

Data Interoperability and Semantics

< Part 1. Encoding base data types >

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

Data Interoperability and Semantics

Outline

- < Part 1. Encoding base data types >
 - Part 1.1. Reminders: binary and hexadecimal strings
 - Part 1.2. Endianness
 - Example: MCF88 LoRa temperature, humidity and pressure sensor payload
 - Part 1.3. Computer number formats
 - Part 1.4. Character encoding
 - Part 1.5. Base32 and Base64 encoding
 - Part 1.6. Date and time
 - Part 1.7. XML Schema Datatypes
 - Part 1.8. Codes: countries, languages, ...
 - Part 1.9. Quantities and Units of measure
 - Part 1.10. Colors

ICM – Computer Science Major – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.1. Reminders: binary and hexadecimal strings

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

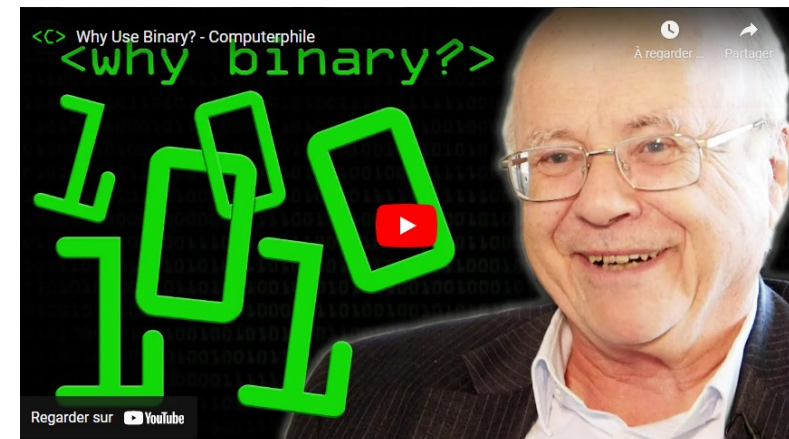
M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

Numbering systems for computers

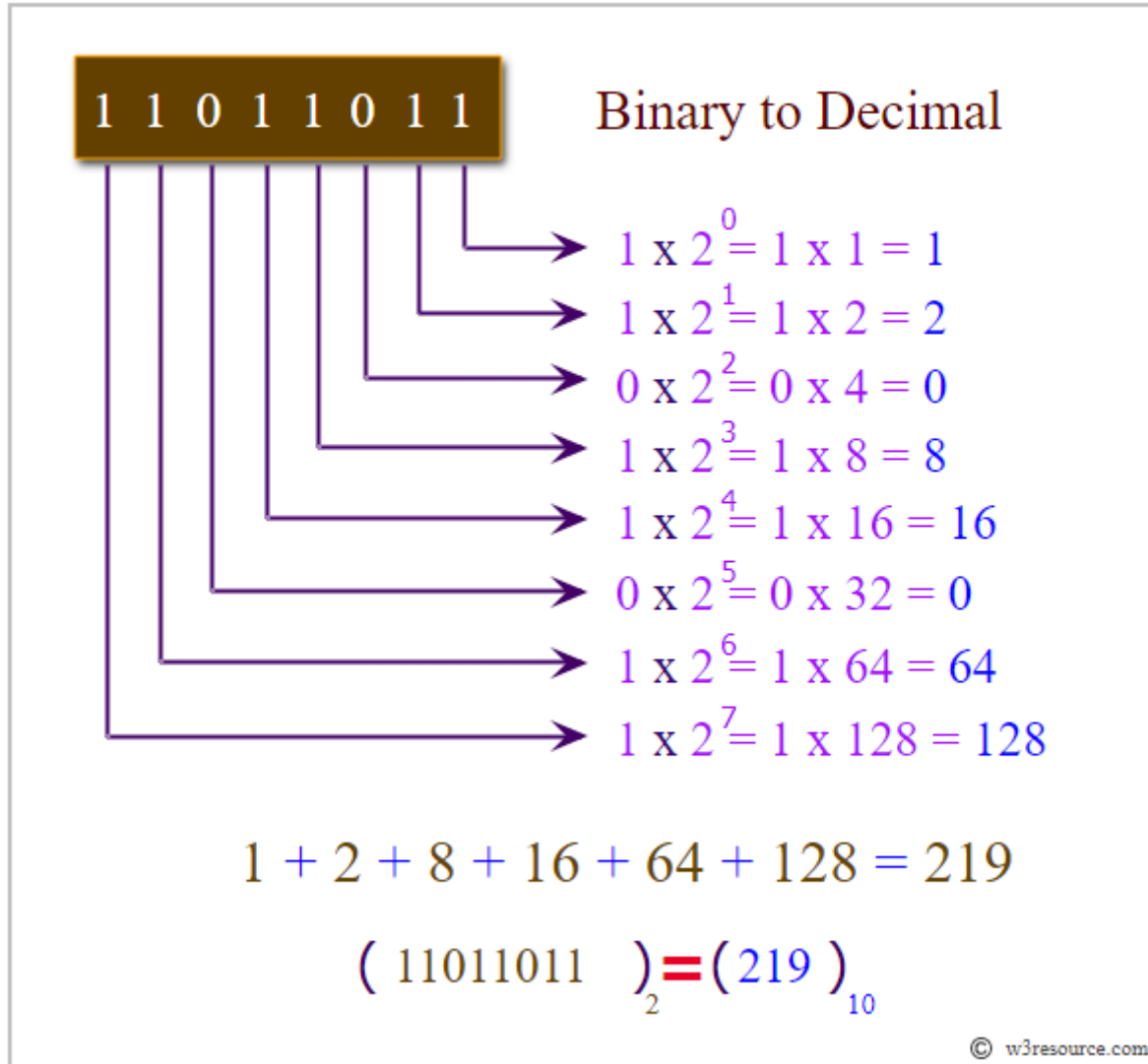
System	Base	Digits	ex python
Binary	2	0,1	0b"01111011"
Octal	8	0,1,2,3,4,5,6,7	0o"173"
Decimal	10	0,1,2,3,4,5,6,7,8,9	123
Hexadecimal	16	0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F	0x"7B"



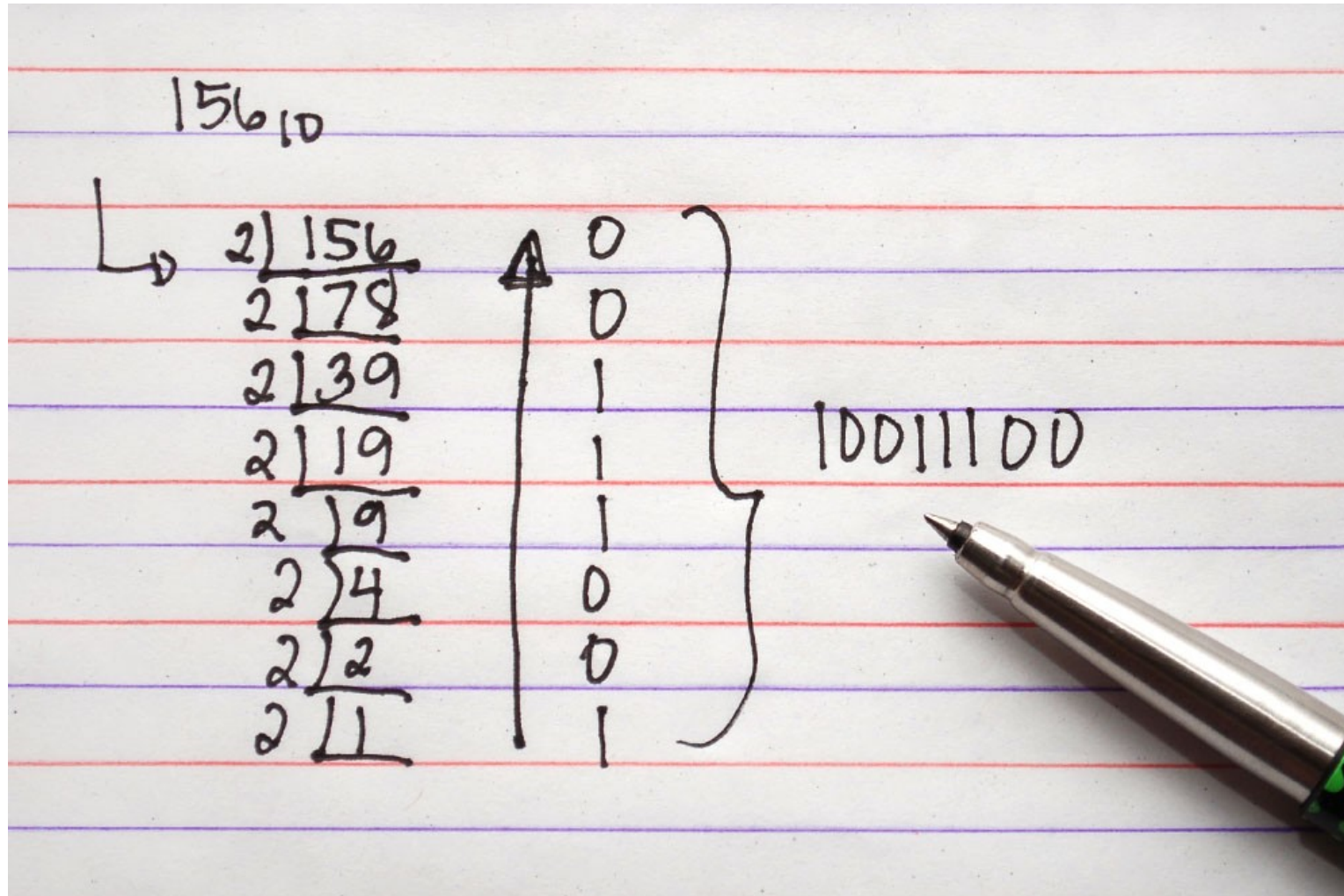
It has been debated a lot at the beginning

<https://www.youtube.com/watch?v=thrx3SBEPt8>

Tips: binary to decimal



Tips: decimal to binary

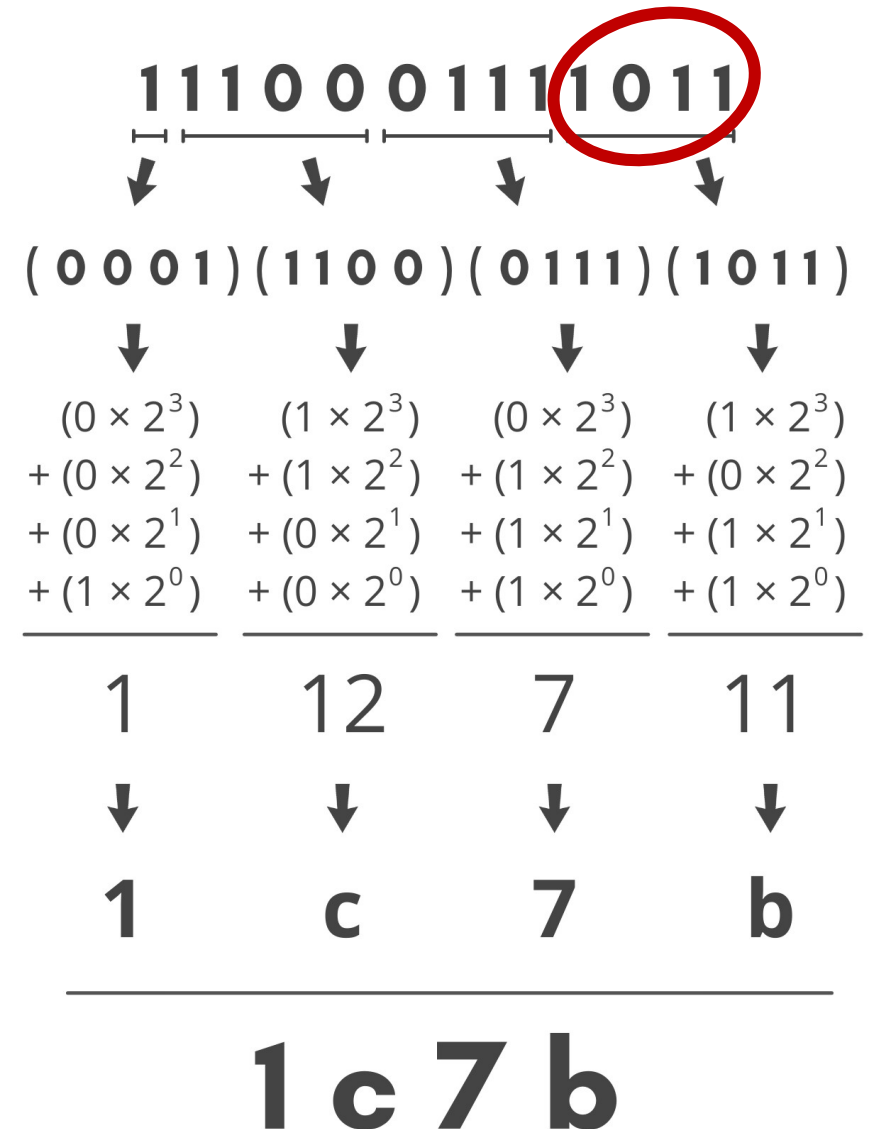


Tips: binary to hexadecimal

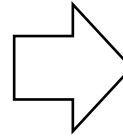
Nibble - In computing, a nibble is a four-bit aggregation, or half an octet. It is also known as half-byte or tetrad. In a networking or telecommunication context, the nibble is often called a semi-octet, quadbit, or quartet.

