Syllabus and Course Organization

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-lefrançois.info

Course unit URL: https://ci.mines-stetlenne.fr/cps2/course/data

Grading policy

Practical Work x2 (PW) Written Exam x2 (WE) Grade = (PW + WE)/2

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-lefrançois.info Course unit URL: https://ch.mines-stetienne-fr/cps2/course/data

Course objectives

This lecture aims at ensuring you are familiar with data and information, at the core of every information technology (IT). The course focuses on four main topics:

- 1. Encoding base data types
- 2. Data formats
- Data schemas and semantics
- 4. The value of data
- 5. The European strategy for data

You will practice with tutorials, and gain experience through a project

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-lefrançois.info
Course unit URL: https://ci.mines-stetienne.fr/cps2/course/data

Syllabus and Course Organization

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-lefrançois.info

Course unit URL: https://ci.mines-stetlenne.fr/cps2/course/data

Grading policy

Practical Work x2 (PW) Written Exam x2 (WE) Grade = (PW + WE)/2

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-lefrançois.info Course unit URL: https://ch.mines-stetienne-fr/cps2/course/data

Course objectives

This lecture aims at ensuring you are familiar with data and information, at the core of every information technology (IT). The course focuses on four main topics:

- 1. Encoding base data types
- 2. Data formats
- Data schemas and semantics
- 4. The value of data
- 5. The European strategy for data

You will practice with tutorials, and gain experience through a project

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-lefrançois.info
Course unit URL: https://ci.mines-stetienne.fr/cps2/course/data

< 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://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

Data Interoperability and SemanticsOutline

- < 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

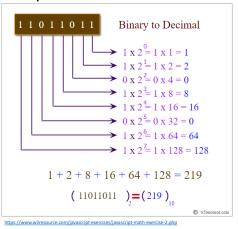
Numbering systems for computers

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"

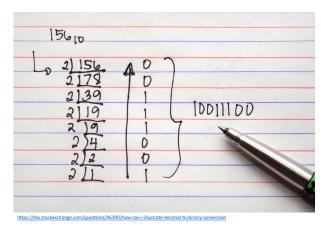


It has been debated a lot at the beginning

Tips: binary to decimal

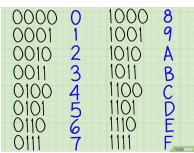


Tips: decimal to binary

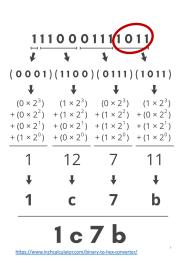


Tips: binary to hexadecimal

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



Binary nibble to hexadecimal digit



Programming with binary strings

```
$ gcc main.c
$ ./a.out
a = 0X65, b = 0X09
a&b = 0X61
a|b = 0X6D
a^b = 0X6C
a << 1 = 0XC2
a >> 1 = 0X32
~a = 0XFFFFF9A
```

(a >> 2) & 0x7 = 0X1

8

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

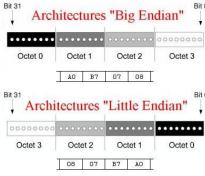
Endianness

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

10

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



Example: WS2812B color leds





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



Created: 29/09/2016 mcf88 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
	10 byte	XX XX	Measure 1, refer to Note1
Data	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

The 10 hutes 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

 - 5 bit for day of the month
 - 5 hits 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
- > 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801 Remove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 me The 3 measurements will be:

- dd7e3721b10a43608801
- e07e3721b20a425d8801



mcf88 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
Jplink ID	1 byte	04	Temperature/Pressure/Humidity
	10 byte	XX XX	Measure 1, refer to Note1
Data	10 byte	XX XX	Measure 2, refer to Note1
	10 byte	XX XX	Measure 3, refer to Note1
3att %	1 byte (optional)	XX	Battery percentage
RFU	4 byte (optional)	XX XX XX XX	Optional RFU byte

The 10 bytes for each measurement are divided as follows:

- 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 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.
- > 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is exp

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

Example: MCF88 LoRa sensors



Created: 29/09/2016 mcf88 Modified: 30/11/2018 MCF88 DATA FRAME FORMAT 1.17

1.2 TEMPERATURE/PRESSURE/HUMIDITY

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

The 10 hutes 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

 - 5 bit for day of the month 5 hits 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
- > 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expressed in

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801 Remove the first byte and divide the other 30 into 3 parts by 10 byte that corresp

- The 3 measurements will be: dc7e3721b40a47608801
- dd7e3721b10a43608801
- e07e3721b20a425d8801 Decipher the first measurement dividing it by groups and applying the necessary transfe
- Measurement date: dc 7e 37 21





Created: 29/09/2016 Modified: 30/11/2018

Author: Colognato Stefano

MCF88 DATA FRAME FORMAT 1.17

1.2 TEMPERATURE/PRESSURE/HUMIDITY

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
	10 byte	XX XX	Measure 1, refer to Note1
Data	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
REU	4 byte (optional)	XX XX XX XX	Optional REU byte

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.

 - 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.
- > 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expres

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

move the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 measure

- dc7e3721b40a47608801
- e07e3721b20a425d8801

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

Measurement date: do 7e 37 21

- Byte swapping, result: 21 37 7e do
- re result in bits will be: 00100001 00110111 01111110 11011100
- The bits are divided as explained above
- > 2000+16 = 2016 Month:
 Result:
- Result:
- Minutes: 110110
- Result: Result:
- > 28*2 = 56
- The date of the measurament will be: 23/09/2016 15:54:56.

17



Modified: 30/11/2018 MCF88 DATA FRAME FORMAT 1.17

1.2 TEMPERATURE/PRESSURE/HUMIDITY

HOME

Created: 29/09/2016

name	size [byte]	hex value	mean
Uplink ID	1 byte	04	Temperature/Pressure/Humidity
	10 byte	XX XX	Measure 1, refer to Note1
Data	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 REU byte

The 10 hutes 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

 - 5 bit for day of the month
 - 5 hits 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
- > 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is expre

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

emove the first byte and divide the other 30 into 3 parts by 10 byte that correspond to 3 mea The 3 measurements will be:

- dd7e3721b10a43608801
- e07e372tb20a425d8801
 Decipher the first measurement dividing it by groups and applying the necessary transformatic. 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:
 Result:
- > 2000+16 = 2016 Result Day:
 Result:
- Minutes: 110110
- Result: 28
- The date of the measurament will be: 23/09/2016 15:54:56.
- Byte swapping, result: 0ab4
- The result (with sign) will be +2740 with two decimal places, then + 27.40 °C.
- 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]	nex value	mean
Jplink ID	1 byte		Temperature/Pressure/Humidity
	10 byte	XX XX	Measure 1, refer to Note1
Data			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

mcf88

Author: Colognato Stefano

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 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.
- > 3 bytes for pressure. Pressure is an unsigned integer with the least significant byte first; it is exp

Created: 29/09/2016

Modified: 30/11/2018

MCF88 DATA FRAME FORMAT 1.17

04dc7e3721b40a47608801dd7e3721b10a43608801e07e3721b20a425d8801

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

- dc7e3721b40a47608801
- e07e3721b20a425d8801

Decipher the first measurement dividi

Measurement date: dc 7e 37 21 nent dividing it by groups and applying the necessary transformations

Byte swapping, result: 21 37 7e dc
 The result in bits will be: 00100001 00110111 01111110 11011100

The bits are divided as explained above

- > 2000+16 = 2016 Month:
 Result:
- Hour:
- Result:
- Minutes: 110110 • Result: 54 Seconds:
 Result:
- > 28*2 = 56 The date of the measurament will be: 23/09/2016 15:54:56.
- o Temperature: b40a
- Byte swapping, result: 0ab4
 The result (with sign) will be +2740 with two decimal places, then + 27.40 °C.
- 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 <<pre>or <<pre>or <<pre>c

• 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

Two's Complement

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

$$B2T(X) = -x_{w-1} \cdot 2^{w-1} + \sum_{i=0}^{w-2} x_i \cdot 2^i$$
short int x = 15213;
short int y = -15213;

C short 2 bytes long

	Decimal	Hex	Binary
Х	15213	3B 6D	00111011 01101101
У	-15213	C4 93	11000100 10010011

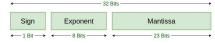
Sign Bit

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

https://slideplayer.com/slide/5048952

IEEE 754 single/double encoding

ttps://www.geeksforgeeks.org/ieee-standard-754-floating-point-numbers/



Single Precision IEEE 754 Floating-Point Standard

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

exceptional cases

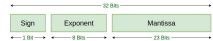
- If E = 255 and F is nonzero, then x = NaN ("Not a number").
- If E = 255, F is zero, and S is 1, then x = Infinity
- If E = 255, F is zero, and S is 0, then x = + Infinity
- If 0 < E < 255, then x = (-1)⁸ × (1. F) × 2^{E-127}, where 1. F represents the binary number created by prefixing F with an implicit 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 "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

- 0 10000001 101000000000000000000000 = $(-1)^0 \times (1.101_2) \times 2^{129-127} = 6.5$
- $1\ 10000001\ 101000000000000000000000 = (-1)^1 \times (1.101_2) \times 2^{129-127} = -6.5.$

IEEE 754 single/double encoding

https://www.geeksforgeeks.org/ieee-standard-754-floating-point-number

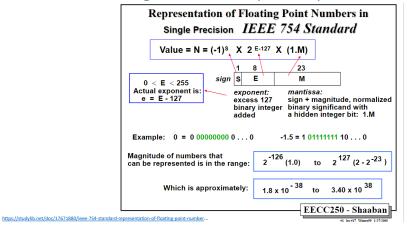


Single Precision IEEE 754 Floating-Point Standard



Double Precision
IEEE 754 Floating-Point Standard

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

- ISO 646 • EBCDIC

Common character encodings [edt]

- ISO 8859-2 Western and Central Europe ISO 8859-3 Western Europe and South European (Turkish, Maltese plus
- . ISO 8859.4 Western Furnne 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 8859-10 Western Furnne with rationalised character set for Nordic
- ISO 8859-13 Baltic languages plus Polish
- . 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 CP730 CP737 CP850 CP852 CP855 CP857 CP858 CP860 CP861 MS-Windows character sets: . Windows-1250 for Central European languages that use Latin script. (Polish.
- . Windows-1251 for Cyrillic alphabets
- Windows-1253 for Greek
- Windows-1255 for Hebrey Windows-1256 for Arabic
- Windows, 1258 for Vietnamese
- KOI8-R KOI8-U KOI7
- ISCII

- . JIS X 0208 is a widely deployed standard for Japanese character encoding that has
- Shift JIS (Microsoft Code page 932 is a dialect of Shift JIS) • EUC-JP
- . JIS X 0213 is an extended version of JIS X 0208. Shift JIS-2004
- EUC-JIS-2004 ISO,2022, ID,2004 • GB 2312
- GBK (Microsoft Code page 936)
- Taiwan Bin5 (a more famous variant is Microsoft Code page 950).
- KS X 1001 is a Korean double-byte character encoding standard • ISO-2022-KR
- UTF-16
- ANSEL or ISO/IEC 6937

Let's focus on the main standards

- ASCII
- UTF-8

Common character encodings [edit]

- ISO 646
 ASCII
- ISO 8859
- . ISO 8859-2 Western and Central Europe
- . ISO 8859-4 Western Furnne and Baltic countries (Lithuania, Estonia, Latvia and Lapp)

 ISO 8859-5 Cyrillic alphabet
- ISO 8859-6 Arabic
- ISO 8859.8 Hebrew ISO 8859-10 Western Furnne with rationalised character set for Nordic
- ISO 8859-13 Baltic languages plus Polish
- 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)
- . Windows-1250 for Central European languages that use Latin script. (Polish. . Windows, 1251 for Cyrillic alphabets
- Windows_1253 for Greek
- Windows-1255 for Hebrew
- Windows, 1256 for Arabic Windows, 1258 for Vietnamese
- KOI8-R KOI8-U KOI7

MS-Windows character sets:

• ISCII

- CP437 CP720 CP737 CP850 CP852 CP855 CP857 CP858 CP860 CP861 . JIS X 0208 is a widely deployed standard for Japanese character encoding that has
 - Shift JIS (Microsoft Code page 932 is a dialect of Shift JIS) • EUC-JP

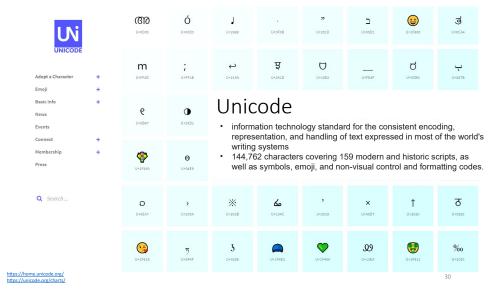
 - . JIS X 0213 is an extended version of JIS X 0208 Shift_JIS-2004
 EUC-JIS-2004
 - ISO.2022, ID.2004
 - GBK (Microsoft Code page 936)
 - . Taiwan Bin5 (a more famous variant is Microsoft Code page 950)
 - KS X 1001 is a Korean double-byte character encoding standard
 - - ANSEL or ISO/IEC 6937

• GB 2312

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

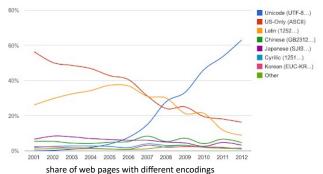
ASCII (7 bits)

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	1	65	41	A	97	61	a
2	2	[START OF TEXT]	34	22		66	42	В	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	1	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	H	104	68	ĥ
9	9	[HORIZONTAL TAB]	41	29)	73	49	1	105	69	i
10	Α	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	В	[VERTICAL TAB]	43	2B	+	75	4B	K	107	6B	k
12	C	[FORM FEED]	44	2C		76	4C	L	108	6C	1
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	1	79	4F	0	111	6F	0
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	P	112	70	р
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	Υ	121	79	У
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	1	124	7C	
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D	1	125	7D	}
30	1E	[RECORD SEPARATOR]	62	3E	>	94	5E	^	126	7E	~
31	1F	[UNIT SEPARATOR]	63	3F	?	95	5F	_	127	7F	[DEL]



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

- 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:
- é (decimal) ⇒ é é (hex) ⇒ é
- o é (mnemonic) ⇒ 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.

http://hcibib.org/multilingual/badchars.htm

Encoding errors

What Could Possibly Go Wrong?

