

Character encoding

Victor Eijkhout

August 2004

Have you ever wondered what goes on between the ‘A’ you hit on your keyboard, the ‘A’ stored in your file, and the ‘A’ that comes out of your printer? Why that letter still comes out of the printer if the file is printed by your friend in Greece who doesn’t use the letter ‘A’? Maybe you know that ‘A’ is character 41 in ascii; if you put it on a web page, and it’s watched by someone in Japan, why don’t they get character 41 in the Kanji alphabet? Do you remember the DOS days when your Mac owning colleague would send you a file and what were supposed to be accented characters would turn into smiley faces? Have you ever pasted text from MS-Word into Emacs, and Emacs wanted to save the document as UTF-8? Just what is that about?

All this, and more, will be explained in this article.

1 History in one byte

Somewhere in the depths of prehistory, people in the western world got to agree on a standard for character codes under 127, ASCII, the American Standard Code for Information Interchange. This standard declares that the letter ‘A’ is character number 41, so if your file contains the bit pattern for 41 (which is 00101001), it will produce an ‘A’ when sent to the printer.

ASCII has a few nice properties, some of which were not shared by another encoding scheme, EBCDIC (which was used almost exclusively by IBM):

- All letters are consecutive, making a test ‘is this a letter’ easy to perform.
- Uppercase and lowercase letters are at a distance of 32; this means that the Shift key on your keyboard simply toggles the sixth bit in the pattern of whatever key you are holding down.
- The first 31 codes, everything below the space character, as well as position 127, are ‘unprintable’, and can be used for such purposes as terminal cursor control.

The ISO 646 standard codified 7-bit ASCII, but it left certain character positions (or ‘code points’) open for national variation. For instance, British usage put a pound sign (£) in the position of the dollar. The ASCII character set was originally accepted as ANSI X3.4 in 1968.

dec

CHAR

hex oct

ASCII CONTROL CODES

<div> <div>b7</div> <div>b6</div> <div>b5</div> <div>BITS</div> <div>b4 b3 b2 b1</div> </div>	<div> <div>0</div> <div>0</div> <div>0</div> <div>0</div> <div>1</div> <div>1</div> <div>1</div> <div>1</div> </div>	<div> <div>0</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>1</div> </div>	<div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> </div>	<div> <div>1</div> <div>0</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> </div>	<div> <div>1</div> <div>0</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> </div>	<div> <div>1</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> </div>	<div> <div>1</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> <div>0</div> <div>1</div> </div>	
	CONTROL		SYMBOLS NUMBERS		UPPERCASE		LOWERCASE	
0 0 0 0	0 NUL	16 DLE	32 SP	48 0	64 @	80 P	96 ,	112 p
0 0 0 1	1 SOH	17 DC1	33 !	49 1	65 A	81 Q	97 a	113 q
0 0 1 0	2 STX	18 DC2	34 "	50 2	66 B	82 R	98 b	114 r
0 0 1 1	3 ETX	19 DC3	35 #	51 3	67 C	83 S	99 c	115 s
0 1 0 0	4 EOT	20 DC4	36 \$	52 4	68 D	84 T	100 d	116 t
0 1 0 1	5 ENQ	21 NAK	37 %	53 5	69 E	85 U	101 e	117 u
0 1 1 0	6 ACK	22 SYN	38 &	54 6	70 F	86 V	102 f	118 v
0 1 1 1	7 BEL	23 ETB	39 ,	55 7	71 G	87 W	103 g	119 w
1 0 0 0	8 BS	24 CAN	40 (56 8	72 H	88 X	104 h	120 x
1 0 0 1	9 HT	25 EM	41)	57 9	73 I	89 Y	105 i	121 y
1 0 1 0	10 LF	26 SUB	42 *	58 :	74 J	90 Z	106 j	122 z
1 0 1 1	11 VT	27 ESC	43 +	59 ;	75 K	91 [107 k	123 {
1 1 0 0	12 FF	28 FS	44 ,	60 <	76 L	92 \ 	108 l	124
1 1 0 1	13 CR	29 GS	45 —	61 =	77 M	93]	109 m	125 }
1 1 1 0	14 SO	30 RS	46 .	62 >	78 N	94 ^	110 n	126 ~
1 1 1 1	15 SI	31 US	47 /	63 ?	79 O	95 .	111 o	127 DEL

Table 1: The ASCII table

Since a computer organizes its bits in 8-bit bytes, and ASCII only codified the codes under 128, this left the codes with the high bit set ('extended ASCII') undefined, and different manufacturers of computer equipment came up with their own way of filling them in. These standards were called 'code pages', and IBM gave a standard numbering to them. For instance, code page 437 is the MS-DOS code page with accented characters for most European languages, 862 is DOS in Israel, 737 is DOS for Greeks.