0000	0	1000	8
0001	1	1001	9
0010	2	1010	A
0011	3	1011	B
0100	4	1100	C
0101	5	1101	D
0110	6	1110	E
0111	7	1111	F

Binary nibble to hexadecimal digit



Programming with binary strings



```
1  #include <stdio.h>
2  int main()
3  {
4      unsigned char a = 0x65;
5      unsigned char b = 0x09;
6
7      // a = 0x65(0110 0101), b = 0x09(0000 1001)
8      printf("a = %#.2X, b = %#.2x\n", a, b);
9
10     // & bitwise AND operator
11     // the result is 0x01(0000 0001)
12     printf("a&b = %#.2X\n", a&b);
13
14     // | bitwise OR operator
15     // the result is 0x6D(0110 1101)
16     printf("a|b = %#.2X\n", a|b);
17
18     // ^ bitwise exclusive OR operator
19     // the result is 0x6C(0110 1100)
20     printf("a^b = %#.2X\n", a^b);
21
22     // << Left shift operator
23     // the result is 0xCA(1100 1010)
24     printf("a << 1 = %#.2X\n", a << 1);
25
26     // >> Right shift operator
27     // the result is 0x32(0011 0010)
28     printf("a >> 1 = %#.2X\n", a >> 1);
29
30     // ~ bitwise One's Complement operator
31     // the result is 0X9A (1001 1010)
32     // or more precisely 0XFFFFFF9A (~ promotes to int)
33     printf("~a = %#.2X\n", ~a);
34
35     // Get 3 bits starting at position 2 (start index 0)
36     // the result is 1(0001)
37     printf("(a >> 2) & 0x7 = %#.1X\n", (a>>2)&7);
38
39     return 0;
40 }
```

```
$ gcc main.c
$ ./a.out
a = 0X65, b = 0x09
a&b = 0X01
a|b = 0X6D
a^b = 0X6C
a << 1 = 0XCA
a >> 1 = 0X32
~a = 0XFFFFFF9A
(a >> 2) & 0x7 = 0X1
```


Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.2. Endianness

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

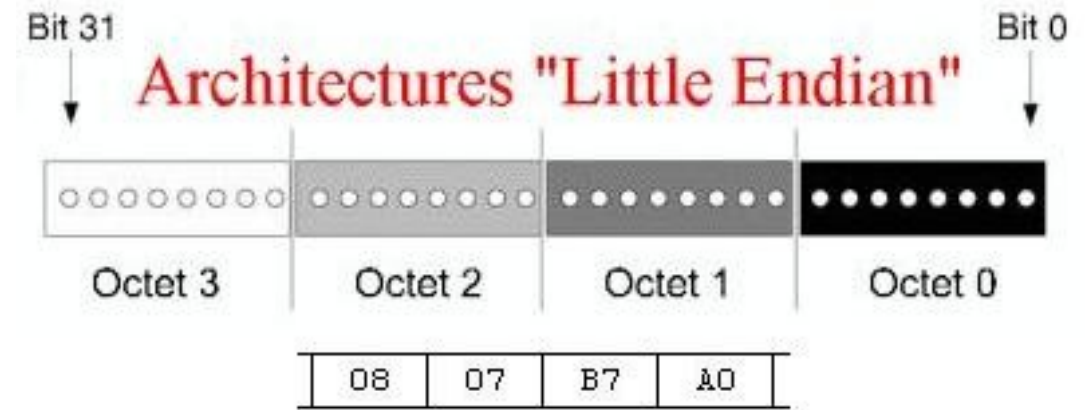
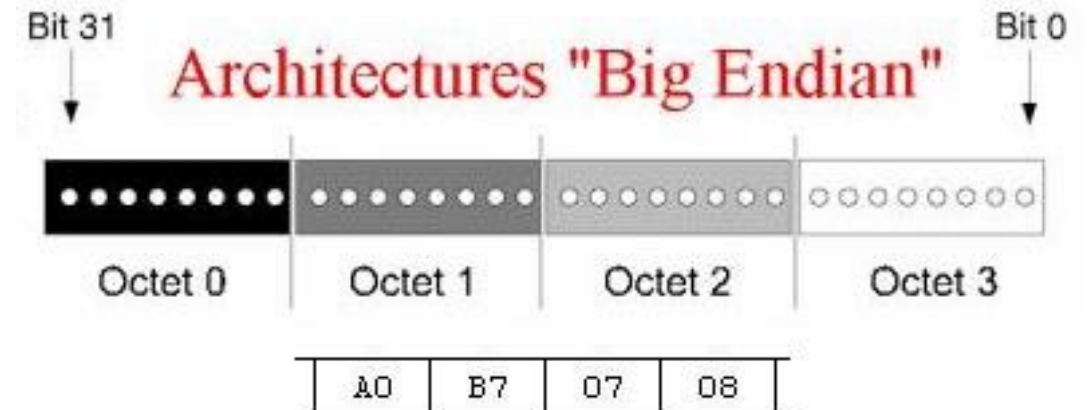
Endianness

*In computing, endianness is the order or sequence of bytes of a word of digital data in computer memory. Endianness is primarily expressed as **big-endian (BE)** or **little-endian (LE)**.*

Acronyms

- LSB - Least significant byte
- MSB - Most significant byte

0XA0B70708



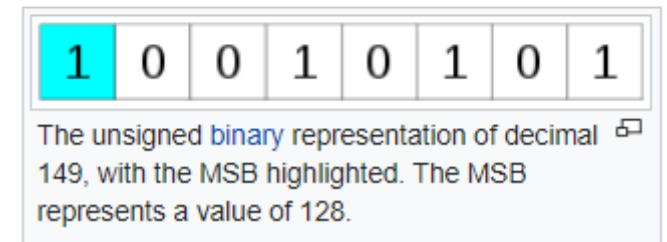
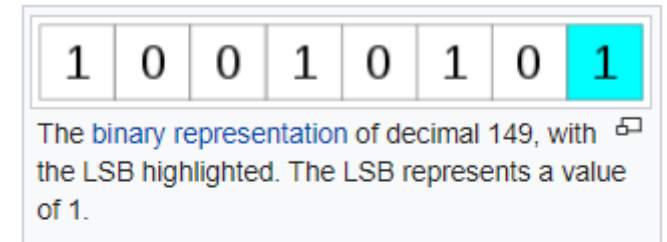
<https://www.sqlpac.com/fr/documents/sybase-ase-12.5.3-dump-load-cross-platforms.html>

Bit endianness or bit-level endianness

Bit endianness or bit-level endianness refers to the transmission order of bits over a serial medium

See course Programming Connected Devices:

- Least significant bit first: used in RS-232, Ethernet, USB...
- Most significant bit first: used in I²C



https://en.wikipedia.org/wiki/Bit_numbering

Example: WS2812B color leds



WS2812B

Intelligent control LED
integrated light source

Composition of 24bit data:

G7	G6	G5	G4	G3	G2	G1	G0	R7	R6	R5	R4	R3	R2	R1	R0	B7	B6	B5	B4	B3	B2	B1	B0
----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

Note: Follow the order of GRB to sent data and the high bit sent at first.



Data Interoperability and Semantics

Part 1. Encoding base data types

Example: MCF88 LoRa temperature, humidity and pressure sensor payload



ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>



Author: Colognato Stefano

Created: 29/09/2016

Modified: 30/11/2018

MCF88 DATA FRAME FORMAT 1.17

1.2 TEMPERATURE/PRESSURE/HUMIDITY

[HOME](#)

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
Data	10 byte	XX XX	Measure 1, refer to Note1
	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
Batt %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

Note1:

The 10 bytes for each measurement are divided as follows:

- 4 bytes are for the date and time. The MSB (most significant byte) is on the right so they must be read from the right. The 4 byte in reverse order are as follows:
 - 7 bit for the offset of the year, starting from the year 2000
 - 4 bit per month
 - 5 bit for day of the month
 - 5 bits for hour
 - 6 bits for minutes
 - 5 bits for half the seconds. The seconds range is from 0 to 31, so the result should be multiplied by 2 to find the actual seconds of the measurement.
- 2 bytes for temperature. The temperature is represented by a signed integer with the least significant byte first. The temperature is expressed in hundreds of a °C degree.
- 1 byte for humidity. Relative humidity is an unsigned integer corresponding to twice the percentage of humidity.
- 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in Pascal.

Example

Sample payload:

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Example: MCF88 LoRa sensors





Author: Colognato Stefano

Created: 29/09/2016

Modified: 30/11/2018

MCF88 DATA FRAME FORMAT 1.17

1.2 TEMPERATURE/PRESSURE/HUMIDITY

[HOME](#)

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
Data	10 byte	XX XX	Measure 1, refer to Note1
	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
Batt %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

Note1:

The 10 bytes for each measurement are divided as follows:

- 4 bytes are for the date and time. The MSB (most significant byte) is on the right so they must be read from the right. The 4 byte in reverse order are as follows:
 - 7 bit for the offset of the year, starting from the year 2000
 - 4 bit per month
 - 5 bit for day of the month
 - 5 bits for hour
 - 6 bits for minutes
 - 5 bits for half the seconds. The seconds range is from 0 to 31, so the result should be multiplied by 2 to find the actual seconds of the measurement.
- 2 bytes for temperature. The temperature is represented by a signed integer with the least significant byte first. The temperature is expressed in hundreds of a °C degree.
- 1 byte for humidity. Relative humidity is an unsigned integer corresponding to twice the percentage of humidity.
- 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in Pascal.

Example

Sample payload:

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Remove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 measurements.

The 3 measurements will be:

- dc7e3721b40a47608801
- dd7e3721b10a43608801
- e07e3721b20a425d8801



1.2 TEMPERATURE/PRESSURE/HUMIDITY

[HOME](#)

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
Data	10 byte	XX XX	Measure 1, refer to Note1
	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
Batt %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

Note1:

The 10 bytes for each measurement are divided as follows:

- 4 bytes are for the date and time. The MSB (most significant byte) is on the right so they must be read from the right. The 4 byte in reverse order are as follows:
 - 7 bit for the offset of the year, starting from the year 2000
 - 4 bit per month
 - 5 bit for day of the month
 - 5 bits for hour
 - 6 bits for minutes
 - 5 bits for half the seconds. The seconds range is from 0 to 31, so the result should be multiplied by 2 to find the actual seconds of the measurement.
- 2 bytes for temperature. The temperature is represented by a signed integer with the least significant byte first. The temperature is expressed in hundreds of a °C degree.
- 1 byte for humidity. Relative humidity is an unsigned integer corresponding to twice the percentage of humidity.
- 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in Pascal.

Example

Sample payload:

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Remove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 measurements.

The 3 measurements will be:

- dc7e3721b40a47608801
- dd7e3721b10a43608801
- e07e3721b20a425d8801

Decipher the first measurement dividing it by groups and applying the necessary transformations:

- Measurement date: dc 7e 37 21
 - Byte swapping, result: 21 37 7e dc



1.2 TEMPERATURE/PRESSURE/HUMIDITY

[HOME](#)

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
Data	10 byte	XX XX	Measure 1, refer to Note1
	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
Batt %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

Note1:

The 10 bytes for each measurement are divided as follows:

- 4 bytes are for the date and time. The MSB (most significant byte) is on the right so they must be read from the right. The 4 byte in reverse order are as follows:
 - 7 bit for the offset of the year, starting from the year 2000
 - 4 bit per month
 - 5 bit for day of the month
 - 5 bits for hour
 - 6 bits for minutes
 - 5 bits for half the seconds. The seconds range is from 0 to 31, so the result should be multiplied by 2 to find the actual seconds of the measurement.
- 2 bytes for temperature. The temperature is represented by a signed integer with the least significant byte first. The temperature is expressed in hundreds of a °C degree.
- 1 byte for humidity. Relative humidity is an unsigned integer corresponding to twice the percentage of humidity.
- 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in Pascal.

Example

Sample payload:

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Remove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 measurements.

The 3 measurements will be:

- dc7e3721b40a47608801
- dd7e3721b10a43608801
- e07e3721b20a425d8801

Decipher the first measurement dividing it by groups and applying the necessary transformations:

- Measurement date: dc 7e 37 21
 - Byte swapping, result: 21 37 7e dc
 - The result in bits will be: 00100001 00110111 01111110 11011100
 - The bits are divided as explained above
 - Year: 0010000
 - ✦ Result: 16
 - $2000+16 = 2016$
 - Month: 1001
 - ✦ Result: 9
 - Day: 10111
 - ✦ Result: 23
 - Hour: 01111
 - ✦ Result: 15
 - Minutes: 110110
 - ✦ Result: 54
 - Seconds: 11100
 - ✦ Result: 28
 - $28*2 = 56$
- The date of the measurement will be: 23/09/2016 15:54:56.

1.2 TEMPERATURE/PRESSURE/HUMIDITY

[HOME](#)

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
Data	10 byte	XX XX	Measure 1, refer to Note1
	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
Batt %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

Note1:

The 10 bytes for each measurement are divided as follows:

- 4 bytes are for the date and time. The MSB (most significant byte) is on the right so they must be read from the right. The 4 byte in reverse order are as follows:
 - 7 bit for the offset of the year, starting from the year 2000
 - 4 bit per month
 - 5 bit for day of the month
 - 5 bits for hour
 - 6 bits for minutes
 - 5 bits for half the seconds. The seconds range is from 0 to 31, so the result should be multiplied by 2 to find the actual seconds of the measurement.
- 2 bytes for temperature. The temperature is represented by a signed integer with the least significant byte first. The temperature is expressed in hundreds of a °C degree.
- 1 byte for humidity. Relative humidity is an unsigned integer corresponding to twice the percentage of humidity.
- 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in Pascal.

Example

Sample payload:

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Remove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 measurements.

The 3 measurements will be:

- dc7e3721b40a47608801
- dd7e3721b10a43608801
- e07e3721b20a425d8801

Decipher the first measurement dividing it by groups and applying the necessary transformations:

- Measurement date: dc 7e 37 21
 - Byte swapping, result: 21 37 7e dc
 - The result in bits will be: 00100001 00110111 01111110 11011100
 - The bits are divided as explained above
 - Year: 0010000
 - ♦ Result: 16
 - $2000+16 = 2016$
 - Month: 1001
 - ♦ Result: 9
 - Day: 10111
 - ♦ Result: 23
 - Hour: 01111
 - ♦ Result: 15
 - Minutes: 110110
 - ♦ Result: 54
 - Seconds: 11100
 - ♦ Result: 28
 - $28*2 = 56$
 - The date of the measurement will be: 23/09/2016 15:54:56.
 - Temperature: b40a
 - Byte swapping, result: 0ab4
 - The result (with sign) will be +2740 with two decimal places, then + 27.40 °C.
 - Humidity: 47
 - In decimal is 71, the humidity is $71/2 = 35.5\%$ rH.
 - Pressure: 608801
 - Byte swapping, result: 018860
 - In decimal, the result is 100448, with two decimal places the pressure is 1004.48 hPa.