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

A©

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), i, is

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, it you see Å or Å where you do not expect it, there are probably too many UTF-8 encodings. Multiple extra encodings have a pattern to them:

1 é 2 é 3 ÃÃ,© 4 ÃÆ,Ã,Æ,Ã∫'Ã,© 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.

http://hcibib.org/multilingual/badchars.htm

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

Encoding errors

Diagnostic Reference

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

Symptom	Diagnosis
é	no problems
é	too much UTF-8 encoding, or viewing UTF-8 encoded text with Latin-1 encoding
Ãf©	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.

http://hcibib.org/multilingual/badchars.htm

Binary to text encoding

Base64

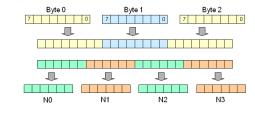


		Table	e 1: The B	ase 64	Alphabet		
.,,							
	Encoding						
0	A	17	R	34	i	51	Z
1	В	18	S	35	j	52	0
2	C	19	T	36	k	53	1
3	D	20	U	37	1	54	2
4	E	21	V	38	m	55	3
5	F	22	W	39	n	56	4
6	G	23	X	40	0	57	5
7	H	24	Y	41	р	58	6
8	I	25	Z	42	9	59	7
9	J	26	a	43	r	60	8
10	K	27	b	44	s	61	9
11	L	28	C	45	t	62	+
12	M	29	d	46	u	63	/
13	N	30	e	47	v		
14	0	31	f	48	W	(pad)	-
15	P	32	g	49	×		
16	Q	33	h	50	у		

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

(underline)

https://datatracker.ietf.org/doc/html/rfc4648

Binary to text encoding

• Base64

8-bit: 00010100 11111011 10011100 00000011 11011001 01: 0-becimal: 5 15 46 28 0 61 37 Output: F P u c A 9 1 1-becimal: 5 15 46 28 0 61 37 Output: F P u c A 9 1 1-becimal: 5 15 46 28 0 61 37 Output: F P u c A 9 1 1-becimal: 0x14fb9c03d9 Hex: 1 4 f b 9 c 0 3 d 9 8-bit: 00010100 11111011 10011100 00000011 11011001 9-ad with 00 6-bit: 000101 001111 101110 011100 000000 11110 100100 Decimal: 5 15 46 28 0 61 36 Output: F P u c A 9 k Input data: 0x14fb9c03			14fb9c0							
6-bit: 000101 001111 101110 011100 000000 111101 100101 Ductmal: 5	Hex:	1 4	f	b 9	C	П	0 3	d	9 7	e
Decimal: 5 15 46 28 0 61 37 Output: F P P u c A 9 1 Input data: 0x14fb9c03d9 Hex: 1 4 f b 9 c 0 3 d 9 8-bit: 00010100 1111011 10011100 00000011 11011001 pad with 00 6-bit: 000101 001111 101110 011100 000000 11101 100100 Decimal: 5 15 46 28 0 61 36 pad with 00 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 00000011 pad with 0000 6-bit: 000101 001111 101110 011100 000000 110000 Decimal: 5 15 46 28 0 48 pad with =	8-bit:	000101	00 1111	1011 10	011100	П	000000	11 1101	1001 0	11111
Output: F P u c A 9 1 Input data: 0x14fb9c03d9 Hex: 1 4 f b 9 c 0 3 d 9 8-bit: 0001010 0111110111 10011100 00000011 11011001 pad with 00 6-bit: 000101 001111 101110 011100 000000 11101 1001100 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 =	6-bit:	000101	001111	101110	011100	П	000000	111101	10010	1 111
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 1111011 10011100 0000001 110000 6-bit: 000101 001111 101110 011100 0000000 110000 Decimal: 5 15 46 28 0 48 pad with =	Decimal:	5	15	46	28		0	61	37	62
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 0 000000 111101 100100 6-bit: 000101 001111 101110 0 000000 111101 100100 0 61 36 pad with 00utput: F P u c A 9 k Input data: 0x14fb9c03 Hex: 1 4 f b 9 c 0 3 8-bit: 00010100 1111011 10011100 00000011 pad with 0000 6-bit: 000101 001111 101110 0 000000 110000 00ccimal: 5 15 46 28 0 48 pad with =	Output:	F	P	u	C		A	9	1	+
8-bit: 00010100 11111011 10011100 00000011 11011001 pad with 00 6-bit: 000101 001111 101110 011100 000000 111101 100100 Decimal: 5	Input da	ta: 0x	14fb9c0	3d9						
6-bit: 000101 001111 101110 011100 000000 111101 100100 Decimal: 5 15 46 28 0 61 36 Dutput: F P u c A 9 k Input data: 0x14fb9c03 Hex: 1 4 f b 9 c 0 3 8-bit: 00010100 1111011 10011100 00000011 pad with 0000 6-bit: 000101 001111 101110 011100 000000 110000 Decimal: 5 15 46 28 0 48 pad with =	Hex:	1 4	f	b 9	C	Т	0 3	d	9	
6-bit: 000101 001111 101110 011100 000000 111101 100100 Decimal: 5 15 46 28 0 61 36 Dutput: 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 =	8-bit:	000101	00 1111	1011 10	011100		000000			а
Decimal: 5 15 46 28 0 61 36 pad with Output: F P u c A 9 k Input data: 0x14fb9c03	6-bit:	000101	001111	101110	011100	ī	000000			
Output: F P u c A 9 k Input data: 0x14fb9c03 Hex: 1 4 f b 9 c 0 3 -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 =	Decimal:	5	15	46	28	ı	0	61	36	
Input data: 0x14fb9c03									ad wit	h =
Hex: 1 4 f b 9 c 0 3 8-bit: 00010100 11111011 10011100 00000011	Output:	F	P	u	C		Α	9	k	=
8-bit: 00010100 11111011 10011100 00000011	Input da	ta: 0x	14fb9c0	3						
pad with 0000 6-bit: 000101 001111 101110 011100 000000 110000 Decimal: 5 15 46 28 0 48 pad with =	Hex:	1 4	f	b 9	C	Т	0 3			
Decimal: 5	8-bit:	000101	00 1111	1011 10	011100	İ)	
pad with =	6-bit:	000101	001111	101110	011100	Т	000000	110000)	
	Decimal:	5	15	46	28					
Output: F P u c A w =							р	ad with	=	=
	Output:	F	P	u	c		Α .	W	=	=

https://datatracker.ietf.org/doc/html/rfc4648 https://www.base64decode.org/

Data Interoperability and Semantics

Part 1. Encoding base data types
Part 1.6. Date and time

Binary to text encoding

- Base64
- Base32
- Base16

		Table	3: The Ba	se 32 Alphabet		
Value	Encoding	Value	Encoding	Value Encoding	Value	Encoding
0	A	9	3	18 S	27	3
1	В	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	0	23 X		
6	G	15	P	24 Y	(pad)	=
7	H	16	0	25 Z		
8	I	17	R	26 2		

			Table	5: The Ba	se 16 /	Alphabet		
ı	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	В	15	F

https://datatracker.ietf.org/doc/html/rfc4648

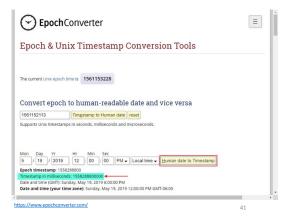
Representing Date and Time is not simple

- 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/

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://cl.mines-stetienne.fr/cps2/course/data

Unix time stamp

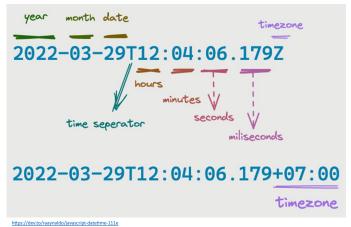
The unix time stamp is a way to track time as a running total of seconds. This count starts at the Unix Epoch on January 1st, 1970 at UTC. Therefore, the unix time stamp is merely the number of seconds between a particular date and the Unix Epoch.



Too many formats



ISO 8601



Data Interoperability and Semantics

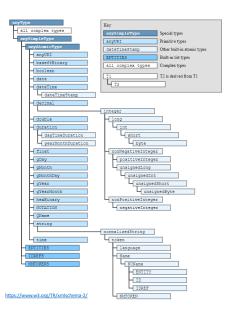
Part 1. Encoding base data types
Part 1.7. XML Schema Datatypes

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://cl.mines-stetienne.fr/cps2/course/data

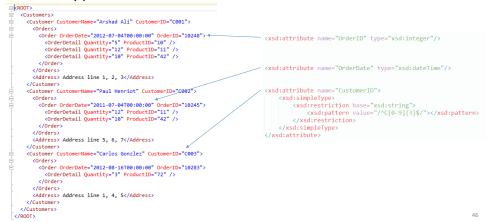
XML Schema Datatypes

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

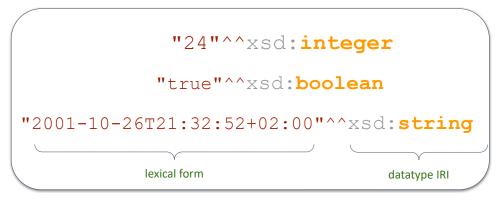
- A value space, which is a set of values.
- A lexical space, which is a set of literals used to denote the values.
- A lexical-to-value mapping
- A small collection of functions, relations, and procedures associated with the datatype.



Goal: type elements and attributes in XML



Anatomy of a XSD literal



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

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

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

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

https://studylib.net/doc/17671880/ieee-754-standard-representation-of-floating-point-number...

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

Country codes

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

ISO 3166 ^[1]		ISO 3166-1 ^[2]				ISO 3166-2 ^[3]		
Country name ^[5]	Official state name ^[6]	Sovereignty ^{[6][7][8]} •	Alpha- 2 ♦ code ^[5]	Alpha- 3 ♦ code ^[5]	Numeric code ^[5] ◆	Subdivision code ♦ links ^[3]	Internet ccTLD ^[9]	
France ^[I]	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	АТ	AUT	040	ISO 3166-2:AT	.at	

Example of country codes
Source: https://en.wikipedia.org/wiki/List_of_ISO_3166_country_code

Language codes

Maxime Lefrançois https://maxime-lefrancois.info

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

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

Data Interoperability

and Semantics

Part 1. Encoding base data types

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

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 ar
Bengali		bn	ben	ben	ben	

Example of language names for the most spoken languages

51

IETF BCP47: Tags for Identifying Languages

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



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.



Example of currency codes

Data Interoperability and Semantics

Part 1. Encoding base data types
Part 1.9. Quantities and Units of measure

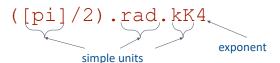
Units of measures: no consensus

- 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





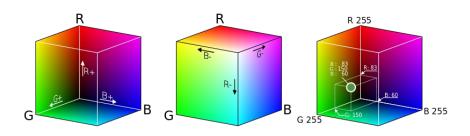
UCUM Code of Unit Concept	EN Unit	EN Symbol	EN Dimension ^a	NCI Concept Code	NCI Term	NCI Abbreviation	SNOMED CT Identifier ^b
[IU]			[arb]	C70497	Anti-Xa Activity International Unit	anti-Xa activity	258997004
Bq	becquerel	Bq	T-1	C42562	Becquerel	Bq	282141004
Bq/g			M-1T-1	C70522	Becquerel per gram	Bq/g	
10^9.[CFU]			[arb]	C68897	Billion Colony Forming Units	Billion CFU	
10^9			1	C71189	Billion Organisms		
m3	cubic metre	m ³	L ³	C42570	Cubic Meter	m ³	396154006
Ci/ml			L-3T-1	C71172	Curie per Millilitre	Ci/ml	
d	day	d	T	C25301	Day	d	258703001
[drp]			L ³	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-3T-1	C71167	Kilobecquerel per Liter	kBq/L	
mmol/l			L-3N	C64387	Millimole per Liter	mmol/L	258813002
[ppm]	part per million	ppm	1	C48523	Part Per Million	ppm	258731005
Pa	pascal	Pa	L-1MT-2	C42547	Pascal	Р	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 ⁻³ M	C48527	Percent Weight Volume	%M/V	396169007
%			1	C48528	Percent Weight Weight	%W/W	118582008
[PFU]			[arb]	C73575	Plaque Forming Unit Equivalent 1000 Mouse LD50	PFU Equivalent 1000 Mouse LD50	
[lb_av]	pound	lb	M	C48531	Pound	LB	258693003
/min	revolution per minute	r.p.m., rev/min, r/min	T-1	C70469	Revolution per Minute	rpm	286549009
[tb'U]			[arb]	C65132	Tuberculin Unit		415758003
[arb'U]{ELISA}				C68875	Enzyme-Linked Immunosorbent Assay Unit	EL. U	

57

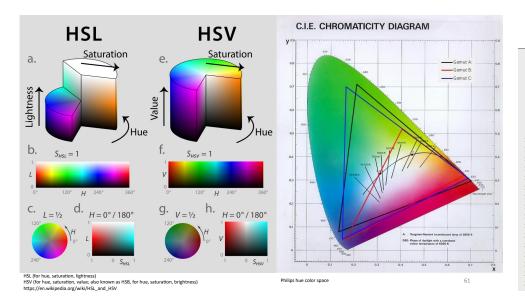
Data Interoperability and Semantics

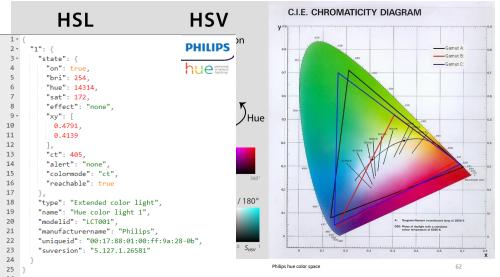
Part 1. Encoding base data types
Part 1.10. Colors

RGB color cube



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://cl.mines-stetienne.fr/cps2/course/data





$HSV \leftrightarrow RGB$ conversion

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

$$C = V \times S$$

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

$$m = V - C$$

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

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

$$R' = R / 255$$

 $G' = G / 255$
 $B' = B / 255$

Cmax = max(R', G', B') Cmin = min(R', G', B) Δ = Cmax – Cmin

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

$$S = \left\{ \begin{array}{l} 0, \ Cmax = 0 \\ \frac{\Delta}{Cmax}, \ Cmax \neq 0 \end{array} \right\}$$

V = Cmax

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-stetlenne.fr/cps2/course/data

< Part 2. Data 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

Data Interoperability and Semantics

Part 2. Data Formats
Part 2.1 Generalities

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 Lefrancois into Systems—Lefrancois Into Systems—Lefra

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

Data Interoperability and SemanticsOutline

- < Part 2. Data formats >
 - Part 2.1 Generalities
 - Part 2.2. Delimiter-separated values
 - Part 2.3. Extensible Markup Language (XML)
 - Part 2.4. JavaScript Object Notation (JSON)
 - Part 2.5. Configuration file formats
 - Part 2.6. YAML Ain't Markup Language (YAML)
 - Part 2.7. Lightweight markup languages
 - Part 2.8. Compressed formats
 - Part 2.9 Multimedia formats

• Part 2.10.3D models 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

File format

standard way that information is encoded for transfer or storage in a computer file

A format is what enables an application to interpret the raw data contained in a file. It is the mode of representation of these data

File formats are marked in the extension of the file name



Contents [hide] 31.1 Video editing production 1 Archive and compressed 20.2 Meteorology 32 Video game data 1.1 Physical recordable media archiving 20.3 Chamietry 33 Video game storage media 2 Computer-aided design 20.4 Mathematics 34 Virtual machines 2.1 Computer-aided design (CAD) 20.5 Biology 34.1 Microsoft Virtual PC, Virtual Server 2.2 Electronic design automation (EDA) 20.6 Biomedical imaging 34.2 EMC VMware ESX, GSX, Workstation, Player 2.3 Test technology 20.7 Biomedical signals (time series) 34.3 VirtualBox 3 Database 20.8 Other biomedical formats 34.4 Parallels Workstation 4 Big Data (Distributed) 20.9 Biometric formats 34.5 QEMU 5 Desktop publishing 21 Programming languages and scripts 35 Web page 6 Document 22 Security 36 Markup languages and other web standards-based formats 7 Financial records 22.1 Certificates and keys 37 Other 7.1 Financial data transfer formats 22.1.1 X.509 37.1 Curso 8 Font file 22.2 Encrypted files 38 Generalized files 9 Geographic information system 22.3 Password files 38.1 General data format: 10 Graphical information organizers 23 Signal data (non-audio) 38.1.1 Text-based 11 Graphics 24 Sound and music 38.2 Generic file extensions 11.1 Color palettes 24.1 Lossless audio 38.2.1 Binary files 11.2 Color management 24.1.1 Uncompressed 38.2.2 Text files 11.3 Raster graphics 24.1.2 Compressed 38.3 Partial files 11.4 Vector graphics 24.2 Lossy audio 38 3.1 Differences and natches 11.5 3D graphics 24.3 Tracker modules and related 38.3.2 Incomplete transfers 12 Links and shortcuts 24.4 Sheet music files 38.4 Temporary files 13 Mathematical 24.5 Other file formats pertaining to audio 39 See also 14 Object code, executable files, shared and dynamically linked libraries 25 Playlist formats 40 Reference 15 Page description language 26 Audio editing and music production 41 External links 16 Personal information manager 27 Recorded television formats 17 Presentation 28 Source code for computer programs 18 Project management software 29 Spreadsheet 19 Reference management software 30 Tabulated data

Open vs unpublished file formats

Open file format

published specification usually maintained by a standards organization.

Linux Information Project: "any format that is published for anyone to read and study but which may or may not be encumbered by patents, copyrights or other restrictions on use"

US government: "An open format is one that is platform independent, machine readable, and made available to the public without restrictions that would impede the re-use of that information"

FR government: loi n° 2004-575 du 21 juin 2004 pour « la confiance dans l'économie numérique »: « On entend par standard ouvert tout protocole de communication, d'interconnexion ou d'échange et tout format de données interopérable et dont les spécifications techniques sont publiques et sans restriction d'accès ni de mise en œuvre. »

Proprietary file formats vs free file formats

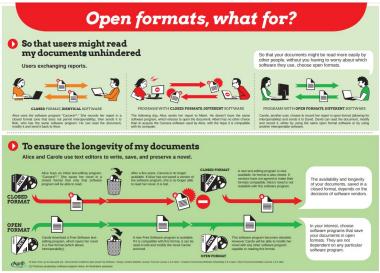
Proprietary file format

Data encoding, s.t. one needs the company's software to decode and interpret it. Either the data encoding specification is secret, or restricted through license.

