

Codificação de caracteres

- Representação binária (ou hexa) de cada símbolo
- ASCII (anos 1960)
American Standard Code for Information Interchange
 - Representação de 128 símbolos (7 bits)
 - $100\ 0001 = 65 = 0x41 = 'A'$
 - O 8º bit (mais significativo) era usado como bit de paridade nas comunicações
 - Faltavam os caracteres latinos acentuados, os chineses, os cirílicos (russos), ...

ASCII

0 null	16 data link escape	32 space	48 0	64 @	80 P	96 `	112 p
1 start of heading	17 device control 1	33 !	49 1	65 A	81 Q	97 a	113 q
2 start of text	18 device control 2	34 "	50 2	66 B	82 R	98 b	114 r
3 end of text	19 device control 3	35 #	51 3	67 C	83 S	99 c	115 s
4 end of transmission	20 device control 4	36 \$	52 4	68 D	84 T	100 d	116 t
5 enquiry	21 negative acknowledge	37 %	53 5	69 E	85 U	101 e	117 u
6 acknowledge	22 synchoronous idle	38 &	54 6	70 F	86 V	102 f	118 v
7 bell	23 end of trans. block	39 '	55 7	71 G	87 W	103 g	119 w
8 backspace	24 cancel	40 (56 8	72 H	88 X	104 h	120 x
9 horizontal tab	25 end of medium	41)	57 9	73 I	89 Y	105 i	121 y
10 line feed	26 substitute	42 *	58 :	74 J	90 Z	106 j	122 z
11 vertical tab	27 escape	43 +	59 ;	75 K	91 [107 k	123 {
12 form feed	28 file separator	44 ,	60 <	76 L	92 \	108 l	124
13 carriage return	29 group separator	45 -	61 =	77 M	93]	109 m	125 }
14 shift out	30 record separator	46 .	62 >	78 N	94 ^	110 n	126 ~
15 shift in	31 unit separator	47 /	63 ?	79 O	95 _	111 o	127 del

Codificações de 8 bits

ISO-8859-1 (caracteres latinos)

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
8_	U+0080	U+0081	U+0082	U+0083	U+0084	U+0085	U+0086	U+0087	U+0088	U+0089	U+008A	U+008B	U+008C	U+008D	U+008E	U+008F
	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA	HTS	HTJ	LTS	PLD	PLU	RI	SS2	SS3
	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
9_	U+0090	U+0091	U+0092	U+0093	U+0094	U+0095	U+0096	U+0097	U+0098	U+0099	U+009A	U+009B	U+009C	U+009D	U+009E	U+009F
	DCS	PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
A_	U+00A0	U+00A1	U+00A2	U+00A3	U+00A4	U+00A5	U+00A6	U+00A7	U+00A8	U+00A9	U+00AA	U+00AB	U+00AC	U+00AD	U+00AE	U+00AF
	NBSP	¡	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	®	¯	
	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
B_	U+00B0	U+00B1	U+00B2	U+00B3	U+00B4	U+00B5	U+00B6	U+00B7	U+00B8	U+00B9	U+00BA	U+00BB	U+00BC	U+00BD	U+00BE	U+00BF
	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
C_	U+00C0	U+00C1	U+00C2	U+00C3	U+00C4	U+00C5	U+00C6	U+00C7	U+00C8	U+00C9	U+00CA	U+00CB	U+00CC	U+00CD	U+00CE	U+00CF
	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï
	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
D_	U+00D0	U+00D1	U+00D2	U+00D3	U+00D4	U+00D5	U+00D6	U+00D7	U+00D8	U+00D9	U+00DA	U+00DB	U+00DC	U+00DD	U+00DE	U+00DF
	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223
E_	U+00E0	U+00E1	U+00E2	U+00E3	U+00E4	U+00E5	U+00E6	U+00E7	U+00E8	U+00E9	U+00EA	U+00EB	U+00EC	U+00ED	U+00EE	U+00EF
	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
F_	U+00F0	U+00F1	U+00F2	U+00F3	U+00F4	U+00F5	U+00F6	U+00F7	U+00F8	U+00F9	U+00FA	U+00FB	U+00FC	U+00FD	U+00FE	U+00FF
	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ
	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255

Codificações de 8 bits

ISO-8859-7 (caracteres gregos)