1.2 TEMPERATURE/PRESSURE/HUMIDITY

[HOME](#)

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
Data	10 byte	XX XX	Measure 1, refer to Note1
	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
Batt %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

Note1:

The 10 bytes for each measurement are divided as follows:

- 4 bytes are for the date and time. The MSB (most significant byte) is on the right so they must be read from the right. The 4 byte in reverse order are as follows:
 - 7 bit for the offset of the year, starting from the year 2000
 - 4 bit per month
 - 5 bit for day of the month
 - 5 bits for hour
 - 6 bits for minutes
 - 5 bits for half the seconds. The seconds range is from 0 to 31, so the result should be multiplied by 2 to find the actual seconds of the measurement.
- 2 bytes for temperature. The temperature is represented by a signed integer with the least significant byte first. The temperature is expressed in hundreds of a °C degree.
- 1 byte for humidity. Relative humidity is an unsigned integer corresponding to twice the percentage of humidity.
- 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in Pascal.

Example

Sample payload:

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Remove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 measurements.

The 3 measurements will be:

- dc7e3721b40a47608801
- dd7e3721b10a43608801
- e07e3721b20a425d8801

Decipher the first measurement dividing it by groups and applying the necessary transformations:

- Measurement date: dc 7e 37 21
 - Byte swapping, result: 21 37 7e dc
 - The result in bits will be: 00100001 00110111 01111110 11011100
 - The bits are divided as explained above
 - Year: 0010000
 - ♦ Result: 16
 - $2000+16 = 2016$
 - Month: 1001
 - ♦ Result: 9
 - Day: 10111
 - ♦ Result: 23
 - Hour: 01111
 - ♦ Result: 15
 - Minutes: 110110
 - ♦ Result: 54
 - Seconds: 11100
 - ♦ Result: 28
 - $28*2 = 56$
 - The date of the measurement will be: 23/09/2016 15:54:56.
 - Temperature: b40a
 - Byte swapping, result: 0ab4
 - The result (with sign) will be +2740 with two decimal places, then + 27.40 °C.
 - Humidity: 47
 - In decimal is 71, the humidity is $71/2 = 35.5\%$ rH.
 - Pressure: 608801
 - Byte swapping, result: 018860
 - In decimal, the result is 100448, with two decimal places the pressure is 1004.48 hPa.

Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.3. Computer number formats

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

C number data types

- char (8 bits) — [0, 255]
- short/int (16 bits) - [-32,767, +32,767] or [0, 65,535]
- long (32 bits) - [-2,147,483,647, +2,147,483,647] or [0, 4,294,967,295]
- long long (64 bits) [-9,223,372,036,854,775,807, +9,223,372,036,854,775,807] or <<positive>>
- float - IEEE 754 single-precision binary floating-point format (32 bits)
- double - IEEE 754 double-precision binary floating-point format (64 bits)

Integer encoding

Unsigned

$$B2U(X) = \sum_{i=0}^{w-1} x_i \cdot 2^i$$

```
short int x = 15213;  
short int y = -15213;
```

- C short 2 bytes long

	Decimal	Hex	Binary
x	15213	3B 6D	00111011 01101101
y	-15213	C4 93	11000100 10010011

Two's Complement

$$B2T(X) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$$

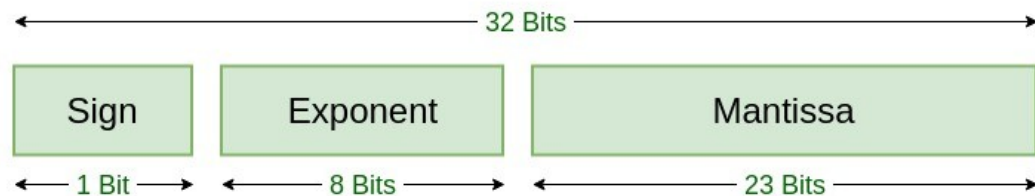
Sign
Bit

Sign Bit

- For 2's complement, most significant bit indicates sign
 - 0 for nonnegative
 - 1 for negative

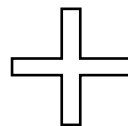
IEEE 754 single/double encoding

<https://www.geeksforgeeks.org/ieee-standard-754-floating-point-numbers/>



Single Precision
IEEE 754 Floating-Point Standard

$$-1^s \times 2^{(\text{exp}-127)} \times 1.\text{frac}$$



exceptional cases:

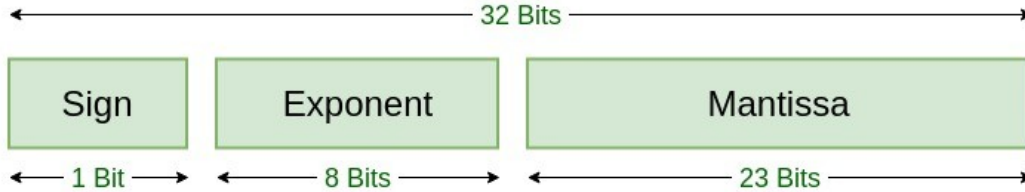
- If $E = 255$ and F is nonzero, then $x = \text{NaN}$ (“Not a number”).
- If $E = 255$, F is zero, and S is 1, then $x = -\text{Infinity}$.
- If $E = 255$, F is zero, and S is 0, then $x = +\text{Infinity}$.
- If $0 < E < 255$, then $x = (-1)^s \times (1.F) \times 2^{E-127}$, where $1.F$ represents the binary number created by prefixing F with an implicit leading 1 and a binary point.
- If $E = 0$ and F is nonzero, then $x = (-1)^s \times (0.F) \times 2^{-126}$. This is an “unnormalized” value.
- If $E = 0$, F is zero, and S is 1, then $x = -0$.
- If $E = 0$, F is zero, and S is 0, then $x = 0$.

Simple examples of conversions

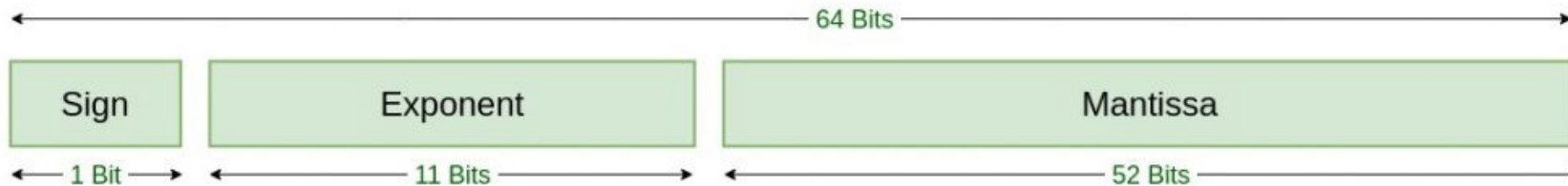
$$\begin{aligned} 0\ 10000000\ 000000000000000000000000 &= (-1)^0 \times (1.0_2) \times 2^{128-127} = 2.0 \\ 0\ 10000001\ 101000000000000000000000 &= (-1)^0 \times (1.101_2) \times 2^{129-127} = 6.5 \\ 1\ 10000001\ 101000000000000000000000 &= (-1)^1 \times (1.101_2) \times 2^{129-127} = -6.5. \end{aligned}$$

IEEE 754 single/double encoding

<https://www.geeksforgeeks.org/ieee-standard-754-floating-point-numbers/>



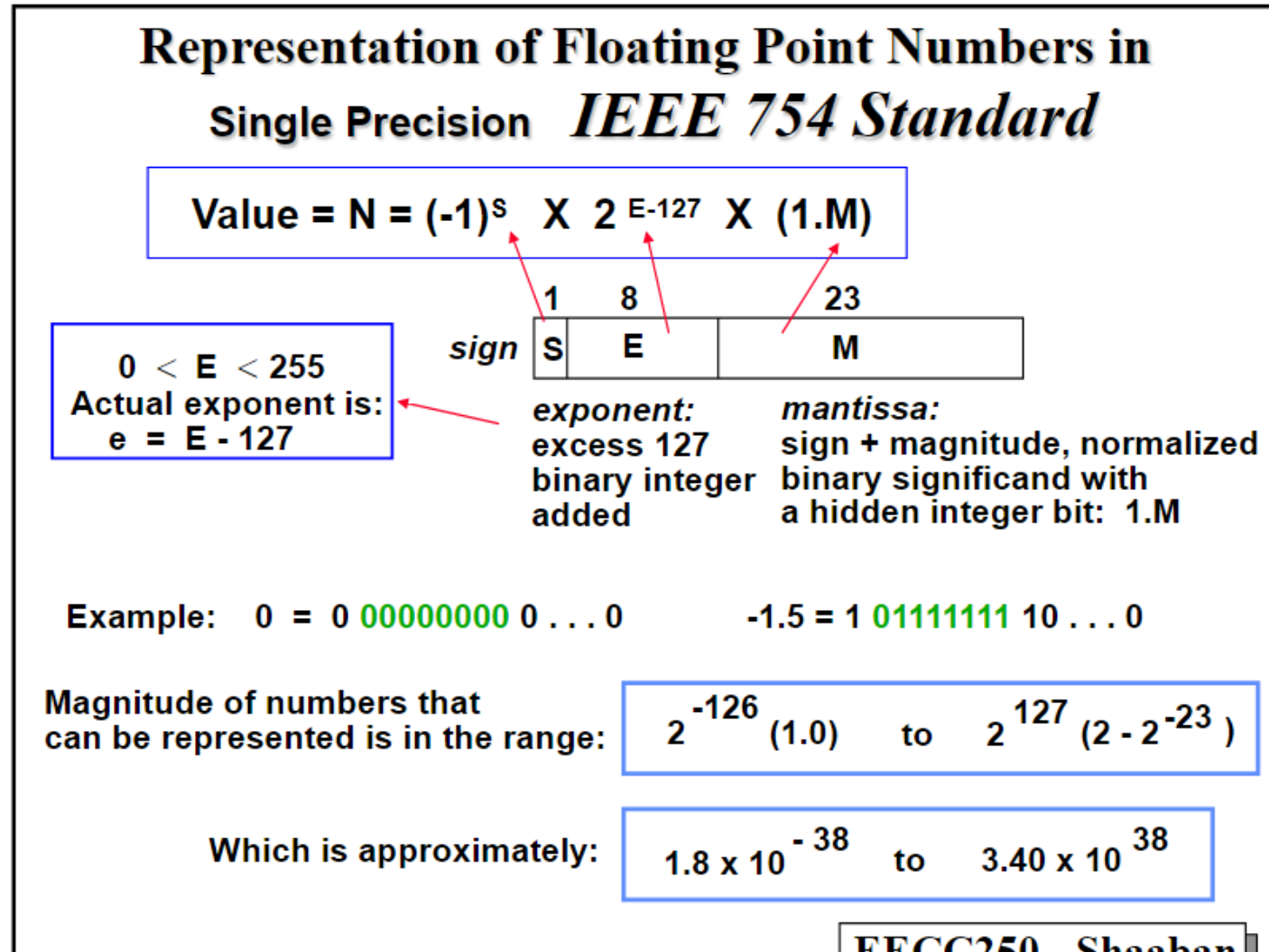
Single Precision
IEEE 754 Floating-Point Standard



Double Precision
IEEE 754 Floating-Point Standard

<https://www.geeksforgeeks.org/ieee-standard-754-floating-point-numbers/>

IEEE 754 single/double/quadruple encoding



Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.4. Character encoding

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

Character encoding

In computing, data storage, and data transmission, character encoding is used to represent a repertoire of characters by some kind of encoding system that assigns a number to each character for digital representation

— https://en.wikipedia.org/wiki/Character_encoding

Common character encodings [\[edit \]](#)

- [ISO 646](#)
 - [ASCII](#)
- [EBCDIC](#)
- [ISO 8859](#):
 - [ISO 8859-1](#) Western Europe
 - [ISO 8859-2](#) Western and Central Europe
 - [ISO 8859-3](#) Western Europe and South European (Turkish, Maltese plus Esperanto)
 - [ISO 8859-4](#) Western Europe and Baltic countries (Lithuania, Estonia, Latvia and Lapp)
 - [ISO 8859-5](#) Cyrillic alphabet
 - [ISO 8859-6](#) Arabic
 - [ISO 8859-7](#) Greek
 - [ISO 8859-8](#) Hebrew
 - [ISO 8859-9](#) Western Europe with amended Turkish character set
 - [ISO 8859-10](#) Western Europe with rationalised character set for Nordic languages, including complete Icelandic set
 - [ISO 8859-11](#) Thai
 - [ISO 8859-13](#) Baltic languages plus Polish
 - [ISO 8859-14](#) Celtic languages (Irish Gaelic, Scottish, Welsh)
 - [ISO 8859-15](#) Added the Euro sign and other rationalisations to ISO 8859-1
 - [ISO 8859-16](#) Central, Eastern and Southern European languages (Albanian, Bosnian, Croatian, Hungarian, Polish, Romanian, Serbian and Slovenian, but also French, German, Italian and Irish Gaelic)
- [CP437](#), [CP720](#), [CP737](#), [CP850](#), [CP852](#), [CP855](#), [CP857](#), [CP858](#), [CP860](#), [CP861](#), [CP862](#), [CP863](#), [CP865](#), [CP866](#), [CP869](#), [CP872](#)
- [MS-Windows character sets](#):
 - [Windows-1250](#) for Central European languages that use Latin script, (Polish, Czech, Slovak, Hungarian, Slovene, Serbian, Croatian, Bosnian, Romanian and Albanian)
 - [Windows-1251](#) for Cyrillic alphabets
 - [Windows-1252](#) for Western languages
 - [Windows-1253](#) for Greek
 - [Windows-1254](#) for Turkish
 - [Windows-1255](#) for Hebrew
 - [Windows-1256](#) for Arabic
 - [Windows-1257](#) for Baltic languages
 - [Windows-1258](#) for Vietnamese
- [Mac OS Roman](#)
- [KOI8-R](#), [KOI8-U](#), [KOI7](#)
- [MIK](#)
- [ISCII](#)
- [TSCII](#)
- [VISCII](#)
- [JIS X 0208](#) is a widely deployed standard for Japanese character encoding that has several encoding forms.
 - [Shift JIS](#) ([Microsoft Code page 932](#) is a dialect of Shift_JIS)
 - [EUC-JP](#)
 - [ISO-2022-JP](#)
- [JIS X 0213](#) is an extended version of JIS X 0208.
 - [Shift_JIS-2004](#)
 - [EUC-JIS-2004](#)
 - [ISO-2022-JP-2004](#)
- Chinese [Guobiao](#)
 - [GB 2312](#)
 - [GBK](#) ([Microsoft Code page 936](#))
 - [GB 18030](#)
- Taiwan [Big5](#) (a more famous variant is [Microsoft Code page 950](#))
 - [Hong Kong HKSCS](#)
- Korean
 - [KS X 1001](#) is a Korean double-byte character encoding standard
 - [EUC-KR](#)
 - [ISO-2022-KR](#)
- [Unicode](#) (and subsets thereof, such as the 16-bit 'Basic Multilingual Plane')
 - [UTF-8](#)
 - [UTF-16](#)
 - [UTF-32](#)
- [ANSEL](#) or [ISO/IEC 6937](#)