Examples of proprietary file formats

- mp3: open standard, but subject to patents in some countries
- dwg: non documented, AutoCAD
- psd: documented, Adobe Photoshop's native image format
- · rar: partially documented, archive and compression file format owned by Alexander L. Roshal
- **zip** (newest versions have patented features)
- gif: CompuServe's Graphics Interchange Format, patent expired in 2004
- pdf: open since 2008 ISO 320000-1. Still some features proprietary by Adobe (forms, scripts)
- doc, xls, ppt: formerly closed/undocumented, now Microsoft Open Specification Promise

https://en.wikipedia.org/wiki/Proprietary format

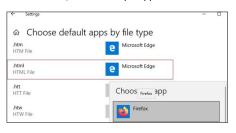


https://www.april.org/en/open-format

How to identify a file type?

Filename extensions

- See https://en.wikipedia.org/wiki/List_of_filename_extensions
- Look up or browse categories https://www.file-extension.info/
- Windows / linux desktops: applications association to filename extension





How to identify a file type?

Filename extensions

- See https://en.wikipedia.org/wiki/List_of_filename_extensions
- Look up or browse categories https://www.file-extension.info/
- Windows / linux desktops: applications association to filename extension
- reason why .htm vs .html ? the 8.3 filename format
- OS hide the extensions + associate applications = creates security issues



How to identify a file type?

File header

may contain metadata about the file and its content

- ex., Exif metadata: image format, size, resolution color space, ...
- ex., AVI header format: https://www.filefix.org/format/avi.html#header

Camera manufacturer	Canon
Camera model	Canon EOS 1200D
Author	Praveen. P
Exposure time	1/60 sec (0.016666666666667)
F-number	f/11
ISO speed rating	200
Date and time of data generation	22:29, 22 November 2018
Lens focal length	41 mm

ex., Exif metadata

Character-based documents may be opened in text editors, and their header interpreted

ex., head of a XML file <?xml version="1.0" encoding="UTF-8"?> <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd"> ex., head of a iCalendar file

VERSION:2.0 PRODID:-//hacksw/handcal//NONSGML_v1.0//EN BEGIN: VEVENT

ex., head of a ISO 10303-21 STEP file

ISO-10303-21; FILE_DESCRIPTION(/* description */ ('A minimal AP214 example wit /* implementation_level */ '2;1'); /* time_stamp */ '2003-12-27T11:57:53', /* author */ ('Lothar Klein'), /* organization */ ('LKSoft'),

How to identify a file type?

File's first bytes: Magic numbers

the first few bytes may be distinctive enough

Hex signature •	ISO 8859-1 •	Offset •	Extension •	Description
23 21	#!	0		Script or data to be passed to the program following the shebang (#I)
52 49 46 46 ?? ?? ?? ?? 41 56 49 20	RIFF????AVIsp	0	avi	Audio Video Interleave video format
FF FB FF F3 FF F2	90 96 98	0	mp3	MPEG-1 Layer 3 file without an ID3 tag or with an ID3v1 tag (which is appended at the end of the file)
49 44 33	ID3	0	mp3	MP3 file with an ID3v2 container
21 3C 61 72 63 68 3E 0A	! <arch>ur</arch>	0	deb	linux deb file
37 7A BC AF 27 1C	7z¼* ' rs	0	7z	7-Zip File Format
1F 88	us <	0	gz tar.gz	GZIP compressed file ^[1]
FD 37 7A 58 5A 00	ý7zXZ _{NIL}	0	xz tar.xz	XZ compression utility using LZMA2 compression
4E 45 53 1A	NESsue	0	nes	Nintendo Entertainment System ROM file[2]

Example of file signatures - https://en.wikipedia.org/wiki/List_of_file_signatures

How to identify a file type?

File's first bytes: Magic numbers

the first few bytes may be distinctive enough

ex., a file that starts with a Byte Order Mark (BOM) tells the system:

- · Highly probable that the text stream is Unicode
- Describes the byte order, or endianness, of the text stream: Big-endian (BE) vs Little-endian (LE)
- Distinguish between Unicode character encoding (UTF-8, UTF-16, UTF-32, ...)

Encoding	Representation (hexadecimal)	Representation (decimal)	Bytes as CP1252 characters
UTF-8 ^[a]	EF BB BF	239 187 191	
UTF-16 (BE)	FE FF	254 255	þÿ
UTF-16 (LE)	FF FE	255 254	ÿþ
UTF-32 (BE)	00 00 FE FF	0 0 254 255	^@^@þÿ (^@ is the null character)
UTF-32 (LE)	FF FE 00 00	255 254 0 0	yp^@^@ (^@ is the null character)

https://en.wikipedia.org/wiki/Byte order mark

How to identify a file type?

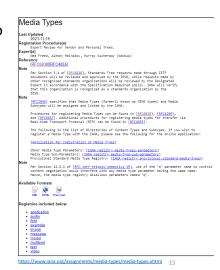
Example: the gzip file format starts with:

- · a 10-byte header, containing:
 - magic number (0x1f8b),
 - the compression method (0x08 for DEFLATE),
 - · 1-byte of header flags,
 - a 4-byte timestamp,
 - · compression flags and
 - · the operating system ID.
- · optional extra headers as allowed by the header flags, including the original filename, and a comment field,

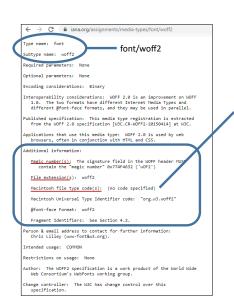
How to identify a file type?

MIME type = Media types

- managed by Internet Assigned Numbers Authority (IANA)
- example: text/html, image/gif, font/woff2, ...
- used by many internet protocols
- demo with your browser: with the Network tab of the developer tools open, see HTTP response headers of requests when accessing https://fonts.google.com/



16



Macintosh file type code(s)

Apple Computer, Inc., "Mac OS: File Type and Creator Codes, and File Formats", Apple Knowledge Article 55381, June 1993, last accessed

https://web.archive.org/web/20080517021842/http://www.info.apple.com/kbnum/n55381

PICT	picture	Stores a PICT image contained in the file
PREF	preference	Stores the environment settings for an application
snd	sound	Stores a sound used in the file
STR	string	Stores a string or hexadecimal data used in the file
STR#	string list	Stores multiple strings used in the file
styl	style	Defines style information, such as the font, color and size of text
TEXT	text	Stores text

https://en.wikipedia.org/wiki/Resource fork#Major resource type



Macintosh Universal Type Identifier code

- · reverse-DNS naming structure
- public UTIs, and organization-specific UTIs
- · multi-heritance



Data Interoperability and Semantics

Part 2. Data Formats

Part 2.2. Delimiter-separated values

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

Delimiter-separated values

separator character

comma

Comma-separated values (CSV) ex: Microsoft Excel csv export1

semicolon

ex: Microsoft Excel csv export2

tab

Tab-separated values (TSV)

colon

ex: Linux \$ cat /etc/passwd

¹ unless the decimal point of the locale is a comma (ex., France)

² when the decimal point of the locale is a comma

persons.csv - Notepad File Edit Format View Help Family Name, Given Name, VIAF ID Ackersdijck, Willem Cornelis, 17959345 Adelung, Friedrich von, 22963658 Afzelius, Arvid August, 49972119 Amerling, Karel, 13331054 Anton, Karl Gottlob von, 183632821 Arwidsson,Adolf Ivar,8184878 Asbjørnsen, Peter Christen, 116587918 Attems, Heinrich, 37665468 Atterbom, Per Daniel Amadeus, 46819248 Balabin, Viktor Petrovich, 44473845 Banks, Joseph, 46830189 Beck, Friedrich, 44338671 Becker, Reinhold von, 42101066 Bernhart, Johann Baptist, 69674335 Bertram, Johann, 32890043 Bilderdijk, Willem, 14882166 Boisserée, Sulpiz, 7483155 Bopp, Franz, 61614118 Borovský, Karel Havlíček, 100277614 Bosković, Jovan, 161354270 Buslaev, Fyodor, 10074560 Cenowa, Florian Stanislaw, 44466031 Chomiakov.Aleksei.66492873

Delimiter-separated values

RFC 4180

separator character

records separated by CRLF

aaa,bbb,ccc CRLF zzz,yyy,xxx CRLF or

aaa,bbb,ccc CRLF zzz,yyy,xxx

optional header file

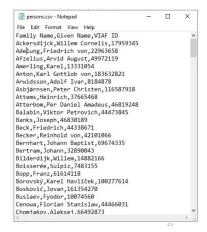
field_name,field_name,field_name CRLF
aaa,bbb,ccc CRLF
zzz,yyy,xxx CRLF

fields enclosed in double quote

"aaa","b CRLF bb","ccc" CRLF zzz,yyy,xxx

escaping cell character

"aaa","b""bb","ccc"



Delimiter-separated values

• Comparison Java libraries:

https://github.com/uniVocity/csv-parsers-comparison

- . CPU: AMD Ryzen 7 1700 Eight-Core Processor @ 4.0 GHz
- RAM: 32 GB
- Storage: 1TB SSD drive
- OS: Arch Linux 64-bit
- JDK: 9.0.4 64-bit (Linux)
- JDK: 1.8.0_144 64-bit (Linux)
- JDK: 1.7.0_80 64-bit (Linux)
- JDK: 1.6.0_45 64-bit (Linux)

Processing 3,173,958 rows of non RFC 4180 compliant input. No quoted values. DK 9								
Parser	Average time	% Slower than best	Best time	Worst time				
uniVocity CSV parser	739 ms	Best time!	707 ms	768 ms				
SimpleFlatMapper CSV parser	861 ms	16%	848 ms	901 ms				
Jackson CSV parser	1212 ms	64%	1169 ms	1238 ms				
Product Collections parser	1409 ms	90%	1389 ms	1451 ms				
Java CSV Parser	1498 ms	102%	1490 ms	1508 ms				
JCSV Parser	1681 ms	127%	1660 ms	1710 ms				
Oster Miller CSV parser	1772 ms	139%	1762 ms	1780 ms				
Gen-Java CSV	1799 ms	143%	1790 ms	1805 ms				
Simple CSV parser	1861 ms	151%	1832 ms	1900 ms				
SuperCSV	1893 ms	156%	1858 ms	1964 ms				
OpenCSV	2022 ms	173%	2007 ms	2037 ms				
Apache Commons CSV	2424 ms	228%	2409 ms	2442 ms				
Way IO Parser	2577 ms	248%	2532 ms	2638 ms				

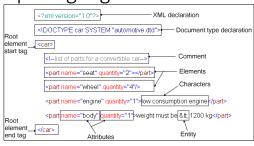
Data Interoperability and Semantics

Part 2. Data Formats

Part 2.3. Extensible Markup Language (XML)

Extensible Markup Language





- v1.0 in 1998, still extensively used in many verticals
- numerous formats based on XML (418 registered on IANA)
 https://en.wikipedia.org/wiki/List of XML markup languages
 application/atom+xml
 application/rdf+xml
- verbosity, complexity and redundancy

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-lefrançois.info_ Course unit URL: https://ci.mines-stetlenne.fr/cps2/course/data

Extensible Markup Language

XML (file format) Filename extension application/xml text/xml [1 Identifier (UTI) UTI conformation <?xml Markup language Extended from Extended to Numerous languages, including XHTML · RSS · Atom · KML Standard (November 26, 2008; 12 years ago) Open format?

Characters and escaping

- unicode implementations: <?xml version="1.0" encoding="UTF-8"?>
- escaping characters: &1t; '<' & '&' ❤ '♥' etc.

Syntactical correctness

- well formed vs ill-formed
- one root tag
- correct nesting
- tag names (approx) start with letter, then alphanumeric or ':'

Schemas and validation

- valid vs invalid
- DTD. or XML Schema

Namespaces

- xmlns:ns1="http://example.org/ns1" xmlns:ns2="http://example.org/ns2"
- allows to use different schemas together: <ns1:Tag> <ns2:Tag>

Type of format Markup language Extended from Extended to Numerous languages, including XHTML · RSS · Atom · KML

Filename

extension

Uniform Type

Identifier (UTI)

Magic number

UTI conformation

Internet

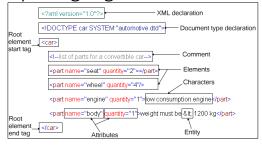
XML (file format) application/xml

(November 26, 2008; 12 years ago) 1.1 (Second Edition)呼 (August 16, 2006; 15 years ago) Open format?

text/xml [1]

<?xml

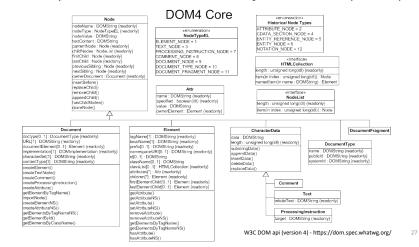
Extensible Markup Language



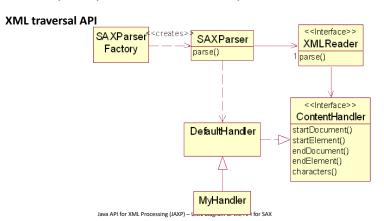
Rules of thumb for modelling with XML:

- · limit the number of tag types and attribute types!
- use attributes for simple datatypes only
- order of elements is important
- group elements when appropriate
- · anticipate how you are going to write XPath queries

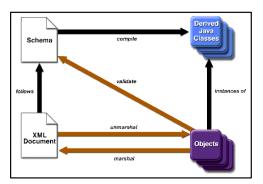
XML manipulation: Document Object Model (DOM)



SAX (Simple API for XML)



XML – OOP data binding



XML-OOP data binding ex Java: https://zetcode.com/java/jaxb/

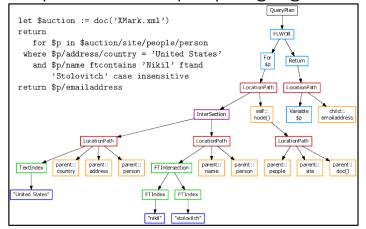


Java Architecture for XML Binding (JAXB) example https://howtodoinjava.com/jaxb/jaxb-annotations/

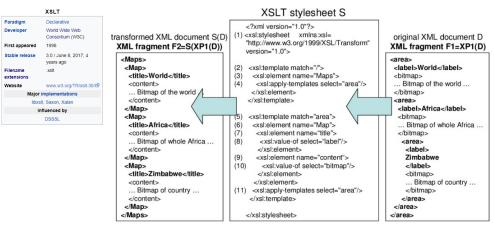
XPath: Declarative tree-traversal language



XQuery: Declarative query language



XSLT: Declarative transformation language



Part 2. Data Formats Part 2.4. JavaScript Object Notation (JSON)

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

Filename extension Internet Type code Uniform Type Identifier (UTI STD 90 (RFC 8259 (₽) 21778:2017 Website

since early 2000s. Now used a lot in software development

JavaScript Object Notation

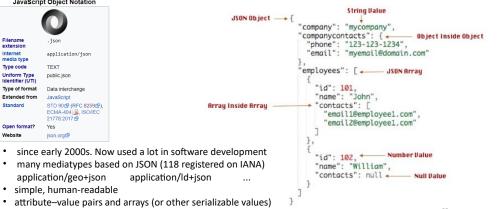
- many mediatypes based on JSON (118 registered on IANA) application/geo+json application/ld+json
- · simple, human-readable
- attribute-value pairs and arrays (or other serializable values)

JavaScript Object Notation



application/geo+json

· simple, human-readable



JavaScript Object Notation

JavaScript Object Notation Internet media type TEXT Uniform Type Identifier (UTI) public.json Type of format Data interchange Extended from JavaScript Standard STD 90r@ (REC 8259r@) 21778:2017 Open format?

