the table, so it is not an encoded color.

This means that only a small part of the colors in the spectrum can be encoded. Numerous nuances (infinite, to be precise) remain left out of the coding. On the other hand, it is unthinkable to have a table of infinite dimensions to accommodate the infinite shades of colors.

This is the problem of the conversion from physical phenomena to "digital" (i.e. made of digits) descriptions: the encoding of physical phenomena always involves a loss of information. Drawing a comparison with mathematics, physical phenomena are characterized by infinitesimal nuances similar to real numbers, while encodings are correspondences with natural numbers, which, despite being infinite, are still much smaller than real numbers and cannot act as references to the aforementioned nuances.

# Physical [Analog] vs. Digital

It is like when people say that looking at a photo can never equal seeing the real thing. In reality, three-dimensionality aside, current color encoding can capture millions of shades, many of which are indistinguishable to the human eye. Therefore, albeit with losses, encoding allows us to have more than good approximations of reality.

We mentioned before that the detailed photo of the sweets results in a much longer encoded sequence of numbers than the very simple drawing of Mario. The length of the encodings naturally depend on what is encoded: the greater the information to be encoded, the longer the result, and in fact the photo of the sweets (with all the details and nuances given by the 12 million pixels) certainly takes up more space in the computer than the Mario drawing (which, in fact, is a sequence of 16 x 12 pixels).

# Compression

There are, however, different ways of describing the sequence of encoded pixels, and, if we choose wisely, we can obtain a shorter description of the same sequence, so that it takes up less space on the computer, and also saves time in transmitting it from one computer to another. We call "compression" a way to obtain sequences of symbols that describe a particular piece of information in a shorter way than the simple list of encodings of each element that composes it.

white pixel, green pixel, green pixel, yellow pixel, yellow pixel, yellow pixel, yellow pixel, green pixel, green pixel, green pixel, green pixel, white pixel

→

1 white pixel,
2 green pixels,
4 yellow pixels,
4 green pixels,
1 white pixel

The description of the Mario image, for example, can be compressed by describing its pixels as shown in the second red box.

Any type of information that can be encoded can be compressed. The compression ratio, i.e. how much we can summarize the description, depends on the content of the information. (Imagine having to describe an all-white 16 x 12 rectangle. How would you do it?)

Let's now see how encoding allows you to have music on your computer.

8

Sounds, and therefore also music, are due to vibrations of the air (or of any medium*) which propagate in the form of waves and which, hitting our eardrums, give rise to what we perceive with our ears. Sound waves are characterized by an amplitude (which determines the loudness of the sound: the wider the wave, the louder the sound), a frequency (which determines the pitch of the note: high notes correspond to high frequencies), and by a shape (which determines the timbre of the sound: my voice, your voice, piano, violin, cymbals, a doorbell, …)

*sounds do not propagate in a vacuum

# Sampling

The encoding of sounds is based on the encoding of the waves that produce them, which is in turn based on a procedure called sampling.

Sampling consists in considering the wave that constitutes the sound only in certain moments of time. Imagine describing the sound wave (with its amplitude, frequency, and shape) in a Cartesian plane, and considering only certain points of this curve.



Measurement of amplitude

Sound wave

Time

1111
1110
1101
1100
1011
1010
1001
1000
0111
0110
0101
0100
0011
0010
0001
0000

Each measurement is assigned a number (byte) according to its amplitude. The end result is a file comprising a string of bytes, eg ...
1001 1110 0001 1010 0111 0100 1111 1101 etc

These points correspond, in the Cartesian plane, to precise coordinates, whose numerical value is used as coding of the "sampled" sound. The encoded sequence of the samples constitutes the encoding of the entire sound data represented by the wave.

Figure 20 - Effect of Increased Resolution and Sampling Rates

Low Resolution and Sampling Rate

Increased Resolution

Increased Resolution and Sampling Rate

**sampling frequency**

There is no universal criterion for establishing the distance between one sample and the next (or its inverse, known as "sampling frequency"). It is easy to imagine that samples that are closer to each other (higher sampling frequency) correspond to a longer encoding, and also a reconstruction that is more faithful to the original wave. Even in the case of sound and music encodings, there are compression methods to synthesize their description. The famous MP3 files are called this way because they take their name from a specific sound encoding technique with compression. Even in this context there are people who say that the sound of an MP3 on a computer will never equal the quality of a live concert: it's all a question of approximations.

Encoding a sound is used to make it processable by a computer (for example, to allow the transfer of a song from the iTunes online store to our computer). However, even in the digital age, we continue to listen with our ears, and to make our listening possible, we need sound waves that propagate in the air. A reconversion from numerical coding to sound waves is therefore necessary: that is, we need speakers which, controlled by the electrical signals produced by the computer according to the numbers contained in the coding of the song, make membranes vibrate which produce waves that we perceive as sounds and music.



**again, physical vs digital**

Finally, let's see how to proceed with video encoding.