Let's focus on the main standards

- ASCII
- UTF-8

Common character encodings [\[edit \]](#)

- [ISO 646](#)
 - [ASCII](#)
- [EBCDIC](#)
- [ISO 8859](#):
 - [ISO 8859-1](#) Western Europe
 - [ISO 8859-2](#) Western and Central Europe
 - [ISO 8859-3](#) Western Europe and South European (Turkish, Maltese plus Esperanto)
 - [ISO 8859-4](#) Western Europe and Baltic countries (Lithuania, Estonia, Latvia and Lapp)
 - [ISO 8859-5](#) Cyrillic alphabet
 - [ISO 8859-6](#) Arabic
 - [ISO 8859-7](#) Greek
 - [ISO 8859-8](#) Hebrew
 - [ISO 8859-9](#) Western Europe with amended Turkish character set
 - [ISO 8859-10](#) Western Europe with rationalised character set for Nordic languages, including complete Icelandic set
 - [ISO 8859-11](#) Thai
 - [ISO 8859-13](#) Baltic languages plus Polish
 - [ISO 8859-14](#) Celtic languages (Irish Gaelic, Scottish, Welsh)
 - [ISO 8859-15](#) Added the Euro sign and other rationalisations to ISO 8859-1
 - [ISO 8859-16](#) Central, Eastern and Southern European languages (Albanian, Bosnian, Croatian, Hungarian, Polish, Romanian, Serbian and Slovenian, but also French, German, Italian and Irish Gaelic)
- [CP437](#), [CP720](#), [CP737](#), [CP850](#), [CP852](#), [CP855](#), [CP857](#), [CP858](#), [CP860](#), [CP861](#), [CP862](#), [CP863](#), [CP865](#), [CP866](#), [CP869](#), [CP872](#)
- [MS-Windows character sets](#):
 - [Windows-1250](#) for Central European languages that use Latin script, (Polish, Czech, Slovak, Hungarian, Slovene, Serbian, Croatian, Bosnian, Romanian and Albanian)
 - [Windows-1251](#) for Cyrillic alphabets
 - [Windows-1252](#) for Western languages
 - [Windows-1253](#) for Greek
 - [Windows-1254](#) for Turkish
 - [Windows-1255](#) for Hebrew
 - [Windows-1256](#) for Arabic
 - [Windows-1257](#) for Baltic languages
 - [Windows-1258](#) for Vietnamese
- [Mac OS Roman](#)
- [KOI8-R](#), [KOI8-U](#), [KOI7](#)
- [MIK](#)
- [ISCII](#)
- [TSCII](#)
- [VISCII](#)
- [JIS X 0208](#) is a widely deployed standard for Japanese character encoding that has several encoding forms.
 - [Shift JIS](#) ([Microsoft Code page 932](#) is a dialect of [Shift_JIS](#))
 - [EUC-JP](#)
 - [ISO-2022-JP](#)
- [JIS X 0213](#) is an extended version of [JIS X 0208](#).
 - [Shift_JIS-2004](#)
 - [EUC-JIS-2004](#)
 - [ISO-2022-JP-2004](#)
- [Chinese Guobiao](#)
 - [GB 2312](#)
 - [GBK](#) ([Microsoft Code page 936](#))
 - [GB 18030](#)
- [Taiwan Big5](#) (a more famous variant is [Microsoft Code page 950](#))
 - [Hong Kong HKSCS](#)
- [Korean](#)
 - [KS X 1001](#) is a Korean double-byte character encoding standard
 - [EUC-KR](#)
 - [ISO-2022-KR](#)
- [Unicode](#) (and subsets thereof, such as the 16-bit 'Basic Multilingual Plane')
 - [UTF-8](#)
 - [UTF-16](#)
 - [UTF-32](#)
- [ANSEL](#) or [ISO/IEC 6937](#)

You will be given this document in appendix of the written exam !

ASCII (7 bits)

Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	@	96	60	`
1	1	[START OF HEADING]	33	21	!	65	41	A	97	61	a
2	2	[START OF TEXT]	34	22	"	66	42	B	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	C	99	63	c
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	e
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	'	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	H	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49	I	105	69	i
10	A	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	B	[VERTICAL TAB]	43	2B	+	75	4B	K	107	6B	k
12	C	[FORM FEED]	44	2C	,	76	4C	L	108	6C	l
13	D	[CARRIAGE RETURN]	45	2D	-	77	4D	M	109	6D	m
14	E	[SHIFT OUT]	46	2E	.	78	4E	N	110	6E	n
15	F	[SHIFT IN]	47	2F	/	79	4F	O	111	6F	o
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	P	112	70	p
17	11	[DEVICE CONTROL 1]	49	31	1	81	51	Q	113	71	q
18	12	[DEVICE CONTROL 2]	50	32	2	82	52	R	114	72	r
19	13	[DEVICE CONTROL 3]	51	33	3	83	53	S	115	73	s
20	14	[DEVICE CONTROL 4]	52	34	4	84	54	T	116	74	t
21	15	[NEGATIVE ACKNOWLEDGE]	53	35	5	85	55	U	117	75	u
22	16	[SYNCHRONOUS IDLE]	54	36	6	86	56	V	118	76	v
23	17	[ENG OF TRANS. BLOCK]	55	37	7	87	57	W	119	77	w
24	18	[CANCEL]	56	38	8	88	58	X	120	78	x
25	19	[END OF MEDIUM]	57	39	9	89	59	Y	121	79	y
26	1A	[SUBSTITUTE]	58	3A	:	90	5A	Z	122	7A	z
27	1B	[ESCAPE]	59	3B	;	91	5B	[123	7B	{
28	1C	[FILE SEPARATOR]	60	3C	<	92	5C	\	124	7C	
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D]	125	7D	}
30	1E	[RECORD SEPARATOR]	62	3E	>	94	5E	^	126	7E	~
31	1F	[UNIT SEPARATOR]	63	3F	?	95	5F	_	127	7F	[DEL]



Adopt a Character



Emoji



Basic Info



News

Events

Connect



Membership



Press

Search ...

ᲀ

U+0D05

Ó

U+00D3

♪

U+2669

.

U+0F0B

”

U+201D

ב

U+05D1

😊

U+1F600

ᲀ

U+0CA4

ᲀ

U+FF4D

;

U+FF1B

↩

U+21A9

ᲀ

U+0A1D

ᲀ

U+15E2

~~~~

U+FE4F

ᲀ

U+0CB0

ᲀ

U+067B

e

U+0B67

◐

U+25D1

🌺

U+1F490

Θ

U+04E9

# Unicode

- information technology standard for the consistent encoding, representation, and handling of text expressed in most of the world's writing systems
- 144,762 characters covering 159 modern and historic scripts, as well as symbols, emoji, and non-visual control and formatting codes.

ᲀ

U+0EA7

›

U+203A

✖

U+203B

ᲀ

U+134C

,

U+2019

×

U+00D7

†

U+2020

ᲀ

U+0920

😘

U+1F618

ᲀ

U+0F4F

ᲀ

U+0296

ᲀ

U+1F9E2

♥

U+1F49A

ᲀ

U+13EA

💰

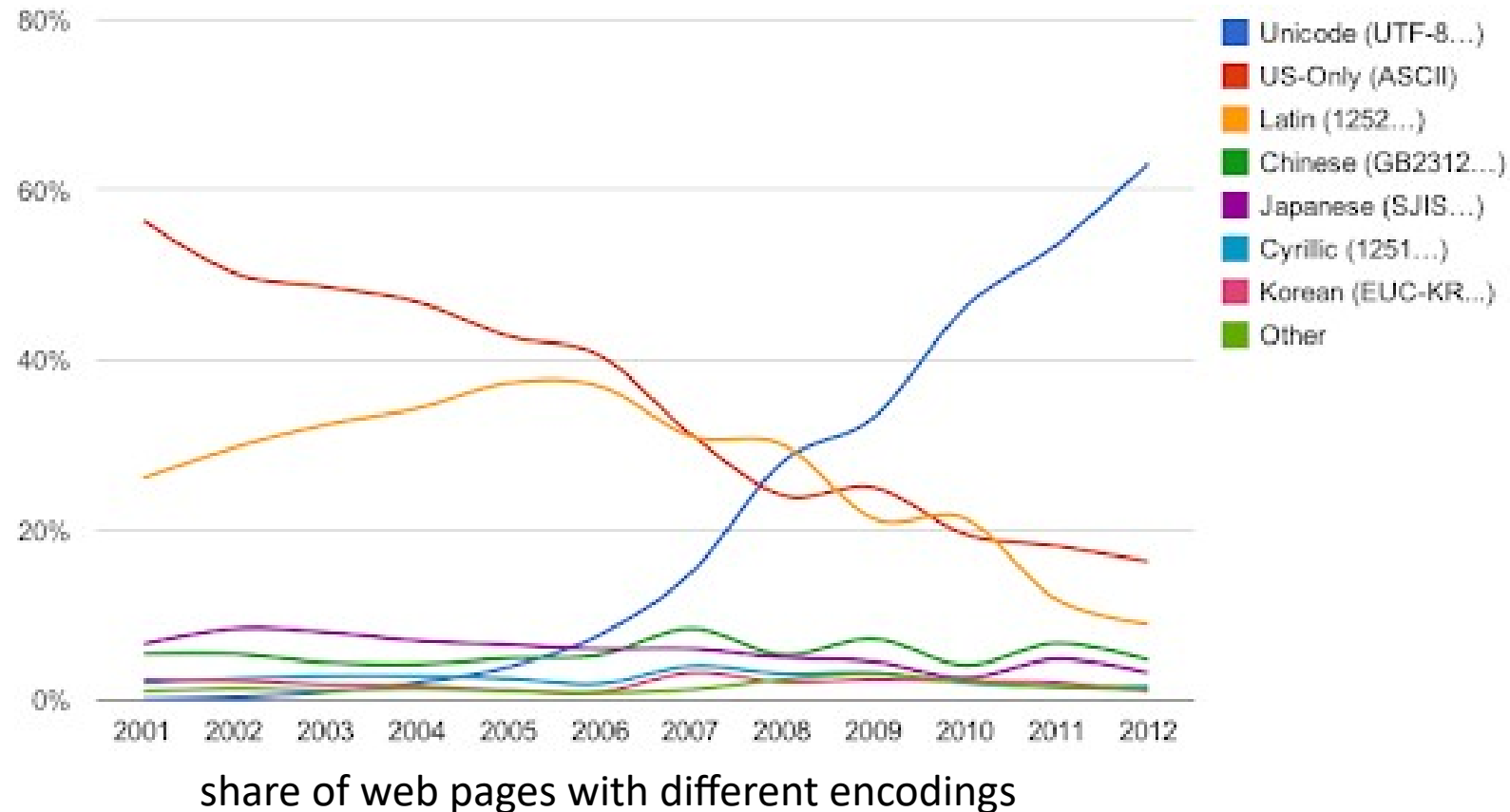
U+1F911

%

U+2030

# UTF-8

- Implementation of Unicode on 1-4 bytes (as little as needed)



# Inserting characters

## Inserting Characters

There are many ways character é ("e" with an acute accent, character code 233 (decimal) in Latin-1 and Unicode), can be inserted into a document:

- On Windows, I hold down the Alt key and type 0233 on the numeric keyboard and release the Alt key. I could use the charmap program, too. Or I could copy and paste it (e.g., é). But entering the code directly is risky because, if the character encoding changes, e.g., from Latin-1 to UTF-8, then the meaning of code 233 changes.
- In an HTML document, I can enter these magical incantations, which are displayed correctly regardless of encoding:
  - `&#233;` (decimal)  $\Rightarrow$  é
  - `&#xE9;` (hex)  $\Rightarrow$  é
  - `&eacute;` (mnemonic)  $\Rightarrow$  é

Note: HTML/XHTML validation programs might not be acquainted with these and complain.

- In Microsoft Word, I type an accent code followed by the accented letter. On Windows, Ctrl+quote, then 'e'. On Mac, Option+quote, then 'e'. Accent codes include: grave=backquote, acute=quote, circumflex=hat, colon=umlaut, comma=cedilla, tilde=tilde, slash=slash, and perhaps others.

# Encoding errors

## What Could Possibly Go Wrong?

If é is UTF-8 encoded, but displayed without decoding, it looks like this:

Ã©

The first 128 characters in the Latin-1 character set (same as ASCII), are simply represented as themselves in UTF-8. The second half of Latin-1 characters are split. The first half of the non-ASCII Latin-1 characters are represented by themselves, preceded by code 194 decimal or C2 hex, so the UTF-8 encoding for character code 191 (decimal), ç, is

Ãç

The second half of the non-ASCII Latin-1 characters are represented by a different character, preceded by code 195 decimal or C3 hex. So, when looking at UTF-8 encodings of Latin-1 characters, if you see Ã or Ä where you do not expect it, there are probably too many UTF-8 encodings. Multiple extra encodings have a pattern to them:

0 é  
1 Ä©  
2 ÄfÄ©  
3 ÄfÄfÄ,Ä©  
4 ÄfÆ'Ä,Æ'Äfä€šÄ,Ä©  
5 you get the idea

Note: If you see boxes in the characters above, it is because the font used is missing that character. There is no way to fix it other than getting a new font or by changing the font. Often, the fonts used in a window title or status bar or JavaScript are more limited than those used elsewhere, so the "alert", "title", and "status" buttons in the [Character Conversion Corner](#) can be used to test characters in those contexts.

Too few encodings can have a bad effect that looks different. When é is not UTF-8 encoded, it can appear like this very high numbered character:



Progressive under-encoding can result in a question mark being displayed.

# Encoding errors

## Diagnostic Reference

You are now ready to diagnose UTF-8 encoding problems (e.g., with é):

| Symptom | Diagnosis                                                                                                                                                                        |
|---------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| é       | no problems                                                                                                                                                                      |
| Ã©      | too much UTF-8 encoding, or viewing UTF-8 encoded text with Latin-1 encoding                                                                                                     |
| ÃfÃ©    | much too much UTF-8 encoding                                                                                                                                                     |
| ◆       | too little UTF-8 encoding                                                                                                                                                        |
| ?       | something bad happened to this character                                                                                                                                         |
|         | wild animals have eaten this character                                                                                                                                           |
| □       | if you see a box, the font in use is missing this character. Firefox 3's boxes contain the hexadecimal value for the missing character, but it's still just a missing character. |

# Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.5. Base32 and Base64 encoding

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

# Binary to text encoding

- Base64

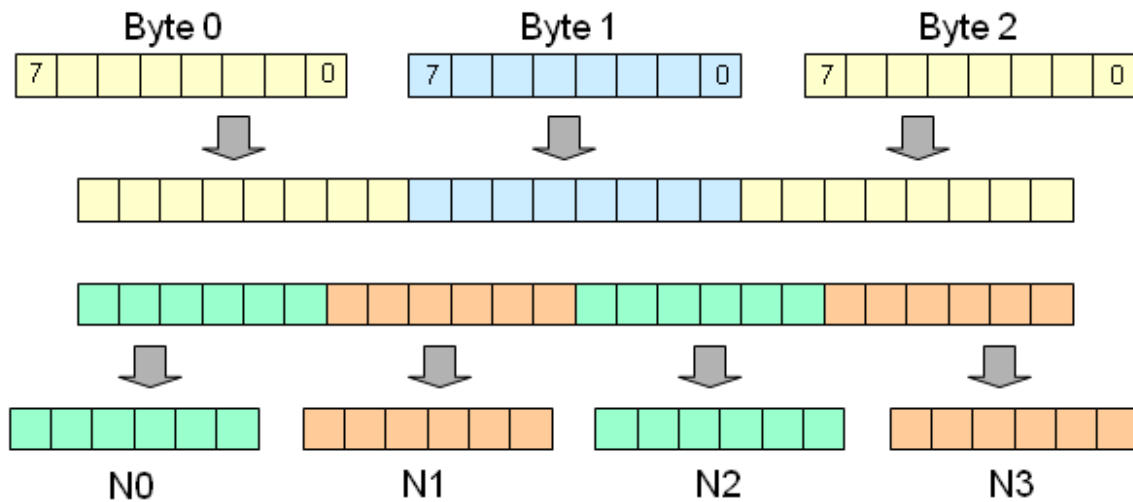


Table 1: The Base 64 Alphabet

| Value | Encoding | Value | Encoding | Value | Encoding | Value | Encoding |
|-------|----------|-------|----------|-------|----------|-------|----------|
| 0     | A        | 17    | R        | 34    | i        | 51    | z        |
| 1     | B        | 18    | S        | 35    | j        | 52    | 0        |
| 2     | C        | 19    | T        | 36    | k        | 53    | 1        |
| 3     | D        | 20    | U        | 37    | l        | 54    | 2        |
| 4     | E        | 21    | V        | 38    | m        | 55    | 3        |
| 5     | F        | 22    | W        | 39    | n        | 56    | 4        |
| 6     | G        | 23    | X        | 40    | o        | 57    | 5        |
| 7     | H        | 24    | Y        | 41    | p        | 58    | 6        |
| 8     | I        | 25    | Z        | 42    | q        | 59    | 7        |
| 9     | J        | 26    | a        | 43    | r        | 60    | 8        |
| 10    | K        | 27    | b        | 44    | s        | 61    | 9        |
| 11    | L        | 28    | c        | 45    | t        | 62    | +        |
| 12    | M        | 29    | d        | 46    | u        | 63    | /        |
| 13    | N        | 30    | e        | 47    | v        |       |          |
| 14    | O        | 31    | f        | 48    | w        | (pad) | =        |
| 15    | P        | 32    | g        | 49    | x        |       |          |
| 16    | Q        | 33    | h        | 50    | y        |       |          |

Table 2: The "URL and Filename safe" Base 64 Alphabet

62 - (minus)  
63 \_  
(underline)



# Binary to text encoding

- Base64

The following example of Base64 data is from [5], with corrections.