Characters and escaping

- · full Unicode character set
- escaping characters \b \f \n \r \t \" \

Data types

- number
- string
- boolean
- array
- object
- null

Semantics

While JSON provides a syntactic framework for data interchange, unambiguous data interchange also requires agreement between producer and consumer on the semantics of specific use of the JSON syntax. One example of where such an agreement is necessary is the serialization of data types defined by the JavaScript syntax that are not part of the JSON standard, e.g., Date, Function, Regular Expression, and "undefined"

JSON Schema



JSONPath: Declarative tree-traversal language

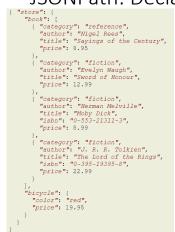
| XPath | JSONPath | Description |
|-------|------------------|---|
| / | \$ | the root object/element |
| | @ | the current object/element |
| / | . or [] | child operator |
| | n/a | parent operator |
| // | | recursive descent. JSONPath borrows this syntax from E4X. |
| * | * | wildcard. All objects/elements regardless their names. |
| @ | n/a | attribute access. JSON structures don't have attributes. |
| [] | 0 | subscript operator. XPath uses it to iterate over element collections and for <u>predicates</u> . In Javascript and JSON it is the native array operator. |
| ı | [.] | Union operator in XPath results in a combination of node sets. JSONPath allows alternate names or array indices as a set. |
| n/a | [start:end:step] | array slice operator borrowed from ES4. |
| [] | ?() | applies a filter (script) expression. |
| n/a | () | script expression, using the underlying script engine. |
| () | n/a | grouping in Xpath |

Examples:

\$.store.book[0].title
\$['store']['book'][0]['title']
\$.store.book[(@.length-1)].title
\$.store.book[?(@.price < 10)].title</pre>

https://goessner.net/articles/JsonPath/

JSONPath: Declarative tree-traversal language



| XPath | JSONPath | Result |
|-------------------------------------|-------------------------------------|--|
| /store/book/author | \$.store.book[*].author | the authors of all books in the store |
| //author | \$author | all authors |
| /store/* | \$.store.* | all things in store, which are some books and a red bicycle. |
| /store//price | \$.storeprice | the price of everything in the store. |
| //book[3] | \$book[2] | the third book |
| //book[last()] | \$book[(@.length-1)]
\$book[-1:] | the last book in order. |
| <pre>//book[position() <3]</pre> | \$book[0,1]
\$book[:2] | the first two books |
| //book[isbn] | \$book[?(@.isbn)] | filter all books with isbn
number |
| //book[price<10] | \$book[?(@.price<10)] | filter all books cheapier than 10 |
| //* | ş* | all Elements in XML
document. All members of
JSON structure. |

https://goessner.net/articles/JsonPath/

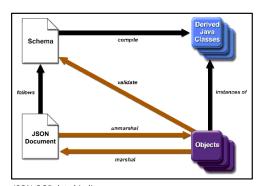
APIs for JSON Processing

example for java: https://javaee.github.io/jsonp/



40

JSON – OOP data binding



JSON-OOP data binding

| 21/2 |
|------|
| ava |
| |





| Parsing Speed | Speed. MB/MS | Parsing Time |
|---------------|--------------|--------------|
| GSON | 100% | 0% |
| Jackson | 58% | 70.87% |
| JSON.simple | 79% | 126.58% |
| JSONP | 44% | 25.49% |

Comparison of Java JSON libraries: Jackson vs. Gson vs. JSON-B vs. JSON-P vs. org.JSON vs. Jsonpath

https://itsallbinary.com/jackson-vs-gson-vs-json-bvs-ison-p-vs-org-ison-vs-isonpath-iava-ison-librari es-features-comparison/

Data Interoperability and Semantics

Part 2. Data Formats Part 2.5. Configuration file 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

.properties files

.properties Filename extension

Mainly for Java-related conf files Ex: i18n1 and l10n2

- key=value,
- key = value,
- key:value,
- key value
- 1i18n: internationalization ² I10n: localization



INI files



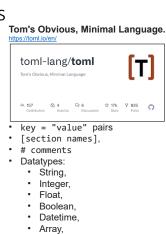
https://en.wikipedia.org/wiki/INI_file

- Developed for Windows for initialization: BOOT.INI. SYSTEM.INI. WIN.INI. ...
- · Similar in Linux systems with .conf, .cfg, ...
- · Syntax used by many programs
 - php.ini
 - · ssh.conf
 - · git.conf

: Last modified 1 April 2001 by John Doe [owner] name = John Doe organization = Acme Widgets Inc. [database] ; use IP address in case network name resolution is not working server = 192.0.2.62 port = 143 file = "payroll.dat

[section] domain = wikipedia.org [section.subsection] foo = bar





· and Table





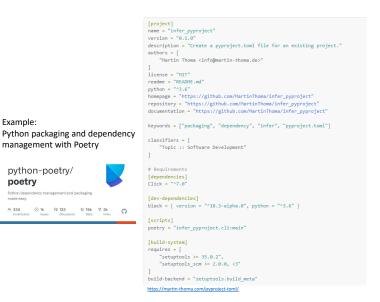
Example:

poetry

management with Poetry

A. 334 ⊙ 1k ♀ 133 ☆ 19k ♀ 2k Contributors Issues Discussions Stars Folia

python-poetry/



Data Interoperability and Semantics

Part 2. Data Formats

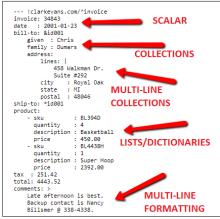
Part 2.6. YAML Ain't Markup Language (YAML)

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

YAML Ain't Markup Language (YAML)



- since early 2000s
- · used a lot for software configuration
- human-readable, superset of JSON



YAML Ain't Markup Language (YAML)



- · Whitespace indentation, tabs forbidden
- end-of-line comments (#)
- lists

```
--- # Favorite movies
- Casablanca
- North by Northwest
- The Man Who Wasn't There
```

```
Or --- # Shopping list [milk, pumpkin pie, eggs, juice]
```

associative arrays

```
--- # Indented Block
name: John Smith
age: 33
--- # Inline Block
{name: John Smith, age: 33}
```

YAML Ain't Markup Language (YAML)



strings

```
data: |
There once was a tall man from Ealing
Who got on a bus to Darjeeling
It said on the door
"Please don't sit on the floor"
So he carefully sat on the ceiling
```

data: >
 Wrapped text
 will be folded
 into a single
 paragraph

Blank lines denote

paragraph breaks

datatypes

```
a: 123 # an integer
b: "123" # a string, disambiguated by quotes
c: 123.0 # a float
d: !!float 123 # also a float via explicit data type prefixed by (!!)
e: !!str 123 # a string, disambiguated by explicit type
f: !!str Yes # a string via explicit type
g: Yes # a boolean True (yamll.1), string "Yes" (yamll.2)
h: Yes we have No bananas # a string, "Yes" and "No" disambiguated by context.
```

50

YAML Ain't Markup Language (YAML)



anchors and references

```
- step: &id001
                                 # defines anchor label &id001
    instrument:
                     Lasik 2000
    pulseEnergy:
                     5.4
                     12
    pulseDuration:
    repetition:
                     1000
    spotSize:
                     1mm
 step: &id002
   instrument:
    pulseEnergy:
    pulseDuration:
                     500
    repetition:
   spotSize:
- step: *id001
                                 # refers to the first step (with anchor &id001)
- step: *id002
                                 # refers to the second step
- step: *id002
```

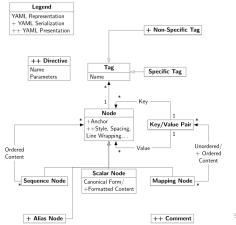
YAML Ain't Markup Language (YAML)

```
https://yaml.org/refcard.html
%YAML 1.1 # Reference card
                                                                                                Misc indicators:
' #' : Throwaway comment indicator.
 Collection indicators:
                                                                                                    '`@' : Both reserved for future use
           Value indicator.
           Nested series entry indicator.
Separate in-line branch entries.
                                                                                                 Special keys:
                                                                                                    '=' : Default "value" mapping key.
'<<' : Merge keys from another mapping.
           Surround in-line series branch.
Surround in-line keyed branch.
                                                                                                Core types: # Default automatic tags.
'!!map' : { Hash table, dictionary, mapping }
 Scalar indicators:
          : Surround in-line unescaped scalar ('' escaped ')
                                                                                                    '!!seq' : { List, array, tuple, vector, sequence }
                                                                                                    '!!str' : Unicode string
        : Surround in-line escaped scalar (see escape codes below).
: Block scalar indicator.
                                                                                                 More types:
                                                                                                     '!!set' : { cherries, plums, apples }
         · Folded scalar indicator
        : Strip chomp modifier ('|-' or '>-')
: Keep chomp modifier ('|+' or '>+').
                                                                                                 '!!omap': [ one: 1, two: 2 ]
Language Independent Scalar types
  1-9 : Explicit indentation modifier ('|1' or '>2').
                                                                                                       ~. null }
                                                                                                                                    : Null (no value)
                                                                                                      1234, 0x4D2, 02333 ]
                                                                                                                                      [ Decimal int, Hexadecimal int, Octal int ]
            # Modifiers can be combined ('|2-', '>+1').
                                                                                                      1 230.15, 12.3015e+02 ]: [ Fixed float, Exponential float ] .inf, -.Inf, .NAN ] : [ Infinity (float), Negative, Not a
 Alias indicators:
        : Anchor property
                                                                                                      Y, true, Yes, ON
                                                                                                                                      Boolean true
                                                                                                       n, FALSE, No, off }
                                                                                                                                   : Boolean false
 Tag property: # Usually unspecified.
             rty: # usually unspectried.
: Unspecified tag (automatically resolved by application).
: Non-specific tag (by default, "!!map"/"!!seq"/"!!str").
: Primary (by convention, means a local "!foo" tag).
                                                                                                    ? !!binary >
R01G...BADS=
                                                                                                        Base 64 binary value
   '!!foo': Secondary (by convention, means "tag:yaml.org,2002:foo").
'!h!foo': Requires "%TAG !h! <prefix>" (and then means "<prefix>foo").
                                                                                                 Escape codes:
                                                                                                 '!<foo>': Verbatim tag (always means "foo").
Document indicators:
  '%' : Directive indicator
          : Document header
  '...': Document terminator.
```

51

YAML Ain't Markup Language (YAML)





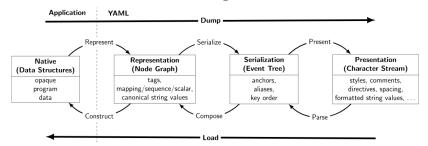
Data Interoperability and Semantics

Part 2. Data Formats

Part 2.7. Lightweight markup 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-stethenne.fr/cps2/course/data

APIs for YAML Processing



Tutorial for Java:

https://www.baeldung.com/java-snake-yaml

Package for JavaScript:

https://www.npmjs.com/package/yaml

Module for Python:

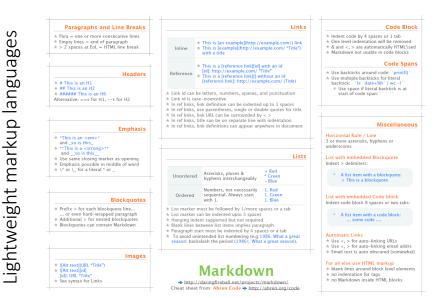
https://pyyaml.org/wiki/PyYAMLDocumentation

5

Comparing language features + HTML export tool + HTML import tool + Tables + Link titles + class attribute + id attribute + Release date 4

| Ų, | Language | TITIBLE EXPORT TOOL + | H TWL IIIport tool • | lables + | LIIIK UUES ¥ | Class attribute • | Tu attribute ▼ | Release date ¥ |
|--------------------------------------|--------------------------|-----------------------|----------------------|---------------------|--------------|-------------------|----------------|------------------------------|
| ള | AsciiDoc | Yes | Yes | Yes | Yes | Yes | Yes | 2002-11-25 ^[1] |
| ag | BBCode | No | No | Yes | No | No | No | 1998 |
| ĭ | Creole | No | No | Yes | No | No | No | 2007-07-04 ^[2] |
| <u> </u> | Gemtext | Yes | ? | No | Yes | No | No | 2020 |
| | GitHub Flavored Markdown | Yes | No | Yes | Yes | No | No | 2011-04-28+ |
| ש | Jira Formatting Notation | Yes | No | Yes | Yes | No | No | 2002+[3] |
| _ | Markdown | Yes | Yes | No | Yes | Yes/No | Yes/No | 2004-03-19 ^{[4][5]} |
| $\stackrel{\hookrightarrow}{\vdash}$ | Markdown Extra | Yes | Yes | Yes ^[6] | Yes | Yes | Yes | 2013-04-11 ^[7] |
| rkup | Marked Text® | Yes | No | Yes | No | No | No | 2021-01 |
| | MediaWiki | Yes | Yes | Yes | Yes | Yes | Yes | 2002 ^[8] |
| Ja | MultiMarkdown | Yes | No | Yes | Yes | No | No | 2009-07-13 |
| Ξ | Org-mode | Yes | Yes ^[9] | Yes | Yes | Yes | Yes | 2003 ^[10] |
| $\boldsymbol{\downarrow}$ | PmWiki | Yes ^[11] | Yes | Yes | Yes | Yes | Yes | 2002-01 |
| | POD | Yes | ? | No | Yes | ? | ? | 1994 |
| . <u>∞</u> 0 | re Structured Text | Yes | Yes ^[9] | Yes | Yes | Yes | auto | 2002-04-02 ^[12] |
| ھ | Slack | No | No | No | Yes | No | No | 2013+[13][14] |
| ≥ | TiddlyWiki | Yes | No | Yes | Yes | Yes | No | 2004-09 ^[15] |
| <u> </u> | Textile | Yes | No | Yes | Yes | Yes | Yes | 2002-12-26 ^[16] |
| Lightwe | Теху | Yes | Yes | Yes | Yes | Yes | Yes | 2004 ^[17] |
| ٠ <u>ـ</u> | txt2tags | Yes | Yes ^[18] | Yes ^[19] | Yes | Yes/No | Yes/No | 2001-07-26 ^[20] |
| _ | WhatsApp | No | No | No | No | No | No | 2016-03-16 ^[21] |

| HTML output | <pre>strongly emphasized</pre> | emphasized text | | <code>code</code> | semantic | |
|-----------------------------|---|------------------------------------|-----|-----------------------------------|---|--|
| ri ini. output | b>bold text | <pre><i>i>italic text</i></pre> | ٠ | <tt>monospace text</tt> | presentational | |
| AsciiDoc | "bold text" | 'italic text' | | +monospace text+ | Can double operators to apply formatting where there is no word | |
| AsciiDoc | -both text- | _italic text_ | | 'monospace text' | example **b**old t**ex**t yields bold lext). | |
| ATX | *bold text* | _italic text_ | | monospace text | email style | |
| BBCode | [b]bold text[/b] | [i]italic text[/i] | | [code]monospace text[/code] | Formatting works across line breaks. | |
| Creole | **bold text** | //italic text// | | {{{monospace text}}} | Triple curly braces are for nowiki which is optionally monospace. | |
| Gemtext | N/A | N/A | | ""alt text
monospace text | Text immediately following the first three backticks is all-text. | |
| Jira Formatting
Notation | *bold text* | _italic text_ | | {{monospace text}} | | |
| Markdown ^[42] | **bold text** | *italic text* | [] | | | |
| Markdown | _bold text_ | _italic text_ | | 'monospace text' | semantic HTML tags | |
| Marked Texts9 | **bold text** | //italic text// | | ;;monospace text;; | semantic HTML tags | |
| MediaWiki | ""bold text"" | ''italic text'' | | <code>monospace text</code> | mostly resorts to inline HTML | |
| | | | | -code- | | |
| Org-mode | *bold text* | /italic text/ | | ~verbatim~ | | |
| PmWiki | '''bold text''' | ''italic text'' | | @@monospace text@@ | | |
| reST | **bold text** | *italic text* | | ''monospace text'' | | |
| Setext | **bold text** | ~italic text~ | | N/A | | |
| Textile ^[43] | *strong* | _emphasis_ | | | semantic HTML tags | |
| lexule | **bold text** | italic text | | @monospace text@ | presentational HTML tags | |
| | **bold text** | *italic text* | | | semantic HTML tags by default, optional support for presentation | |
| Texy! | | //italic text// | | `monospace text` | | |
| TiddlyWiki | | | | 'monospace text' | | |
| Hadiywiki | ''bold text'' | //italic text// | | ''monospace text'' | | |
| txt2tags | **bold text** | //italic text// | | ''monospace text'' | | |
| POD | B <bold text=""></bold> | I <italic text=""></italic> | | C <monospace text=""></monospace> | Indented text is also shown as monospaced code. | |
| Slack | *bold text* | _italic text_ | | 'monospace text' | ""block of monospaced text"" | |
| WhatsApp | *bold text* | _italic text_ | | ""monospace text"" | | |



Part 2. Data Formats Part 2.8. Compressed formats

Compression

S

نة

 σ

A compressed file is an archive that contains one or more file that have been reduced in size.

many different file compression types: zip, arc, arj, rar, cab, tar.gz, ...

lossless file compression

reduce redundancy without loosing data ex. AAABBBBBCC -> A3B5C2

lossy file compression

ok to loose some data ex. video, audio, images

How compression works



Compression - Computerphile

Computerphile ❷ 384K views • 8 years ago

Most of us deal with data compression on a daily basis, but what is it and how does it work? Professor David Brailsford introduces compression with regards to text and pictures. http://www.faceboo.



