

## Introduction to Data Engineering 1

Paul Blondel

UTSEUS, Shanghai

March 30st, 2021

*'You can have data without information, but you cannot have information without data.'*  
Daniel Keys Moran

### 1 Introduction

### 2 Memory

- Different memories for different usages
- How it works?

### 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

### 4 File formats

### 5 Databases

### 6 Data compression

### 1 Introduction

### 2 Memory

- Different memories for different usages
- How it works?

### 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

### 4 File formats

### 5 Databases

### 6 Data compression

Introduction

Introduction to Data Engineering 1 1 / 63

Since the beginning of computers, data is a key component:

- Operating systems are stored as data in the hard-drive and loaded in the main memory
- Computer users' data are stored on hard-drive, optical discs, etc

Because, with computers, data can not only be read but also easily written and erased, one can...

- ...install/uninstall a software
- ...update the OS
- ...write a word document
- ...have a custom system and a personal experience

#### The software era

Compared to an old-school calculator: the flexibility to access and edit the data made a great difference and opened the era of software creation.

However nowadays, data is more than just your operating system, all your software and some word documents...

There have been a series of recent discoveries and breakthroughs making possible the **extraction, the storage, the pre-processing and the processing, the analysis and the building of statistical and Machine Learning models on bigger and bigger amounts of data.**

#### The era of data

Brand new applications are nowadays unlocked thanks to that: making the engineering of data a top priority in most companies worldwide.

#### First thing first!

Before manipulating petabytes of data, let's understand how to change information into data, and first, how to store data into... the memory!

## 1 Introduction

## 2 Memory

- Different memories for different usages
- How it works?

## 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

## 4 File formats

## 5 Databases

## 6 Data compression

## Different types of memory:

- CPU registers
  - Stores data directly used by the CPU (instructions, memory addresses, program counter, accumulator, etc.)
- Cache Memory
  - Temporary stores copies of the data from frequently used and/or nearby main memory locations
- Main Memory (Random-access memory or RAM)
  - Stores running softwares and opened documents
- Mass Storage Memory
  - Stores installed softwares and all documents
  - Two main technologies on the market:
    - ★ Hard-disk drive (HDD) (invented in 1954)
    - ★ Solid-state drive (SSD) (invented in 1978)
- Optical discs
- GPU memory units

## How it works?

## • CPU Registers

- Y1, Y2, PSW, Z: registers used to store operation values, store computed value and remain
- RI, CO: registers storing instructions
- RAD: register storing memory address
- RDO, general registers: registers storing data

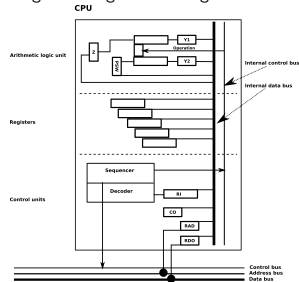


Figure: Simple CPU architecture

## How it works?

## • Cache Memory

- Located between the Computer Processing Unit (CPU) and the main memory
- Help reducing the number of accesses to the main memory, reducing the overall memory access time
  - ★ By temporary storing copies of the data from frequently used and/or nearby main memory locations

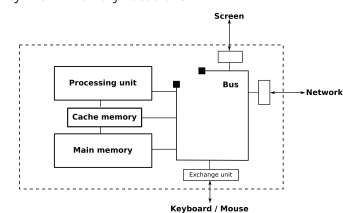


Figure: Von Neumann computer architecture

## How it works?

## • RAM (Main Memory)

- All the data is stored in an array of storage cells
- 1 bit = 1 storage cell (DRAM: bit stored in a capacitor)
- Volatile memory: need electric power for storage

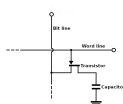


Figure: DRAM memory cell

## • HDD (Mass Storage Memory)

- All the data is stored in disks of ferromagnetic material
- 1 bit = direction of magnetic field lines

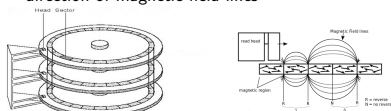


Figure: Left: disks of ferromagnetic material right: magnetic field lines

## How it works?

## • SSD (Mass Storage Memory)

- All the data is stored in an array of floating-gate transistors
- No mechanical moving component = faster
- Non-volatile memory: electrons trapped in the floating-gate (surrounded by oxide layers)

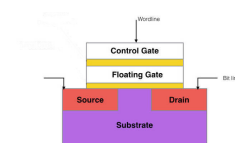


Figure: Floating gate transistor cell

## Optical Discs

- Mechanical moving component, disc moving and laser head reading
- Non-volatile memory: bits are "physically" created ("holes" or chemical reactions)

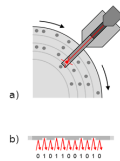


Figure: Reading optical discs

- The laser arm and the disc move alike the HDD mechanism (a)
- Bits are represented by pits etched on the disc itself (b)

Summary of the main properties of computer memories:

|                  | CPU Registers                    | Cache                         | RAM                    | HDD / SSD           | Optical Discs              |
|------------------|----------------------------------|-------------------------------|------------------------|---------------------|----------------------------|
| Access speed     | Ultra fast (1 or few CPU cycles) | Very fast (20000-300000 MO/s) | Fast (2000-26600 MO/s) | Slow (50-5000 MO/s) | Very Slow (0.150-3 MO/s)   |
| Storage capacity | Ultra low (32/64 bits)           | Very low (512KB-8MB)          | Low (4-32GB)           | High (256GB - 1TB)  | Moderate High (700MO-25GB) |
| Storage property | Volatile                         | Not Volatile                  | Volatile               | Not Volatile        | Not Volatile               |

## Computer Memory Hierarchy

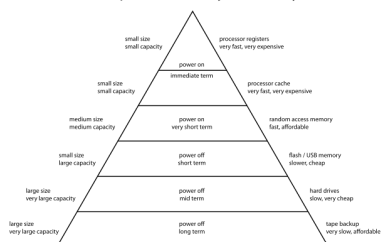


Figure: The Memory Hierarchy: a trade-off between storage capacity and speed

### Notes

To optimize the speed of execution: top of the pyramid memories should be use as much as possible! (Require to design your program accordingly).

