

Outline

- DNS
- Charsets
- Email
- Web
- **HTML**

Web elements

- Protocol
 - **HTTP** (HyperText Transfer Protocol)
- Information (format)
 - **HTML** (HyperText Markup Language)
- LINK to information
 - **URI** (Uniform Resource Identifier):
URN (Name), **URL** (Locator)

HTML – Hyper-Text Markup Language, HTML

- In 1986 ISO standardized the Standard Generalized Markup Language (**SGML**). SGML introduced the `<>` syntax, and has been used in large documentation projects.
- Tim Berners-Lee defined **HTML** in 1989 inspired in SGML. HTML design goal was **displaying formatted** text documents with **hyperlinks** (including links to other documents) in **web browsers**.
- Based on **tags** e.g. `<head> data </head>`
- **Example:**

```
<html>
<head>
  <title>Basic html document</title>
</head>
<body>
  <h1><font color="red">First Heading</font></h1>
  <p>first paragraph.</p>
</body>
</html>
```

First Heading

first paragraph.

Terminology:

- **element**
- **attribute**
- **text**

HTML – Hyper-Text Markup Language, HTML

- HTML features (1):
 - **Forms**: The document accept user inputs that are sent to the server
 - **Scripting**: Allow adding programs. The program executes on the client's machine when the document loads, or at some other time such as when a link is activated.
- **javascript example**:

```
<html>
<head>
<script type="text/javascript">
  function displaymessage() {
    alert("Hello World!");
  }
</script>
</head>
<body>
  <form>
    <input type="button"
      value="Click me!" onclick="displaymessage()" />
  </form>
</body>
</html>
```

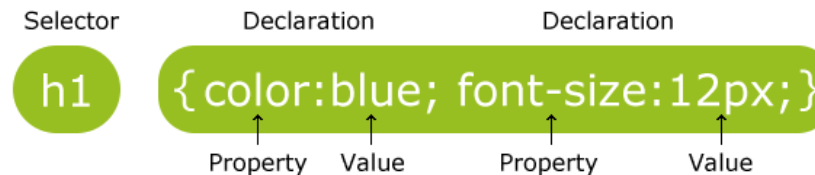


HTML – Hyper-Text Markup Language, HTML

- HTML features (2):

- Cascading Style Sheets, CSS:** Allows describing the *physical layout* in a separate document. E.g. thousand of HTML pages can use the same CSS. If the style must be changed, only the CSS need to be updated.

- CSS Syntax**



Source: <http://www.w3schools.com/xml/>

- CSS example**

- Content of the file “**mystyle.css**”:

```
h1 {color:red; font-size:20px;}
p {margin-left:20px; color:blue; font-size:18px;}
```

```
<html>
<head>
<link rel="stylesheet" type="text/css" href="mystyle.css" />
</head>
<body>
  <h1>First Heading</h1>
  <p>first paragraph.</p>
</body>
</html>
```

First Heading

first paragraph.

Outline

- DNS
- **Charsets**
- Email
- Web
- HTML

Languages, cultures, alphabets

7400 million people (2016)