Elegant Compression in Text (The LZ 77 Method) - Computerphile

Computerphile ❷ 436K views • 8 years ago

Text compression methods such as LZ can reduce file sizes by up to 80%. Professor Brailsford explains the nuts and bolts of how it is done. Original Compression film: http://youtu.be/Lto-ajuqW3w...

How Huffman Trees Work - Computerphile

Computerphile **②** 225K views • 8 years ago

How do we derive the most compact codes for a situation? Huffman Trees can help. Professor Brailsford explains how computer scientists like their trees to be upside down. "Entropy in Compression.

Inflate gzip Filename .gz extension ${\it application/gzip}^{[2]}$ Internet media type Uniform Type org.gnu.gnu-zip-archive Identifier (UTI) Magic number Developed by Jean-loup Gailly and Mark Adler Type of format Data compression Open format?

Website

Deflate

ZIP file format Filename .zip .zipx application/zip [1] media type com.pkware.zip-archive Identifier (UTI Magic number PK\x03\x04 PK\x05\x06 (empty) PK\x07\x08 (spanned) PKWARE Inc Developed by 14 February 1989; 32 Initial release years ago Latest release 6.3.9 (15 July 2020; 15 months ago) Data compression Extended to JAR (EAR, RAR (Java), WAR) Office Open XML (Microsoft) Open Packaging Conventions OpenDocument (ODF) XPI (Mozilla extensions) APPNOTEr₽ from Standard PKWARE ISO/IEC 21320-1:2015 (a subset of ZIP file format

Open format?



Filename application/x-tar media type public.tar-archive Identifier (UTI

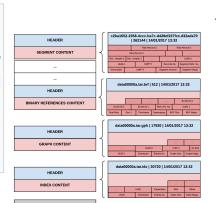
ustar\000 at byte offset 257 (for POSIX versions)

ustar \040 \040 \0 (for old GNU tar format)[1] absent in pre-POSIX versions

Latest release Type of

various (various) File archiver POSIX since POSIX.1, presently in the definition of

pax[1]r



EMPTY

gzip Filename Internet application/gzip^[2] Uniform Type org.gnu.gnu-zip-archive Identifier (UTI) Magic number Developed by Jean-loup Gailly and Mark Adler Type of format Data compression gzip.org (obsolete)

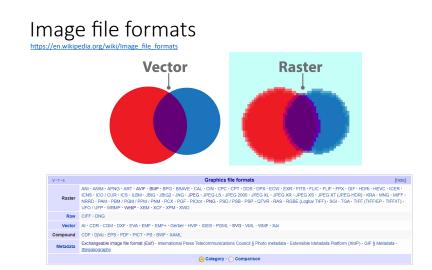
Data Interoperability and Semantics

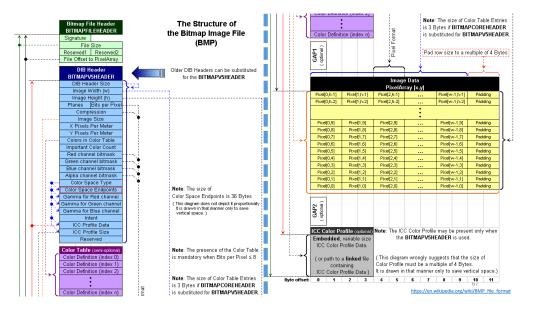
gzip.org ☐ (obsolete)

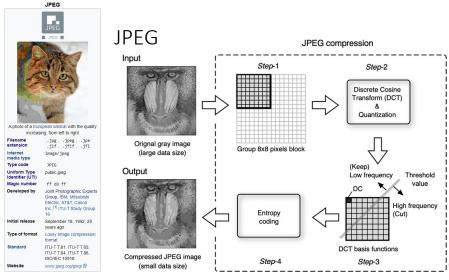
Part 2. Data Formats Part 2.9 Multimedia formats

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

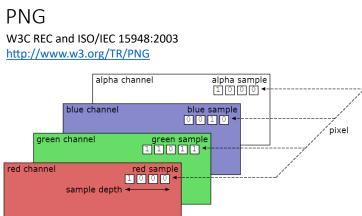
| ·T•E | | Multimedia compression and container formats | [hide] | | | |
|----------------------|--|---|------------------|--|--|--|
| Video
compression | ISO, IEC,
MPEG | | art 2 / HEVC) · | | | |
| | ITU-T, VCEG | H.120 · DCT (H.261 · H.262 · H.263 · H.264 / AVC · H.265 / HEVC · H.266 / VVC · DV) | | | | |
| | SMPTE | VC-1 · VC-2 · VC-3 · VC-5 · VC-6 | | | | |
| | TrueMotion | TrueMotion S · DCT (VP3 · VP6 · VP7 · VP8 · VP9 · AV1) | | | | |
| | Others | Apple Video - AVS - Bink - Cinepak - Daala - DVI - FFV1 - Huffyor - Indeo - Lagarth - Microsoft Video 1 - MSU Lossless - OMS V
ProRes (422 - 4444) - Lightfilme (Animation - Graphics) - RealVideo - Strucker - Sorenson Video Spart
VMV - XEB - YULS | | | | |
| | | MPEG-1 Layer II (Multichannel) · MPEG-1 Layer I · MPEG-1 Layer III (MP3) · AAC (HE-AAC · AAC-LD) · MPEG Surround · MPEG-4 SLS · MPEG-4 DST · MPEG-4 HVXC · MPEG-4 CELP · MPEG-D USAC · MPEG-H 3D Audio | 3-4 ALS · | | | |
| | ITU-T | G.711 (A-law · μ-law) · G.718 · G.719 · G.722 · G.722.1 · G.722.2 · G.723 · G.723.1 · G.726 · G.728 · G.729 · G.729.1 | | | | |
| Audio | IETF | Opus · iLBC · Speex · Vorbis | | | | |
| compression | 3GPP | AMR · AMR-WB · AMR-WB+ · EVRC · EVRC-B · EVS · GSM-HR · GSM-FR · GSM-EFR | | | | |
| | ETSI | AC-3 · AC-4 · DTS | | | | |
| | Others | ACELP - ALAC - Asao - ATRAC - AVS - CELT - Codec 2 - DRA - FLAC - ISAC - MELP - Monkey's Audio - MT9 - Musepack - Optin QCELP - RCELP - RealAudio - RTAudio - SBC - SD2 - SHN - SILK - Siren - SMV - SVOPC - TTA (True Audio) - TwinVQ - VMR-W WavPack - VMA - MDA - aptX - aptX + D0 - aptX tow Latency - aptX Adaptive - LDAC - LHDC - LLAC | | | | |
| Image | | IETF, CCITT Group 4 - DCT (HEIC - HEVC - JPEG · JPEG XL - JPEG XR - JPEG XT - TIFF/EP) - Arithmetic (JBIG - JBIG2) - JPEG | EG-LS - JPEG XS | | | |
| compression | 01 | thers APNG · BPG · DCT (AVIF · AV1) · DJVu · EXR · FLIF · ICER · MNG · PGF · QTVR · WBMP · WebP | | | | |
| | | MPEG-ES (MPEG-PES) · MPEG-PS · MPEG-TS · ISO/IEC base media file format · MPEG-4 Part 14 (MP4) · Motion JPEG 2000 · MPEG media transport | MPEG-21 Part 9 - | | | |
| | ITU-T | H.222.0 · T.802 | | | | |
| Containers | IETF | RTP · Ogg | | | | |
| | SMPTE | GXF · MXF | | | | |
| | Others | 3GP and 3G2 · AMV · ASF · AIFF · AVI · AU · BPG · Bink (Smacker) · BMP · DivX Media Format · EVO · Flash Video · HEIF · IFF Matroska (WebM) · QuickTime File Format · RatDVD · RealMedia · RIFF (WAV) · MOD and TOD · VOB, IFO and BUP | · M2TS · | | | |
| ollaborations | NETVC · MPEC | G LA • HEVC Advance • Alliance for Open Media | | | | |
| Methods | Discrete cosine transform (DCT - MDCT) - Entropy (Arithmetic - Huffman - Modified) - FFT - LPC (ACELP - CELP - LSP - WLPC) - Lossiess - Lossy - LZ (DEFLATE - LZW) - PCM (A-law - µ-law - ADPCM - DPCM) - Transform - Wavelet (Daubechies - DWT - Transform) | | | | | |
| Lists | Comparison of | audio coding formats - Comparison of video codecs - List of codecs | | | | |
| | | See Compression methods for techniques and Compression software for codecs | | | | |
| uthority contro | l: National librar | ries / France면 (data)면 https://en.wikipedia.org/wiki/1 | Template:Compre | | | |

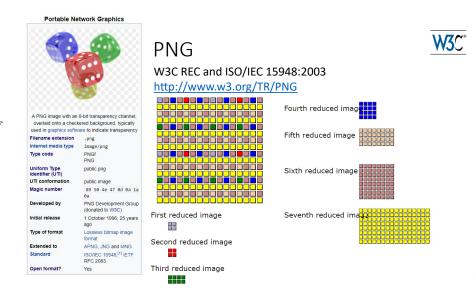














Modern image file formats: WebP

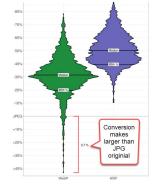
- WebP lossless images are 26% smaller in size compared to PNGs.
- WebP lossy images are 25-34% smaller than comparable JPEG images at equivalent structural similarity index measure quality index.

| Image | File Name | Original Size | Compressed JPG | WebP Format |
|-------|-------------------|---------------|----------------|-------------|
| - | jpg-to-webp-1.jpg | 480 KB | 407 KB | 43 KB |
| | jpg-to-webp-2.jpg | 659 KB | 578 KB | 113 KB |
| | jpg-to-webp-3.jpg | 787 KB | 715 KB | 127 KB |
| | jpg-to-webp-4.jpg | 617 KB | 543 KB | 61 KB |

Modern image file formats: AVIF



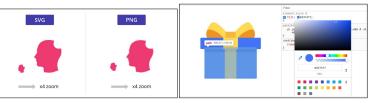
- In most cases, AVIF generates a smaller image payload than WebP.
- In a few cases, WebP generates a larger version than the original JPEG.
- ImageEngine will recognize and deliver the most format with the smallest possible payload.



https://imageengine.io/blog/how-efficient-is-avif/







What is an SVG File Used For and Why Developers Should be Using Them # Published Jan 19, 2021 - https://deliciousbrains.com/svg-advantages-developers/

Recommendations for images in the browser

https://developer.mozilla.org/en-US/docs/Web/Media/Formats

Photographs

AVIF*, WebP, or JPEG

Icons

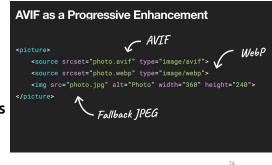
SVG, Lossless WebP, or PNG

Screenshots

Lossless WebP or PNG

Diagrams, drawings, and charts

SVG, PNG



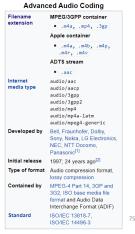
* Ladded AVIE here

Audio file formats

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







Video file formats



Open format?



Recommendations for audio and video files

Audio-only files

| and only and | | | | | |
|---|--------------------------------------|--|--|--|--|
| If you need | Consider using this container format | | | | |
| Compressed files for general-purpose playback | MP3 (MPEG-1 Audio Layer III) | | | | |
| Losslessly compressed files | FLAC with ALAC fallback | | | | |
| Uncompressed files | WAV | | | | |

Video files

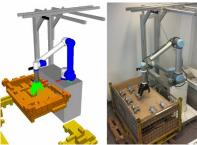
| If you need | Consider using this container format | | | |
|---|---|--|--|--|
| General purpose video, preferably in an open format | WebM (ideally with MP4 fallback) | | | |
| General purpose video | MP4 (ideally with WebM or Ogg fallback) | | | |
| High compression optimized for slow connections | 3GP (ideally with MP4 fallback) | | | |
| Compatibility with older devices/browsers | QuickTime (ideally with AVI and/or MPEG-2 fallback) | | | |

Data Interoperability and Semantics

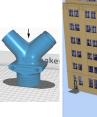
Part 2. Data Formats Part 2.10 3D models

(focusing on Web and Cyber-Physical Systems application domains)

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

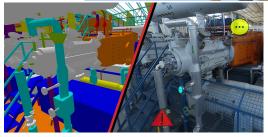












3D models

Many file formats













Key features of a 3D file:

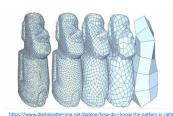
- geometry,
- surface texture,
- scene details,
- animation of the model
- element properties,

https://all3dp.com/2/most-common-3d-file-formats-model/

Shape Geometry

a) Approximate Mesh Encoding b) Precise Mesh Encoding « tesselations »: decompose a surface in polygons (usually triangles)

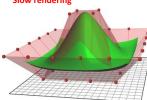
Good for 3D printing Files may be huge.



https://all3dp.com/2/most-common-3d-file-formats-model/

Parametric surfaces made of control points and knots. Example: non-uniform rational basis spline (NURBS)

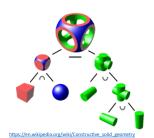
Exact at any resolution Slow rendering



c) Constructive Solid Geometry

Primitive shapes that are combined using **Boolean operations**

Good for CAD



Surface Textures

a) Texture Mapping

every point in the 3D model's surface (or the polygonal mesh) is mapped to a twodimensional image.



https://metalbyexample.com/textures-and-samplers/



b) Face Attributes

each face of the mesh has a set of attributes Color, texture, material type, reflection, refraction,...



Scene detail

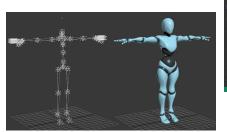
• layout of the 3D model in terms of cameras, light sources, and other nearby 3D models

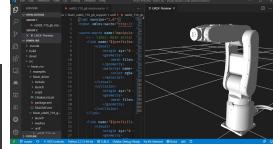


https://www.creativeshrimp.com/light-texture-tutorial-1-lighting-book.html

Animation

• skeletal animation: skeleton + joints



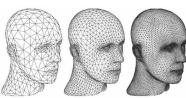


Unified Robot Description Format (URDF) in Visual Studio Code

 $\underline{\text{http://www.downloads.redway3d.com/downloads/public/documentation/wf_skeletal_animation.html}}$

Popular 3D file formats: **STL** ("stereolithography")





- · Popular for 3D printing
- · Triangular mesh,
- · No color information
- · ASCII or Binary representation
- Superseded by OBJ, 3MF, AMF, ...

Popular 3D file formats: Wavefront OBJ file format



Popular for 3D printing and 3D graphics

Triangular mesh

· + other kinds of interpolation: Tayolr, B-splines...