GPU memories: a complex assembly:

- Registers: very fast access, thread, few space
- Local Memory: slow access, thread, more space
- Shared Memory: fast access, block, more space
- Global Memory: slow access, all blocks, more space
- Constant Memory: fast access, all blocks, more space
- Texture Memory: fast access, all blocks, read-only

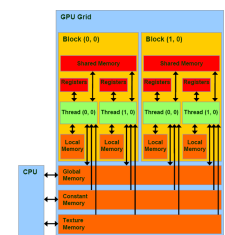


Figure: Memory units of Nvidia GPUs

GPU memories: A more complex trade-off

With GPU memory units: a trade-off storage capacity VS access speed VS accessibility. That's why coding optimal code for GPU is hard.

## 1 Introduction

## 2 Memory

- Different memories for different usages
- How it works?

## 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

## 4 File formats

## 5 Databases

## 6 Data compression

## What is data?

Data is the result of a measurement

Examples: The diameter of the sun is **1.39 million kilometers**, the temperature of the room is **30°C**, and so on.

A data ...

- ... is **simple** or **complex** (composed of simple data)
- ... has a **type** (integer, character, and so on)
- ... has a set of authorized **instructions**

Examples: 'hello' is a complex data (composed of simple character data), 30°C is a simple integer data. 'hello' \* 30 is not authorized, but 2 \* 30 is.

The most common units are:

- 1 bit (0 or 1)
- 1 byte (or octet) = 8 bits
- 1 KB (or KO) = 1000 bytes<sup>1</sup>
- 1 MB (or MO) = 10<sup>6</sup> bytes = 1000 KB
- 1 GB (or GO) = 10<sup>9</sup> bytes = 10<sup>6</sup> KB = 1000 MB
- 1 TB (or TO) = 10<sup>12</sup> bytes

Example: 'hello' is stored on 5 bytes (each character is stored on 1 byte, see later in this course why it is the case...).

<sup>1</sup> 1 KB is sometimes defined as equal to 1024 bytes, but this definition is no longer valid. In 1998, the International Electrotechnical Commission stated that 1 KB = 1000 bytes and that 1024 bytes = 1 KiB (kibibyte).

Data types in C/C++:

| Type          | Typical Bit Width | Typical Range                   |
|---------------|-------------------|---------------------------------|
| char          | 1byte             | -127 to 127 or 0 to 255         |
| unsigned char | 1byte             | 0 to 255                        |
| signed char   | 1byte             | -127 to 127                     |
| int           | 4bytes            | -2147483648 to 2147483647       |
| unsigned int  | 4bytes            | 0 to 4294967295                 |
| signed int    | 4bytes            | -2147483648 to 2147483647       |
| long int      | 8bytes            | -2,147,483,648 to 2,147,483,647 |
| float         | 4bytes            | 1.2E-38 to 3.4E+38              |
| double        | 8bytes            | 2.3E-308 to 1.7E+308            |

```
#include <iostream>
int main(void) {
    std::cout << sizeof(char) << std::endl;
}
```

**C/C++ Structure** (aggregation of variables of different types)

Size of a structure = sum of variables sizes (1 byte: C++ and empty)

Example of a 8 bytes structure:

```
struct {
    int a; // = 4 bytes
    float b; // = 4 bytes
};
struct test; // 8 bytes
```

**C Union** (only one variable is used among many)

Size of an union = size of the biggest element of the union.

Example of a 8 bytes union:

```
union {
    int a; // = 4 bytes
    double b; // = 8 bytes
};
union test; // 8 bytes
```

**C++ Class** (aggregation of variables + methods)

Size of an instanced class = size of its variables

(+ N \* (size of its variables) if N inheritances)

Example of a 16 bytes object:

```
class A{
    int a; // = 4 bytes
    float b; // = 4 bytes
    A() {}
};
class B : A{
    int ap; // = 4 bytes
    float bp; // = 4 bytes
    B(): A() {}
};
B btest; // 16 bytes object
```

Other language may have different data type sizes ...

Data types in Java:

| Type    | Typical Bit Width | Typical Range   |
|---------|-------------------|---|
| byte    | 1 byte            | -128 to 127   |
| short   | 2 bytes           | -32,768 to 32,767                                       |
| int     | 4 bytes           | -2,147,483,648 to 2,147,483,647                         |
| long    | 8 bytes           | -9,223,372,036,854,775,808 to 9,223,372,036,854,775,807 |
| float   | 4 bytes           | Sufficient for storing 6 to 7 decimal digits            |
| double  | 8 bytes           | Sufficient for storing 15 decimal digits                |
| boolean | 1 bit             | True or false values                                    |
| char    | 2 bytes           | Stores a single character/letter or ASCII values        |

```
public class example {
    public static void main(String [] args) {
        System.out.println(ObjectSizeCalculator.getObjectSize(3)); // 16
    }
}
```

Data types in Python (3.9.2 on Linux x86-64)

| Type        | Typical Bit Width |
|-------------|-------------------|
| int         | 28 bytes          |
| float       | 24 bytes          |
| empty str   | 49 bytes          |
| empty list  | 56 bytes          |
| empty tuple | 40 bytes          |
| empty dict  | 232 bytes         |
| set         | 216 bytes         |
| frozenset   | 216 bytes         |
| None        | 16 bytes          |

```
import sys
a = sys.getsizeof(12)
print(a)
b=sys.getsizeof('geeks')
print(b)
```

Storing integer values is relatively straightforward.

Examples:

- You all know that  $8 = 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 0 \times 2^0$ . Thus, 8 is equal to 1000 in binary code (we just need 4 bits to store 8). We can use a "char" type to store this value in C++.
- Another example:  $11 = 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$ . The binary form of 11 is just: 1011. Again we just need 4 bits, we can use a "char" type to store this value in C++.
- A last example:  $888 = 1101111000$  in binary code. In this case we need to store 10 bits.  $10\text{bits} > 1\text{bytes}$  and  $10\text{bits} < 2\text{bytes}$ : we can use a short type to store this value in C++.

#### Rule

For unsigned values: you need to count the numbers of bits required to store the integer, and you know which type you can use. For signed values: beware! You also need to have a bit for the sign of the integer!

Storing floating numbers is little bit more complicated.

Floating point numbers usually follow **IEEE-754 representations** (mostly the **IEEE-754 Single Precision** and the **IEEE-754 Double Precision** representations).