22% speak Chinese, 11% English, 7% Spanish, 0,1% Catalan

Apart from languages, there are cultures and alphabets

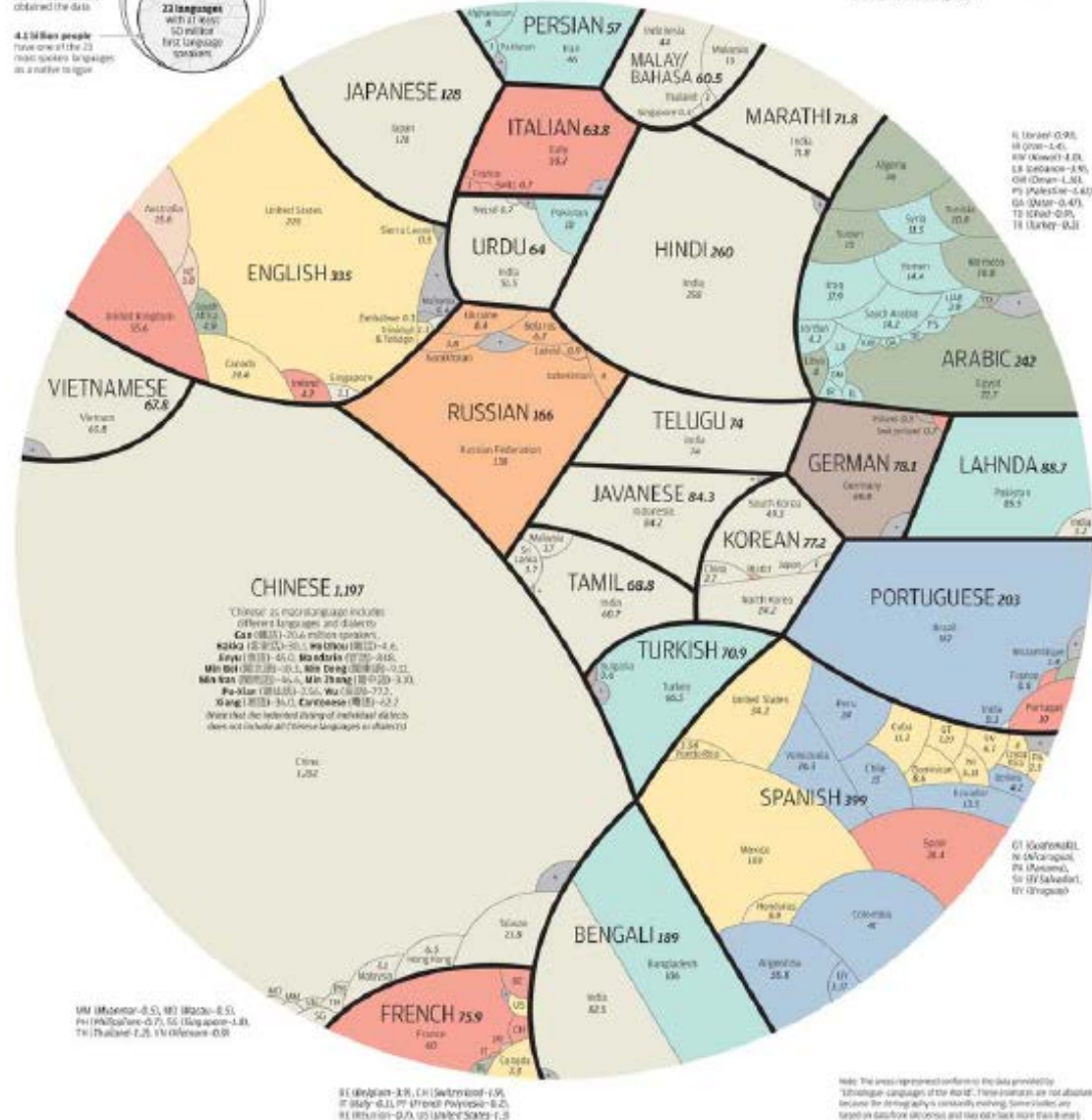
- Language with several cultures: es_ES, es_CO ("locale")
- Alphabet shared by several languages (e.g. català & français)

Culture:

- Messages, character sets, transliteration, ordering, search in strings, hours and dates, numbers and currency, pronunciation, ...

Interaction between agents in different languages and cultures:
alphabets and character sets

There are at least 7,000 known languages alive in the world today. Twenty-three of these languages are a mother tongue for more than 30 million people. The 25 languages make up the native tongue of 4.1 billion people. We represent each language within black borders and then provide the numbers of native speakers (in millions) by country. The colour of these countries shows how languages have taken root in many different regions.



Languages, cultures, alphabets

Internacionalization (i18n), Localization (l10n)

Alphabets

- "base": ascii
- National: e.g.: latin-1 (includes ascii), kanji
- International: e.g.: unicode (includes latin-1 and “all” languages)

Expression or language negotiation (in HTTP):

Accept-Language: es, ca, en-gb, en
Accept-Charset: iso-8859-15, unicode-9-0
...



Content-Language: ca
Content-Type: text/html; charset=utf-8
...

English is the default ...




Character sets

Characters are encoded following several conventions:

- **repertoire**: a set of characters (name and representation (glyph))
- **code**: correspondence between repertoire and natural numbers.
- **encoding**: method (algorithm) to convert code numbers into a sequence of octets (> 256 characters)
- US-ASCII: 95 characters + control=128: 7 bits (1 octet sent)