Here is cp473:

	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
10	▶	◀	‡	!!	¶	§	■	‡	↑	↓	→	←	↖	↗	↘
20	!	!"	#	\$	%	&	'	()	*	+	,	-	.	/
30	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>
40	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N
50	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^
60	€	a	b	c	d	e	f	g	h	i	j	k	l	m	n
70	p	q	r	s	t	u	v	w	x	y	z	{		}	~
80	Ç	ü	ë	â	ä	à	â	ç	è	ë	ï	î	ï	Ä	Å
90	É	æ	œ	ô	ö	ò	û	ü	ö	ü	φ	£	¥	℔	f
A0	á	í	ó	ú	ñ	Ñ	æ	ø	ó	í	ó	ú	ñ	«	»
B0	☼	☼	☼	☼	☼	☼	☼	☼	☼	☼	☼	☼	☼	☼	☼
C0	L	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥
D0	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥	⊥
E0	α	β	γ	π	Σ	σ	μ	τ	Φ	Θ	Ω	δ	∞	∅	∈
F0	≡	±	≥	≤	∫	J	÷	∞	°	°	°	√	n	ε	■

MacRoman:

80	À	Á	Ç	È	É	Ë	Ü	ä	ä	â	ä	ä	Ç	é	è
90	ê	ë	î	ï	î	ï	ñ	ô	ö	ö	ö	ü	ü	ü	ü
A0	†	°	φ	£	§	•	¶	§	©	©	™	~	~	~	~
B0	∞	±	≤	≥	¥	μ	∂	Σ	Π	Π	∫	≈	Ω	æ	∅
C0	ó	í	í	√	f	≈	Δ	«	»	...	~	~	~	~	~
D0	-	-	«	»	«	»	÷	◊	ÿ	ÿ	/	℥	<	>	fi
E0	‡	•	,	,	%	^	^	^	^	^	^	^	^	^	^
F0	♥	ö	ü	ü	ü	ü	ü	ü	ü	ü	ü	ü	ü	ü	ü

E0	€	E2	f	E4	...	E6	†	E8	§	EA	S	EC	œ	EE	ž
	91	92	93	94	95	96	97	98	99	9A	9B	9C	9D	9E	9F
A0	i	A2	ƒ	A4	¥	A6	£	A8	...	BA	®	BC	«	BD	®
B0	±	B2	2	B4	...	B6	¶	B8	1	BB	»	BC	%	BD	®
C0	À	Ā	Ĉ	Ċ	Ċ	C6	Ĉ	C8	Ē	CA	Ĉ	CC	Ĉ	CD	Ĉ
D0	Đ	Ñ	Ŋ	Ŋ	Ŋ	D6	×	D8	Ŋ	D8	Ŋ	DC	Ŋ	DD	Ŋ
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF

The international variants were standardized as ISO 646-DE (German), 646-DK (Danish), et cetera. Originally, the dollar sign could still be replaced by the currency symbol, but after a 1991 revision the dollar is now the only possibility.

A0	A1	A2	A3	A4	A5	A6	A7	A8	A9	AA	AB	AC	AD	AE	AF
B0	B1	B2	B3	B4	B5	B6	B7	B8	B9	BA	BB	BC	BD	BE	BF
C0	C1	C2	C3	C4	C5	C6	C7	C8	C9	CA	CB	CC	CD	CE	CF
D0	D1	D2	D3	D4	D5	D6	D7	D8	D9	DA	DB	DC	DD	DE	DF
E0	E1	E2	E3	E4	E5	E6	E7	E8	E9	EA	EB	EC	ED	EE	EF
F0	F1	F2	F3	F4	F5	F6	F7	F8	F9	FA	FB	FC	FD	FE	FF