```
Input data: 0x14fb9c03d97e
Hex:      1  4  f  b  9  c      |  0  3  d  9  7  e
8-bit:    00010100 11111011 10011100 | 00000011 11011001 01111110
6-bit:    000101 001111 101110 011100 | 000000 111101 100101 111110
Decimal:  5      15      46      28      |  0      61      37      62
Output:   F      P      u      c      A      9      l      +
```

```
Input data: 0x14fb9c03d9
Hex:      1  4  f  b  9  c      |  0  3  d  9
8-bit:    00010100 11111011 10011100 | 00000011 11011001
                                           pad with 00
6-bit:    000101 001111 101110 011100 | 000000 111101 100100
Decimal:  5      15      46      28      |  0      61      36
                                           pad with =
Output:   F      P      u      c      A      9      k      =
```

```
Input data: 0x14fb9c03
Hex:      1  4  f  b  9  c      |  0  3
8-bit:    00010100 11111011 10011100 | 00000011
                                           pad with 0000
6-bit:    000101 001111 101110 011100 | 000000 110000
Decimal:  5      15      46      28      |  0      48
                                           pad with =
Output:   F      P      u      c      A      w      =      =
```

# Binary to text encoding

- Base64
- Base32
- Base16

Table 3: The Base 32 Alphabet

| Value | Encoding | Value | Encoding | Value | Encoding | Value | Encoding |
|-------|----------|-------|----------|-------|----------|-------|----------|
| 0     | A        | 9     | J        | 18    | S        | 27    | 3        |
| 1     | B        | 10    | K        | 19    | T        | 28    | 4        |
| 2     | C        | 11    | L        | 20    | U        | 29    | 5        |
| 3     | D        | 12    | M        | 21    | V        | 30    | 6        |
| 4     | E        | 13    | N        | 22    | W        | 31    | 7        |
| 5     | F        | 14    | O        | 23    | X        |       |          |
| 6     | G        | 15    | P        | 24    | Y        | (pad) | =        |
| 7     | H        | 16    | Q        | 25    | Z        |       |          |
| 8     | I        | 17    | R        | 26    | 2        |       |          |

Table 5: The Base 16 Alphabet

| Value | Encoding | Value | Encoding | Value | Encoding | Value | Encoding |
|-------|----------|-------|----------|-------|----------|-------|----------|
| 0     | 0        | 4     | 4        | 8     | 8        | 12    | C        |
| 1     | 1        | 5     | 5        | 9     | 9        | 13    | D        |
| 2     | 2        | 6     | 6        | 10    | A        | 14    | E        |
| 3     | 3        | 7     | 7        | 11    | B        | 15    | F        |

# Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.6. Date and time

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

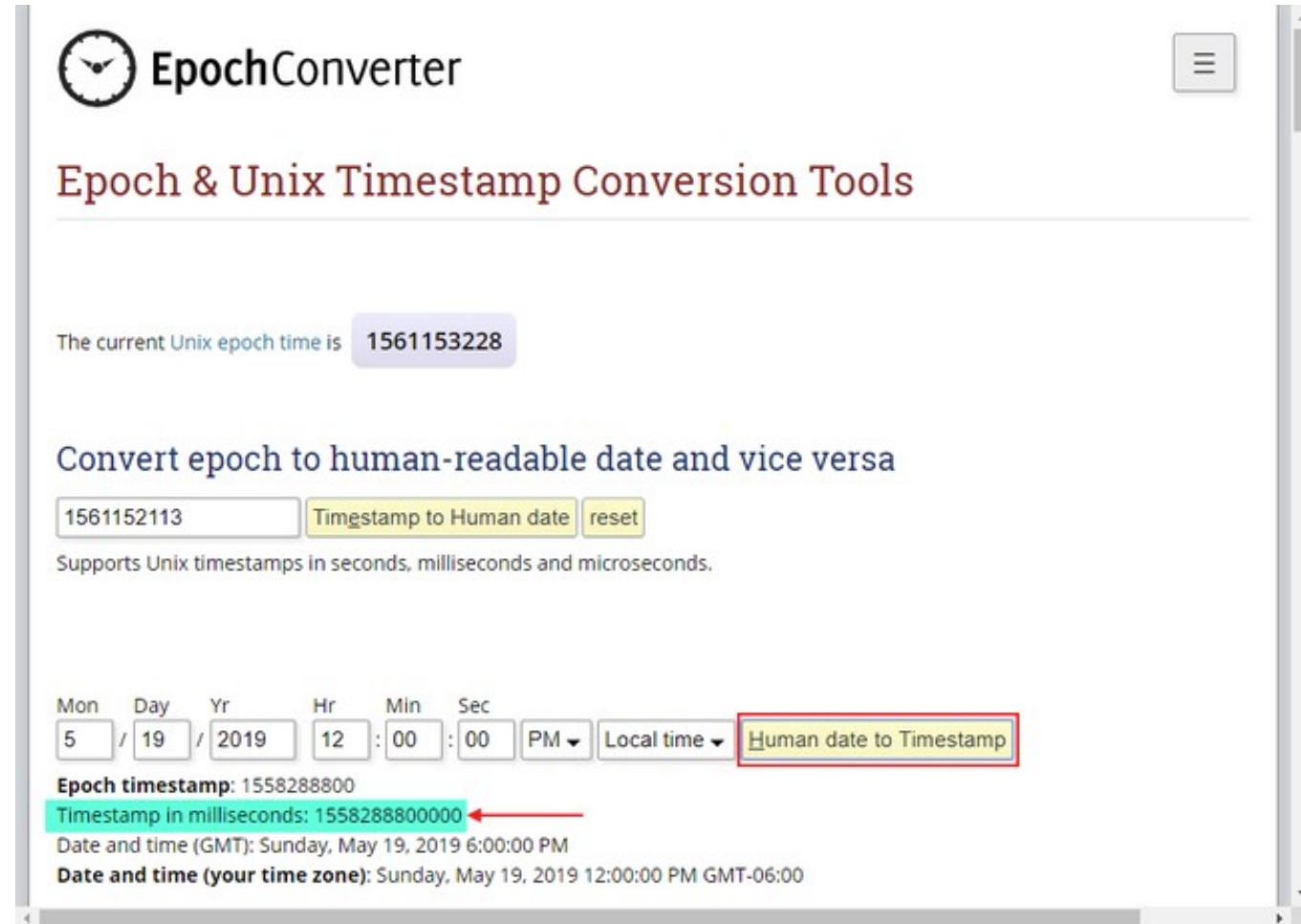
Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

# Representing Date and Time is not simple

- Leap year / bissextile years
  - *a calendar year that contains an additional day added to keep the calendar year synchronized with the astronomical year or seasonal year.*
- Leap seconds
  - *a one-second adjustment that is occasionally applied to Coordinated Universal Time (UTC), to accommodate the difference between precise time (International Atomic Time (TAI), as measured by atomic clocks) and imprecise observed solar time (UT1), which varies due to irregularities and long-term slowdown in the Earth's rotation*
- Timezones
  - UTC, GMT, CET, CEST, ... <https://www.timeanddate.com/time/zones/>

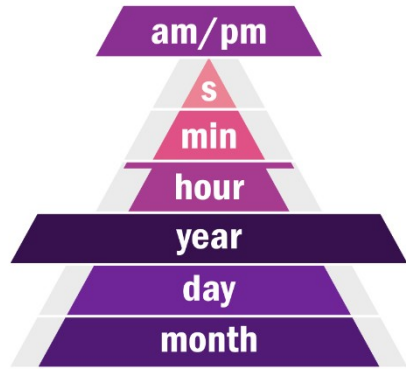
# Unix time stamp

*The unix time stamp is a way to track time as a running total of seconds. This count starts at the Unix Epoch on January 1st, 1970 at UTC. Therefore, the unix time stamp is merely the number of seconds between a particular date and the Unix Epoch.*

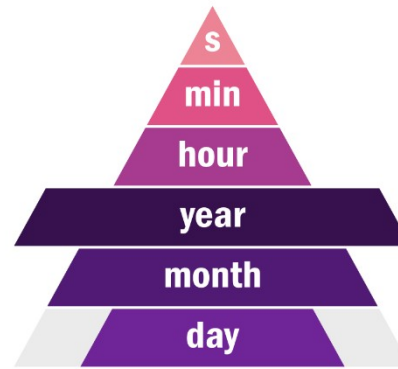


The screenshot shows the EpochConverter website. At the top, there's a logo with a clock icon and the text "EpochConverter". Below it, the title "Epoch & Unix Timestamp Conversion Tools" is displayed. A section indicates "The current Unix epoch time is 1561153228". The main heading is "Convert epoch to human-readable date and vice versa". There are two input fields: one for a timestamp (1561152113) and a button "Timestamp to Human date", and another for a human date (5/19/2019 12:00:00 PM) with a button "Human date to Timestamp". Below these, it says "Supports Unix timestamps in seconds, milliseconds and microseconds." The date and time fields are set to "5 / 19 / 2019 12 : 00 : 00 PM" with a dropdown for "Local time". The "Human date to Timestamp" button is highlighted with a red box. Below the date fields, the "Epoch timestamp" is 1558288800. The "Timestamp in milliseconds" is 1558288800000, highlighted with a green box and a red arrow. The "Date and time (GMT)" is "Sunday, May 19, 2019 6:00:00 PM". The "Date and time (your time zone)" is "Sunday, May 19, 2019 12:00:00 PM GMT-06:00".

# Too many formats



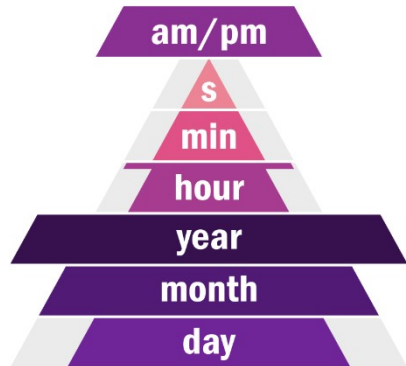
USA/Philippines  
**12/31/2021, 11:59:59 PM**  
12:00 AM – 11:59 PM



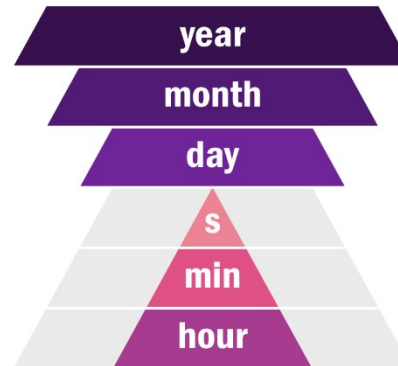
Worldwide  
**31.12.2021, 23:59:59**  
00:00 – 23:59



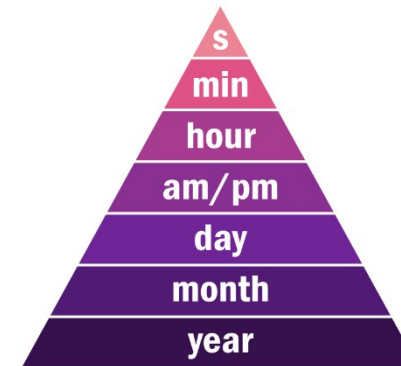
ISO 8601  
**2021-12-31T23:59:59**  
00:00 – 23:59



Australia/India  
**31/12/2021, 11:59:59 PM**  
12:00 AM – 11:59 PM

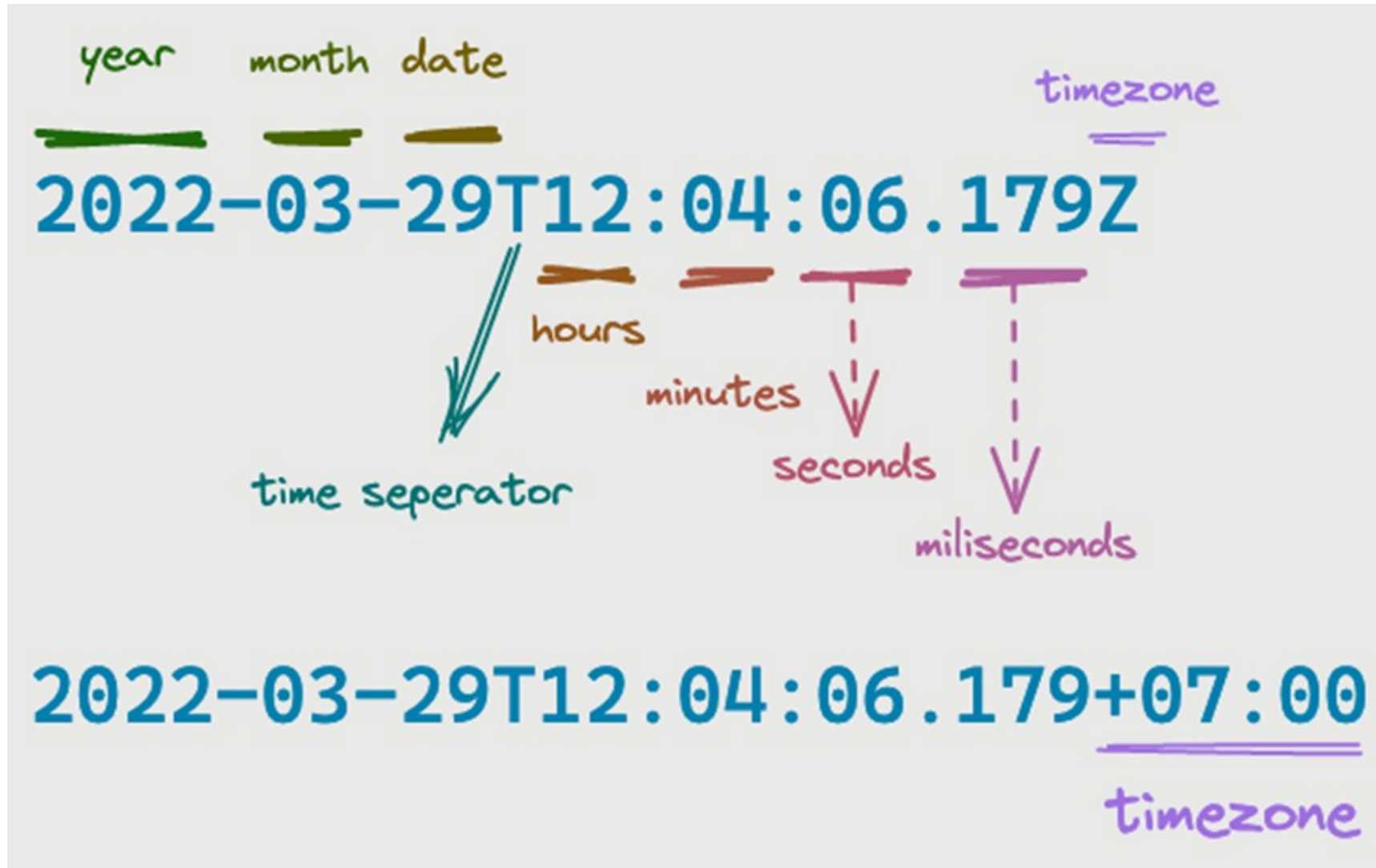


Parts of Spain  
**23:59:59 of 31/12/2021**  
00:00 – 23:59



East Asia  
**2021/12/31 PM 11:59:59**  
AM 0:00 – PM 11:59

# ISO 8601





# Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.7. XML Schema Datatypes

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

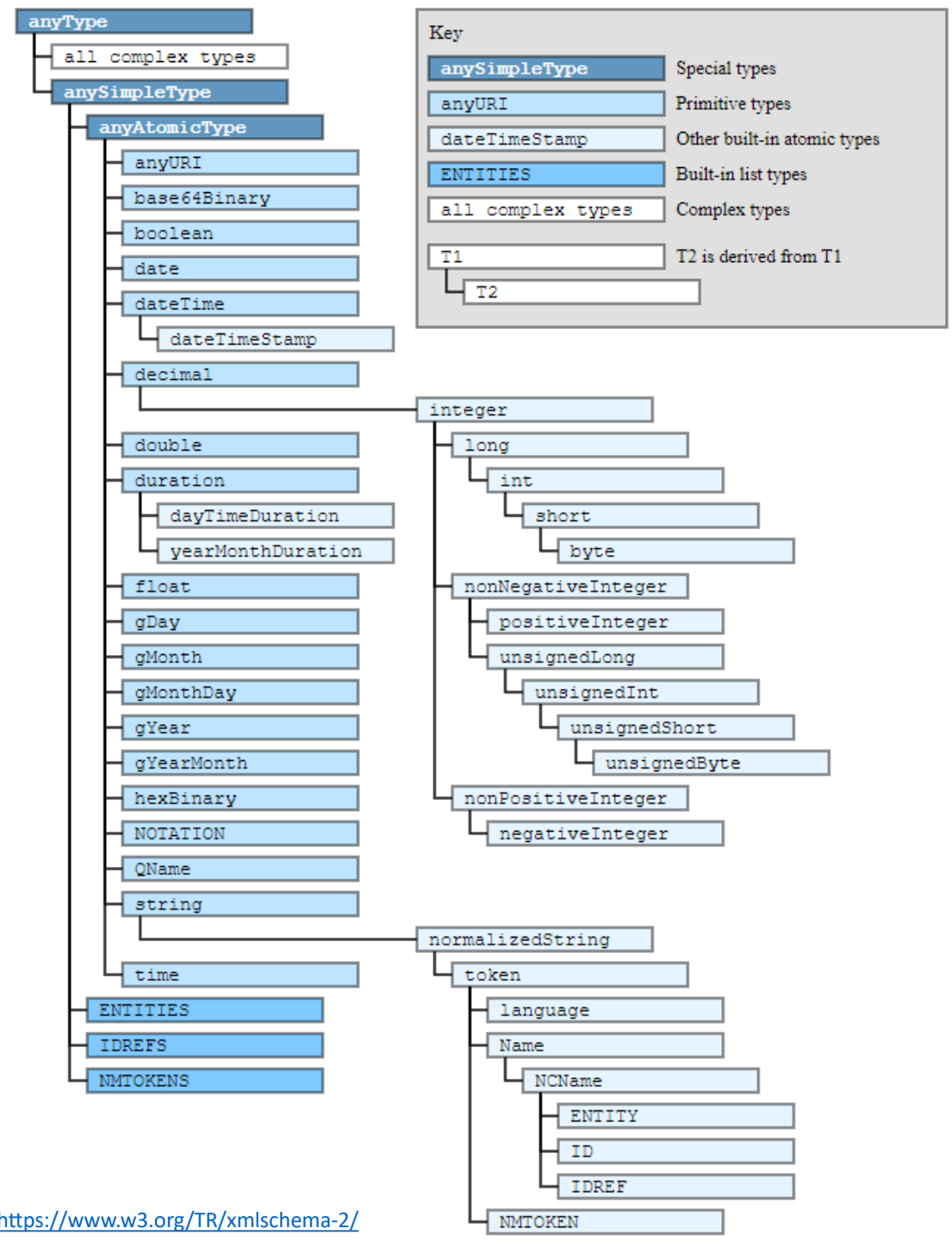
Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

# XML Schema Datatypes

A **datatype** is denoted by a IRI and has three properties:

- A **value space**, which is a set of values.
- A **lexical space**, which is a set of **literals** used to denote the values.
- A lexical-to-value mapping
- A small collection of *functions, relations, and procedures* associated with the datatype.



# Goal: type elements and attributes in XML

The diagram illustrates the relationship between an XML document and its XSD schema. On the left, an XML tree structure is shown with three customers, each having orders and an address. On the right, the XSD schema is displayed, defining the structure and data types for the XML elements and attributes. Arrows point from specific XML elements to their corresponding XSD definitions.

**XML Document Structure:**