- IEEE-754 Single Precision:

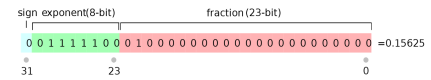


Figure: 8 bits for exponent part, 23 bits for fraction part and 1 bit for sign

- exponent part = real exponent + 127 (in binary)
  - ★ 127 is a bias to handle negative exponents
- fraction part = fraction part (in binary)
- sign part = 0 (if positive) or 1 (if negative)
- number (in binary) =  $(-1)^{\text{sign}} \times 1 \times \text{fraction (in binary)} \times 2^{\text{real exponent}}$

- Example of a real in IEEE-754 Single Precision binary representation:
  - $23.375 = (10111.011)_2 = (-1)^0 \times 1.0111011 \times 2^{-4}$
  - sign part =  $0_2$
  - fraction part =  $(011101100000000000000000)_2$
  - exponential part =  $(00000100)_2 + (10000000)_2 = (10000100)_2$
  - Final bit storage =  $(01000010001101100000000000000000)_2$
- IEEE-754 Double Precision:

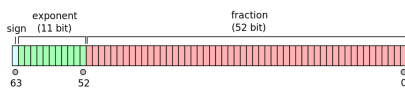


Figure: 11 bits for exponent part, 52 bits for fraction part and 1 bit for sign

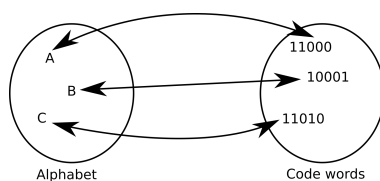
- exponent part = real exponent + 1023 (in binary)
  - ★ 1023 is a bias to handle negative exponents
- fraction part = fraction part (in binary)
- sign part = 0 (if positive) or 1 (if negative)
- number (in binary) =  $(-1)^{\text{sign}} \times 1 \times \text{fraction (in binary)} \times 2^{\text{real exponent}}$

Dealing with encoding characters:

- We know how integers and floating numbers are stored, but what about characters?
- Indeed, characters cannot be directly transformed into bits (such as integers or floating number).
- And all information need to be transformed into bits to be processed by a computer...
- We need an **encoding mechanism** to transform characters into bits.

What is encoding? Encoding the information consists in matching symbols to binary representations (code words) using an encoding map.

In the given example below the information is made of alphanumerical characters<sup>2</sup>:



<sup>2</sup>it is a bijection

How to get an encoding map of code words?

We can consider the character space as a system and each character as a state in this system. And we can use a fixed-state coding approach to get an encoding map:

- Each state of the system is coded by a certain number of bits (called length of code) with:
  - 1 bit: we can code 2 states (0 or 1)
  - 2 bits: we can code 4 states (00, 01, 10, 11)
  - N bits: we can code  $2^N$  states

The number of states required to code  $P$  states is  $n$  such as:  $2^{(n-1)} < P \leq 2^n$

Example: We have 26 letters in the alphabets, how many bits are necessary to code  $P$  values? We need 5 bits, because:  $2^4 < 26 \leq 2^5$

Examples of character encoding systems:

- The **Baudot code** is encoded on 5 bits. With the Baudot code we can code up to 32 characters (and thus the 26 letters of the alphabet).
- The **ASCII table** (7 bits encoding). With the ASCII table uppercase letters, system and special characters can also be encoded (we can code 128 characters).

