

# Text File Coding

Leonardo Araújo

UFSJ

# Introduction to Text Encoding

## ■ Why Encoding Matters:

- Human-Readable Format: Text encoding allows computers to store and exchange text in a format that is both human-readable and machine-processable.
- Standardization: Encodings like UTF-8 provide a standard way for computers and systems to interpret text, avoiding confusion that could arise from regional or proprietary encodings.
- Ensures data integrity across different systems.
  - Proper encoding preserves the accuracy of text when shared across different platforms, preventing data corruption and ensuring consistent display of characters.

## ■ Overview:

- Evolution from simple to complex encoding systems.

## Character vs. Glyph vs. Font

- **Character:** Abstract unit in encoding (e.g., 'A' in Unicode).
- **Glyph:** Visual form of a character (how 'A' looks in Arial vs. Times).
- **Font:** Collection of glyphs sharing a design style.



Figure 1: Glyph vs font.

# Jacquard

- Joseph Marie Jacquard in Lyon in 1801.

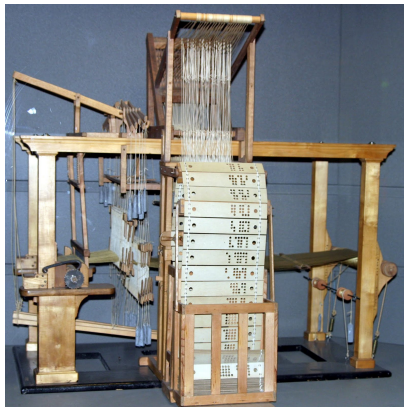


Figure 2: Jacquard's loom.

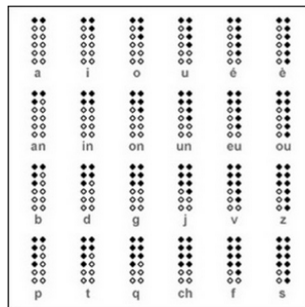
# Night Writing to Braille: Evolution in Tactile Encoding

## ■ Night Writing:

- Inventor: Charles Barbier, 1815, for silent military communication.
- Structure: 12-dot cells for phonetic sounds, complex for practical use.

## ■ Braille:

- Creator: Louis Braille, 1824 (first published 1829), adapted from Night Writing.
- Innovation: 6-dot cell, simpler, and more accessible for the blind.
- First Binary Encoding: Each cell represents binary combinations.
- Universal Adoption: Braille became the standard for blind communication worldwide.
- Expansion: Adapted for various languages, math, music, and more.



	1	2	3	4	5	6
1	a	i	o	u	é	è
2	an	in	on	un	eu	ou
3	b	d	g	j	v	z
4	p	t	q	ch	f	s
5	l	m	n	r	gn	ll
6	oi	oin	ian	ien	ion	ieu

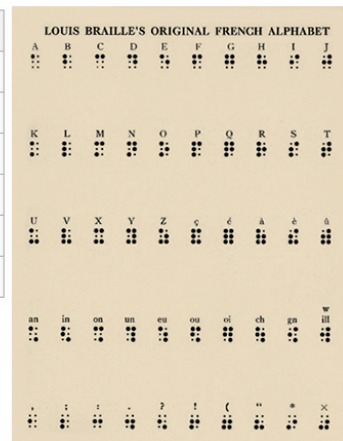


Figure 3: Night writing and Braille.

## Morse Code

- **History:** Developed by Samuel Morse and Alfred Vail in the early 1840s.
- **Mechanism:** Uses dots, dashes, and spaces for letters, numbers, and punctuation.
- **Usage:** Primarily telegraphy, but also in radio communication.

A ·-·	N -·	1 ·-·-·-·
B -·-·	O -·-·	2 ··-·-·
C -·-·-	P ·-·-·	3 ··-·-·
D -··	Q -·-·-	4 ·-·-·
E ·	R ·-·	5 ·-·-·
F ·-·-·	S ···	6 -·-·-·
G -·-·	T -	7 -·-·-·
H ·-·-·	U ·-·	8 -·-·-·
I ··	V ·-·-·	9 -·-·-·
J ·-·-·-	W ·-·-	0 -·-·-·-
K -·-·	X -·-·-	· ·-·-·-·
L ·-·-·	Y -·-·-·	, -·-·-·-
M -·-	Z -·-·	? ·-·-·-

Figure 4: Morse Code

# Baudot and Murray Code

## ■ Baudot Code:

- Invented by Émile Baudot, 5-bit code for telegraphy.
- Limited characters, used shift for numbers/letters.

