

# Storing text in binary

Created by Pamela Fox.

Computers store more than just numbers in binary. But how can binary numbers represent non-numbers such as letters and symbols?

As it turns out, all it requires is a bit of human cooperation. We must agree on **encodings**, mappings from a character to a binary number.

## ASCII encoding

ASCII was one of the first standardized encodings. It was invented back in the 1960s when telegraphy was the primary form of long-distance communication, but is still in use today on modern computing systems.

Teletypists would type messages on teleprinters such as this one:



Photo of a teletype machine, composed of a mechanical keyboard, a piece of paper coming out with typed letters, and a mechanism for reading input paper strips.

An ASR 33 teletype machine. Image source: [Marcin Wichary](#)

The teleprinter would then use the ASCII standard to encode each typed character into binary and then store or transmit the binary data.

This page from a 1972 teleprinter manual shows the 128 ASCII codes:

USASCII code chart

<div><div>b7b6b5</div><div>b4b3b2b1</div><div>Bits</div><div>Column</div><div>Row</div></div>					000	001	010	011	100	101	110	111	
					0	1	2	3	4	5	6	7	
00000					0	NUL	DLE	SP	0	@	P	\	p
00001					1	SOH	DC1	!	1	A	Q	a	q
00100					2	STX	DC2	"	2	B	R	b	r
00101					3	ETX	DC3	#	3	C	S	c	s
01000					4	EOT	DC4	\$	4	D	T	d	t
01001					5	ENQ	NAK	%	5	E	U	e	u
01100					6	ACK	SYN	&	6	F	V	f	v
01101					7	BEL	ETB	'	7	G	W	g	w
10000					8	BS	CAN	(	8	H	X	h	x
10001					9	HT	EM	)	9	I	Y	i	y
10100					10	LF	SUB	*	:	J	Z	j	z
10101					11	VT	ESC	+	;	K	[	k	{
11000					12	FF	FS	,	<	L	\	l	
11001					13	CR	GS	-	=	M	]	m	}
11100					14	SO	RS	.	>	N	^	n	~
11101					15	SI	US	/	?	O	_	o	DEL

The first 32 codes represent "control characters," characters which cause some effect besides printing a letter. "BEL" (encoded in binary as 0000111) caused an audible bell or beep. "ENQ" (encoded as 0000101) represented an enquiry, a request for the receiving station to identify themselves.

The control characters were originally designed for teleprinters and telegraphy, but many have been re-purposed for modern computers and the Internet—especially "CR" and "LF". "CR" (0001101) represented a "carriage return" on teleprinters, moving the printing head to the start of the line. "LF" (0001010) represented a "line feed", moving the printing head down one line. Modern Internet protocols, such as HTTP, FTP, and SMTP, use a combination of "CR" + "LF" to represent the end of lines.

The remaining 96 ASCII characters look much more familiar.

Here are the first 8 uppercase letters:

Binary	Character
1000001	A
1000010	B
1000011	C
1000100	D
1000101	E
1000110	F
1000111	G
1001000	H

There are several problems with the ASCII encoding, however.

The first big problem is that ASCII only includes letters from the English alphabet and a limited set of symbols.

A language that uses less than 128 characters could come up with their own version of ASCII to encode text in just their language, but what about a text file with characters from multiple languages? ASCII couldn't encode a string like: "Hello, José, would you care for Glühwein? It costs 10 €".

And what about languages with thousands of logograms? ASCII could not encode enough logograms to cover a Chinese sentence like

"你好，想要一盘饺子吗？十块钱。"

The other problem with the ASCII encoding is that it uses **7** bits to represent each character, whereas computers typically store information in bytes—units of **8** bits—and programmers don't like to waste memory.

When the earliest computers first started using ASCII to encode characters, different computers would come up with various ways to utilize the final bit. For example, HP computers used the eighth bit to represent characters used in European countries (e.g. "£" and "Ü"), TRS-80 computers used the bit for colored graphics, and Atari computers used the bit for inverted white-on-black versions of the first 128 characters.

The result? An "ASCII" file created in one application might look like gobbledy gook when opened in another "ASCII"-compatible application.

Computers needed a new encoding, an encoding based on 8-bit bytes that could represent all the languages of the world.

# Unicode

But first, how many characters do you even need to represent the world's languages? Which characters are basically the same across languages, even if they have different sounds?

In 1987, a group of computer engineers attempted to answer those questions. They eventually came up with Unicode, a universal character set which assigns each a "code point" (a hexadecimal number) and a name to each character. ^33cubed

For example, the character "ą" is assigned to "U+0105" and named "Latin Small Letter A with Ogonek". There's a character that looks like "ą" in 13 languages, such as Polish and Lithuanian. Thus, according to Unicode, the "ą" in the Polish word "robią" and the "ą" in the Lithuanian word "aslą" are both the same character. Unicode saves space by unifying characters across languages.

But there are still quite a few characters to encode. The Unicode character set started with 7,129 named characters in 1991 and has grown to 137,929 named characters in 2019. The majority of those characters describe logograms from Chinese, Japanese, and Korean, such as "U+6728" which refers to "木". It also includes over 1,200 emoji symbols ("U+1F389" = "🍷").

Unicode is a character set, but it is *not* an encoding. Fortunately, another group of engineers tackled the problem of efficiently encoding Unicode into binary.

## UTF-8

In 1992, computer scientists invented UTF-8, an encoding that is compatible with ASCII encoding but also solves its problems. <sup>5</sup>start superscript, 5, end superscript

UTF-8 can describe every character from the Unicode standard using either 1, 2, 3, or 4 bytes.

When a computer program is reading a UTF-8 text file, it knows how many bytes represent the next character based on how many <sub>1</sub> bits it finds at the beginning of the byte.

Number of bytes	Byte 1	Byte 2	Byte 3	Byte 4
1	0xxxxxxx			
2	110xxxxx	10xxxxxx		
3	1110xxxx	10xxxxxx	10xxxxxx	
4	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

If there are no <sub>1</sub> bits in the prefix (if the first bit is a <sub>0</sub>), that indicates a character represented by a single byte. The remaining 7 bits of the byte are used to represent the original 128 ASCII characters. That means a sequence of 8-bit ASCII characters is also a valid UTF-8 sequence.

Two bytes beginning with <sub>110</sub> are used to encode the rest of the characters from Latin-script languages (e.g. Spanish, German) plus other languages such as Greek, Hebrew, and Arabic. Three bytes starting with <sub>1110</sub> encode most of the characters for Asian languages (e.g. Chinese, Japanese, Korean). Four bytes starting with <sub>11110</sub> encode everything else,



from rarely used historical scripts to the increasingly commonly used emoji symbols.

Most modern programming languages have built-in support for UTF-8, so most programmers never need to know exactly how to convert from characters to binary.

The UTF-8 encoding standard is now the dominant encoding of HTML files on the web, accounting for 94.5% of webpages as of December 2019.

🔍 If you right click and select "view page source" on this webpage right now, you can search for the string "utf-8" and see that this webpage is encoded as UTF-8.

Generally, a good encoding is one that can represent the maximum amount of information with the least number of bits. UTF-8 is a great example of that, since it can encode common English letters with just 1 byte but is flexible enough to encode thousands of letters with additional bytes.

UTF-8 is only one possible encoding, however. UTF-16 and UTF-32 are alternative encodings that are also capable of representing all Unicode characters. There are also language specific encodings such as Shift-JIS for Japanese. Computer programs can use the encoding that best suits their needs and constraints.