```
<ROOT>
  <Customers>
    <Customer CustomerName="Arshad Ali" CustomerID="C001">
      <Orders>
        <Order OrderDate="2012-07-04T00:00:00" OrderID="10248">
          <OrderDetail Quantity="5" ProductID="10" />
          <OrderDetail Quantity="12" ProductID="11" />
          <OrderDetail Quantity="10" ProductID="42" />
        </Order>
      </Orders>
      <Address> Address line 1, 2, 3</Address>
    </Customer>
    <Customer CustomerName="Paul Henriot" CustomerID="C002">
      <Orders>
        <Order OrderDate="2011-07-04T00:00:00" OrderID="10245">
          <OrderDetail Quantity="12" ProductID="11" />
          <OrderDetail Quantity="10" ProductID="42" />
        </Order>
      </Orders>
      <Address> Address line 5, 6, 7</Address>
    </Customer>
    <Customer CustomerName="Carlos Gonzlez" CustomerID="C003">
      <Orders>
        <Order OrderDate="2012-08-16T00:00:00" OrderID="10283">
          <OrderDetail Quantity="3" ProductID="72" />
        </Order>
      </Orders>
      <Address> Address line 1, 4, 5</Address>
    </Customer>
  </Customers>
</ROOT>
```

**XSD Schema Definitions:**