USASCII code chart



				<div><div>0 0 0 0</div><div>0 0 0 1</div><div>0 0 1 0</div><div>0 0 1 1</div><div>0 1 0 0</div><div>0 1 0 1</div><div>0 1 1 0</div><div>0 1 1 1</div><div>1 0 0 0</div><div>1 0 0 1</div><div>1 0 1 0</div><div>1 0 1 1</div><div>1 1 0 0</div><div>1 1 0 1</div><div>1 1 1 0</div><div>1 1 1 1</div></div>								
				0	1	2	3	4	5	6	7	
0	0	0	0	0	NUL	DLE	SP	0	@	P	`	p
0	0	0	1	1	SOH	DC1	!	1	A	Q	a	q
0	0	1	0	2	STX	DC2	"	2	B	R	b	r
0	0	1	1	3	ETX	DC3	#	3	C	S	c	s
0	1	0	0	4	EOT	DC4	\$	4	D	T	d	t
0	1	0	1	5	ENQ	NAK	%	5	E	U	e	u
0	1	1	0	6	ACK	SYN	&	6	F	V	f	v
0	1	1	1	7	BEL	ETB	'	7	G	W	g	w
1	0	0	0	8	BS	CAN	(8	H	X	h	x
1	0	0	1	9	HT	EM)	9	I	Y	i	y
1	0	1	0	10	LF	SUB	*	:	J	Z	j	z
1	0	1	1	11	VT	ESC	+	;	K	[k	{
1	1	0	0	12	FF	FS	,	<	L	\	l	
1	1	0	1	13	CR	GS	-	=	M]	m	}
1	1	1	0	14	SO	RS	.	>	N	^	n	~
1	1	1	1	15	SI	US	/	?	O	_	o	DEL

ISO 8859

- ISO 8859-1 (ISO Latin 1): 190 + control = 256: 1 octet
Western European, default for HTTP

- More variants

ISO 8859-15 extends -1 + ÿ, €

ISO 8859-2 (Central European)

ISO 8859-4 (North European)

ISO 8859-5 (Cyrillic)

ISO 8859-6 (Arabic) — Most common Arabic glyphs

ISO 8859-7 (Greek)

ISO 8859-8 (Hebrew) — modern Hebrew.

ISO 8859-9 (Turkish, Kurdish)

ISO 8859-11 (Thai) — Contains most glyphs needed

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
	í	í	í	í	í	í	í	í	í	í	í	í	í	í	í
B0	B1	B2	B3	B4	B5	B6	B7	B8	B9	BA	BB	BC	BD	BE	BF
°	±	²	³	¼	½	¾	¸	¹	º	»	¼	½	¾	¸	¹
C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA	DB	DC	DD	DE	DF
Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF
ä	ñ	ö	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

Universal Coded Character Set Unicode

All characters from all written languages + math + emoticons
= Universal Character Set (UCS)

Encoding: UCS-4 bytes (fixed length)

Proportional spacing, language independent

Unicode consortium: synchronized with ISO,

- Unicode 9.0.0 (7/2016): 128,172 symbols 🤪 🍲
- U+hex code: U+0020 = ' '

Character Encodings: Universal Transformation Format (UTF)

- Difficulty or impossibility to transport 8 or 16 bits data in Internet protocols:
- UTF-7, **UTF-8**, UTF-16, UTF-32 (variable length)



Universal Coded Character Set Unicode



• UTF-8 Encoding

- Determine high-order bits from the number of octets
- Fill in the bits marked x

Char. number range (hexadecimal)	UTF-8 octet sequence (binary)
0000 0000-0000 007F	0xxxxxxx
0000 0080-0000 07FF	110xxxxx 10xxxxxx
0000 0800-0000 FFFF	1110xxxx 10xxxxxx 10xxxxxx
0001 0000-0010 FFFF	11110xxx 10xxxxxx 10xxxxxx 10xxxxxx

• Example

- character: €
- code point: U+20AC
- code point in binary (12 bits): 10 0000 1010 1100
- 3 code units required:
- UTF-8: 11100010 10000010 10101100
- UTF-8 in hex: E282AC

UTF-8

Unicode (or Universal Coded Character Set) Transformation Format – 8-bit

This table shows UTF-8 as it is since 2003 (the `x` characters are replaced by the bits of the code point):

UTF-8 (2003)

Number of bytes	Bits for code point	First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4
1	7	U+0000	U+007F	0xxxxxxx			
2	11	U+0080	U+07FF	110xxxxx	10xxxxxx		
3	16	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
4	21	U+10000	U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

Character	Octal code point	Binary code point	Binary UTF-8	Octal UTF-8	Hexadecimal UTF-8
\$ U+0024	044	010 0100	00100100	044	24
¢ U+00A2	0242	000 1010 0010	11000010 10100010	302 242	C2 A2
€ U+20AC	020254	0010 0000 1010 1100	11100010 10000010 10101100	342 202 254	E2 82 AC
Ⓞ U+10348	0201510	0 0001 0000 0011 0100 1000	11110000 10010000 10001101 10001000	360 220 215 210	F0 90 8D 88

Source: Wikipedia