А0	В1	Г2	Д3	Е4	Ж5	З6	И7	Й8	Ј9	Љ10	Њ11	Ћ12	Ќ13	Ў14	Ф15	Х16	Ц17	Ч18	Ш19	Щ20	Ъ21	Ь22	Э23	Ю24	Я25
а0	б1	в2	г3	д4	е5	ж6	з7	и8	й9	ј10	љ11	њ12	ќ13	ў14	ф15	х16	ц17	ч18	ш19	щ20	ъ21	ь22	э23	ю24	я25
А0	В1	Г2	Д3	Е4	Ж5	З6	И7	Й8	Ј9	Љ10	Њ11	Ћ12	Ќ13	Ў14	Ф15	Х16	Ц17	Ч18	Ш19	Щ20	Ъ21	Ь22	Э23	Ю24	Я25
а0	б1	в2	г3	д4	е5	ж6	з7	и8	й9	ј10	љ11	њ12	ќ13	ў14	ф15	х16	ц17	ч18	ш19	щ20	ъ21	ь22	э23	ю24	я25
А0	В1	Г2	Д3	Е4	Ж5	З6	И7	Й8	Ј9	Љ10	Њ11	Ћ12	Ќ13	Ў14	Ф15	Х16	Ц17	Ч18	Ш19	Щ20	Ъ21	Ь22	Э23	Ю24	Я25
а0	б1	в2	г3	д4	е5	ж6	з7	и8	й9	ј10	љ11	њ12	ќ13	ў14	ф15	х16	ц17	ч18	ш19	щ20	ъ21	ь22	э23	ю24	я25
А0	В1	Г2	Д3	Е4	Ж5	З6	И7	Й8	Ј9	Љ10	Њ11	Ћ12	Ќ13	Ў14	Ф15	Х16	Ц17	Ч18	Ш19	Щ20	Ъ21	Ь22	Э23	Ю24	Я25
а0	б1	в2	г3	д4	е5	ж6	з7	и8	й9	ј10	љ11	њ12	ќ13	ў14	ф15	х16	ц17	ч18	ш19	щ20	ъ21	ь22	э23	ю24	я25
А0	В1	Г2	Д3	Е4	Ж5	З6	И7	Й8	Ј9	Љ10	Њ11	Ћ12	Ќ13	Ў14	Ф15	Х16	Ц17	Ч18	Ш19	Щ20	Ъ21	Ь22	Э23	Ю24	Я25
а0	б1	в2	г3	д4	е5	ж6	з7	и8	й9	ј10	љ11	њ12	ќ13	ў14	ф15	х16	ц17	ч18	ш19	щ20	ъ21	ь22	э23	ю24	я25

Reading material: The history of ASCII out of telegraph codes [1]. A history, paying attention to multilingual use [4]. Bob Bemer, the ‘father of ascii’ [2]. A detailed discussion of ISO 8859, Latin-1 [11].

4

Informally, the term ‘character set’ (also ‘character code’ or ‘code’) used to mean something like ‘a table of bytes, each with a character shape’. With only the English alphabet to deal with that is a good enough definition. These days, much more general cases are handled, mapping one octet into several characters, or several octets into one character. The definition has changed accordingly:

A *charset* is a method of converting a sequence of octets into a sequence of characters. This conversion may also optionally produce additional control information such as directionality indicators.

(From RFC 2978) A conversion the other way may not exist, since different octet combinations may map to the same character. Another complicating factor is the possibility of switching between character sets; for instance, ISO 2022-JP is the standard ASCII character set, but the escape sequence ESC \$ @ switches to JIS X 0208-1978.

To disentangle the concepts behind encoding, we need to introduce a couple of levels:

ACR Abstract Character Repertoire: the set of characters to be encoded; for example, some alphabet or symbol set. This is an unordered set of characters, which can be fixed (the contents of ISO 8859-1), or open (the contents of Unicode).

CCS Coded Character Set: a mapping from an abstract character repertoire to a set of nonnegative integers. This is what is meant by ‘encoding’, ‘character set definition’, or ‘code page’; the integer assigned to a character is its ‘code point’.

There used to be a drive towards unambiguous abstract character names across repertoires and encodings, but Unicode ended this, as it provides (or aims to provide) more or less a complete list of every character on earth.