## ■ Murray Code (ITA2):

- Extension of Baudot, improved by Donald Murray.
- Added lower case, more symbols.

LETTERS FIGURES		A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	CARRIAGE RETURN	LINE FEED	LETTERS	FIGURES	SPACE	ALL-SPACE NOT IN USE
CODE ELEMENTS	1	●	●		●	●	●				●	●						●		●		●		●	●	●	●			●	●		
	2	●		●				●		●	●	●	●				●	●	●			●	●	●	●				●	●	●	●	
	3	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○
	4		●	●			●	●	●		●	●		●	●	●	●	●	●			●	●	●	●	●	●	●	●	●	●	●	●
	5		●		●		●	●	●			●	●	●	●	●	●	●			●		●	●	●	●	●	●	●		●	●	○

● INDICATES A MARK ELEMENT (A HOLE PUNCHED IN THE TAPE)

○ INDICATES POSITION OF A SPROCKET HOLE IN THE TAPE

**The International Telegraph Alphabet**

Figure 5: ITA2 Baudot-Murray code



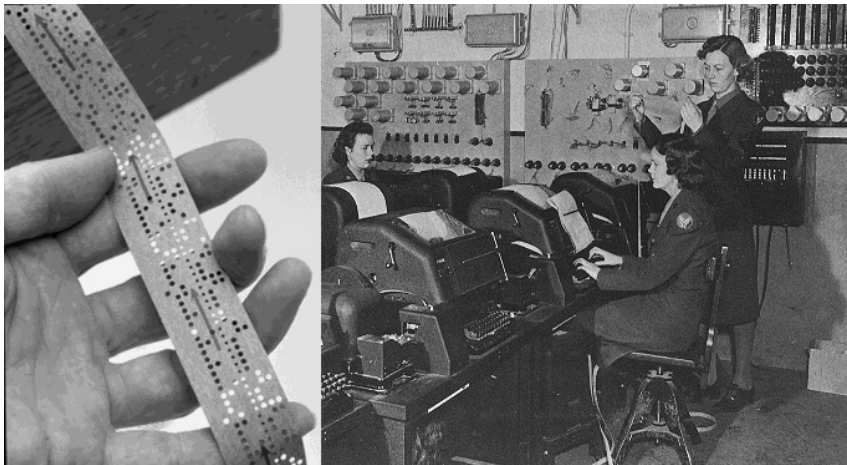
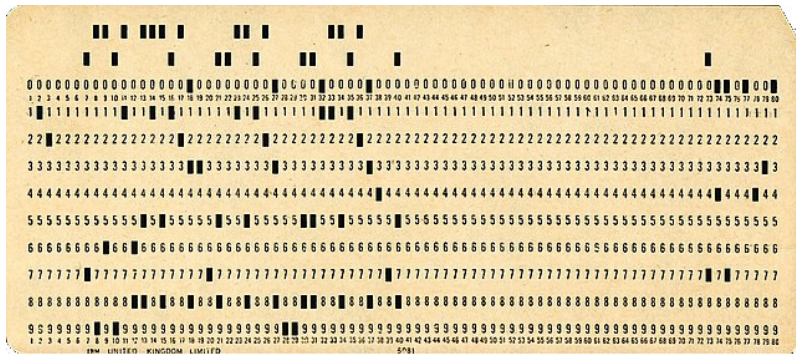


Figure 6: Teletype and perforated strip.

## EBCDIC (Extended Binary Coded Decimal Interchange Code)

- **History:** Developed by IBM for mainframe computers.
- **Characteristics:** the first character encodings created for data processing on large-scale systems.
  - Used in legacy systems (IBM 1401, 7090, System/360).
- EBCD, a subset of EBCDIC.



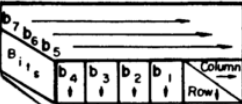
	-	1	2	3	4	5	6	7	8	9	2-8	3-8	4-8	5-8	6-8	7-8
-		1	2	3	4	5	6	7	8	9	:	#	@	'	=	"
Y	&	A	B	C	D	E	F	G	H	I	[	.	<	(	+	!
X	-	J	K	L	M	N	O	P	Q	R	]	\$	*	)	;	^
0	0	/	S	T	U	V	W	X	Y	Z	\	,	%	_	>	?

Figure 7: 12-row/80-column IBM punched card and EBCDIC table.

## ASCII and Extended ASCII

- **ASCII (American Standard Code for Information Interchange):**
  - 7-bit code, 128 characters including 33 non-printing control codes.
  - Standardized in 1963 (ANSI).
  - Backward Compatibility: Despite its age, ASCII remains widely used today for compatibility reasons.
- **Extended ASCII:**
  - 8-bit code, 256 characters, allowing for additional symbols and characters.