• No animation, no deformation

Reference to companion file MTL (Material Template Library)

https://people.sc.fsu.edu/~jburkardt/data/mtl/mtl.html

Popular 3D file formats: Autodesk 3DS Max file format



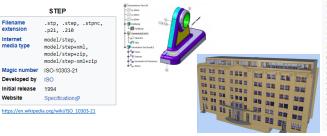


- Autodesk 3D Studio MAX 1996
- · Popular for architecture, engineering, education, manufacturing
- geometry, appearance, scene, and animation
- Triangular mesh, total 65536 triangles
- Directional light sources are not supported

Reference: http://paulbourke.net/dataformats/3ds/ Tree structure made of chunks

- 6 byte Chunk header: 1-2 = chunk ID
- 3-6 = little-endian length of the chunk
- Next bytes: chunk's data. including sub-chunks

Popular 3D file formats: ISO 10303-21 STEP-files



#16-PRODUCT('A8081','Test Part 1','',(#18)); #17-PRODUCT_RELATED_PRODUCT_CATEGORY('part',\$,(#16)); #20-ORGANIZATION ROLE('id owner

- · Popular for in many engineering domains that rely on CAD
 - Mechanical engineering
 - AECOO industry (Architecture, Engineering, Construction, Owner Operator)
- ASCII. not storage-efficient
- Schema defined separately in the EXPRESS data modeling language ISO 10303-11

</ Part 2. Data 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

< Part 3. Data schemas and semantics >

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

< Part 3. Data schemas and semantics > Part 3.1. Data schemas

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 Lef

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

Data Interoperability and SemanticsOutline

- < Part 3. Data schemas and semantics >
 - Part 3.1. Data schemas
 - Part 3.1.1. XML Schema
 - Part 3.1.2. JSON Schema
 - Part 3.2. Semantics
 - Part 3.2.1. Heterogeneities and data conflicts
 - · Part 3.2.2. Controlled vocabularies and ontologies
 - Part 3.2.3. Resource Description Framework
 - Part 3.2.4. RDFa: Rich structured data markup for web documents
 - Part 3.2.5. JSON-LD: JSON for Linking Data

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-lefrançois.info
Course unit URL: https://ci.mines-stetienne.fr/cps2/course/data

Data Interoperability and Semantics

< Part 3. Data schemas and semantics >
Part 3.1. Data schemas
Part 3.1.1. XML Schema

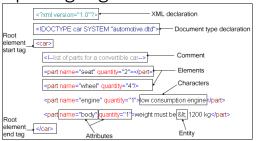
ICM – Toolbox Engineering and Interoperability of Software Systems – Course unit on Data Interoperability and Semantics M1 Cyber Physical and Social Systems – Course unit on Data Interoperability and Semantics Maxime Lefrançois https://maxime-lefrancois.info

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

reminder reminder

Extensible Markup Language

XML (file format) Filename extension application/xml text/xml [1 Identifier (UTI) UTI conformation <?xml Markup language Extended from Extended to Numerous languages, including XHTML · RSS · Atom · KML Standard (November 26, 2008; 12 years ago) (August 16, 2006; 15 years ago) Open format?



- v1.0 in 1998, still extensively used in many verticals
- numerous formats based on XML (418 registered on IANA) https://en.wikipedia.org/wiki/List of XML markup languages application/atom+xml application/rdf+xml
- verbosity, complexity and redundancy

Extensible Markup Language

XML (file format) Filename extension application/xml Internet text/xml [1 Identifier (UTI) UTI conformation Type of format Markup language Extended from Extended to Numerous languages, including XHTML · RSS · Atom · KML (November 26, 2008; 12 years ago) 1.1 (Second Edition)呼 (August 16, 2006; 15 years ago)

Characters and escaping

- unicode implementations: <?xml version="1.0" encoding="UTF-8"?>
- escaping characters: < '<' & '&' ❤ '♥' etc.

Syntactical correctness

- · well formed vs ill-formed
- · one root tag
- · correct nesting
- tag names (approx) start with letter, then alphanumeric or ':'

Schemas and validation

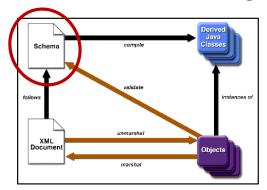
- · valid vs invalid
- · DTD. or XML Schema

Namespaces

- xmlns:ns1="http://example.org/ns1" xmlns:ns2="http://example.org/ns2"
- allows to use different schemas together: <ns1:Tag> <ns2:Tag>

reminder

XML – OOP data binding



XML-OOP data binding ex Java: https://zetcode.com/java/jaxb/



Java Architecture for XML Binding (JAXB) example

XMI Schema

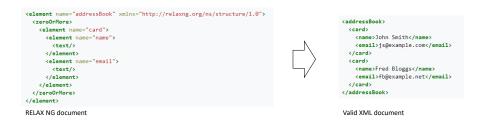
Document Type Definition (DTD, 2008)

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN"</pre>
 "http://www.w3.org/TR/xhtml1/DTD/xhtml1-transitional.dtd">
link to some .dtd document
     [...]
     <!ELEMENT html (head, body)>
     <!ELEMENT p (#PCDATA | p | ul | dl | table | h1|h2|h3)*>
     [...]
      <!ATTLIST img
        src CDATA
                              #REQUIRED
        id
              ID
                             #IMPLIED
        sort CDATA
                             #FIXED "true
        print (yes | no) "yes"
```

```
<?xml version="1.0" encoding="utf-8"?>
<!DOCTYPE html>
 <!-- the XHTML document body starts here-->
 <html xmlns="http://www.w3.org/1999/xhtml">
 c/h+m1>
HTML 5 : no more link to a DTD
```

XML Schema

• Regular Language for XML Next Generation (RELAX NG, 2001) https://relaxng.org/



XML Schema

• XML Schema (W3C, 2003) The only one you should use (if you need to)