#### ASCII TABLE

| Decimal Hex Char | Decimal Hex Char | Decimal Hex Char | Decimal Hex Char |
|------------------|------------------|------------------|------------------|
| 0 0 NUL          | 16 10 P          | 32 20 D          | 48 30 0          |
| 1 1 SOH          | 17 11 Q          | 33 21 E          | 49 31 1          |
| 2 2 STX          | 18 12 R          | 34 22 F          | 50 32 2          |
| 3 3 ETX          | 19 13 S          | 35 23 G          | 51 33 3          |
| 4 4 EOT          | 20 14 T          | 36 24 H          | 52 34 4          |
| 5 5 ENQ          | 21 15 U          | 37 25 I          | 53 35 5          |
| 6 6 ACK          | 22 16 V          | 38 26 J          | 54 36 6          |
| 7 7 BEL          | 23 17 W          | 39 27 K          | 55 37 7          |
| 8 8 BS           | 24 18 X          | 40 28 L          | 56 38 8          |
| 9 9 HT           | 25 19 Y          | 41 29 M          | 57 39 9          |
| 10 A LF          | 26 1A Z          | 42 2A N          | 58 3A 0          |
| 11 B VT          | 27 1B [          | 43 2B O          | 59 3B 1          |
| 12 C FF          | 28 1C \          | 44 2C P          | 60 3C 2          |
| 13 D SO          | 29 1D ^          | 45 2D Q          | 61 3D 3          |
| 14 E SH          | 30 1E _          | 46 2E R          | 62 3E 4          |
| 15 F OH          | 31 1F `          | 47 2F S          | 63 3F 5          |
| 16 G NL          | 32 20 A          | 48 30 T          | 64 40 6          |
| 17 H VT          | 33 21 B          | 49 31 U          | 65 41 7          |
| 18 I SO          | 34 22 C          | 50 32 V          | 66 42 8          |
| 19 J SH          | 35 23 D          | 51 33 W          | 67 43 9          |
| 20 K ETX         | 36 24 E          | 52 34 X          | 68 44 0          |
| 21 L EOT         | 37 25 F          | 53 35 Y          | 69 45 1          |
| 22 M ENQ         | 38 26 G          | 54 36 Z          | 70 46 2          |
| 23 N ACK         | 39 27 H          | 55 37 [          | 71 47 3          |
| 24 O BEL         | 40 28 I          | 56 38 \          | 72 48 4          |
| 25 P BS          | 41 29 J          | 57 39 ^          | 73 49 5          |
| 26 Q HT          | 42 2A K          | 58 3A _          | 74 4A 6          |
| 27 R LF          | 43 2B L          | 59 3B `          | 75 4B 7          |
| 28 S VT          | 44 2C M          | 60 3C A          | 76 4C 8          |
| 29 T SO          | 45 2D N          | 61 3D B          | 77 4D 9          |
| 30 U SH          | 46 2E O          | 62 3E C          | 78 4E 0          |
| 31 V OH          | 47 2F P          | 63 3F D          | 79 4F 1          |
| 32 W NL          | 48 30 Q          | 64 40 E          | 80 50 2          |
| 33 X VT          | 49 31 R          | 65 41 F          | 81 51 3          |
| 34 Y SO          | 50 32 S          | 66 42 G          | 82 52 4          |
| 35 Z SH          | 51 33 T          | 67 43 H          | 83 53 5          |
| 36 [ ETX         | 52 34 U          | 68 44 I          | 84 54 6          |
| 37 \ EOT         | 53 35 V          | 69 45 J          | 85 55 7          |
| 38 ] ENQ         | 54 36 W          | 70 46 K          | 86 56 8          |
| 39 ^ ACK         | 55 37 X          | 71 47 L          | 87 57 9          |
| 40 _ BEL         | 56 38 Y          | 72 48 M          | 88 58 0          |
| 41 ` BS          | 57 39 Z          | 73 49 N          | 89 59 1          |
| 42 A HT          | 58 3A [          | 74 4A O          | 90 5A 2          |
| 43 B LF          | 59 3B \          | 75 4B P          | 91 5B 3          |
| 44 C VT          | 60 3C ^          | 76 4C Q          | 92 5C 4          |
| 45 D SO          | 61 3D _          | 77 4D R          | 93 5D 5          |
| 46 E SH          | 62 3E A          | 78 4E S          | 94 5E 6          |
| 47 F OH          | 63 3F B          | 79 4F T          | 95 5F 7          |
| 48 G NL          | 64 40 C          | 80 60 U          | 96 60 8          |
| 49 H VT          | 65 41 D          | 81 61 V          | 97 61 9          |
| 50 I SO          | 66 42 E          | 82 62 W          | 98 62 0          |
| 51 J SH          | 67 43 F          | 83 63 X          | 99 63 1          |
| 52 K ETX         | 68 44 G          | 84 64 Y          | 100 64 2         |
| 53 L EOT         | 69 45 H          | 85 65 Z          | 101 65 3         |
| 54 M ENQ         | 70 46 I          | 86 66 [          | 102 66 4         |
| 55 N ACK         | 71 47 J          | 87 67 \          | 103 67 5         |
| 56 O BEL         | 72 48 K          | 88 68 ^          | 104 68 6         |
| 57 P BS          | 73 49 L          | 89 69 _          | 105 69 7         |
| 58 Q HT          | 74 4A M          | 90 6A `          | 106 6A 8         |
| 59 R LF          | 75 4B N          | 91 6B A          | 107 6A 9         |
| 60 S VT          | 76 4C O          | 92 6B B          | 108 6B 0         |
| 61 T SO          | 77 4D P          | 93 6C C          | 109 6B 1         |
| 62 U SH          | 78 4E Q          | 94 6D D          | 110 6B 2         |
| 63 V OH          | 79 4F R          | 95 6E E          | 111 6B 3         |
| 64 W NL          | 80 50 S          | 96 6F F          | 112 6B 4         |
| 65 X VT          | 81 51 T          | 97 70 G          | 113 6B 5         |
| 66 Y SO          | 82 52 U          | 98 71 H          | 114 6B 6         |
| 67 Z SH          | 83 53 V          | 99 72 I          | 115 6B 7         |
| 68 [ ETX         | 84 54 W          | 100 73 J         | 116 6B 8         |
| 69 \ EOT         | 85 55 X          | 101 74 K         | 117 6B 9         |
| 70 ] ENQ         | 86 56 Y          | 102 75 L         | 118 6C 0         |
| 71 ^ ACK         | 87 57 Z          | 103 76 M         | 119 6C 1         |
| 72 _ BEL         | 88 58 [          | 104 77 N         | 120 6C 2         |
| 73 ` BS          | 89 59 \          | 105 78 O         | 121 6C 3         |
| 74 A HT          | 90 5A ^          | 106 79 P         | 122 6C 4         |
| 75 B LF          | 91 5B _          | 107 7A Q         | 123 6C 5         |
| 76 C VT          | 92 5C A          | 108 7B R         | 124 6C 6         |
| 77 D SO          | 93 5D B          | 109 7C S         | 125 6C 7         |
| 78 E SH          | 94 5E C          | 110 7D T         | 126 6C 8         |
| 79 F OH          | 95 5F D          | 111 7E U         | 127 6C 9         |
| 80 G NL          | 96 60 E          | 112 7F V         | 128 6D 0         |
| 81 H VT          | 97 61 F          | 113 80 W         | 129 6D 1         |
| 82 I SO          | 98 62 G          | 114 81 X         | 130 6D 2         |
| 83 J SH          | 99 63 H          | 115 82 Y         | 131 6D 3         |
| 84 K ETX         | 100 64 I         | 116 83 Z         | 132 6D 4         |
| 85 L EOT         | 101 65 J         | 117 84 [         | 133 6D 5         |
| 86 M ENQ         | 102 66 K         | 118 85 \         | 134 6D 6         |
| 87 N ACK         | 103 67 L         | 119 86 ^         | 135 6D 7         |
| 88 O BEL         | 104 68 M         | 120 87 _         | 136 6D 8         |
| 89 P BS          | 105 69 N         | 121 88 A         | 137 6D 9         |
| 90 Q HT          | 106 6A O         | 122 89 B         | 138 6E 0         |
| 91 R LF          | 107 6B P         | 123 8A C         | 139 6E 1         |
| 92 S VT          | 108 6C Q         | 124 8B D         | 140 6E 2         |
| 93 T SO          | 109 6D R         | 125 8C E         | 141 6E 3         |
| 94 U SH          | 110 6E S         | 126 8D F         | 142 6E 4         |
| 95 V OH          | 111 6F T         | 127 8E G         | 143 6E 5         |
| 96 W NL          | 112 70 U         | 128 8F H         | 144 6E 6         |
| 97 X VT          | 113 71 V         | 129 90 I         | 145 6E 7         |
| 98 Y SO          | 114 72 W         | 130 91 J         | 146 6E 8         |
| 99 Z SH          | 115 73 X         | 131 92 K         | 147 6E 9         |
| 100 [ ETX        | 116 74 Y         | 132 93 L         | 148 6F 0         |
| 101 \ EOT        | 117 75 Z         | 133 94 M         | 149 6F 1         |
| 102 ] ENQ        | 118 76 [         | 134 95 N         | 150 6F 2         |
| 103 ^ ACK        | 119 77 \         | 135 96 O         | 151 6F 3         |
| 104 _ BEL        | 120 78 ^         | 136 97 P         | 152 6F 4         |
| 105 ` BS         | 121 79 _         | 137 98 Q         | 153 6F 5         |
| 106 A HT         | 122 7A A         | 138 99 R         | 154 6F 6         |
| 107 B LF         | 123 7B B         | 139 100 S        | 155 6F 7         |
| 108 C VT         | 124 7C C         | 140 101 T        | 156 6F 8         |
| 109 D SO         | 125 7D D         | 141 102 U        | 157 6F 9         |
| 110 E SH         | 126 7E E         | 142 103 V        | 158 70 0         |
| 111 F OH         | 127 7F F         | 143 104 W        | 159 70 1         |
| 112 G NL         | 128 80 G         | 144 105 X        | 160 70 2         |
| 113 H VT         | 129 81 H         | 145 106 Y        | 161 70 3         |
| 114 I SO         | 130 82 I         | 146 107 Z        | 162 70 4         |