USASCII code chart



Column Row	0	1	2	3	4	5	6	7
0	NUL	DLE	SP	@	P	\	p	
1	SOH	DC1	!	A	Q	a	q	
2	STX	DC2	"	B	R	b	r	
3	ETX	DC3	#	C	S	c	s	
4	EOT	DC4	\$	D	T	d	t	
5	ENQ	NAK	%	E	U	e	u	
6	ACK	SYN	&	F	V	f	v	
7	BEL	ETB	'	G	W	g	w	
8	BS	CAN	(	H	X	h	x	
9	HT	EM	)	I	Y	i	y	
10	LF	SUB	*	J	Z	j	z	
11	VT	ESC	+	K	[	k	{	
12	FF	FS	,	L	\	l		
13	CR	GS	-	M	]	m	}	
14	SO	RS	.	N	^	n	~	
15	SI	US	/	O	_	o	DEL	

Figure 8: ASCII Table

Character	Binary (Uppercase)	Binary (Lowercase)	Character
A	01000001	01100001	a
B	01000010	01100010	b
C	01000011	01100011	c
...	...	...	...
Z	01011010	01111010	z
1	00110001	00100001	!
2	00110010	01000000	"
3	00110011	00100011	#
4	00110100	00100100	\$
...	...	...	...

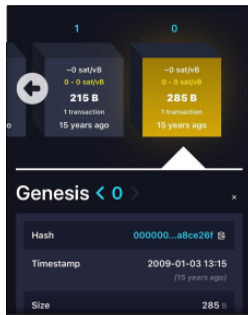
# ASCII Art

```

| .-----|
||                                     || | | | | |
||           -----                 ||
||           . ; ; ; ; ; ; ; .      ||
||           / ; ; ; ; ; ; ; ; \    ||
||           / ; / ^      ^ - ; ; ; ; . . ||
||           | ; | _ _      \ ; ; ; |   ||
|| | . - . | ; | e ^ / e ^   | ; ; ; |   ||
||           | ; |   |       | ; ; ; | ' -- ||
||           | ; |   ' _      | ; ; ; |   ||
||           | ; ; \  -- '    / | ; ; ; |   ||
||           | ; ; ; ; ; --- ' \ | ; ; ; |   ||
||           | ; ; ; ; |       | ; ; ; |   ||
||           | ; ; . - '       | ; ; ; |   ||
|| | ' -- | / ^                | ; ; ; | -- . ||
|| | ; ; ; ;      .           ; ; ; ; \ ; ; ; |
|| | ; ; ; ; - . ; _      / . - ; ; ; ; ; |
|| | ; ; ; ; ; ; ; ; ; ; ; ; ; ; ; ; ; ; ; |
|-----|

```

# BTC Genesis Block



**Bitcoin Genesis Block**

**Raw Hex Version**

```

00000000 01 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
00000010 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
00000020 00 00 00 00 3B A3 ED FD 7A 7B 12 B2 7A C7 2C 3E .....;fig{.*&,>
00000030 67 76 8F 61 7F C8 1B C3 88 8A 51 32 3A 9F 88 AA .....gv.a.B.A"8Q2;V,*
00000040 4B 1E 5E 4A 29 AB 5F 49 FF FF 00 1D 1D AC 2B 7C .....K."j)*_Iyy...~+|
00000050 01 01 00 00 01 00 00 00 00 00 00 00 00 00 00 00 .....
00000060 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
00000070 00 00 00 00 00 00 FF FF FF FF 4D 04 FF FF 00 1D .....yyyym.yy..
00000080 01 04 45 54 68 65 20 54 69 6D 65 73 20 30 33 2F .....The Times 03/
00000090 4A 61 6E 2F 32 30 30 39 20 43 68 61 6E 63 65 6C .....lor on brink of
000000A0 6C 6F 72 20 6F 6E 20 62 72 69 68 6B 20 6F 66 20 .....second bailout f
000000B0 73 65 63 6F 6E 64 20 62 61 69 6C 6F 75 74 20 66 .....or banksffff..b.
000000C0 6F 72 20 62 61 6E 6B 73 FF FF FF FF 01 00 F2 05 .....*....CA.g5?psh"
000000D0 2A 01 00 00 00 43 41 04 67 8A FD 80 FE 55 48 27 ....._gn/q0-.\'0' (A9.;
000000E0 19 67 F1 A6 71 30 37 10 5C DE A8 28 80 39 09 A6 .....ybab.ab?10K7L18K
000000F0 79 62 E0 EA 1F 61 DE B6 49 F6 BC 3F 4C EF 38 C4 .....dU.A.A.b\BN+*.W
00000100 F3 55 04 K5 1E C1 12 DE 5C 38 4D P7 BA 0B 8D 57 .....8!p+kh._'....
00000110 8A 4C 70 2B 6B F1 1D 5F AC 00 00 00 00 .....
  
