

HTML CHARACTER ENCODINGS

http://www.tutorialspoint.com/html/html_character_encodings.htm

Copyright © tutorialspoint.com

Character encoding is a method of converting bytes into characters. To validate or display an HTML document properly, a program must choose a proper character encoding.

The most common character set or character encoding in use on computers is ASCII **The American Standard Code for Information Interchange**, and this is probably the most widely used character set for encoding text electronically.

ASCII encoding supports only the upper- and lowercase Latin alphabet, the numbers 0-9, and some extra characters which make a total of 128 characters in all. You can have a look at complete set of [Printable ASCII Characters](#)

However, many languages use either accented Latin characters or completely different alphabets. ASCII does not address these characters; therefore you need to learn about character encodings if you want to use any non-ASCII characters.

The International Standards Organization created a range of character sets to deal with different national characters. For the documents in English and most other Western European languages, the widely supported encoding ISO-8859-1 is used.

Here is the list of Character Set being used around the world along with their description.

Character Set	Description
ISO-8859-1	Latin alphabet part 1 Covering North America, Western Europe, Latin America, the Caribbean, Canada, Africa
ISO-8859-2	Latin alphabet part 2 Covering Eastern Europe
ISO-8859-3	Latin alphabet part 3 Covering SE Europe, Esperanto, miscellaneous others
ISO-8859-4	Latin alphabet part 4 Covering Scandinavia/Baltics <i>and others not in ISO – 8859 – 1</i>
ISO-8859-5	Latin/Cyrillic alphabet part 5
ISO-8859-6	Latin/Arabic alphabet part 6
ISO-8859-7	Latin/Greek alphabet part 7
ISO-8859-8	Latin/Hebrew alphabet part 8
ISO-8859-9	Latin 5 alphabet part 9 Same as ISO-8859-1 except Turkish characters replace Icelandic ones
ISO-8859-10	Latin 6 Latin 6 Lappish, Nordic, and Eskimo
ISO-8859-15	The same as ISO-8859-1 but with more characters added
ISO-2022-JP	Latin/Japanese alphabet part 1
ISO-2022-JP-2	Latin/Japanese alphabet part 2
ISO-2022-KR	Latin/Korean alphabet part 1

The Unicode Consortium was then set up to devise a way to show all characters of different languages, rather than have these different incompatible character codes for different languages.

Therefore, if you want to create documents that use characters from multiple character sets, you will be able to do so using the single Unicode character encodings.

Unicode therefore specifies encodings that can deal with a string in special ways so as to make enough space for the huge character set it encompasses. These are known as UTF-8, UTF-16, and UTF-32.

Character Set	Description
UTF-8	A Unicode Translation Format that comes in 8-bit units that is, it comes in bytes. A character in UTF8 can be from 1 to 4 bytes long, making UTF8 variable width.
UTF-16	A Unicode Translation Format that comes in 16-bit units that is, it comes in shorts. It can be 1 or 2 shorts long, making UTF16 variable width.
UTF-32	A Unicode Translation Format that comes in 32-bit units that is, it comes in longs. It is a fixed-width format and is always 1 "long" in length.

The first 256 characters of Unicode character sets correspond to the 256 characters of ISO-8859-1.

By default, HTML 4 processors should support UTF-8, and XML processors are supposed to support UTF-8 and UTF-16; therefore all XHTML-compliant processors should also support UTF-16.

Loading [Mathjax]/jax/output/HTML-CSS/fonts/TeX/fontdata.js