**8**

Actually, having image and sound encoding techniques available, it is easy to imagine that they can be combined to create video encoding.



=

Additional techniques are needed to take into account the synchronization between images and sounds, and compression techniques based on the idea of not describing all the pixels of each image, but only describing the initial one and then focusing on the differences between an image and the next (better compression with very similar frames, worse compression with scene changes).

 +  + ...

3.14159265358979323846264338327950288419716939937510
5820974944592307816406286208998628034825342117067982
1480865132823066470938446095505822317253594081284811
1745028410270193852110555964462294895493038196442881
0975665933446128475648233786783165271201909145648566
9234603486104543266482133936072602491412737245870066
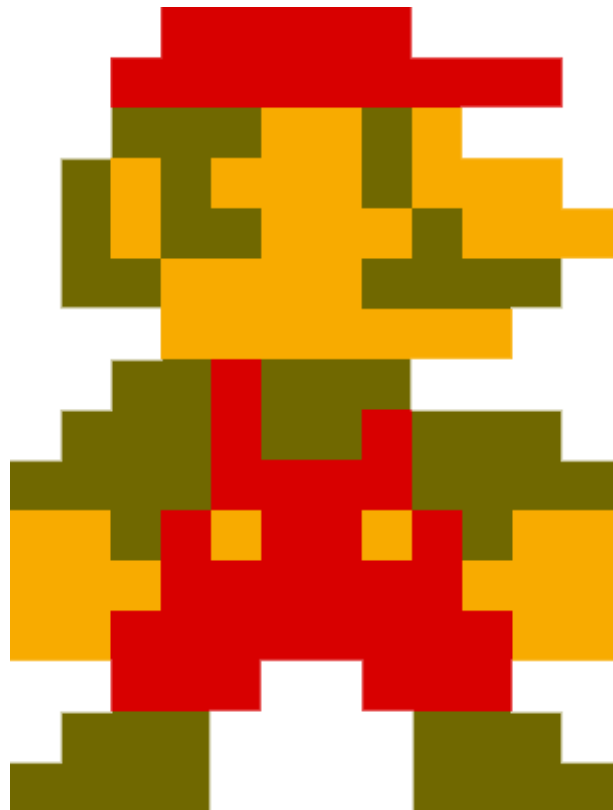0631558817488152092096282925409171536436789259036001
1330530548820466521384146951941511609433057270365759
5919530921861173819326117931051185480744623799627495
6735188575272489122793818301194912983367336244065664
3086021394946395224737190702179860943702770539217176
2931767523846748184676694051320005681271452635608277
7857713427577896091736371787214684409012249534301465
4958537105079227968925892354201995611212902196086403
4418159813629774771309960518707211349999998372978049
9510597317328160963185950244594553469083026425223082
5334468503526193118817101000313783875288658753320838
1420617177669147303598253490428755468731159562863882
3537875937519577818577805321712268066130019278766111
9590921642019899....

As already mentioned, the approximation that occurs when passing from physical phenomena such as colors and sound waves to numerical encoding suitable for a computer is reminiscent of the relationship between real numbers and natural numbers in mathematics. Perhaps the use of natural numbers to talk about the concept of encoding may seem limiting and may seem to invite the involvement of rational numbers, but the underlying argument does not change: even with rational numbers, the encodings are unable to faithfully reproduce a physical phenomenon, but they bring with them inevitable approximations. The number in the previous slide, for example, represents only a part of the decimal expansion of the number $\pi$. We can approximate it as we like, but part of the description of $\pi$ will still be left out and cannot be included, since no calculator is capable of containing an infinite sequence of digits. The important thing is that the approximation reaches such detail that the reconstruction of the physical phenomenon is indistinguishable from the original to the eyes (or ears) of the human being who uses the computer to process it.

Let's consider again the concept of encoding with the example of the image of Mario. To be processed by a computer, it must be transformed into a sequence of digits.

5678728729830939487489749874987484671011876278628767835463541376761291236213524136512438736149835023857467547790909394238472938472834782374923847238947283472983742836471651551562546546354635437647365847598475982475767151101982492849384934028391809174037547478347473647836473463784394010913481340946307456038475603847650183746508731465013874650783456038475603748657834657465873450618374560183746501837645675840387657348018736457108857346578310458713465317845110387465
0

Actually, these digits must themselves be encoded to be processed by a computer.
This further encoding transforms them into a sequence of 0s and 1s.