```



Figure 9: The message embedded by Satoshi Nakamoto in Bitcoin's first block (the Genesis Block). The message reads, "The Times 03/Jan/2009 Chancellor on brink of second bailout for banks," which was a headline from The Times newspaper on that date.



# Base64

Base64 is a binary-to-text encoding scheme that represents binary data in an ASCII string format. Each Base64 digit represents exactly 6 bits of data, providing a way to encode binary data as text.

Base64 is used in:

- email attachments,
- embedding binary data in XML, JSON, or HTML, and
- data exchange in APIs.

```
$ base64 /tmp/tux.png
iVBORw0KGgoAAAANSUHEUgAAADIAAAA7CAYAAAA5MNl5AAAAxHpUWHRSYXcgCHJvZmlsZSB0eXB1
IGV4aWYAAHjabVBbDsMgDPvnFDsCiVMix6GPSbvBjr8AaVW2WcINceSahOP9eoZHA5MEWbKmkLI0
SJHC1QqNA7UzRencAXaNs7A6gJbC21yXDX5/Nmny2B8qLXLzUg3F9ZZKOL++mXkidAstXp3o7Jd
kbtAbLDHs2Iqmu9PWI84Q8cJjUtN2D/3bNvbF/sPmA8QojGgIWdakyBqAhszkg0SxGpB7p3TzBby
b08nwgcmkllHdJ9h5QAAAYRpQ0NQSUNDIHByb2ZpbGUAAHicfZE9SMNAHMFU0WRiogdpDhkqLrY
RUUcaxWkUCHUCq06mFz6BU0akhQXR8G140DHYtXBxVlXB1dBEPwAcXZwUnsREv+XFFrEeHDcj3f3
HnfvAKFRYZrVFQc03TbTyYSYZa2KPa8IIYJBjCMoM8uYk6QUfMfXPQJ8vYvxLP9zf45+NW8xICAS
x5lh2sQbxDObtsF5nzjMSrJKfE48YdIFiR+5rnj8xrnossAzW2YmPU8cJhaLHax0MCuZGvE0cVTV
dMoXsh6rnLc4a5Uaa92TvzCU1leWuU5zBEksYgkSRCiooYwKbMRo1UmxxKb9hI8/4volcinkKoOR
YwFVaJBdP/gf/07WkKxNekmhBND94jgfo0DPLtCs0873seM0T4DgM3ClT/3VBjD7SXq9rUWPgIFt
40K6rSl7w0UOMPxkyKbsSkGaQqEAvJ/RN+WAoVugb83rrbWP0wcgQ12lboCDQ2CsSNnrPu/u7ezt
3z0t/n4AkJJysizOae0AAA14aVRydfHNTDpj20uYWRvYmUueG1wAAAAAAA8P3hwYWNrZXQgYmVn
aw49Iu+7vyIgaWQ9Iic1TTBNcENlaGlIenJlU3p0VGNa2M5ZCI/Pgo8eDp4bXBtZXRhIHhtbG5z
Ong9ImFkb2JlOm5zOm1ldGEvIiB40nhtcHRrPSJYTVAgQ29yZSA0LjQuMC1FeGl2MiI+CIA8cmRm
OLJERiB4bWxuc2pyZGY9Imh0dHA6Ly93d3cudzMub3JnLzE5OTkvdMDIvMjItcmRmLXN5bnRheC1u
```



Figure 10: Encoding of Tux image into base64.

## Keys and Addresses Encoding in Bitcoin

- All keys and addresses are encoded using appropriate methods:
  - **Base58Check**: For legacy addresses and private keys.
  - **Bech32**: For SegWit addresses.

### Prefix Summary Table

Data Type	Prefix	Example
Legacy Address	0x00	1PMycacnJa...UAs
SegWit Address	bc1	bc1qw508d6q...
Testnet Address	0x6F	mhPo5P2RVu5...rEo
Private Key (WIF)	0x80	5J76fRXQYWk...U6q

## Base58Check

- Char set: 1 2 3 4 5 6 7 8 9 A B C D E F G H J K L M N P Q R S T U V W X Y Z a b c d e f g h i j k m n o p q r s t u v w x y z.
  - a-z, A-Z, and 0-9, with visually ambiguous characters (0, O, l, I) removed.