| 115 J SH         | 131 83 J         | 147 108 [        | 163 70 5         |
| 116 K ETX        | 132 84 K         | 148 109 \        | 164 70 6         |
| 117 L EOT        | 133 85 L         | 149 110 ^        | 165 70 7         |
| 118 M ENQ        | 134 86 M         | 150 111 _        | 166 70 8         |
| 119 N ACK        | 135 87 N         | 151 112 A        | 167 70 9         |
| 120 O BEL        | 136 88 O         | 152 113 B        | 168 71 0         |
| 121 P BS         | 137 89 P         | 153 114 C        | 169 71 1         |
| 122 Q HT         | 138 8A Q         | 154 115 D        | 170 71 2         |
| 123 R LF         | 139 8B R         | 155 116 E        | 171 71 3         |
| 124 S VT         | 140 8C S         | 156 117 F        | 172 71 4         |
| 125 T SO         | 141 8D T         | 157 118 G        | 173 71 5         |
| 126 U SH         | 142 8E U         | 158 119 H        | 174 71 6         |
| 127 V OH         | 143 8F V         | 159 120 I        | 175 71 7         |
| 128 W NL         | 144 90 W         | 160 121 J        | 176 71 8         |
| 129 X VT         | 145 91 X         | 161 122 K        | 177 71 9         |
| 130 Y SO         | 146 92 Y         | 162 123 L        | 178 72 0         |
| 131 Z SH         | 147 93 Z         | 163 124 M        | 179 72 1         |
| 132 [ ETX        | 148 94 [         | 164 125 N        | 180 72 2         |
| 133 \ EOT        | 149 95 \         | 165 126 O        | 181 72 3         |
| 134 ] ENQ        | 150 96 ^         | 166 127 P        | 182 72 4         |
| 135 ^ ACK        | 151 97 _         | 167 128 Q        | 183 72 5         |
| 136 _ BEL        | 152 98 A         | 168 129 R        | 184 72 6         |
| 137 ` BS         | 153 99 B         | 169 130 S        | 185 72 7         |
| 138 A HT         | 154 100 C        | 170 131 T        | 186 72 8         |
| 139 B LF         | 155 101 D        | 171 132 U        | 187 72 9         |
| 140 C VT         | 156 102 E        | 172 133 V        | 188 73 0         |
| 141 D SO         | 157 103 F        | 173 134 W        | 189 73 1         |
| 142 E SH         | 158 104 G        | 174 135 X        | 190 73 2         |
| 143 F OH         | 159 105 H        | 175 136 Y        | 191 73 3         |
| 144 G NL         | 160 106 I        | 176 137 Z        | 192 73 4         |
| 145 H VT         | 161 107 J        | 177 138 [        | 193 73 5         |
| 146 I SO         | 162 108 K        | 178 139 \        | 194 73 6         |
| 147 J SH         | 163 109 L        | 179 140 ^        | 195 73 7         |
| 148 K ETX        | 164 110 M        | 180 141 _        | 196 73 8         |
| 149 L EOT        | 165 111 N        | 181 142 A        | 197 73 9         |
| 150 M ENQ        | 166 112 O        | 182 143 B        | 198 74 0         |
| 151 N ACK        | 167 113 P        | 183 144 C        | 199 74 1         |
| 152 O BEL        | 168 114 Q        | 184 145 D        | 200 74 2         |
| 153 P BS         | 169 115 R        | 185 146 E        | 201 74 3         |
| 154 Q HT         | 170 116 S        | 186 147 F        | 202 74 4         |
| 155 R LF         | 171 117 T        | 187 148 G        | 203 74 5         |
| 156 S VT         | 172 118 U        | 188 149 H        | 204 74 6         |
| 157 T SO         | 173 119 V        | 189 150 I        | 205 74 7         |
| 158 U SH         | 174 120 W        | 190 151 J        | 206 74 8         |
| 159 V OH         | 175 121 X        | 191 152 K        | 207 74 9         |
| 160 W NL         | 176 122 Y        | 192 153 L        | 208 75 0         |
| 161 X VT         | 177 123 Z        | 193 154 M        | 209 75 1         |
| 162 Y SO         | 178 124 [        | 194 155 N        | 210 75 2         |
| 163 Z SH         | 179 125 \        | 195 156 O        | 211 75 3         |
| 164 [ ETX        | 180 126 ^        | 196 157 P        | 212 75 4         |
| 165 \ EOT        | 181 127 _        | 197 158 Q        | 213 75 5         |
| 166 ] ENQ        | 182 128 A        | 198 159 R        | 214 75 6         |
| 167 ^ ACK        | 183 129 B        | 199 160 S        | 215 75 7         |
| 168 _ BEL        | 184 130 C        | 200 161 T        | 216 75 8         |
| 169 ` BS         | 185 131 D        | 201 162 U        | 217 75 9         |
| 170 A HT         | 186 132 E        | 202 163 V        | 218 76 0         |
| 171 B LF         | 187 133 F        | 203 164 W        | 219 76 1         |
| 172 C VT         | 188 134 G        | 204 165 X        | 220 76 2         |
| 173 D SO         | 189 135 H        | 205 166 Y        | 221 76 3         |
| 174 E SH         | 190 136 I        | 206 167 Z        | 222 76 4         |
| 175 F OH         | 191 137 J        | 207 168 [        | 223 76 5         |
| 176 G NL         | 192 138 K        | 208 169 \        | 224 76 6         |
| 177 H VT         | 193 139 L        | 209 170 ^        | 225 76 7         |
| 178 I SO         | 194 140 M        | 210 171 _        | 226 76 8         |
| 179 J SH         | 195 141 N        | 211 172 A        | 227 76 9         |
| 180 K ETX        | 196 142 O        | 212 173 B        | 228 77 0         |
| 181 L EOT        | 197 143 P        | 213 174 C        | 229 77 1         |
| 182 M ENQ        | 198 144 Q        | 214 175 D        | 230 77 2         |
| 183 N ACK        | 199 145 R        | 215 176 E        | 231 77 3         |
| 184 O BEL        | 200 146 S        | 216 177 F        | 232 77 4         |
| 185 P BS         | 201 147 T        | 217 178 G        | 233 77 5         |
| 186 Q HT         | 202 148 U        | 218 179 H        | 234 77 6         |
| 187 R LF         | 203 149 V        | 219 180 I        | 235 77 7         |
| 188 S VT         | 204 150 W        | 220 181 J        | 236 77 8         |
| 189 T SO         | 205 151 X        | 221 182 K        | 237 77 9         |
| 190 U SH         | 206 152 Y        | 222 183 L        | 238 78 0         |
| 191 V OH         | 207 153 Z        | 223 184 M        | 239 78 1         |
| 192 W NL         | 208 154 [        | 224 185 N        | 240 78 2         |
| 193 X VT         | 209 155 \        | 225 186 O        | 241 78 3         |
| 194 Y SO         | 210 156 ^        | 226 187 P        | 242 78 4         |
| 195 Z SH         | 211 157 _        | 227 188 Q        | 243 78 5         |
| 196 [ ETX        | 212 158 A        | 228 189 R        | 244 78 6         |
| 197 \ EOT        | 213 159 B        | 229 190 S        | 245 78 7         |
| 198 ] ENQ        | 214 160 C        | 230 191 T        | 246 78 8         |
| 199 ^ ACK        | 215 161 D        | 231 192 U        | 247 78 9         |
| 200 _ BEL        | 216 162 E        | 232 193 V        | 248 79 0         |
| 201 ` BS         | 217 163 F        | 233 194 W        | 249 79 1         |
| 202 A HT         | 218 164 G        | 234 195 X        | 250 79 2         |
| 203 B LF         | 219 165 H        | 235 196 Y        | 251 79 3         |
| 204 C VT         | 220 166 I        | 236 197 Z        | 252 79 4         |
| 205 D SO         | 221 167 J        | 237 198 [        | 253 79 5         |
| 206 E SH         | 222 168 K        | 238 199 \        | 254 79 6         |
| 207 F OH         | 223 169 L        | 239 200 ^        | 255 79 7         |
| 208 G NL         | 224 170 M        | 240 201 _        | 256 79 8         |
| 209 H VT         | 225 171 N        | 241 202 A        | 257 79 9         |
| 210 I SO         | 226 172 O        | 242 203 B        | 258 80 0         |
| 211 J SH         | 227 173 P        | 243 204 C        | 259 80 1         |
| 212 K ETX        | 228 174 Q        | 244 205 D        | 260 80 2         |
| 213 L EOT        | 229 175 R        | 245 206 E        | 261 80 3         |
| 214 M ENQ        | 230 176 S        | 246 207 F        | 262 80 4         |
| 215 N ACK        | 231 177 T        | 247 208 G        | 263 80 5         |
| 216 O BEL        | 232 178 U        | 248 209 H        | 264 80 6         |
| 217 P BS         | 233 179 V        | 249 210 I        | 265 80 7         |
| 218 Q HT         | 234 180 W        | 250 211 J        | 266 80 8         |
| 219 R LF         | 235 181 X        | 251 212 K        | 267 80 9         |
| 220 S VT         | 236 182 Y        | 252 213 L        | 268 81 0         |
| 221 T SO         | 237 183 Z        | 253 214 M        | 269 81 1         |
| 222 U SH         | 238 184 [        | 254 215 N        | 270 81 2         |
| 223 V OH         | 239 185 \        | 255 216 O        | 271 81 3         |
| 224 W NL         | 240 186 ^        | 256 217 P        | 272 81 4         |
| 225 X VT         | 241 187 _        | 257 218 Q        | 273 81 5         |
| 226 Y SO         | 242 188 A        | 258 219 R        | 274 81 6         |
| 227 Z SH         | 243 189 B        | 259 220 S        | 275 81 7         |
| 228 [ ETX        | 244 190 C        | 260 221 T        | 276 81 8         |
| 229 \ EOT        | 245 191 D        | 261 222 U        | 277 81 9         |
| 230 ] ENQ        | 246 192 E        | 262 223 V        | 278 82 0         |
| 231 ^ ACK        | 247 193 F        | 263 224 W        | 279 82 1         |
| 232 _ BEL        | 248 194 G        | 264 225 X        | 28               |

