

Character Encodings and Unicode

This chapter describes how characters are encoded into bits.

After this chapter you should be able to convert between characters, Unicode code points and their UTF-8 encodings.

Unicode Character Encodings

Data often represents printable characters.

A standard called **Unicode** assigns a unique number, called a “code point,” to over 100,000 of the characters used by more than 100 languages and scripts. Unicode is used by most operating systems and Internet applications.

Exercise 1: How many bits would be required to uniquely identify 100,000 different characters? (Hint: $2^{16} = 65536$).

Each character could be represented with a 32-bit (4 byte) number. However, the UTF-8 (Unicode Transformation Format - 8 bits) format is widely used since it allows many documents to be stored or transmitted using fewer bytes. It also has some practical advantages¹

ASCII (American Standard Code for Information Interchange) was an earlier character encoding that used 7 bits to encode letters from the English alphabet, numbers, and the most common punctuation symbols. UTF-8 is exactly the same as ASCII for these first 128 (2^7) characters. This means that ASCII documents are already encoded as UTF-8.

The table below shows the ASCII table which is the first “code chart” from the Unicode standard. The columns are labelled with the most significant (first) hex digit and the rows with the least-significant (second) hex digit of the numerical value of each character.

ASCII also includes some non-printable control codes (values 0 to 31) that were used to control printers. For example the line feed (LF) character would move the paper in the printer up one line.

Other Unicode characters require between 2 and 4 bytes according to the rules summarized in Table 3-6 of the Unicode standard shown below. Unicode values between 128 and 2047 include most characters from European languages and can be encoded in two bytes. Values from 2048 to 65535 include most

¹The UTF-8 encoding has no zero bytes and all values $<0x80$ represent ASCII characters.

CJK (Chinese, Japanese and Korean) characters and require three bytes. Some rarely-used symbols (e.g. emoticons or Mahjong tiles) have four-byte encodings.

Encoding to UTF-8

Step 1 From the value of the character’s code point choose a sequence of bytes (all numbers in hexadecimal):

code point	prefixes
0 – 7F	00
80 – 07FF	C0 80
800 – FFFF	E0 80 80
> FFFF	F0 80 80 80

Step 2 Convert the code point to binary and divide the value into groups of 6 bits.

Step 3 Starting at the right, add each group of 6 bits to the corresponding byte, starting at the right.

Example

The codepoint for the CJK character for potato (薯) is U+85F7. From the table above, this must be encoded into the three bytes E0, 80 and 80. The code point in binary is 1000 0101 1111 0111. The groups of 6 bits, starting on the right are thus 11 0111 (37), 01 01 11 (17), and 00 1000 (08). Adding these to the prefixes the bytes in the UTF-8 encoding are (E0+08=E8), (80+17=97), and (80+37=B7).

Exercise 2: The Chinese character for “Rice” (the grain) is 米 with Unicode value (code point) U+7C73. What is the UTF-8 encoding for this character?

Decoding from UTF-8

The most significant nybble of each byte in a UTF-8 sequence identifies its purpose:

Table 3-6. UTF-8 Bit Distribution

Scalar Value	First Byte	Second Byte	Third Byte	Fourth Byte
00000000 0xxxxxxx	0xxxxxxx			
00000yyy yyxxxxxx	110yyyyy	10xxxxxx		
zzzyyyyy yyxxxxxx	1110zzzz	10yyyyyy	10xxxxxx	
000uuuuu zzzzyyyy yyxxxxxx	11110uuu	10uuzzzz	10yyyyyy	10xxxxxx

nybble	meaning
0 – 7	ASCII, one byte
8, 9, A, or B	the second, third or fourth byte of an encoding
C or D	start of a 2-byte encoding
E	start of a 3-byte encoding
F	start of a 4-byte encoding

Step 1 Scan for a byte whose most significant nybble is not 8, 9, A or B. This will be the first byte of a UTF-8 sequence.

Step 2 Extract the appropriate number of bytes (1 through 4).

Step 3 Delete the bits indicated by 0 or 1 in Table 3-6 above.

Step 4 Concatenate the remaining bits. This is the code point in binary. Group every 4 bits starting on the right and convert the nybbles to hex.

Example

Find the first Unicode character in the UTF sequence BC, D0, BE. First we ignore the byte BC because the initial nybble is B. Then we extract 2 bytes because the most significant nybble of the byte D0 indicates a 2-byte encoding. Then we convert the two bytes to binary: D0 = 1101 0000 and BE = 1011 1110. Deleting the leading three bits from the first byte and the leading two bits from the second we are left with 1 0000 11 1110. Grouping into nybbles this is 0100 0011 1110 and converting to hexadecimal this is 43F. Looking this up in the Unicode charts (or on unicode.org) show this is the character π, “CYRILLIC SMALL LETTER PE.”

Exercise 3: Find the codepoint of the first Unicode character in the sequence of bytes: A0 88 EB 8C 80 EC.

Text versus Binary Number Representations

It’s important to understand the difference between text that represents a number and binary data. For example, the character ‘1’ would be transmitted with a UTF-8 encoding as the byte 0x31 while a byte with the value 1 could be transmitted as 0x01.

Numbers can be stored in files or transmitted over communication systems in either binary format (e.g. one 8-bit value per byte) or in text format (as a sequence of numeric characters). Numbers in text format can be more easily interpreted by humans since they are sequences of printable characters.

Exercise 4: Four numbers are transmitted as the following CSV file:
 2, 1
 9, 3
 How many bytes are required to transmit these four numbers formatted this way? Note that a “line feed” character is required at the end of each line and that spaces and commas also need to be transmitted.

How many bytes are required to transmit these four numbers if they are transmitted, one after another, if each is encoded as a 16-bit number? What if each number was encoded as a 32-bit number?

	000	001	002	003	004	005	006	007
0	NUL 0000	DLE 0010	SP 0020	0 0030	@ 0040	P 0050	` 0060	p 0070
1	SOH 0001	DC1 0011	! 0021	1 0031	A 0041	Q 0051	a 0061	q 0071
2	STX 0002	DC2 0012	" 0022	2 0032	B 0042	R 0052	b 0062	r 0072
3	ETX 0003	DC3 0013	# 0023	3 0033	C 0043	S 0053	c 0063	s 0073
4	EOT 0004	DC4 0014	\$ 0024	4 0034	D 0044	T 0054	d 0064	t 0074
5	ENQ 0005	NAK 0015	% 0025	5 0035	E 0045	U 0055	e 0065	u 0075
6	ACK 0006	SYN 0016	& 0026	6 0036	F 0046	V 0056	f 0066	v 0076
7	BEL 0007	ETB 0017	' 0027	7 0037	G 0047	W 0057	g 0067	w 0077
8	BS 0008	CAN 0018	(0028	8 0038	H 0048	X 0058	h 0068	x 0078
9	HT 0009	EM 0019) 0029	9 0039	I 0049	Y 0059	i 0069	y 0079
A	LF 000A	SUB 001A	* 002A	: 003A	J 004A	Z 005A	j 006A	z 007A
B	VT 000B	ESC 001B	+ 002B	; 003B	K 004B	[005B	k 006B	{ 007B
C	FF 000C	FS 001C	, 002C	< 003C	L 004C	\ 005C	l 006C	 007C
D	CR 000D	GS 001D	- 002D	= 003D	M 004D] 005D	m 006D	} 007D
E	SO 000E	RS 001E	. 002E	> 003E	N 004E	^ 005E	n 006E	~ 007E
F	SI 000F	US 001F	/ 002F	? 003F	O 004F	_ 005F	o 006F	DEL 007F