010101010101010001001110101010101010101
011101010000101010101011101010101010110
101010010011111110101010101010100010101
010101010111111001010001101010100011101
010101011110101101010001010100111010011
010100111000101011101010100010101011010
100011011010101110101001010001010101000
101110101010101110010101010100010101010
101010101110101010100010101010111010101
010110101011111010111000001110001110011
100110101011100111000110001010101011111
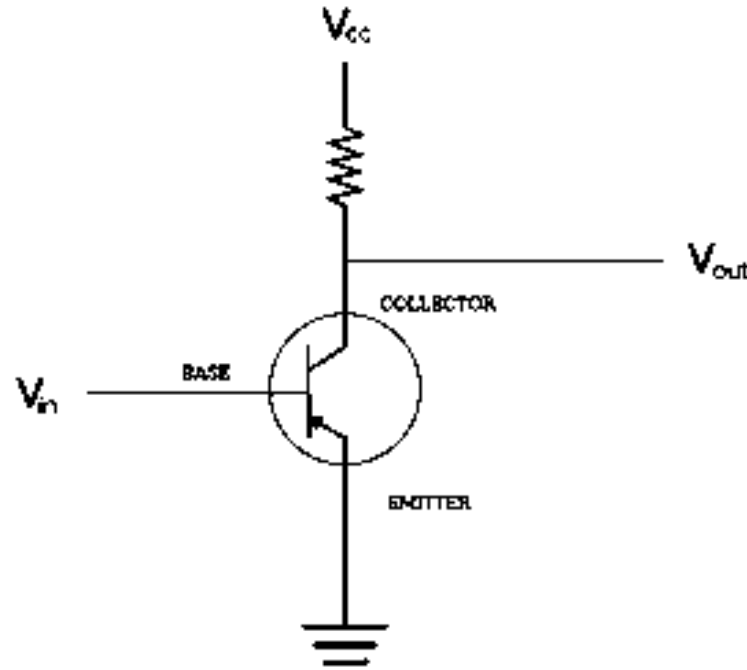000011100001010101010110010101011111000

# What is a bit?

It is the basic data unit processed by a computer: it is 0 or 1.

bit = binary + digit.

8 bit (b) = 1 Byte (B).

The encoding that uses only 0 and 1
is called binary encoding
(because it only uses 2 digits).

# Why 0 and 1?



High electric voltage: 1
Low electric voltage: 0

The electronic circuits that make up a computer are built to respond to high or low voltage electrical signals. High voltage is interpreted as a "1", while low voltage is interpreted as a "0". The restriction to only two values has no physical reason: it is possible to build circuits that respond with numerous voltage levels. For example, we could have 10 different ones, interpretable as the 10 digits (from "0" to "9") that are usually used in mathematics.

However, the advantages of having only two signals are numerous:
the circuits are simpler to make and cost less; furthermore, the output signals, even in the presence of disturbances due to natural causes, are easier to interpret, with less possibility of error. (e.g. $0.8 \approx 1$)

# Multiples

| |
|---|
| 1 KB (kilo) = 1000 B |
| 1 MB (mega) = 1 million B |
| 1 GB (giga) = 1 billion B |
| 1 TB (tera) = 1000 billion B |

# Binary encoding of numbers

Numbers with base 10:
$215 = 2 \times 10^2 + 1 \times 10^1 + 5 \times 10^0$

Numbers with base 2:
$110010111 = 1 \times 2^8 + 1 \times 2^7 + 0 \times 2^6 + 0 \times 2^5 + 1 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$

## Exercise

Convert the following numbers from base 2 to base 10: 101, 1000, 11011.

**Exercice**

Convert the following numbers from base 10 to base 2: 8, 23, 144, 201.

While the definition of a binary system helps us solve the first exercise, we must find new methods to solve the second. There are two.

First method:

given the number n in base 10, we look for the largest power of 2 that is less than or equal to n.

If it is the same, we have solved the problem: we write a 1 in the position corresponding to that power of 2, followed by zeros.

For example, 8 is a power of 2: $2^3$ to be precise, so its binary encoding will be 1000.

However, if the largest power of 2 that is less than or equal to n is less than n (let's call it k), let's set it aside and we calculate the difference n-k.

We repeat the same procedure with n-k, and look for the largest power of 2 that is less than or equal to it.

We continue until we are able to express n as a sum of powers of 2.

We take the list of powers and write a 1 in the corresponding positions, 0 in the others.

For example, the largest power of 2 contained in 23 is 16 ($2^4$). Their difference is 7, in which 4 ($2^2$) is contained. The difference is 3 where there is 2 ($2^1$), after which only 1 ($2^0$) remains.

Writing the powers of 2 present in order we get 10111.

**Second method:**

we divide the number by 2 and obtain quotient and remainder.
As long as the quotient is not 0, we take it as the new dividend and continue
dividing. When we get zero quotient, we have to write the list of remainders in
reverse order to get the binary encoding of the initial number.
Let's take 144 as an example:

$$144 : 2 = 72 \text{ with remainder } 0$$
$$72 : 2 = 36 \text{ with remainder } 0$$
$$36 : 2 = 18 \text{ with remainder } 0$$
$$18 : 2 = 9 \text{ with remainder } 0$$
$$9 : 2 = 4 \text{ with remainder } 1$$
$$4 : 2 = 2 \text{ with remainder } 0$$
$$2 : 2 = 1 \text{ with remainder } 0$$
$$1 : 2 = 0 \text{ with remainder } 1$$

144 in base 2 is 10010000 (these bits are the remainders in the reverse order
of writing)