CEF Character Encoding Form: a mapping from a set of nonnegative integers that are elements of a CCS to a set of sequences of particular code units. A ‘code unit’ is an integer of a specific binary width, for instance 8 or 16 bits. A CEF then maps the code points of a coded character set into sequences of code point, and these sequences can be of different lengths inside one code page. For instance ASCII uses a single 7-bit unit; UTF-8 uses one to four 8-bit units. We will discuss the UTF encodings below.

CES Character Encoding Scheme: a reversible transformation from a set of sequences of code units (from one or more CEFs to a serialized sequence of bytes. In single-byte cases such as ASCII and UTF-8 this mapping is trivial. With the two-byte scheme UCS-2 there is a single ‘byte order mark’, after which the code units are trivially mapped to bytes. On the other hand, ISO 2022, which uses escape sequences to switch between different encodings, is a complicated CES.

Additionally, there are the concepts of

CM Character Map: a mapping from sequences of members of an abstract character repertoire to serialized sequences of bytes bridging all four levels in a single operation. These maps are what gets assigned MIBenum values by IANA; see section 4.1.

TES Transfer Encoding Syntax: a reversible transform of encoded data. This data may or may not contain textual data. Examples of a TES are base64, uuencode, and quoted-printable, which all transform a byte stream to avoid certain values.

3 Unicode and UTF encodings

The systems above functioned quite well as long as you stuck to one language or writing system. Poor dictionary makers. More or less simultaneously two efforts started that aimed to incorporate all the world's character sets in one standard: Unicode standard (originally 2-byte), and ISO 10646 (originally 4-byte). Unicode was extended further, so that it has all code points up to 10FFFFFF, which is slightly over a million.

Two international standards organizations, the Unicode Consortium and ISO/IEC JTC1/SC2, started designing a universal standard that was to be a superset of all existing character sets. These standards are now synchronized. Unicode has elements that are not in 10646, but they are compatible where it concerns straight character encoding.

ISO 10646 defines UCS, the 'Universal Character Set'. This is in essence a table of official names and code numbers for characters. Unicode adds to this rules for hyphenation, bi-directional writing, and more.

The full Unicode list of code points can be found online, broken down by blocks [14], and downloadable [17].

3.1 BMP and earlier Unicode subplanes

Characters in Unicode are mostly denoted hexadecimally as U+wx yz, for instance U+0041 is 'Latin Capital Letter A'. The range U+0000–U+007F (0–127) is identical to US-ASCII (ISO 646 IRV), and U+0000–U+00FF (0–255) is identical to Latin 1 (ISO 8859-1).

The original 2-byte subset is now called 'BMP' for Basic Multilingual Plane, or plane 0. These are the Unicode code points that are nonzero in the last two bytes. Other 'planes' have been defined that have one or more bits set outside the last two bytes.

BMP (Basic Multilingual Plane) The first plane defined in Unicode/ISO 10646, designed to include all scripts in active modern use. The BMP currently includes the Latin, Greek, Cyrillic, Devangari, hiragana, katakana, and Cherokee scripts, among others, and a large body of mathematical, APL-related, and other miscellaneous characters. Most of the Han ideographs in current use are present in the BMP, but due to the large number of ideographs, many were placed in the Supplementary Ideographic Plane.

SMP (Supplementary Multilingual Plane; plane 1) This contains mostly ancient writing systems. Some of these you'll have heard of, such as Linear B, cuneiform, Aztec, and Maya; others are fairly obscure, such as Tangut, a language used in Central China between 1000 and 1500.

SIP (Supplementary Ideographic Plane) The third plane (plane 2) defined in Unicode/ISO 10646, designed to hold all the ideographs descended from Chinese writing (mainly found in Vietnamese, Korean, Japanese and Chinese) that aren't found in the Basic Multilingual Plane. The BMP was supposed to hold all ideographs in modern use; unfortunately, many Chinese dialects (like Cantonese and Hong Kong Chinese) were overlooked; to write these, characters from the SIP are necessary. This is one reason even non-academic software must support characters outside the BMP.

3.2 Unicode encodings

Unicode is basically a numbered list of characters. When they are used in a file, their numbers can be encoded in a number of ways. To name the obvious example: if only the first 128 positions are used, the long Unicode code point can be truncated to just one byte. Here are a few encodings:

UTF-32 Little used: this is a four-byte encoding. (UTF stands for ‘UCS Transformation Format’.)