## Example 1:

The max frequency of the human voice is between 300Hz and 3400Hz.

Let's choose  $F_{max} = 4000\text{Hz}$

According to the Shannon law,  $F_{sampling}$  must be  $\geq 2 \times F_{max}$

So let's choose  $F_{sampling} = 8000\text{Hz}$  (8000 samples per seconds)

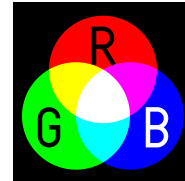
Besides, to get a correct audio reproduction each sample should be encoded on 12 bits.

So, in order to store in real time the human voice we should have a Bit Rate of  $8000 \times 12 = 96000 \text{ bits/second}$

## Example 2:

As you know, the color can be decomposed into three different primary colors (red, green and blue).

The weighted sum of these three primary colors give the color (this is the additive synthesis).



Each image point in an image is represented by two values: the luminance (luminous intensity) and the chrominance (conveying color information).

These two values are bounded together by this relation:

$$Y = 0.3 \times R + 0.59 \times G + 0.11 \times B$$

Y: luminance R, G and B: chrominance components

For example, with the lowest quality of European Digital TV standard we have:

- 625 lines per image (576 lines are really used).
- 720 points per line
- 25 images/second

Thus we can store:

- 720 points per line for the luminance (encoded on 8bits)
- 360 points each for the blue and red chrominance (encoded on 8bits)

So we have 1440 points per line.

Number of bits to store for one image =  $(1440 \times 8 \times 576) = 6635520 \text{ bits}$

With a framerate of 25 images/second: The Bit Rate must be:  $6635520 \times 25 = 166\text{Mbit/s}$

Actually this Bit Rate is difficult to realize, one solution is to compress the data to lower the required Bit Rate.

This is the goal of the "Motion Picture Expert Group" or MPEG: they define the compression algorithm for sound and video (see after in the course for more info about compression)...

## 1 Introduction

## 2 Memory

- Different memories for different usages
- How it works?

## 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

## 4 File formats

## 5 Databases

## 6 Data compression

On computers, data is mostly stored on files.

