

# Character encodings

Low-level view: Byte streams

1 Byte = 8 bit, enough to encode  $2^8 = 256$  characters

Standards: • ASCII, ~ 1967, 7-bit

• various 8-bit extensions, most important ISO 8859

16 versions 1987-2001, e.g. latin-1 "Western European"

Issues: • Barely works for single latin-like alphabets, difficult to mix languages

• 8 bit not enough for many languages (Chinese, arabic, ...)

Data model:

(1) Byte array ("C-style strings")

• A string is an array of bytes, direct indexing

• Convention: Termination with 0-character (0x00)

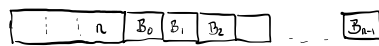
• Simple things are simple and fast, low memory overhead

• Everything else hard to get 100% right

• Memory safety hard to get right → security implementation

• Still used in system-level programming, small embedded systems

(2) Length-prefixed ("Pascal-style strings") (Python: bytes, bytearray)



length prefix

• Makes a lot of string algorithms easier, faster, and more robust

• No special handling of 0x00

• Prefix-size is a compromise between maximal string length and memory requirements for small strings (rarely an issue nowadays)

• Strings are different from standard arrays

## High-level view:

Character streams

"Character" = Unicode code point

### Standard:

- Unicode 1.0, 1981

16-bit

Targets "modern" languages,

$2^{14} = 16384$  "should be enough for all modern usage" (Becker '88)

rest are "private code points"

- Unicode 2.0, 1996 beyond 16 bit

17 planes, first plane (BMP) is, more or less, old unicode

Codepoints  $U+0000 \dots U+10FFFF$  ( $\approx 1$  million)

- Now at version 15.0, keeps changing/adding characters

### Encodings:

UTF-32 : 32-bit per character

often internal format, direct indexing possible!

UTF-16 : BMP via 16 bits, other planes as "surrogate characters"

UTF-8 : ASCII as 8-bit, BMP 1-3 byte surrogates, other planes 4 bytes

empirically most efficient (even for Chinese web pages!)

ASCII is valid UTF-8! ✓

### Usage:

- Unix operating systems, WWW : UTF-8

- Windows : UTF-16

- In high-level programming languages (Python, ...), encoding details are hidden from programmer

	0	1	2	3	4	5	6	7	8	9	A	B	C	D	E	F
0	0 ^@ NUL NULL	1 ^A SOH START OF HEADING	2 ^B STX START OF TEXT	3 ^C ETX END OF TEXT	4 ^D EOT END OF TRANSM.	5 ^E ENQ ENQUIRY	6 ^F ACK ACKNOWL- EDGE	7 ^G BEL BELL	8 ^H BS BACKSP.	9 ^I HT CHARACT. TAB'TION	10 ^J LF LINE FEED	11 ^K VT LINE TAB'TION	12 ^L FF FORM FEED	13 ^M CR CARRIAGE RETURN	14 ^N SO SHIFT OUT	15 ^O SI SHIFT IN
1	16 ^P DLE DATALINK ESCAPE	17 ^Q DC1 DEVICE CONTROL1	18 ^R DC2 DEVICE CONTROL2	19 ^S DC3 DEVICE CONTROL3	20 ^T DC4 DEVICE CONTROL4	21 ^U NAK NEG. ACK- NOWLEDGE	22 ^V SYN SYNCHR. IDLE	23 ^W ETB END OF TRANS.	24 ^X CAN CANCEL	25 ^Y EM END OF MEDIUM	26 ^Z SUB SUBS- TITUTE	27 ^[ ESC ESCAPE	28 ^\ FS INFO. SEP. 4	29 ^] GS INFO. SEP. 3	30 ^^ RS INFO. SEP. 2	31 ^_ US INFO. SEP. 1
2	#32 SPACE	&#33; ! EXCLAM. MARK	&#34; " QUOT. MARK	&#35; # NUMBER SIGN	&#36; \$ DOLLAR SIGN	&#37; % PERCENT SIGN	&#38; & AMPERS- SAND	&#39; ' APOS- TROPHE	&#40; ( LEFT PAREN.	&#41; ) RIGHT PAREN.	&#42; * ASTERISK	&#43; + PLUS SIGN	&#44; , COMMA	&#45; - HYPHEN- MINUS	&#46; . FULL STOP	&#47; / SOLIDUS
3	#48 DIGIT ZERO	#49 1 DIGIT ONE	#50 2 DIGIT TWO	#51 3 DIGIT THREE	#52 4 DIGIT FOUR	#53 5 DIGIT FIVE	#54 6 DIGIT SIX	#55 7 DIGIT SEVEN	#56 8 DIGIT EIGHT	#57 9 DIGIT NINE	#58 : COLON	#59 ; SEMI- COLON	#60 < LS.-THAN SIGN	#61 = EQUALS SIGN	#62 > GR.-THAN SIGN	#63 ? QUEST- ION MARK
4	#64 @ COMM'IAL AT	#65 A	#66 B	#67 C	#68 D	#69 E	#70 F	#71 G	#72 H	#73 I	#74 J	#75 K	#76 L	#77 M	#78 N	#79 O
5	#80 P	#81 Q	#82 R	#83 S	#84 T	#85 U	#86 V	#87 W	#88 X	#89 Y	#90 Z	#91 [ LEFT SQ. BRACKET	#92 \ REVERSE SOLIDUS	#93 ] RT. SQ. BRACKET	#94 ^ CIRCUM'X ACCENT	#95 _ LOW LINE
6	#96 ` GRAVE ACCENT	#97 a	#98 b	#99 c	#100 d	#101 e	#102 f	#103 g	#104 h	#105 i	#106 j	#107 k	#108 l	#109 m	#110 n	#111 o
7	#112 p	#113 q	#114 r	#115 s	#116 t	#117 u	#118 v	#119 w	#120 x	#121 y	#122 z	#123 { L. CURLY BRACKET	#124   VERTICAL LINE	#125 } R. CURLY BRACKET	#126 ~ TILDE	127 ^? DEL DELETE