UTF-16 A two-byte encoding. Its precursor, UCS-2, encoded the BMP; UTF-16 has a way of going beyond that to encode the whole of Unicode.

UTF-8 A one-byte scheme; details below.

UTF-7 Another one-byte scheme, but now the high bit is always off. Certain byte values act as ‘escape’, so that higher values can be encoded. Like UTF-1 and SCSU, this encoding is only of historical interest.

There is an important practical reason for a one-byte encoding such as UTF-8. Multi-byte encodings such as UCS-2 are wasteful of space, if only traditional ASCII is needed. Furthermore, they would break software that is expecting to walk through a file with `s++` and such. Also, they would introduce many zero bytes in a file, which would play havoc with Unix software that uses null-termination for strings.

Then there would be the problem of whether two bytes are stored in low-endian or high-endian order. For this reason it was suggested to store `FE FF` or `FF FE` at the beginning of each file as the ‘Unicode Byte Order Mark’. Formally, `FEFF` is the Unicode ‘zero width nobreak space’ character, which can innocently be inserted anywhere. Conversely `FFEF` is defined to be illegal, so encountering this is a sign that bytes should be interpreted little-endian. Of course this plays havoc with files such as shell scripts which expect to find `# !` at the beginning of the file.

3.3 UTF-8

UTF-8, standardized as RFC 3629, is an encoding where the positions up to 127 are encoded ‘as such’; higher numbers are encoded in groups of 2 to 6 bytes. (Tim Bray describes this as ‘kind of racist’ [3]: the further east a language comes from, the more overhead is involved in its encoding.) In a multi-byte group, the first byte is in the range `0xC0–0xFD` (192–252). The next up to 5 bytes are in the range `0x80–0xBF` (128–191, bit pattern starting with 10). Note that `8 = 1000` and `B = 1011`, so the highest two bits are always 10, leaving six bits for encoding).

U-00000000 – U-0000007F	7 bits	0xxxxxxx	
U-00000080 – U-000007FF	11 = 5 + 6	110xxxxx	10xxxxxx
U-00000800 – U-0000FFFF	16 = 4 + 2 × 6	1110xxxx	10xxxxxx 10xxxxxx
U-00010000 – U-001FFFFF	21 = 3 + 3 × 6	11110xxx	10xxxxxx (3 times)
U-00200000 – U-03FFFFFF	26 = 2 + 4 × 6	111110xx	10xxxxxx (4 times)
U-04000000 – U-7FFFFFFF	31 = 1 + 5 × 6	1111110x	10xxxxxx (5 times)

All bites in a multi-byte sequence have their high bit set.

IETF documents such as RFC 2277 require support for this encoding in internet software. Readable introductions can be found all over the internet [19]; see also the history of UTF-8 [20].

3.4 Unicode tidbits

3.4.1 Line breaking

The Unicode standard describes line breaking: it has a mechanism for specifying tables of character pairs between which line breaks are allowed or forbidden [15, 18].

3.4.2 Bi-directional writing

Most scripts are left-to-right, but Arabic and Hebrew run right-to-left. Characters in a file are stored in ‘logical order’, and usually it is clear in which direction to render them, even if they are used mixed. Letters have a ‘strong’ directionality: unless overridden, they will be displayed in their natural direction. The first letter of a paragraph with strong direction determines the main direction of that paragraph [16].

أوروبا، برمجيات الحاسوب + انترنت :

تصبح عالميا مع يونيكود

تسجل الآن لحضور المؤتمر الدولي العاشر ليونيكود، الذي سيعقد في 10-12 آذار 1997 بمدينة ماينتس ألمانيا. وسيجمع المؤتمر بين خبراء من كافة قطاعات الصناعة على الشبكة العالمية انترنت ويونيكود. حيث سيتم على الصعيدين الدولي والمحلي على حد سواء مناقشة سبل استخدام يونيكود في النظم القائمة وفيما يخص التطبيقات الحاسوبية. الخطوط تصميم النصوص والحوسبة متعددة اللغات.

عندما يريد العالم أن يتكلم فهو يتحدث بلغة يونيكود.

However, when differently directional texts are embedded, some explicit help is needed. The problem arises with letters that have only weak directionality. The following is a sketch of a problematic case.