```
<xsd:attribute name="OrderID" type="xsd:integer"/>

<xsd:attribute name="OrderDate" type="xsd:dateTime"/>

<xsd:attribute name="CustomerID">
  <xsd:simpleType>
    <xsd:restriction base="xsd:string">
      <xsd:pattern value="/^[0-9]{3}$/"></xsd:pattern>
    </xsd:restriction>
  </xsd:simpleType>
</xsd:attribute>
```

**Arrows indicating mappings:**

- From `OrderID="10248"` in the first XML `<Order>` element to `<xsd:attribute name="OrderID" type="xsd:integer"/>`.
- From `OrderDate="2011-07-04T00:00:00"` in the second XML `<Order>` element to `<xsd:attribute name="OrderDate" type="xsd:dateTime"/>`.
- From `CustomerID="C002"` in the second XML `<Customer>` element to the `<xsd:attribute name="CustomerID">` block.

# Anatomy of a XSD literal

"24"^^xsd:**integer**

"true"^^xsd:**boolean**

"2001-10-26T21:32:52+02:00"^^xsd:**string**

lexical form

datatype IRI

@prefix xsd: <http://www.w3.org/2001/XMLSchema#>.

# xsd:int vs xsd:integer

A **xsd:int** represents a signed 32-bit integer

A **xsd:integer** is an integer unbounded value

# xsd:float vs xsd:double vs xsd:decimal

A **xsd:float** is patterned after the IEEE single-precision 32-bit floating point datatype

$(\backslash+|-)?([0-9]+(\backslash.[0-9]*)?|\backslash.[0-9]+)([Ee](\backslash+|-)?[0-9]+)?|(\backslash+|-)?INF|NaN$

A **xsd:double** is patterned after the IEEE double-precision 64-bit floating point datatype

$(\backslash+|-)?([0-9]+(\backslash.[0-9]*)?|\backslash.[0-9]+)([Ee](\backslash+|-)?[0-9]+)?|(\backslash+|-)?INF|NaN$

A **xsd:integer** represents a subset of the real numbers, which can be represented by decimal numerals

$(\backslash+|-)?([0-9]+(\backslash.[0-9]*)?|\backslash.[0-9]+)$

# Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.8. Codes: countries, languages, ...

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics





Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>



# Country codes

## ISO 3166-1 – Codes for the representation of names of countries and their subdivisions – Part 1: Country codes

| ISO 3166 <sup>[1]</sup>                                                                                          |                                      |                                    | ISO 3166-1 <sup>[2]</sup>     |                               |                               | ISO 3166-2 <sup>[3]</sup>               |                                 |
|------------------------------------------------------------------------------------------------------------------|--------------------------------------|------------------------------------|-------------------------------|-------------------------------|-------------------------------|-----------------------------------------|---------------------------------|
| Country name <sup>[5]</sup> ♦                                                                                    | Official state name <sup>[6]</sup> ♦ | Sovereignty <sup>[6][7][8]</sup> ♦ | Alpha-2 code <sup>[5]</sup> ♦ | Alpha-3 code <sup>[5]</sup> ♦ | Numeric code <sup>[5]</sup> ♦ | Subdivision code links <sup>[3]</sup> ♦ | Internet ccTLD <sup>[9]</sup> ♦ |
|  France <sup>[1]</sup>          | The French Republic                  | UN member state                    | FR                            | FRA                           | 250                           | ISO 3166-2:FR                           | .fr                             |
|  United States of America (the) | The United States of America         | UN member state                    | US                            | USA                           | 840                           | ISO 3166-2:US                           | .us                             |
|  China                        | The People's Republic of China       | UN member state                    | CN                            | CHN                           | 156                           | ISO 3166-2:CN                           | .cn                             |
|  Austria                      | The Republic of Austria              | UN member state                    | AT                            | AUT                           | 040                           | ISO 3166-2:AT                           | .at                             |

### Example of country codes

Source: [https://en.wikipedia.org/wiki/List\\_of\\_ISO\\_3166\\_country\\_codes](https://en.wikipedia.org/wiki/List_of_ISO_3166_country_codes)

# Language codes

**ISO 639** is a standardized nomenclature used to classify languages. Each language is assigned a two-letter (639-1) and three-letter (639-2 and 639-3) lowercase abbreviation

| ISO language name  | 639-1 | 639-2/T | 639-2/B | 639-3    |                                       |
|--------------------|-------|---------|---------|----------|---------------------------------------|
| English            | en    | eng     | eng     | eng      |                                       |
| Chinese            | zh    | zho     | chi     | zho + 16 | macrolanguage                         |
| Hindi              | hi    | hin     | hin     | hin      |                                       |
| Spanish, Castilian | es    | spa     | spa     | spa      |                                       |
| French             | fr    | fra     | fre     | fra      |                                       |
| Arabic             | ar    | ara     | ara     | ara + 29 | macrolanguage, Standard Arabic is arb |
| Bengali            | bn    | ben     | ben     | ben      |                                       |

Example of language names for the most spoken languages

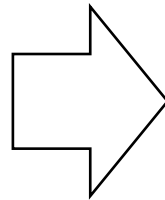
Source: [https://en.wikipedia.org/wiki/List\\_of\\_ISO\\_639-1\\_codes](https://en.wikipedia.org/wiki/List_of_ISO_639-1_codes)

# IETF BCP47: Tags for Identifying Languages

Internet Engineering Task Force (IETF) « Best Current Practice » (BCP)

Region subtags

en-GB  
es-005  
zh-Hant-HK



| HTML Demo: lang                                                                   |     | RESET                                                                    |
|-----------------------------------------------------------------------------------|-----|--------------------------------------------------------------------------|
| HTML                                                                              | CSS | OUTPUT                                                                   |
| 1 <p>This paragraph is English, but the language is not specifically defined.</p> |     | This paragraph is English, but the language is not specifically defined. |
| 2                                                                                 |     |                                                                          |
| 3 <p lang="en-GB">This paragraph is defined as British English.</p>               |     | (In British English) This paragraph is defined as British English.       |
| 4                                                                                 |     |                                                                          |
| 5 <p lang="fr">Ce paragraphe est défini en français.</p>                          |     | (In French) Ce paragraphe est défini en français.                        |
| 6                                                                                 |     |                                                                          |

Read more in the BCP 47 spec:

2.2.4 Region Subtag

4.1 Choice of Language Tag

[https://developer.mozilla.org/en-US/docs/Web/HTML/Global\\_attributes/lang](https://developer.mozilla.org/en-US/docs/Web/HTML/Global_attributes/lang)

# Currency codes

**ISO 4217** defines alpha codes and numeric codes for the representation of currencies and provides information about the relationships between individual currencies and their minor units.

Active ISO 4217 currency codes<sup>[1]</sup> [hide]

| Code ↕ | Num ↕ | D <sup>[a]</sup> ↕ | Currency ↕                  | Locations listed for this currency <sup>[b]</sup> ↕                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |
|--------|-------|--------------------|-----------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| AED    | 784   | 2                  | United Arab Emirates dirham |  United Arab Emirates                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
| EUR    | 978   | 2                  | Euro                        |  Åland Islands (AX),  European Union (EU),  Andorra (AD),  Austria (AT),  Belgium (BE),  Cyprus (CY),  Estonia (EE),  Finland (FI),  France (FR),  French Southern and Antarctic Lands (TF),  Germany (DE),  Greece (GR),  Guadeloupe (GP),  Ireland (IE),  Italy (IT),  Latvia (LV),  Lithuania (LT),  Luxembourg (LU),  Malta (MT),  French Guiana (GF),  Martinique (MQ),  Mayotte (YT),  Monaco (MC),  Montenegro (ME),  Netherlands (NL),  Portugal (PT),  Réunion (RE),  Saint Barthélemy (BL),  Saint Martin (MF),  Saint Pierre and Miquelon (PM),  San Marino (SM),  Slovakia (SK),  Slovenia (SI),  Spain (ES),  Vatican City (VA) |
| KRW    | 410   | 0 <sup>[c]</sup>   | South Korean won            |  South Korea                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| GBP    | 826   | 2                  | Pound sterling              |  United Kingdom,  Isle of Man (IM, see Manx pound),  Jersey (JE, see Jersey pound),  Guernsey (GG, see Guernsey pound),  Tristan da Cunha (SH-TA)                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
| CHF    | 756   | 2                  | Swiss franc                 |  Switzerland,  Liechtenstein (LI)                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |

## Example of currency codes

Source: [https://en.wikipedia.org/wiki/ISO\\_4217](https://en.wikipedia.org/wiki/ISO_4217)

# Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.9. Quantities and Units of measure

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

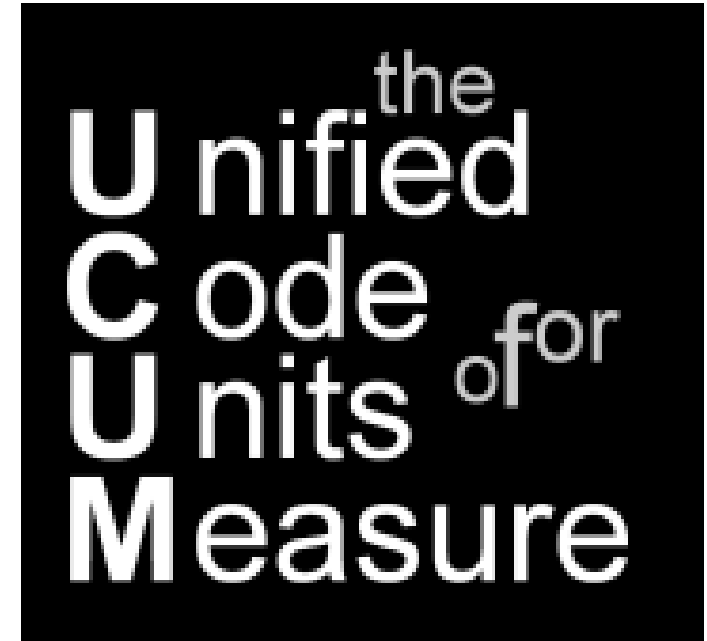
Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

# Units of measures: no consensus

- BIPM (International Bureau of Weights and Measures)
- ISO/IEC 1000 – ISO/IEC 80000
- VIM (international Vocabulary for Measurements)
- UnitsML
- UCUM (Unified Code for Units of Measure)
- UNECE Recommendation 20
- Sweet
- QUDT
- ...

# UCUM: Unified Code for Units of Measure

- A code system intended to include *all units* of measures being contemporarily used in international science, engineering, and business.
- **Used** by international organizations and standards
- **No ambiguity** possible
- **Clear semantics** of units
- **Con: Problematic** custom license



$([\text{pi}]/2) \cdot \text{rad} \cdot \text{kK}^4$

Diagram illustrating the components of the UCUM expression  $([\text{pi}]/2) \cdot \text{rad} \cdot \text{kK}^4$ :

- simple units**: Points to the  $([\text{pi}]/2)$  term.
- simple units**: Points to the  $\text{rad}$  term.
- simple units**: Points to the  $\text{kK}$  term.
- exponent**: Points to the  $4$  term.



| UCUM Code of Unit Concept | EN Unit               | EN Symbol                 | EN Dimension <sup>a</sup>        | NCI Concept Code | NCI Term                                          | NCI Abbreviation                  | SNOMED CT Identifier <sup>b</sup> |
|---------------------------|-----------------------|---------------------------|----------------------------------|------------------|---------------------------------------------------|-----------------------------------|-----------------------------------|
| [IU]                      |                       |                           | [arb]                            | C70497           | Anti-Xa Activity International Unit               | anti-Xa activity                  | 258997004                         |
| Bq                        | becquerel             | Bq                        | T <sup>-1</sup>                  | C42562           | Becquerel                                         | Bq                                | 282141004                         |
| Bq/g                      |                       |                           | M <sup>-1</sup> T <sup>-1</sup>  | C70522           | Becquerel per gram                                | Bq/g                              |                                   |
| 10 <sup>9</sup> .[CFU]    |                       |                           | [arb]                            | C68897           | Billion Colony Forming Units                      | Billion CFU                       |                                   |
| 10 <sup>9</sup>           |                       |                           | 1                                | C71189           | Billion Organisms                                 |                                   |                                   |
| m <sup>3</sup>            | cubic metre           | m <sup>3</sup>            | L <sup>3</sup>                   | C42570           | Cubic Meter                                       | m <sup>3</sup>                    | 396154006                         |
| Ci/ml                     |                       |                           | L <sup>-3</sup> T <sup>-1</sup>  | C71172           | Curie per Millilitre                              | Ci/ml                             |                                   |
| d                         | day                   | d                         | T                                | C25301           | Day                                               | d                                 | 258703001                         |
| [drp]                     |                       |                           | L <sup>3</sup>                   | C48491           | Drop Dosing Unit                                  | Gtt                               | 404218003                         |
| [IU]/ml                   |                       |                           | [arb]                            | C67377           | International Unit per Millilitre                 | IU/mL                             | 259002007                         |
| k[USP'U]                  |                       |                           | [arb]                            | C71202           | Kilo United States Pharmacopoeia Unit             | KUSP'U                            |                                   |
| kBq/l                     |                       |                           | L <sup>-3</sup> T <sup>-1</sup>  | C71167           | Kilobecquerel per Liter                           | kBq/L                             |                                   |
| mmol/l                    |                       |                           | L <sup>-3</sup> N                | C64387           | Millimole per Liter                               | mmol/L                            | 258813002                         |
| [ppm]                     | part per million      | ppm                       | 1                                | C48523           | Part Per Million                                  | ppm                               | 258731005                         |
| Pa                        | pascal                | Pa                        | L <sup>-1</sup> MT <sup>-2</sup> | C42547           | Pascal                                            | P                                 | 259016002                         |
| %                         | per cent              | %                         | 1                                | C48570           | Percent                                           | %                                 | 118582008                         |
| %                         |                       |                           | 1                                | C48571           | Percent Volume per Volume                         | %V/V                              | 419569009                         |
| g/ml                      | per cent (w/v)        | %(w/v)                    | L <sup>-3</sup> M                | C48527           | Percent Weight Volume                             | %M/V                              | 396169007                         |
| %                         |                       |                           | 1                                | C48528           | Percent Weight Weight                             | %W/W                              | 118582008                         |
| [PFU]                     |                       |                           | [arb]                            | C73575           | Plaque Forming Unit Equivalent<br>1000 Mouse LD50 | PFU Equivalent<br>1000 Mouse LD50 |                                   |
| [lb_av]                   | pound                 | lb                        | M                                | C48531           | Pound                                             | LB                                | 258693003                         |
| /min                      | revolution per minute | r.p.m., rev/min,<br>r/min | T <sup>-1</sup>                  | C70469           | Revolution per Minute                             | rpm                               | 286549009                         |
| [tb'U]                    |                       |                           | [arb]                            | C65132           | Tuberculin Unit                                   |                                   | 415758003                         |
| [arb'U]{ELISA}            |                       |                           |                                  | C68875           | Enzyme-Linked Immunosorbent Assay Unit            | EL. U                             |                                   |

# Data Interoperability and Semantics

Part 1. Encoding base data types

Part 1.10. Colors

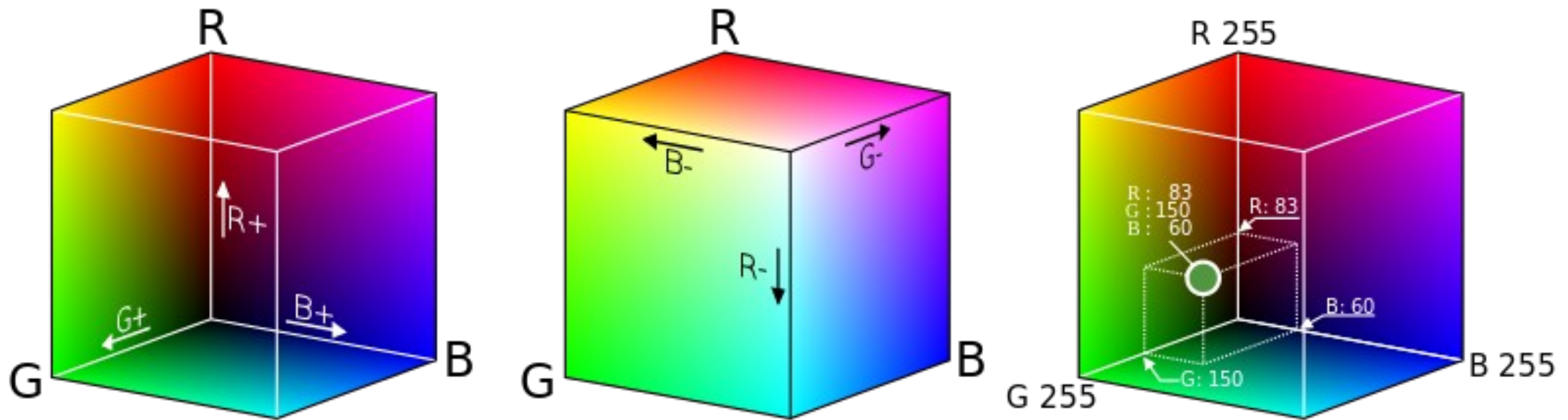
ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

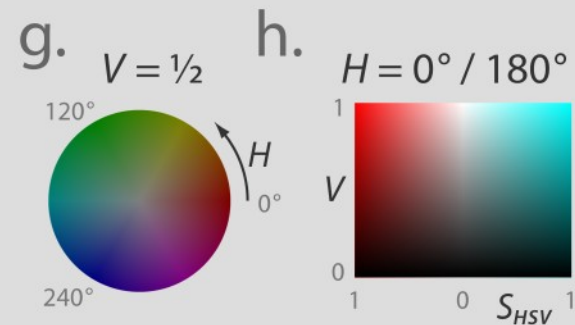
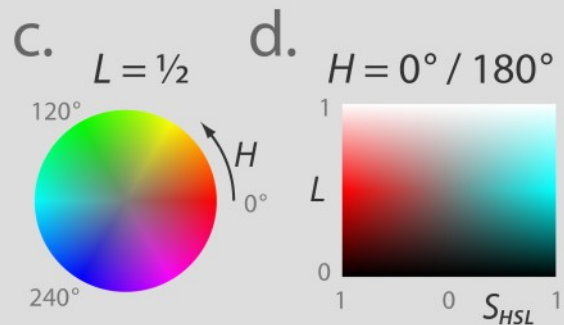
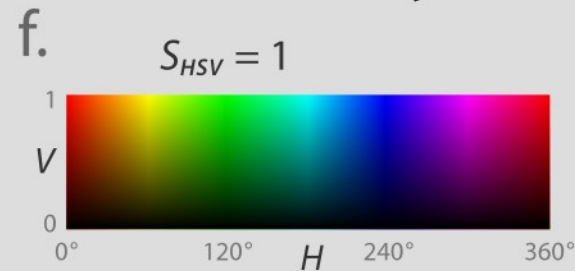
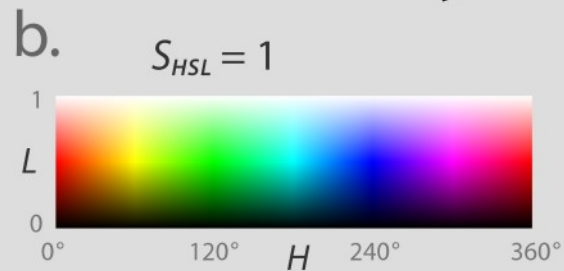
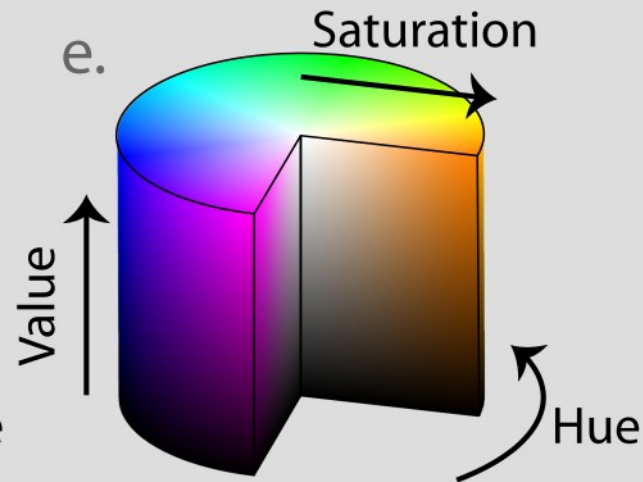
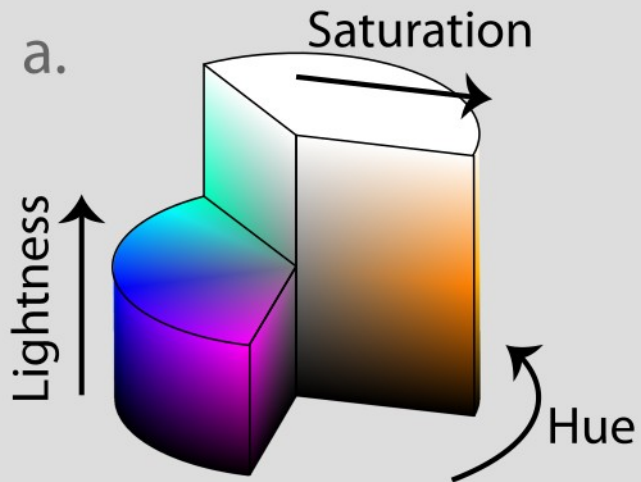
Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>

# RGB color cube

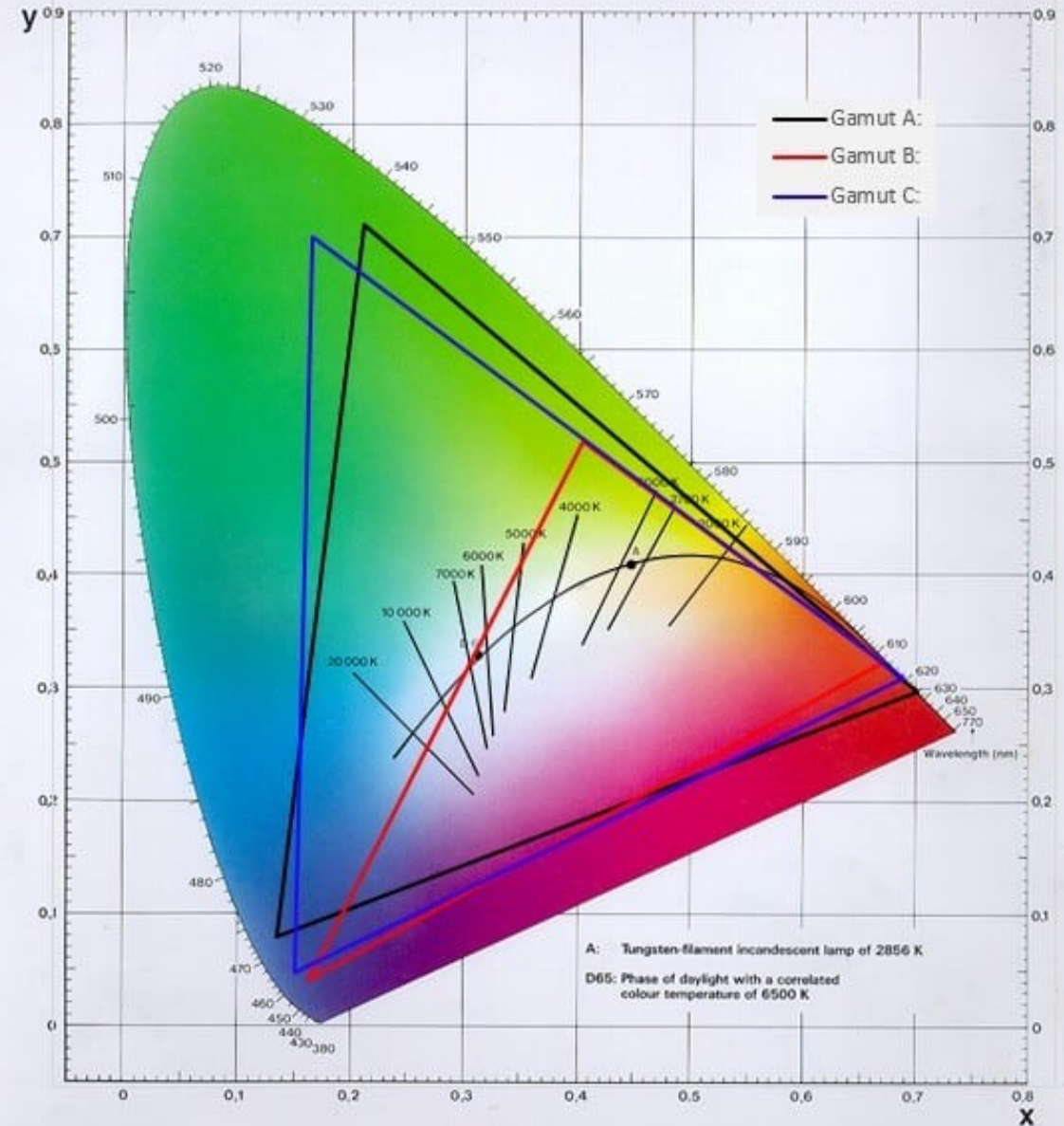


# HSL

# HSV



## C.I.E. CHROMATICITY DIAGRAM



HSL (for hue, saturation, lightness)

HSV (for hue, saturation, value; also known as HSB, for hue, saturation, brightness)

[https://en.wikipedia.org/wiki/HSL\\_and\\_HSV](https://en.wikipedia.org/wiki/HSL_and_HSV)

Philips hue color space



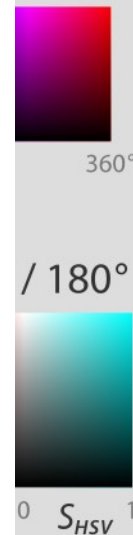
# HSL

# HSV

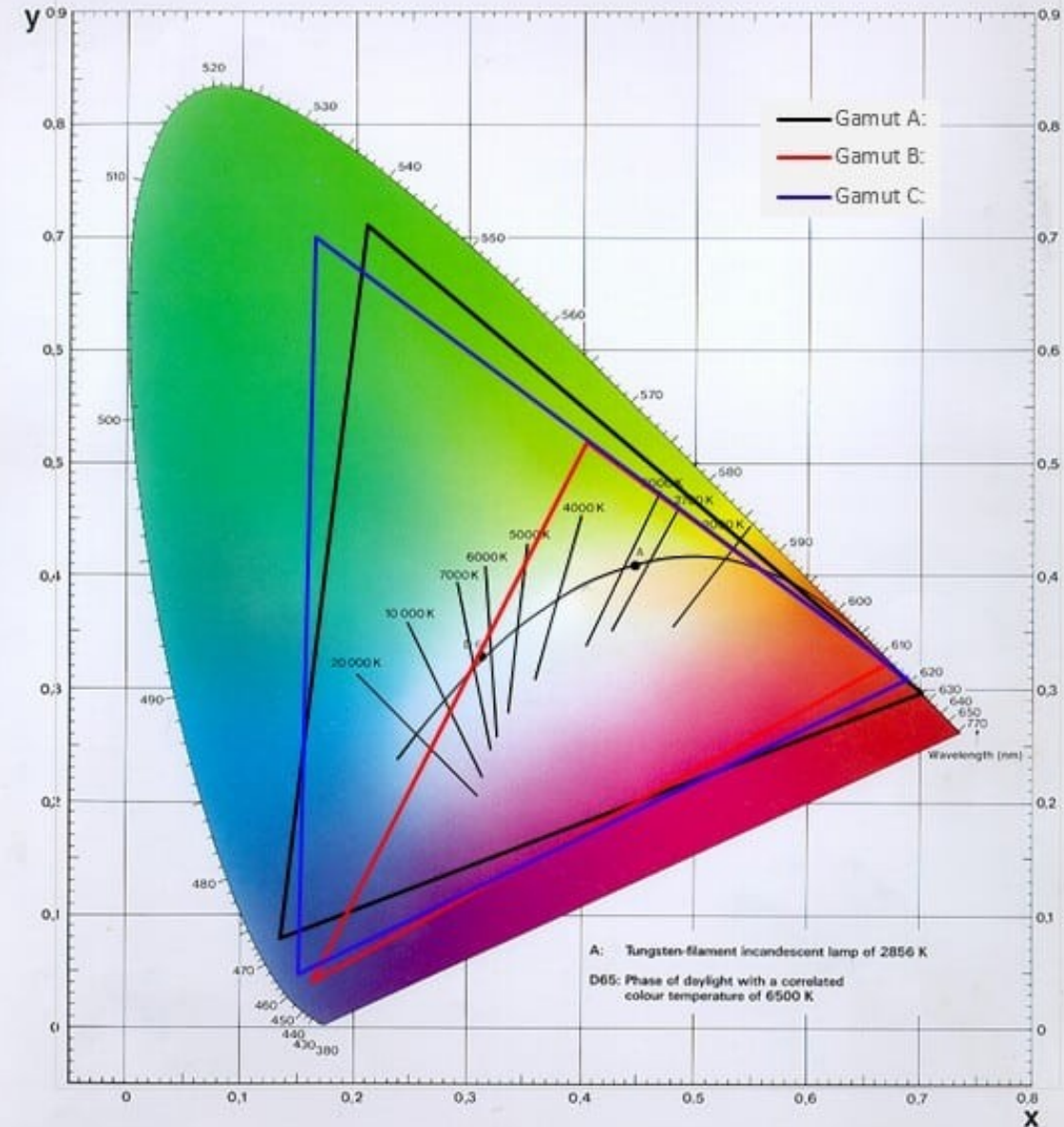
```