We have two categories of file formats: **binary** and **text** file formats:

- Binary file formats:
  - ▶ Can store multiple types of information in a single file (audio, text, video). Example: `xlsx`
  - ▶ Very application dependent (mostly, only Excel can read `xlsx`, you cannot open `avi` file with Excel, etc)
  - ▶ Here 1 byte = 1 custom data (huge byte sequence dependency)
- Text file formats:
  - ▶ Encoding dependent (ASCII, ISO-8859-1, GB18030, etc)
  - ▶ Can be opened on any system by any editor
  - ▶ Easy to read and understand
  - ▶ Less prone to get corrupted, more robust (1 byte error = may be just one error on an isolated and non important character)
  - ▶ 1 byte = 1 readable character

A non-exhaustive list of interesting file formats in Data Engineering:

- CSV, TSV: column and tab-based text file formats
- XLSX XLS: sheets+columns+media binary Excel file format
- JSON: key-value-based structured file format
- XML: tag-based structured file format
- AVRO: Hadoop binary dump file format
- PARQUET: column first binary file format
- JPG, PNG and so on: image binary file formats
- MPEG, OVG and so on: video binary file formats
- MP3 and so on: audio binary file formats
- Many other...

## CSV (Comma Separated Value) - Text file

- A CSV is a file with the values comma-separated, this is used to store column and row data
- CSV files can be used with Microsoft Excel, Google Spreadsheets, etc
- It differs from other spreadsheet format: you can only have one single sheet in a file. You cannot store formulas, images, etc.

CSV file example:

```
column_name_1, column_name_2, column_name_3, ...  
value1, value2, value3, ...
```

Code example (python/pandas):

```
import pandas as pd  
data = pd.read_csv(sep=",", "example.csv")  
print(data['column1'])
```

Code example (awk):

```
awk -F',' '{ print $1; }'
```

## TSV (Tabular Separated Value) - Text file

- Values are tab-separated
- TSV is an alternative to the common comma-separated values (CSV) format
- CSV: difficulties if you use commas in values...

Example:

```
column_name_1 column_name_2 column_name_3 ...  
value1 value2 value3 ...
```

Code (using Python/pandas):

```
import pandas as pd  
data = pd.read_csv(sep="t", "example.csv")  
print(data['column1'])
```

Code (using awk):

```
awk -F'\t' '{ print $1; }'
```

## XLSX, XLS (Excel file format) - Binary format

- Official spreadsheet format of Microsoft Excel
- Store column/raw data in sheets
- Can store more than one sheet, plus images, videos, etc.

Code (using pandas):

```
import pandas  
data = pd.ExcelFile('example.xls')  
df1 = pd.read_excel(data, 'my sheet 1')  
df2 = pd.read_excel(data, 'my sheet 2')
```

## JSON (JavaScript Object Notation) - Text format

- Data are stored and structured as attribute-value pairs
- Returned format of some APIs (Application Programming Interface)
- Open standard file format

Example:

```
{ "sport": {  
    "ball": {  
        "solo": ["tennis", "pingpong"],  
        "team": ["soccer", "rugby"]  
    }  
}}
```

Code (using pandas):

```
import json  
data = json.load(open("example.json"))  
print(data['sport']['ball'])  
print(data['sport']['ball']['solo'])
```

## XML (eXtensive Markup Language) - Text format

- Format used to structure data for storage and transport
- The structure is made of tags, options and clear text
- This is the returned type of some APIs (ex: WeChat API)

Example:

```
<?xml version="1.0" encoding="UTF-8"?>  
<email>  
  <from>Jani</from>  
  <body>Message 1</body>  
</email>
```

Code example (Python):

```
import xml.etree.ElementTree as ET  
tree = ET.parse("example.xml")  
root = tree.getroot()  
item = root.find("email")
```



## Parquet - Binary format

- Designed for large-scale analytical querying
- A flat columnar storage format (different than CSV/TSB that are row-based)
- Designed to bring efficiency
  - ▶ Row-based easier to conceptualize, but slower to use in some case
  - ▶ Column-based permits to skip non-relevant data quickly

| ROW-BASED STORAGE             | COLUMNAR STORAGE                     |
|-------------------------------|--------------------------------------|
| 1 MARG THOMPSON WASHINGTON 27 | ID: 1 2 3                            |
| 2 JIM THOMPSON DENVER 33      | FIRST NAME: MARG, JIM, JACK          |
| 3 JACK RILEY SEATTLE 51       | LAST NAME: THOMPSON, THOMPSON, RILEY |
|                               | CITY: WASHINGTON, DENVER, SEATTLE    |
|                               | AGE: 27 33 51                        |

Code example (Python):

```
import pyarrow.parquet as pq
pq.write_table(table, 'example.parquet')
```

## Pickle (Serialized format) - Binary format

- Serialization refers to the process of converting object from the main memory to the disk
- Pickle is a well-known Python module permitting that
- Very useful to store objects, structures or Machine Learning models

Code example:

```
import pickle
class exampleClass:
    a = 35
    b = [1, 2, 3]
object = exampleClass()
outfile = open("example.bin", 'wb')
pickled_object = pickle.dump(object, outfile)
infile = open("example.bin", 'rb')
loaded_pickled_object = pickle.load(infile)
```

## 1 Introduction

## 2 Memory

- Different memories for different usages
- How it works?

## 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

## 4 File formats

## 5 Databases

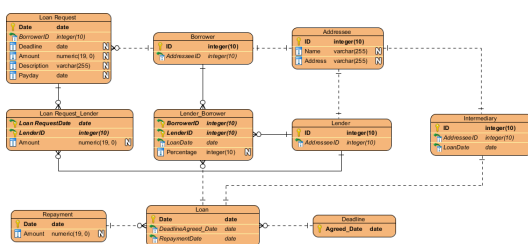
## 6 Data compression

Another popular way to save data is to store it in a **database**.

What is a database?

- Alternative at storing data in the file system
- An elegant way to **centralize all the data in one place**
- Force users and developers to access and store the data in a **unique and well defined way**
- Permit to **structure the data and manage it in a powerful way**
- **Different levels of permission** can be set
  - ▶ Administrator can read, write and delete everything
  - ▶ Users can only access their databases and in a specific way
- Databases usually offer a **set of very powerful functionality to extract and write data**
  - ▶ Accessible through a language (SQL), an API and/or the command line
  - ▶ Can be directly "plugged" to a software

Example of a relationship diagram of a relational database:



Actually, there are two main types of databases:

- 1 Relational Databases (using the SQL language)
  - ▶ MySQL: Relatively easy to use, suitable for most applications
  - ▶ PostgreSQL: Complex to use, but more powerful than MySQL
  - ▶ SQLite: Easy to use
  - ▶ Etc.
- 2 Non-relational Databases called NoSQL
  - ▶ MongoDB
  - ▶ ElasticSearch
  - ▶ Etc.