Memory: he said "I NEED WATER!", and expired.

Display: he said "RETAW DEEN I!", and expired.

If the exclamation mark is to be part of the Arabic quotation, then the user can select the text 'I NEED WATER!' and explicitly mark it as embedded Arabic (<RLE> is Right-Left Embedding; <PDF> Pop Directional Format), which produces the following result:

Memory: he said "<RLE>I NEED WATER!<PDF>", and expired.

Display: he said "!RETAW DEEN I", and expired.

A simpler method of doing this is to place a Right Directional Mark <RLM> after the exclamation mark. Since the exclamation mark is now not on a directional boundary, this produces the correct result.

Memory: he said "I NEED WATER!<RLM>", and expired.

Display: he said "!RETAW DEEN I", and expired.

3.5 Unicode and oriental languages

‘Han unification’ is the Unicode strategy of saving space in the oriental languages (traditional Chinese, simplified Chinese, Japanese, Korean: ‘CJK’) by recognizing common characters. This idea is not uncontroversial [6].

4 Further tidbits

4.1 A bootstrapping problem

In order to know how to interpret a file, you need to know what character set it uses. This problem also occurs in MIME mail encoding (section 4.5), which can use many character sets. Names and numbers for character sets are standardized by IANA: the Internet Assigned Names Authority [8]. However, in what character set do you write this name down?

Fortunately, everyone agrees on (7-bit) ASCII, so that is what is used. A name can be up to 40 characters from us-ascii.

As an example, here is the iana definition of ASCII:

```
name ANSI_X3.4-1968
reference RFC1345,KXS2
MIBenum 3
source ECMA registry
aliases iso-ir-6,ANSI_X3.4-1986,ISO_646.irv:1991,ASCII,ISO646-US,
        US-ASCII (preferred MIME name),us,IBM367,cp367,csASCII
```

The MIBenum (Management Information Base) is a number assigned by IANA¹. The full list of character sets can be found online [9], and RFC 3808 is a memo that describes the IANA Charset MIB.

4.2 Unicode in programming languages

Before Unicode, a system called the ‘Double Byte Character Set’ was invented to accommodate Asian languages, where some characters were stored in one, others in two bytes. This is very messy, since you can not simply write `s++` or `s--` to traverse a string. Instead you have to use functions from some library that understands these encodings. While this system is now only of historical interest, the string handling problem is back in force with UTF-8.

Many modern languages (Python, C99) have support for Unicode. In C99 (which is the new standard for C) this is done through so-called ‘wide characters’. For instance, `L'x'` is a wide character and `L"xyz"` is a string of wide characters. Such strings can be manipulated through equivalents of the normal string library. For instance, `wcscpy` acts like `strcpy` but on wide strings. General Unicode characters can be represented as `\u0000` for 4-byte and `\U00000000` for up to 8-byte characters.

The two-byte UTF-16 encoding is popular among programmers, since it can handle almost any practically encountered character without extensions to longer byte sequences.

4.3 Character codes in HTML

HTML can access unusual characters in several ways:

- With a decimal numerical code: ` ` is a space token. (HTML 4 supports hexadecimal codes.)

1. Apparently these numbers derive from the Printer MIB, RFC 1759.

- With a vaguely symbolic name [12, 7]: `©` is the copyright symbol.
- The more interesting way is to use an encoding such as UTF-8 (section 3.2) for the file. For this it would be nice if the server could state that the file is
`Content-type: text/html; charset=utf-8`
but it is also all right if the file starts with
`<META HTTP-EQUIV="Content-Type" CONTENT="text/html; charset=utf-8">`

Description	Char Code	Entity name
non-breaking space	<code>&#160;</code> -->	<code>&nbsp;</code> -->
inverted exclamation	<code>¡</code> <code>&#161;</code> --> <code>¡</code>	<code>&iexcl;</code> --> <code>¡</code>
cent sign	<code>¢</code> <code>&#162;</code> --> <code>¢</code>	<code>&cent;</code> --> <code>¢</code>
pound sterling	<code>£</code> <code>&#163;</code> --> <code>£</code>	<code>&pound;</code> --> <code>£</code>
general currency sign	<code>¤</code> <code>&#164;</code> --> <code>¤</code>	<code>&curren;</code> --> <code>¤</code>
yen sign	<code>¥</code> <code>&#165;</code> --> <code>¥</code>	<code>&yen;</code> --> <code>¥</code>
broken vertical bar	<code> </code> <code>&#166;</code> --> <code> </code>	<code>&brvbar;</code> --> <code> </code>

It is requirement that user agents can at least parse the `charset` parameter, which means they have to understand us-ascii.