1 {
2   "1": {
3     "state": {
4       "on": true,
5       "bri": 254,
6       "hue": 14314,
7       "sat": 172,
8       "effect": "none",
9       "xy": [
10        0.4791,
11        0.4139
12      ],
13       "ct": 405,
14       "alert": "none",
15       "colormode": "ct",
16       "reachable": true
17     },
18     "type": "Extended color light",
19     "name": "Hue color light 1",
20     "modelid": "LCT001",
21     "manufacturername": "Philips",
22     "uniqueid": "00:17:88:01:00:ff:9a:28-0b",
23     "swversion": "5.127.1.26581"
24   }
25 }

```



## C.I.E. CHROMATICITY DIAGRAM



# HSV $\leftrightarrow$ RGB conversion

When  $0 \leq H < 360$ ,  $0 \leq S \leq 1$  and  $0 \leq V \leq 1$ :

$$C = V \times S$$

$$X = C \times (1 - |(H / 60^\circ) \bmod 2 - 1|)$$

$$m = V - C$$

$$(R', G', B') = \begin{cases} (C, X, 0) & , 0^\circ \leq H < 60^\circ \\ (X, C, 0) & , 60^\circ \leq H < 120^\circ \\ (0, C, X) & , 120^\circ \leq H < 180^\circ \\ (0, X, C) & , 180^\circ \leq H < 240^\circ \\ (X, 0, C) & , 240^\circ \leq H < 300^\circ \\ (C, 0, X) & , 300^\circ \leq H < 360^\circ \end{cases}$$

$$(R, G, B) = ((R' + m) \times 255, (G' + m) \times 255, (B' + m) \times 255)$$

$$R' = R / 255$$

$$G' = G / 255$$

$$B' = B / 255$$

$$C_{\max} = \max(R', G', B')$$

$$C_{\min} = \min(R', G', B')$$

$$\Delta = C_{\max} - C_{\min}$$

$$H = \begin{cases} 0^\circ, \Delta = 0 \\ 60^\circ \times \left( \frac{G' - B'}{\Delta} \bmod 6 \right), C_{\max} = R' \\ 60^\circ \times \left( \frac{B' - R'}{\Delta} + 2 \right), C_{\max} = G' \\ 60^\circ \times \left( \frac{R' - G'}{\Delta} + 4 \right), C_{\max} = B' \end{cases}$$

$$S = \begin{cases} 0, C_{\max} = 0 \\ \frac{\Delta}{C_{\max}}, C_{\max} \neq 0 \end{cases}$$

$$V = C_{\max}$$

# Data Interoperability and Semantics

</ Part 1. Encoding base data types >

ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics

M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics

Maxime Lefrançois <https://maxime-lefrancois.info>

Course unit URL: <https://ci.mines-stetienne.fr/cps2/course/data>