## Comparisons between Relational and Non-Relational databases:

| Property                      | Relational Databases                                       | Non-Relational Databases                                  |
|-------------------------------|--|---|
| Data model                    | Relational   | Non-relational  |
| Structure                     | Table-based, with columns and rows                         | Document based, key-value pairs, graph, or wide-column    |
| Schema                        | A predefined and strict schema in which every record (row) | A dynamic schema, records don't need to be of same nature |
| Query language                | Structured Query Language (SQL)                            | Varies from database to database                          |
| Scalability                   | Vertical   | Horizontal  |
| Ability to add new properties | Need to alter the schema first (migration: risky!)         | Possible without disturbing anything                      |

## Code example (Python and SQLite - a relational database):

```
import sqlite3
conn = sqlite3.connect('mydatabase.db')
cursor = conn.execute("SELECT id, name, address, salary from COMPANY")
for row in cursor:
    print("NAME = {}".format(row[1]))
    print("SALARY = {}".format(row[2]))
conn.close()
```

## Code example (Python and MongoDB - a non-relational database):

```
from pymongo import MongoClient
client = MongoClient(host="localhost", port=27017)
db = client["mydb"]
collection = db.addresses
address = collection.find_one({"NAME": "Jon"})
collection.insert_one({"NAME": "Paul", "PHONE": "12334"}) #JSON format
client.close()
```

## 1 Introduction

## 2 Memory

- Different memories for different usages
- How it works?

## 3 Data

- Data, types and properties
- Units
- Data types
- Encoding characters
- How to deal with images, audio and video

## 4 File formats

## 5 Databases

## 6 Data compression

## Data Compression, what is it and why we need it in Data Engineering?

- Data Compression: process of encoding information using fewer bits than the original representation
- We need it to:
  - ▶ Save transmission capacity
  - ▶ Save transmission time
  - ▶ Reduce storage occupancy
  - ▶ Reduce computation
- So it's all about saving resources and money...

Basically two main types of compression:

- Lossy compression
- Lossless compression

## Examples of popular algorithms used for compressing data:

### Lossless algorithms:

- Run-length
- Huffman (used with: zip, gzip, png, etc.)
- LZW (Lempel-Ziv-Welch) (used with: gif, tiff)
- LZMA (Lempel-Ziv-Markov chain algorithm) (used with 7zip)
- Etc.

### Lossy algorithms:

- JPEG (images)
- MPEG (videos)
- MP3 (audio)

## What is a lossless compression?

- It means there is **no loss of information** when the data is compressed
- **Original data can be exactly recovered** from the compressed data (exact size before and after)
- Usually used with **applications that cannot tolerate a loss of information**, for example:
  - ▶ word documents
  - ▶ computer programs
  - ▶ company's customer dataset
  - ▶ Etc.
- Principle:
  - ▶ **Redundant data is removed during compression**
  - ▶ **Redundant data is added back during decompression...**

What is a lossy compression?

- Lossy compression results in **lost data and quality from the original representation**
- However, because it removes data from original representation it **usually mean a greater compression**
- Remove just what have to be removed:
  - ▶ JPEG may reduce image's size by more than 80% without noticeable effect
  - ▶ MP3 may reduce down to one tenth the size of the original audio file
- Trade-off: **a greater compression will take up less space, but may alter the quality**

In Data Engineering we are mostly interested in used lossless compression. Let's compare<sup>3</sup>these lossless compression file formats then:

- **gzip**: uses Lempel-Ziv coding (LZ77) (tar czf pack.tar.gz rep/)
- **bzip2**: uses the Burrows-Wheeler block sorting text compression algorithm and Huffman coding (tar cjf pack.tar.bz2 rep/)
- **zip**: ZIP MSDOS (zip -r pack.zip rep/)
- **rar**: proprietary archive file format (rar a pack.rar rep/)
- **lha**: based on Lempel-Ziv-Storer-Szymanski-Algorithm (LZSS) and Huffman coding (lha a pack.lha rep/)
- **lzma**: Lempel-Ziv-Markov chain algorithm (tar -lzma -cf pack.tar.lzma rep/)
- **lzop**: Similar to gzip but favors speed over compression ratio (tar -lzop -cf pack.tar.lzop rep/)

<sup>3</sup>Original study: <https://binfalse.de/2011/04/04/comparison-of-compression/>

Comparing the compression of binaries (total: 176.753.125 Bytes):

| Method | Compressed Size | % of original | Time in s   |
|--------|-----------------|---------------|-------------|
| gzip   | 161.999.804     | 91.65         | 10.18       |
| bzip2  | 161.634.685     | 91.45         | 71.76       |
| zip    | 179.273.428     | 101.43        | 13.51       |
| rar    | 175.085.411     | 99.06         | 156.46      |
| lha    | 180.357.628     | 102.04        | 35.82       |
| lzma   | 157.031.052     | <b>88.84</b>  | 129.22      |
| lzop   | 165.533.609     | 93.65         | <b>4.16</b> |

Comparing plain text compression (total: 40.040.854 Bytes):

| Method | Compressed Size | % of original | Time in s   |
|--------|-----------------|---------------|-------------|
| gzip   | 11.363.931      | 28.38         | 1.88        |
| bzip2  | 9.615.929       | 24.02         | 13.63       |
| zip    | 12.986.153      | 32.43         | 1.6         |
| rar    | 11.942.201      | 29.83         | 8.68        |
| lha    | 13.067.746      | 32.64         | 8.86        |
| lzma   | 8.562.968       | <b>21.39</b>  | 30.21       |
| lzop   | 15.384.624      | 38.42         | <b>0.38</b> |

Comparing the compression of a mix of data (plain+media+binary+pictures+random, total: 355.971.825 Bytes):

| Method | Compressed Size | % of original | Time in s   |
|--------|-----------------|---------------|-------------|
| gzip   | 294.793.255     | 82.81         | 20.43       |
| bzip2  | 290.093.007     | 81.49         | 141.89      |
| zip    | 313.670.439     | 88.12         | 23.78       |
| rar    | 305.083.648     | 85.70         | 246.63      |
| lha    | 315.669.631     | 88.68         | 64.81       |
| lzma   | 283.475.568     | <b>79.63</b>  | 258.05      |
| lzop   | 307.644.076     | 86.42         | <b>7.89</b> |

For the next courses:

- Enough of the theory for now!
- Bring your laptop computer!
- You will need Linux on virtual machine:  
<https://itsfoss.com/install-linux-in-virtualbox/>
- Or installed on a USB stick:  
<https://itsfoss.com/create-live-usb-of-ubuntu-in-windows/>