```
<?xml version="1.0"?>
                                                                     <?xml version="1.0"?>
<xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
<xs:element name="note">
                                                                     xmlns="https://www.w3schools.com"
 <xs:complexType>
                                                                     xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
  <xs:seauence>
                                                                     xsi:schemaLocation="https://www.w3schools.com/xml note.xsd">
    <xs:element name="to" type="xs:string"/>
                                                                       <to>Tove</to>
     <xs:element name="from" type="xs:string"/>
                                                                       <from>Jani</from>
    <xs:element name="heading" type="xs:string"/>
                                                                      <heading>Reminder</heading>
    <xs:element name="body" type="xs:string"/>
                                                                      <body>Don't forget me this weekend!</body>
                                                                     </note>
 </xs:complexType>
                                                                    Valid XML document
</xs:element>
</r>
```

XML Schema document

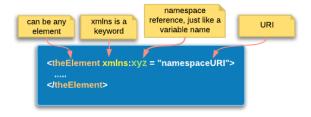
https://www.w3schools.com/xml/schema_intro.asp 10

XML Schema

- Tutorials
 - https://www.tutorialspoint.com/xml/xml schemas.htm
 - https://www.w3schools.com/xml/schema intro.asp
- What to do with XML Schemas
 - validate documents
 - generate forms
 - generate classes (any OOP language)

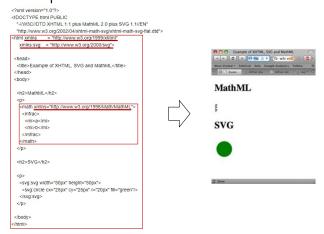
Multiple XML Schemas?

- How to combine different XML schemas in a single file?
- Using XML namespaces



https://www.w3schools.com/xml/schema_intro.asp 11

Example: A document with XHTML + MathML + SVG



Data Interoperability and Semantics

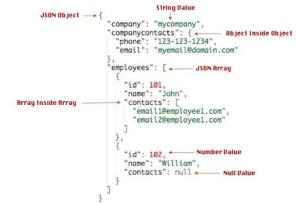
< Part 3. Data schemas and semantics >
Part 3.1. Data schemas
Part 3.1.2. JSON Schema

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/cos/course/data

reminder

JavaScript Object Notation

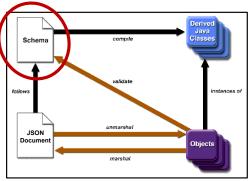




JSON Schema



JSON – OOP data binding









JSON-OOP data binding

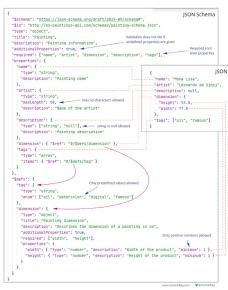
JSON Schema

Examples:

- https://json-schema.org/learn/miscellaneous-examples.html
- https://www.fiware.org/developers/data-models/
- https://oneiota.org/







JSON Schema Specification (latest)

https://json-schema.org/draft/2020-12/json-schema-validation.html

7. Vocabularies for Semantic Content With "format"
7.1. Foreword
7.1. Foreword
7.1. Foreword
7.2. Implementation Requirements
7.2. Implementation Requirements
7.2. Implementation Requirements
7.2. Contents of Semantic Semantic
7.2. Contents of Semantic
7.2. Contents of Semantic
7.2. Email Addresses
8. A Content Email Addresses
8. A Content Email Addresses
8. A Content Email E

Data Interoperability and Semantics

< Part 3. Data schemas and semantics > Part 3.2. Semantics

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 Lef

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

< Part 3. Data schemas and semantics > Part 3.2. Semantics

Part 3.2.1. Heterogeneities and data conflicts

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

DEMO: Compare JSON documents

Different modeling choices were made, which make these two services completely non-interoperable:

- · the lat/long coordinates: string vs number
- the UNIX timestamps vs dates and times
- the choice of keys and the semantics (meaning) of the values
- the units of temperature, pressure, wind speed, ...
- the semantics of wind direction
- the value for "icon": "03n" (if we follow our nose on the website, we may figure out it refers to http://openweathermap.org/img/w/03n.png)
- the country codes ISO 3166-1 ALPHA-2 and ISO 3166-1 ALPHA-3 (example of Australia and Austria)

DEMO: Compare JSON documents

 $\frac{https://samples.openweathermap.org/data/2.5/weather?id=2172797\&appid=b6907d289e10d714a6e88b30761fae22$

and

https://www.prevision-meteo.ch/services/json/lausanne

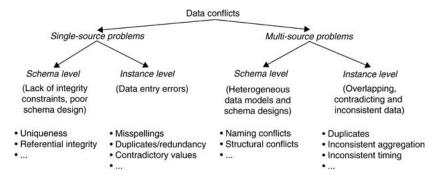
and

https://github.com/smart-data-models/dataModel.Weather/blob/master/WeatherForecast/examples/example-normalized.json

22

Data conflicts

Data conflicts are deviations between data intended to capture the same state of a real-world entity. Data with conflicts are often called "dirty" data and can mislead analysis performed on it.



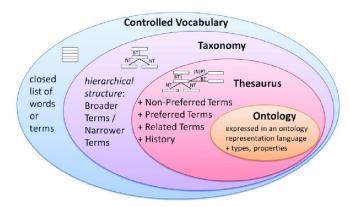
< Part 3. Data schemas and semantics > Part 3.2. Semantics

Part 3.2.2. Controlled vocabularies, thesauri, taxonomies

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

Controlled vocabulary, taxonomy, thesaurus, ontology



Kopácsi, Sándor & Hudak, Rastislav & Ganguly, Raman. (2017). Implementation of a Classification Server to Support Metadata Organization for Long Term Preservation Systems Mitteilungen der Vereinigung Österreichischer Bibliothekarinnen und Bibliothekare. 70. 225. 10.31263/voebm.v70i2.1897.

Controlled vocabularies

... an established list of standardized terminology for use in indexing and retrieval of information

- OFCD

... an organized arrangement of words and phrases used to index content and/or retrieve content through browsing or searching

- Getty institute

... an standardized – yet dynamic – set of terms and phrases authorized for use in an indexing system to describe a subject area or information domain

- SCIP

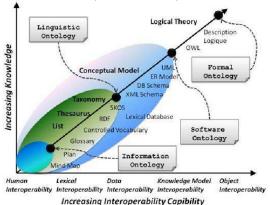
vocabulary for which the entries, i.e. definition/term pairs, are controlled by a Source Authority based on a rulebase and process for addition/deletion of entries

- ISO/IEC 15944-5:2008(en)

Information technology — Business Operational View

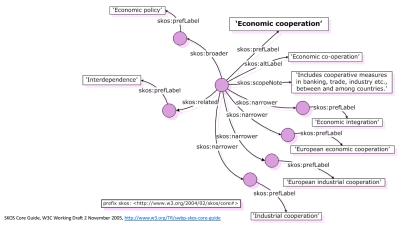
20

Controlled vocabulary, taxonomy, thesaurus, ontology



Roussey, Catherine & Pinet, François & Kang, Myoung-Ah & Corcho, Oscar & Falquet, Gilles & Métral, Claudine & Teller, Jacques & Tweed, Christopher. (2011 Ontologies for Interoperability. 10.1007/978-0-85729-724-2_3.

SKOS Simple Knowledge Organization System



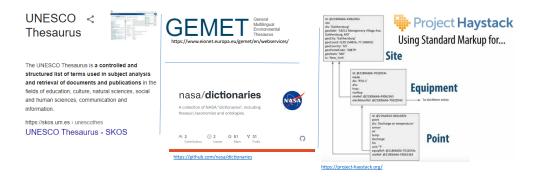
Data Interoperability and Semantics

< Part 3. Data schemas and semantics > Part 3.2. Semantics

Part 3.2.3. Resource Description Framework

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://maxime-lefrancois.info
Course unit URL: https://ci.mines-stetienne.fr/cps2/course/data

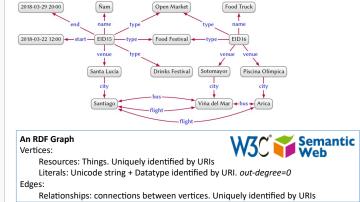
Examples of thesauri, taxonomies, ...



List at https://www.w3.org/2001/sw/wiki/SKOS/Datasets

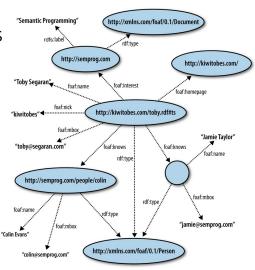
RDF-Resource Description Framework





CURIE = Compact URIs

@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>. @prefix rdfs: <http://www.w3.org/2000/01/rdf-schema#>. @prefix foaf: <http://xmlns.com/foaf/0.1/>. @prefix schema: <http://schema.org/>.



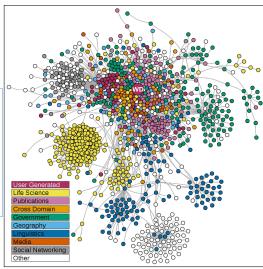
The Linked Data

It's possible to use URIs defined in an external graph

- > to use a reference identification system
- > to augment a graph anyone can say anything about anything

One can specify vertice co-reference with some special edges: owl:sameAs from W3C OWL vocabulary

chile:Santiago owl:sameAs → geo:SantiagoDeChile



Linked Open Data cloud visualisation: each node is a RDF dataset. Links indicate the existence of external identification links. Source: wikipedia

5 ★ OPEN DATA

Tim Berners-Lee, the inventor of the Web and Linked Data initiator, suggested a 5-star deployment scheme for Open Data.



https://5stardata.info/en/

Examples of RDF vocabularies & ontologies



Documentation

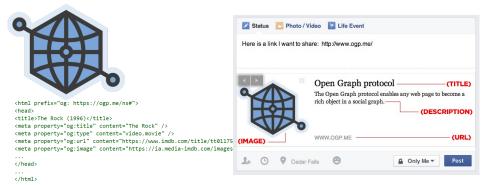
- Getting Started: A simple introduction to microdata and using schema.org for marking up your site
- . Schemas: The actual schemas, arranged in a hierarchy, with a page for each item in the schema.
- The full type hierarchy: The full type hierarchy, in a single file. · Frequently asked questions
- Data model: a brief note on the data model used, etc.
- Extension Mechanism: The extension mechanism that can be used to extend the schema Schema.org Discussion Group: Forum for finding answers to questions, etc.
- . Feedback form: Please give us feedback, report bugs, etc.

markup their web pages and email messages. Many applications from Google, Microsoft, Pinterest, Yandex and others already use these vocabularies to power rich, extensible experiences.

Founded by Google, Microsoft, Yahoo and Yandex, 2011

Test with https://validator.schema.org/ or https://developers.google.com/search/docs/advanced/structured-data

Examples of RDF vocabularies & ontologies



Open Graph Protocol https://ogp.me/ (Facebook, 2010), test with https://developers.facebook.com/tools/debug/

Data Interoperability and Semantics

< Part 3. Data schemas and semantics > Part 3.2. Semantics

Part 3.2.4. RDFa: Rich structured data markup for web documents

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-lefrançois.info

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

RDFa Resource Description Framework in Attributes

See RDFa 1.1 Primer - Rich Structured Data Markup for Web Documents - https://www.w3.org/TR/rdfa-primer/

RDFa Resource Description Framework in Attributes

See RDFa 1.1 Primer - Rich Structured Data Markup for Web Documents - https://www.w3.org/TR/rdfa-primer/







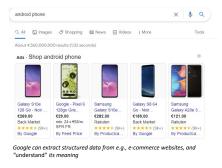
Figure 1: On the left, what browsers see. On the right, what humans see. Can we bridge the gap so that browsers see more of what we see?

RDFa Resource Description Framework in Attributes

See RDFa 1.1 Primer - Rich Structured Data Markup for Web Documents - https://www.w3.org/TR/rdfa-primer/

The Result





Data Interoperability and Semantics

< Part 3. Data schemas and semantics > Part 3.2. Semantics

Part 3.2.5. JSON-LD: JSON for Linking Data

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

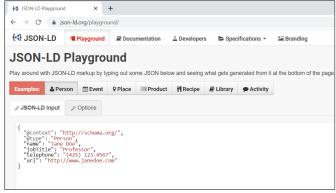
JavaScript Object Notation for Linked Data

See JSON-LD 1.1 - A JSON-based Serialization for Linked Data - https://www.w3.org/TR/json-ld11/



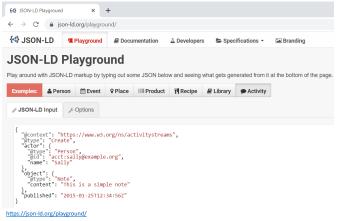


JavaScript Object Notation for Linked Data



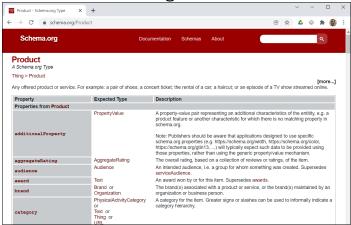
https://json-ld.org/playground/

JavaScript Object Notation for Linked Data



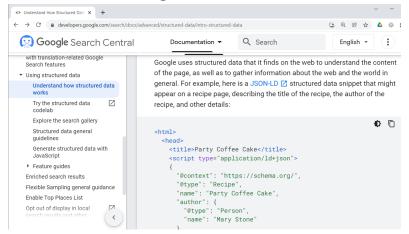
Product - Schema.org Type × + ← → C 🍙 schema.org/Product Examples No Markup Microdata RDFa JSON-LD Structure Example encoded as JSON-LD in a HTML script tag. <script type="application/ld+json"> "@context": "https://schema.org", "@type": "Product", "aggregateRating": { "@type": "AggregateRating", "ratingValue": "3.5", "reviewCount": "11" "description": "0.7 cubic feet countertop microwave. Has six preset cooking categories and co "name": "Kenmore White 17\" Microwave", "image": "Kenmore-microwave-17in.jpg", "offers": {
 "@type": "Offer",
 "availability": "https://schema.org/InStock",
 "price": "55.00", "priceCurrency": "USD" "review": ["@type": "Review",
"author": "Bllie",
"datePublished": "2011-04-01",
"reviewBody": "The lamp burned out and now I have to replace it.",

JSON-LD and schema.org

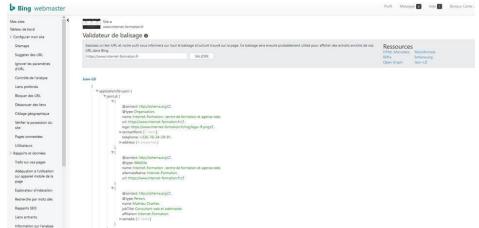


46

JSON-LD and Google



JSON-LD and Bing



Data Interoperability and Semantics

</ Part 3. Data schemas and semantics >

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 Intigos (Maxime Lefrançois Intigos) (Maxime Lefrançois Intigos) (Maxime Lefrançois Intigos) (Maxime Lefrançois Intigos) (Maxime Lefrançois) (Maxime Lefrançois

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

< Part 4. The value of data >

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://maxime-lefrancois.info
Course unit URL: https://ci.mines-stetienne.fr/cps2/course/data

Data Interoperability and Semantics

Part 4. The value of data
Part 4.1. Value as one of the V's of Big Data

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 SemanticsOutline

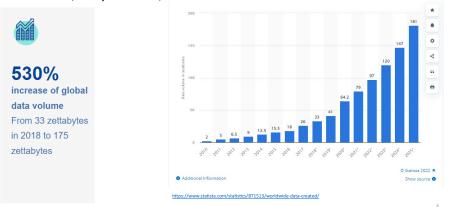
- < Part 4. The value of data >
 - Part 4.1. Value as one of the V's of Big Data
 - Part 4.2. Data, Information, Knowledge, Metadata, ...
 - Part 4.3. Interoperability unlocks the value of data
 - Part 4.4. Open data generates economic and societal value
 - Part 4.5. Machine-actionability of data increases its value

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

Orders of magnitude of the total amount of data created, captured, copied, and consumed globally?

Multiple-byte units									
Decimal			Binary						
Value		Metric	Value		IEC		Legacy		
1000	kB	kilobyte	1024	KiB	kibibyte	KB	kilobyte		
1000 ²	МВ	megabyte	1024 ²	MiB	mebibyte	МВ	megabyte		
1000 ³	GB	gigabyte	1024 ³	GiB	gibibyte	GB	gigabyte		
10004	ТВ	terabyte	1024 ⁴	TiB	tebibyte	ТВ	terabyte		
1000 ⁵	РΒ	petabyte	1024 ⁵	PiB	pebibyte		-		
1000 ⁶	ЕВ	exabyte	1024 ⁶	EiB	exbibyte		-		
1000 ⁷	ZΒ	zettabyte	1024 ⁷	ZiB	zebibyte		_		
1000 ⁸	YΒ	yottabyte	1024 ⁸	YiB	yobibyte		-		
Orders of magnitude of data									

Projected figures 2025 for the total amount of data created, captured, copied, and consumed globally

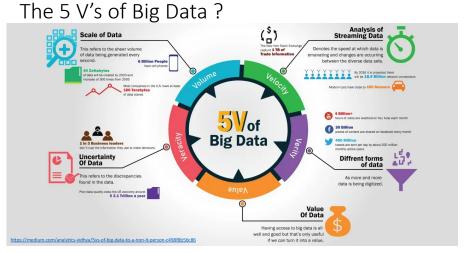


Projected figures 2025



https://ec.europa.eu/info/strategy/priorities-2019-2024/europe-fit-digital-age/european-data-strategy_en#projected-figures-2025

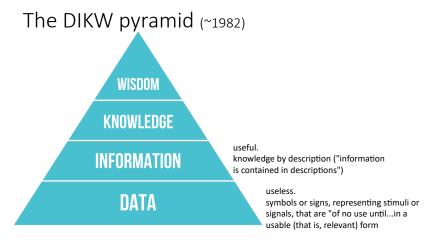


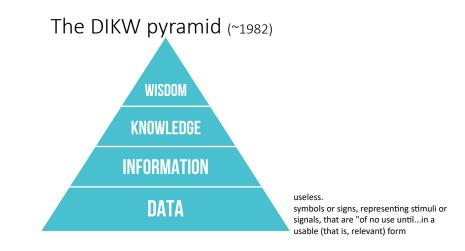


Part 4. The value of data
Part 4.2. Data, Information, Knowledge, Metadata, ...

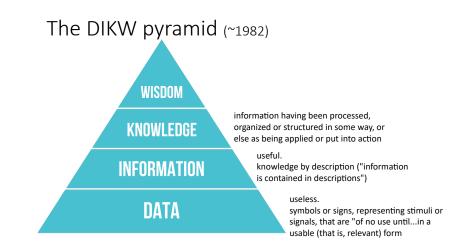
ICM — Toolbox Engineering and Interoperability of Software Systems — Course unit on Data Interoperability and Semantics M1 Cyber Physical and Social Systems — Course unit on Data Interoperability and Semantics Maxime Lefrançois https://maxime-lefrancois.info

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

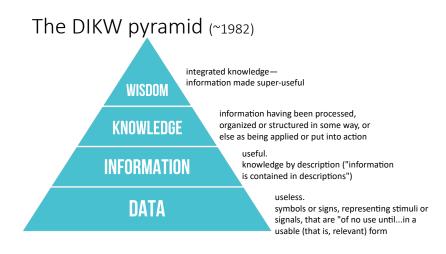




https://en.wikipedia.org/wiki/DIKW_pyramid_, and chosen definitions



11 https://en.wikipedia.org/wiki/DHKW_pyramid, and chosen definitions 1 https://en.wikipedia.org/wiki/DHKW_pyramid, and chosen definitions 1 https://en.wikipedia.org/wiki/DHKW_pyramid, and chosen definitions 1



https://en.wikipedia.org/wiki/DIKW_pyramid_, and chosen definitions

The DIKW pyramid (~1982)

Still...

KNOWLEDGE we often use

INFORMATION "Data "

to generalize

Definitions on Data

- ISO/IEC20546:2019 (Big data - Overview and vocabulary)

Dataset

Identifiable collection of data available for access or download in one or more formats

Data

 $\textit{Re-interpretable representation of information in a formalized manner suitable for communication, interpretation, or processing$

Note 1 to entry: Data can be processed by humans or by automatic means.

Metadata

Data about data or data elements, possibly including their data descriptions and data about data ownership, access paths, access rights and data volatility

Information

Data that are processed, organised and correlated to produce meaning.

Note 1 to entry: Information concerns facts, concepts, objects, events, ideas, processes, etc.

Metadata

"data that provides information about other data"

- Descriptive metadata the descriptive information about a resource.
 - · For discovery and identification.
 - · Ex: title, abstract, author, keywords.
- Structural metadata containers of data, how compound objects are put together
 - Ex, how pages are ordered to form chapters.
 - . Ex: types, versions, relationships, and other characteristics of digital materials.
- Administrative metadata the information to help manage a resource
 - Ex: resource type, permissions, time, when and how it was created.
- Reference metadata contents and quality
 - · For quality assessment of the data
 - Ex: conceptual metadata, quality metadata, methodological metadata.
- · Statistical metadata, also called process data,
 - May describe processes that collect, process, or produce statistical data.
 - · Number of rows, columns, etc.
- Legal metadata –

https://en.wikipedia.org/Exi/Meense, creator, copyright holder

Part 4. The value of data
Part 4.3. Interoperability unlocks the value of data

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 Analytics vs Data Processing

- ISO/IEC20546:2019 (Big data - Overview and vocabulary)

Data Analytics

Composite concept consisting of data acquisition, data collection, data validation, data processing, including data quantification, data visualisation and data interpretation

Data Processing

Systematic performance of operations upon data

Note 1 to entry: Example: Arithmetic or logic operations upon data, merging or sorting of data, or operations on text, such as editing, sorting, merging, storing, retrieving, displaying, or printing.

18

Data silos problem illustrated

Repositories of fixed data that are isolated, incompatible, or not integrated

Marketing Purchasing Sales NAME OF THE **BUDGET RANGE** CLIENT NAME OF THE WHAT BOOKS **BUDGET RANGE** CLIENT HE BUYS NAME_HIS **BIRTHDAY BUDGET RANGE** BIRTHDAY WHAT BOOKS DEVICE BIRTHDAY HE LIKES

Data integration

Simple schematic for a data warehouse.
The Extract, transform, load

(ETL) process extracts information from the source databases, transforms it and then loads it into the data warehouse.

Simple schematic for a dataintegration solution.

A system designer constructs a mediated schema against which users can run queries. The virtual database interfaces with the source databases via wrapper code if required. Data Source A

Data Source A

Wrapper

Wrapper

Data Source B Wrapper Wrapper Mediated Schema "Virtual Database" Data Source C Wrapper

combining **heterogeneous data** and providing users with a unified view of them.

Sub-areas:

- Data warehousing
- Data migration
- Enterprise application/information integration
- Master data management

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

Interoperability vs Portability

— IEEE Standard Computer Dictionary — ISO/IEC19941:2017 (Cloud computing – interoperability and portability)

Interoperability

Ability of two or more systems or components to exchange information and to mutually use the information that has been exchanged

Portability

Ability to easily transfer data from one system to another without being required to re-enter data

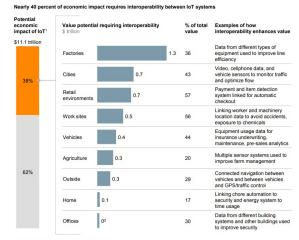