Open this link in your browser, and additionally view the source: <http://www.unicode.org/unicode/iuc10/x-utf8.html>. How well does your software deal with it?

4.4 Keyboards and control characters

Unprintable ascii codes are accessible through the control modifier key; for this reason they are also called ‘control codes’ or control characters. The control key, combined with a regular key, zeros bits 2 and 3 of the ascii code of that key. For instance, you can hit `Ctrl-[` to get `Esc`.

The way key presses generate characters is typically controlled in software. This mapping from keyboard scan codes to 7 or 8-bit characters is called a ‘keyboard’, and can be changed dynamically in most operating systems.

Using the modifier keys, one can generate more keystrokes than can be described in 8 bits, so keyboards can send an ‘escape sequence’: one escape character followed by one or more regular characters. The escape character is mostly ascii NULL or ESC [10].

4.5 Characters in email

The protocols for internet mail are based on ‘7-bit ASCII’, that is, the high bit in every byte transmitted is supposed to be off. This is a problem for any message that has text outside of ASCII, such as when accented characters from the various ISO 8859 character sets are used. It also makes transmitting binary data such as images impossible. For this reason the ‘Multipurpose Internet Mail Extensions’ (MIME) were designed. MIME uses several encoding schemes, such as base64 or printed quotable, to turn arbitrary data into 7-bit ascii.

The email standard RFC 822 states that anything outside 7-bit ascii has to be encoded with `uuencode`. This means that the sender and recipient need decoding program; it is decidedly overkill if a message is plain ASCII apart from a few accented characters.

The MIME protocol (RFC 2045 and 2046) inserts headers in a mail message, stating for each message section the content type and the encoding that is used for that section of the message. These encodings are also used for attachments, in which case the content type should give an indication what application can handle the attachment after its decoding. ‘Helpful’ mail programs that automatically invoke such applications have been a source of Trojans (malicious softwares) in the past.

4.6 FTP

FTP is a very old ARPA protocol for transferring files from one computer to another. It knows ‘binary’ and ‘text’ mode: in binary mode bytes are transferred without further interpretation, but the text mode is concerned with files that contain lines of text. Unfortunately, line ends are different between operating systems, and their transfer in text mode is not well defined. Some ftp programs adjust line ends; others, such as `Fetch` on the Mac, actually do code page translation.

5 Character issues in T_EX / L^AT_EX

5.1 Diacritics

Before 1990, T_EX was a 7-bit system: only characters 0–127 in the input could be recognized, and fonts were also limited to 127 positions. This meant that there was not enough space in fonts for letters with accents, so accents (diacritics) were implemented as things to put on top of characters, even when, as with the cedilla, they are under the letter. This leads to the problem that T_EX can not hyphenate a word with accents, since the accent introduces a space in the word (technically: an explicit kern).

Both problems were remedied to a large extent with the ‘Cork font encoding’, which contains most accented letters as single characters. This means that accents are correctly placed by design, and also that the word can be hyphenated, since the kern has disappeared.

These fonts with accented characters became possible when T_EX version 3 came out around 1990. This introduced full 8-bit compatibility, both in the input side and in the font addressing.

5.2 L^AT_EX input file access to fonts

If an input file for L^AT_EX is allowed to contain all 8-bit octets, we get all the problems of compatibility that plagued regular text files. This is solved by the package `inputenc`:

```
\usepackage[code]{inputenc}
```

where `codes` is `applemac`, `ansinew`, or various other code pages.

This package makes all unprintable ASCII characters, plus the codes over 127, into active characters. The definitions are then dynamically set depending on the code page that is loaded.

5.3 L^AT_EX output encoding

The `inputenc` package does not solve the whole problem of producing a certain font character from certain keyboard input. It only mapped a byte value to the T_EX command for producing a character. To map such commands to actual code point in a font file, the T_EX and L^AT_EX formats contain lines such as

```
\chardef\i="10
```

declaring that the dotless-i is at position 16. However, this position is a convention, and other people – type manufacturers – may put it somewhere else.