	_0	_1	_2	_3	_4	_5	_6	_7	_8	_9	_A	_B	_C	_D	_E	_F
8_	U+0080	U+0081	U+0082	U+0083	U+0084	U+0085	U+0086	U+0087	U+0088	U+0089	U+008A	U+008B	U+008C	U+008D	U+008E	U+008F
	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA	HTS	HTJ	LTS	PLD	PLU	RI	SS2	SS3
	128	129	130	131	132	133	134	135	136	137	138	139	140	141	142	143
9_	U+0090	U+0091	U+0092	U+0093	U+0094	U+0095	U+0096	U+0097	U+0098	U+0099	U+009A	U+009B	U+009C	U+009D	U+009E	U+009F
	DCS	PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
	144	145	146	147	148	149	150	151	152	153	154	155	156	157	158	159
A_	U+00A0	U+2018	U+2019	U+00A3	U+20AC	U+20AF	U+00A6	U+00A7	U+00A8	U+00A9	U+037A	U+00AB	U+00AC	U+00AD		U+2015
	NBSP	‘	’	£	€	ƒ	ı	§	¨	©		«	¬			—
	160	161	162	163	164	165	166	167	168	169	170	171	172	173	174	175
B_	U+00B0	U+00B1	U+00B2	U+00B3	U+0384	U+0385	U+0386	U+00B7	U+0388	U+0389	U+038A	U+00BB	U+038C	U+00BD	U+038E	U+038F
	°	±	2	3	´	ˆ	À	·	È	É	Ì	»	Ó	½	Υ	Ω
	176	177	178	179	180	181	182	183	184	185	186	187	188	189	190	191
C_	U+0390	U+0391	U+0392	U+0393	U+0394	U+0395	U+0396	U+0397	U+0398	U+0399	U+039A	U+039B	U+039C	U+039D	U+039E	U+039F
	ĭ	À	B	Γ	Δ	E	Z	H	Θ	I	K	Λ	M	N	Ξ	O
	192	193	194	195	196	197	198	199	200	201	202	203	204	205	206	207
D_	U+03A0	U+03A1		U+03A3	U+03A4	U+03A5	U+03A6	U+03A7	U+03A8	U+03A9	U+03AA	U+03AB	U+03AC	U+03AD	U+03AE	U+03AF
	Π	P		Σ	T	Υ	Φ	X	Ψ	Ω	İ	Ÿ	á	é	η	í
	208	209	210	211	212	213	214	215	216	217	218	219	220	221	222	223
E_	U+03B0	U+03B1	U+03B2	U+03B3	U+03B4	U+03B5	U+03B6	U+03B7	U+03B8	U+03B9	U+03BA	U+03BB	U+03BC	U+03BD	U+03BE	U+03BF
	Û	α	β	γ	δ	ε	ζ	η	θ	ι	κ	λ	μ	ν	ξ	ο
	224	225	226	227	228	229	230	231	232	233	234	235	236	237	238	239
F_	U+03C0	U+03C1	U+03C2	U+03C3	U+03C4	U+03C5	U+03C6	U+03C7	U+03C8	U+03C9	U+03CA	U+03CB	U+03CC	U+03CD	U+03CE	
	π	ρ	ς	σ	τ	υ	φ	χ	ψ	ω	ϊ	ϋ	ό	ύ	ώ	
	240	241	242	243	244	245	246	247	248	249	250	251	252	253	254	255

Unicode

- Sistema de codificação de caracteres capaz de representar texto de qualquer sistema de escrita existente (mais de 100 mil símbolos)
 - Sistemas latino, árabe, cirílico, chinês, hebraico, ...
 - Símbolos matemáticos, geométricos, musicais, setas, ícones, emojis, ...
 - Escrita cuneiforme, Braille, runas, élfico, ...
 - Próximos: hieróglifos egípcios, alfabeto babilônico, ...
- Referências:
 - <http://unicode.org/>
 - <http://www.unicodetables.com/>



Unicode

- **UTF-8** (*8-bit Unicode Transformation Format*)
 - Codificação Unicode de comprimento variável (1 a 4 bytes), que pode representar qualquer caráter (*code point*) do padrão.
 - Padrão usado na Internet e na Web
- UTF-16
 - Codificação de 2 ou 4 bytes, usada principalmente para escrita em idiomas dos países asiáticos
- UTF-32
 - Codificação de comprimento fixo de 4 bytes

UTF-8

Bits do <i>code point</i>	Primeiro <i>code point</i>	Último <i>code point</i>	Bytes	Byte 1	Byte 2	Byte 3	Byte 4
7	U+0000	U+007F	1	0xxxxxxx			
11	U+0080	U+07FF	2	110xxxxx	10xxxxxx		
16	U+0800	U+FFFF	3	1110xxxx	10xxxxxx	10xxxxxx	
21	U+10000	U+1FFFFF	4	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

- 1 byte - Tabela ASCII
- 2 bytes - Caracteres latinos, hebraicos, gregos, ...
- 3 bytes - Caracteres chineses, japoneses, coreanos, ...
- 4 bytes - Alguns outros caracteres e símbolos

UTF-8

Representação de strings

C	o	n	c	e	i	ç	ã	o
43	6F	6E	63	65	69	C3A7	C3A3	6F

9 caracteres – 11 bytes