2.

Interoperability at different levels: from the hardware to policies

- Policies: Laws, regulations, norms, ...
- Behaviour: Processus, services, operations, ...
- **Semantic**: Knowledge domain vocabularies, taxonomies, ontologies, ...
- Syntactic: Encoding, Formats, Schemas, ...
- Transport: see also Open Systems Interconnections (OSI) network layers
- OS: Kernel, File systems,...
- Hardware: Instruction Set Architectures, ...

Interoperability as an enabler for the potential value of data

ex: Internet of Things



Includes sized applications only; includes consumer surpli
 Less than \$100 billion.

NOTE: Numbers may not sum due to rounding.

SOURCE: Expert interviews; McKinsey Global Institute analysis

 $\underline{\text{https://www.mckinsey.com/business-functions/mckinsey-digital/our-insights/the-internet-of-things-the-value-of-digitizing-the-physical-world and the transfer of the physical formula of the phys$

ISO/IEC definitions of interoperability

Transport interoperability

interoperability where information exchange uses an established communication infrastructure between the participating systems

— ISO/IEC 22123-1:2021, Cloud computing — Part 4: Vocabulary

Syntactic interoperability

interoperability such that the formats of the exchanged information can be understood by the participating systems

- ISO/IEC 22123-1:2021, Cloud computing - Part 4: Vocabulary

Semantic data interoperability

interoperability so that the meaning of the data model within the context of a subject area is understood by the participating systems

- ISO/IEC 22123-1:2021, Cloud computing - Part 4: Vocabulary

Behavioural interoperability

interoperability so that the actual result of the exchange achieves the expected outcome

— ISO/IEC 22123-1:2021, Cloud computing — Part 4: Vocabulary

Policy interoperability

interoperability while complying with the legal, organizational, and policy frameworks applicable to the participating systems

— ISO/IEC 22123-1:2021, Cloud computing — Part 4: Vocabulariy4

Data Interoperability – a definition

Data interoperability addresses the ability of systems and services that create, exchange and consume data to have clear, shared expectations for the contents, context and meaning of that data.

- https://datainteroperability.org/

Syntactic interoperability

interoperability such that the formats of the exchanged information can be understood by the participating systems

- ISO/IEC 22123-1:2021, Cloud computing - Part 4: Vocabulary

Semantic data interoperability

interoperability so that the meaning of the data model within the context of a subject area is understood by the participating systems

- ISO/IEC 22123-1:2021, Cloud computing - Part 4: Vocabulary

Data Interoperability and Semantics

Part 4. The value of data

Part 4.4. Sharing data generates economic and societal value

Norm-based interoperability

- Hardware and/or software standards
- Standard = Detailed set of technical requirements intended to establish a certain uniformity (in a field of hardware or software development)
- Standard ---- Norm: no clear boundary
 - docx: Microsoft
 - HTML, XML: World Wide Web consortium (W3C)
 - SIM cards: European Telecommunication Standard Institute (ETSI)
 - Postscript became standardized after release
 - De facto standards (not produced by a standard development organization). example: csv, json, rar, pdf, java, flash

The Data Value chain

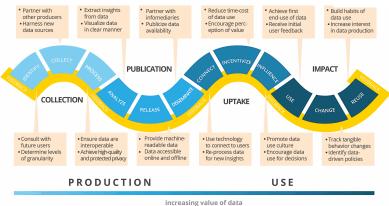




increasing value of data

The Data Value chain





Open Data Watch. "The data value chain: Moving from production to impact." Data2X. https://opendatawatch.com/publications/the-data-value-chain-moving-from-production-to-impact (2018)

The Data Value roadblocks



PRODUCTION

USE

increasing value of data



Roadblocks for **production** include lack of financial, human, and technological resources; low data literacy; lack of trust between users and data collectors; blindspots in data gaps; lack of country ownership; and lack of government desire for transparency.



Roadblocks for use include low political support; lack of data relevance to decisions; poor quality; lack of trust in government data use; no rewards or results of data use; financial constraints; corruption; data silos; and lack of partnerships between infomediaries.

Open Data Watch. "The data value chain: Moving from production to impact." Data2X. https://opendatawatch.com/publications/the-data-value-chain-moving-from-production-to-impact (2018)

Open Data in the US



- Data.gov 2009
- Legal framework:
 - The U.S. Open Government Directive of December 8, 2009, required that all agencies post at least three high-value data sets online and register them on Data.gov within 45 days
 - OPEN Government Data Act, as part of the Foundations for Evidence Based Policymaking Act (2019)

Open Data in France



- France at the forefront of Open Data in Europe:
 - Légifrance 1999
- Legal framework:
 - "The society has the right of requesting account from any public agent of its administration." (Declaration of rights of man and of the citizen of 1789)
 - Law on the liberty of access to administrative documents (1978)
 - Euopean directove 2003 + French Law 2005 + Decree 2011
 - Bill on a Digital Republic (2016)
 - The law on Energy Transition (2015)
- 2014: Chief Data Officer in the French public administration



Open Data in Europe





data.europa.eu/europeandataportal





34

Open Data, Open Content, and Open Knowledge

The Open Definition

The Open Definition sets out principles that define "openness" in relation to data and content.

It makes **precise** the meaning of "open" in the terms "**open data**" and "**open content**" and thereby ensures **quality** and encourages **compatibility** between different pools of open material.

It can be summed up in the statement that:

"Open means anyone can freely access, use, modify, and share for any purpose (subject, at most, to requirements that preserve provenance and openness)."

Put most succinctly:

"Open data and content can be freely used, modified, and shared by anyone for any purpose"

Open Work

1. Open License or Status

The work must be in the public domain or provided under an open license

2. Access

The **work** *must* be provided as a whole and at no more than a reasonable one-time reproduction cost, and *should* be downloadable via the Internet without charge.

3. Machine Readability

The **work** *must* be provided in a form readily processable by a computer and where the individual elements of the work can be easily accessed and modified.

4. Open Format

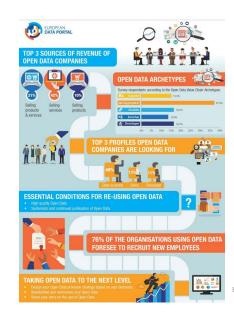
The work must be provided in an open format.

http://opendefinition.org/ 35 http://opendefinition.org/ 36











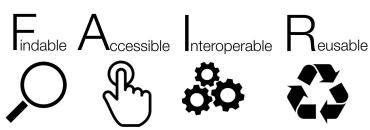


Part 4. The value of data

Part 4.5. Machine-actionability of data increases its value

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-stetlenne.fr/cps2/course/data

FAIR Principles



The FAIR principles emphasize machine-actionability (i.e., the capacity of computational systems to find, access, interoperate, and reuse data with none or minimal human intervention) because humans increasingly rely on computational support to deal with data as a result of the increase in volume, complexity, and creation speed of data.

https://www.go-fair.org/fair-principles/

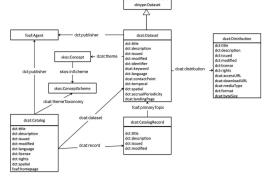
FAIR Guiding Principles for scientific data management and stewardship (2016)

https://www.go-fair.org/fair-principles/

FAIR Guiding Principles for scientific data management and stewardship (2016)

Example of standards for Metadata

Dublin Core



Data Catalog Vocabulary (DCAT) - Version 2 W3C Recommendation 04 February 2020

- Demo Dublin core or schema.org in web pages
- Demo DCAT for open data portal catalogs

schema.org

< Part 4. The value of data >

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

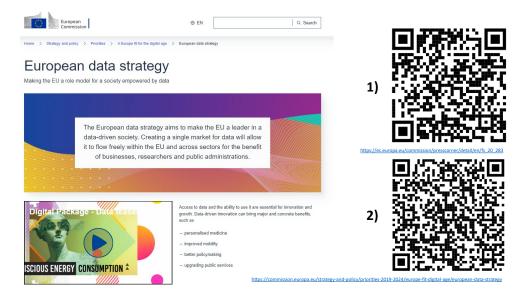
< Part 5. The European Data Strategy >



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 SemanticsOutline

- < Part 5. The European Data Strategy >
 - Part 5.1. The European legislation on open data
 - Part 5.2. The General Data Protection Regulation (GDPR)
 - Part 5.3. The Data Governance Act (DGA)
 - Part 5.4. The Data Act
 - Part 5.5. Data Spaces



Data Interoperability and Semantics

Part 5. The value of data
Part 5.1. The European legislation on open data

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

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



Home > Strategy > Priorities 2019-2024 > A Europe fit for the digital age > European data strategy

European legislation on open data

(adopted 05/2022 - applicable 09/2023)

The Directive on open data and the re-use of public sector information provides common rules for a European market for government-held data.

The "Open Data Directive" (EU) 2019/1024 entered into force on 16 July 2019

It replaced the Public Sector Information (PSI) Directive of 2003.

EU countries had to transpose Directive (EU) 2019/1024 by 16 July 2021.

The Commission Implementing Regulation (EU) 2023/138 adopted a list of specific high-value datasets by way of an implementing act.





Home > Strategy > Priorities 2019-2024 > A Europe fit for the digital age > European data strategy

European legislation on open data

(adopted 05/2022 – applicable 09/2023)

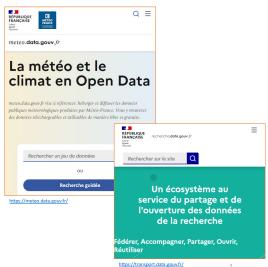


The Commission Implementing Regulation (EU) 2023/138 adopted a list of specific high-value datasets by way of an implementing act.



Actual impact (FR)





Data Interoperability and Semantics

Part 5. The value of data Part 5.2. The General Data Protection Regulation (GDPR)

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

General Data Protection Regulation (GDPR) (2016)

- Requiring the consent of subjects for data processing
- Anonymizing collected data to protect privacy
- Providing data breach notifications
- Safely handling the transfer of data across borders
- Requiring certain companies to appoint a data protection officer to oversee GDPR compliance



https://gdpr.eu/

General Data Protection Regulation (GDPR) (2016)





10

Data Interoperability and Semantics

Part 5. The value of data
Part 5.3. The Data Governance Act (DGA)

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-lefrançois.info
Course unit URL: https://ci.mines-stetienne.fr/cps2/course/data





Home > Strategy > Priorities 2019-2024 > A Europe fit for the digital age > European data strategy

The European Data Governance Act

(adopted 05/2022 – applicable 09/2023)

3.6.2022 EN

Official Journal of the European Union

REGULATION (EU) 2022/868 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

on European data governance and amending Regulation (EU) 2018/1724 (Data Governance Act)

(Text with EEA relevance)

ttp://data.europa.eu/eli/reg/2022/868/oj

Setting up a new European way of data governance will facilitate data sharing across sectors and Member States. It will create wealth for society, and provide control to citizens and trust to companies.



Data Interoperability and Semantics

Part 5. The value of data Part 5.4. The Data Act

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



Home > Strategy > Priorities 2019-2024 > A Europe fit for the digital age > European data strategy

The European Data Act

(adopted 01/2024 - applicable 09/2025)







Home > Strategy > Priorities 2019-2024 > A Europe fit for the digital age > European data strategy

The European Data Act

(adopted 01/2024 - applicable 09/2025)

EN Official Journal Series L 2023/2854 22.12.2023 REGULATION (EU) 2023/2854 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

of 13 December 2023

(Text with EEA relevance)

ttp://data.europa.eu/eli/reg/2023/2854/oj

The Data Act will make more data available for use. It will set up all economic sectors in the EU.





Home > Strategy > Priorities 2019-2024 > A Europe fit for the digital age > European data strategy

The European Data Act

(adopted 01/2024 – applicable 09/2025)

17

Chapter VIII - Interoperability

Article 33 – Essential requirements regarding interoperability of data, of data sharing mechanisms and services, as well as of common European data spaces

- Participants in data spaces that offer data or data services to other participants shall comply with the following essentia requirements to facilitate the interoperability of data, of data sharing mechanisms and services, as well as of common European data spaces which are purpose- or sector-specific or cross-sectoral interoperable frameworks for common standards and practices to share:
- the dataset content, use restrictions, licences, data collection methodology, data quality and uncertainty shall be sufficiently described, where
 applicable, in a machine-readable format, to allow the recipient to find, access and use the data
- the data structures, data formats, vocabularies, classification schemes, taxonomies and code lists, where available, shall be described in a
 publicly available and consistent manner
- the technical means to access the data, such as application programming interfaces, and their terms of use and quality of service shall be sufficiently described to enable automatic access and transmission of data between parties, including continuously, in bulk download or in real time in a machine-readable format where that is technically feasible and does not hamper the good functioning of the connected product
- where applicable, the means to enable the interoperability of tools for automating the execution of data sharing agreements, such as smart contracts shall be provided.

Data Interoperability and Semantics

Part 5. The value of data Part 5.5. Data Spaces

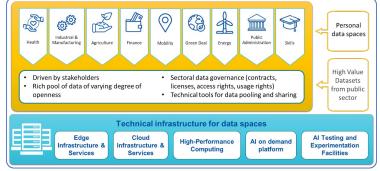
ICM — Toolbox Engineering and Interoperability of Software Systems — Course unit on Data Interoperability and Semantics M1 Cyber Physical and Social Systems — Course unit on Data Interoperability and Semantics Maxime Lefrançois https://maxime-lefrançois.info

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

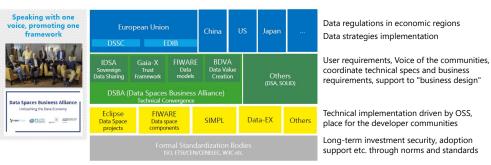
Common European data spaces

"purpose or sector specific or cross-sectoral interoperable frameworks for common standards and practices to share or jointly process data for, inter alia, the development of new products and services, scientific research or civil society initiatives."

[source: Data Act]



Regulatory, business, and technical foundation for Data Spaces within the Edge-Cloud-Continuum



https://data-spaces-business-alliance.eu/download/33968

https://digital-strategy.ec.europa.eu/en/fibrary/building_data-economy-brochure

19 Adapted from a presentation by L. Nagel at the <u>Data Spaces Symposium 2024</u>

20

The Data Spaces Business Alliance

Unleashing the European Data Economy

Accelerating Business transformation in the Data Economy

The Data Spaces Business Alliance (DSBA) accelerates business transformation in the data economy. It's the first initiative of its kind, uniting industry players to realize a data-driven future in which organizations and individuals can unlock the full value of their data. Data spaces are key to achieving sovereign, interoperable and trustworthy data-sharing across businesses and societies — a key step to the data economy of the future. The Alliance embraces this reality, converging the best skills, assets, and experience in Europe into a one-stop-shop for data spaces, from inception to deployment.

The Data Spaces Business Alliance are Gaia-X European Association for Data and Cloud AISBL, the Big Data Value Association (BDVA), FIWARE Foundation, and the International Data Spaces Association (IDSA). Together they represent 1,000+ leading key industry

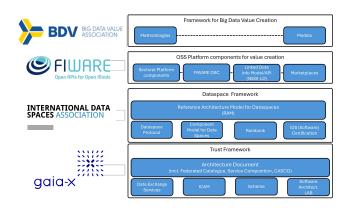
players, associations, research organizations, innovators, and policymakers worldwide. With

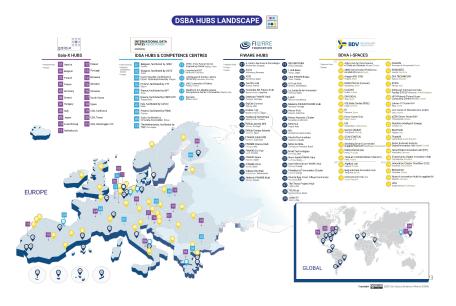
this cross-industry expertise, resources and know-how, the Alliance drives awareness, evangelizes technology, shapes standards, and enables integration across industries.

The Data Spaces Business Alliance

Unleashing the European Data Economy







Our members are the backbone of IDSA

Adapted from a presentation by L. Nagel at the Data Spaces Symposium 2024



People contributing

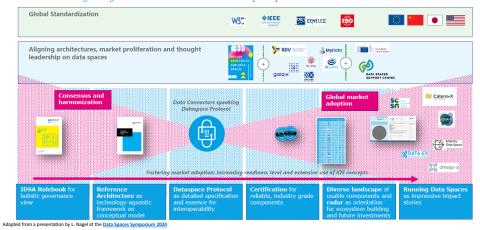
31 Countries

INTERNATIONAL DATA

A holistic approach to bring data spaces to global scale

INTERNATIONAL DATA
SPACES ASSOCIATION

IDSA on its way to a global standard – with the dataspace protocol in its core

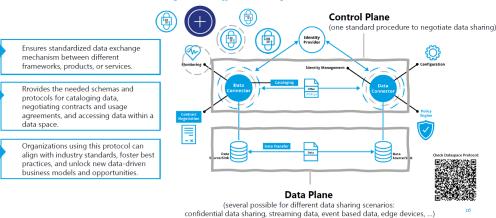




Dataspace Protocol V1.0 → ISO Standard

INTERNATIONAL DATA
SPACES ASSOCIATION

Enables standardized data exchange across different data space instances.

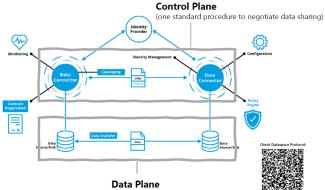


Dataspace Protocol V1.0 → ISO Standard

INTERNATIONAL DATA SPACES ASSOCIATION

Enables standardized data exchange across different data space instances.

Control Plane decides who can access the data and how. Data Plane is where the action (data sharing) happens. Conceptually divided, can be combined practically



(several possible for different data sharing scenarios: confidential data sharing, streaming data, event based data, edge devices, ...)

Make the connection and enable data economy

The key to data spaces is the data connector

- » Connects participants in a data space to share, utilize, benefit from data.
- » Ensures trust through IDS Certification and cyber security assessment.
- » Connects to trust frameworks and identity management
- » Includes identity & policy management, ensures data usage control.
- » Guarantees interoperability.
- » Understands and enforces data usage policies.
- » Master for other connectors of diverse feature sets.



IDSA Storu

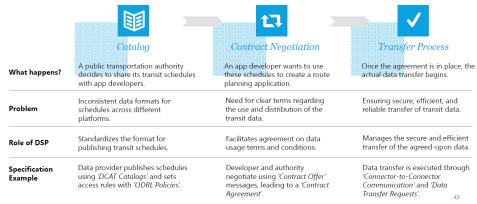
Adapted from a presentation by L. Nagel at the Data Spaces Symposium 2024

INTERNATIONAL DATA
SPACES ASSOCIATION

Standardized Data Exchange

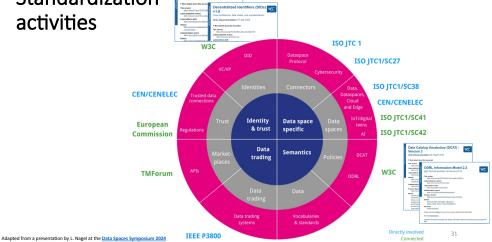
What does this mean? How does Dataspace Protocol ensure that?





Adapted from a presentation by L. Nagel at the Data Spaces Symposium 2024

Standardization activities



Data Interoperability and Semantics

< Part 5. The European Data Strategy >



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