This is handled by the ‘font encoding’ mechanism. The various people working on the L^AT_EX font schemes have devised a number of standard font encodings. For instance, the OT1 encoding corresponds to the original 128-character set. The T1 encoding is a 256-character extension thereof, which includes most accented characters for Latin alphabet languages.

A font encoding is selected with

```
\usepackage[T1]{fontenc}
```

A font encoding definition contains lines such as

```
\DeclareTextSymbol{\AE}{OT1}{29}  
\DeclareTextSymbol{\OE}{OT1}{30}  
\DeclareTextSymbol{\O}{OT1}{31}  
\DeclareTextSymbol{\ae}{OT1}{26}  
\DeclareTextSymbol{\i}{OT1}{16}
```

5.4 T_EX beyond 8 bits

The above L^AT_EX packages allow flexible handling of (8-bit) codepages, essentially the ISO 8859 standard. For handling of other alphabets, a number of styles have been written over the years. However, their continued support is often uncertain. The only project that aims at use of Unicode throughout T_EX’s codebase is Omega [13].

Contents

1	History in one byte	1	4.2	<i>Unicode in programming languages</i>	9
2	Character sets and encodings	4	4.3	<i>Character codes in HTML</i>	9
3	Unicode and UTF encodings	6	4.4	<i>Keyboards and control characters</i>	10
3.1	<i>BMP and earlier Unicode subplanes</i>	6	4.5	<i>Characters in email</i>	10
3.2	<i>Unicode encodings</i>	7	4.6	<i>FTP</i>	11
3.3	<i>UTF-8</i>	7	5	Character issues in \TeX / \LaTeX	11
3.4	<i>Unicode tidbits</i>	8	5.1	<i>Diacritics</i>	11
3.5	<i>Unicode and oriental languages</i>	8	5.2	<i>\LaTeX input file access to fonts</i>	11
4	Further tidbits	9	5.3	<i>\LaTeX output encoding</i>	12
4.1	<i>A bootstrapping problem</i>	9	5.4	<i>\TeX beyond 8 bits</i>	12

References

- [1] Annotated history of ascii. <http://www.wps.com/projects/codes/index.html>.
- [2] Bob berner homepage. <http://www.trailing-edge.com/~bobbemer/>.
- [3] Tim Bray. Characters vs bytes. [urlhttp://www.tbray.org/ongoing/When/200x/2003/04/26/UTF](http://www.tbray.org/ongoing/When/200x/2003/04/26/UTF).
- [4] Brief history of character codes in north american, europe, and east asia. <http://tronweb.super-nova.co.jp/characodehist.html>.
- [5] Codepage & co. <http://aspell.net/charsets/codepages.html>.
- [6] Han unification. http://en.wikipedia.org/wiki/Han_unification.
- [7] Character entity references in HTML4. <http://www.w3.org/TR/html401/sgml/entities.html>.
- [8] <http://www.iana.org/>.
- [9] IANA character set names. <http://www.iana.org/assignments/character-sets>.
- [10] IBM PC keyboard scan codes. <http://jimprice.com/jim-asc.shtml#keycodes>.
- [11] The ISO latin 1 character repertoire. <http://www.cs.tut.fi/~jkorpela/latin1/index.html>.
- [12] Character entities for ISO latin 1. <http://www.cs.tut.fi/~jkorpela/HTML3.2/latin1.html>.
- [13] Omega project home page. <http://omega.enstb.org/>.
- [14] Unicode. <http://www.fileformat.info/info/unicode/index.htm>.
- [15] Unicode standard annex 14, line breaking properties. <http://www.unicode.org/reports/tr14/>.
- [16] Unicode standard annex 9, the bidirectional algorithm. <http://www.unicode.org/reports/tr9/>.
- [17] Unicode code chart and scripts. <http://www.unicode.org/charts/>.

- [18] Unicode line breaking rules: explanations and criticism. <http://www.cs.tut.fi/~jkorpela/unicode/linebr.html>.
- [19] UTF-8 and Unicode FAQ for Unix/Linux. <http://www.cl.cam.ac.uk/~mgk25/unicode.html>.
- [20] UTF-8 history. <http://www.cl.cam.ac.uk/~mgk25/ucs/utf-8-history.txt>.