### example

- 3 bytes: 0xFFFFFFFF
- Base 58: 2UzHL
- Steps:
  - $0xFFFFFFFF = 16777215$
  - $16777215 \bmod 58 = 19 = L$
  - $289262 \bmod 58 = 16 = H$
  - $4987 \bmod 58 = 57 = z$
  - $85 \bmod 58 = 27 = U$
  - $1 \bmod 58 = 1 = 2$

## Bech32

- Char set: q p z r y 9 x 8 g f 2 t v d w 0 s 3 j n 5 4 k h c e 6 m u a 7 l.
  - a-z, and 0-9, without the following characters: 1, b, i, and, o (b, i, and o can easily be confused with 8, 1, and 0, especially in handwriting or certain fonts).
  - Commonly mistaken characters (e.g. 5 vs S, 2 vs Z, p vs q vs g, etc.) are always one bit different.
  - BCH codes, GF(32), polynomial
 
$$g(x) = x^6 + 29x^5 + 22x^4 + 20x^3 + 21x^2 + 29x + 18.$$
  - Error detection of 4 errors in up to 89 characters.

Talk - Pieter Wuille: New Address Type for SegWit Addresses  
 (Some of) the math behind Bech32 addresses

# ASCII Smuggling

ASCII smuggling is a technique that leverages Unicode characters, which are invisible in user interfaces but can be interpreted by large language models (LLMs), to embed hidden instructions or data within text. This method allows attackers to manipulate AI responses or exfiltrate sensitive information without the user's awareness, by embedding these hidden Unicode tags within clickable hyperlinks or documents shared in chats.

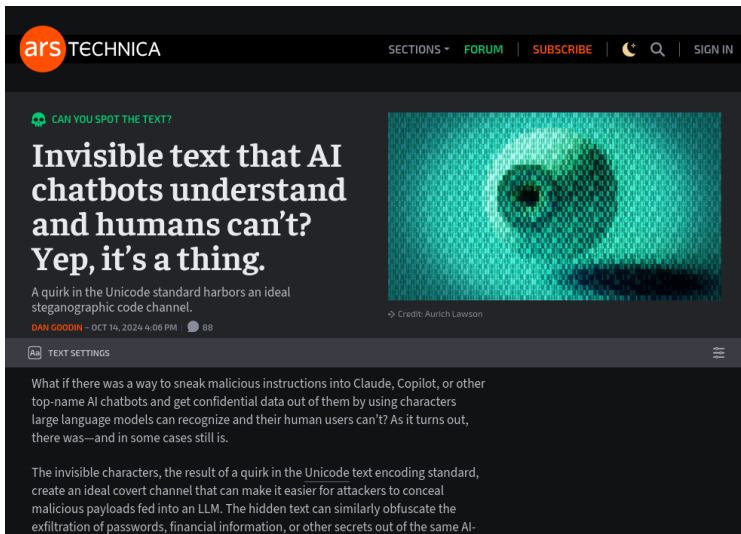


Figure 11: Ars Technica

## ASCII Smuggler

Convert ASCII text to Unicode Tags which are invisible in most UI elements.

Check if a text has hidden Unicode Tags embedded with Decode.

Can you spot the text?

Invisible text that AI chatbots understand and humans can't? Yep, it's a thing.

A quirk in the Unicode standard harbors an ideal steganographic code channel.

Encode

Decode

[Advanced Options](#)

Can you spot the Easter Egg text?

Invisible text that AI chatbots understand and humans can't? Yep, it's a thing.

A quirk in the Unicode standard harbors an ideal steganographic code channel.

**Hidden Unicode Tags discovered.**

Clear

Figure 12: <https://embracethered.com/blog/ascii-smuggler.html>



# ISO/IEC Standards

## ■ ISO/IEC 8859:

- Series for 8-bit character encoding supporting multiple languages.
- ISO-8859-1 (Western Europe), also known as ISO Latin 1.
  - The first 128 characters are identical to ASCII.
  - 0x00 to 1F and 0x80 to 0x9F (hex) used for C0 and C1 control codes.
  - C0 set was originally defined in ISO 646 (ASCII) (e.g., Start of Heading, Start of Text, End of Text, End of Transmission, ...).
  - C1 are additional control codes (e.g., Padding Character, High Octet Preset, Break Permitted Here, No Break Here, ...).

## ■ ISO/IEC 10646:

- Universal character set (UCS) for multilingual text.

	00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
00	<u>MUL</u> 0000	<u>STX</u> 0001	<u>SOT</u> 0002	<u>ETX</u> 0003	<u>EOT</u> 0004	<u>ENQ</u> 0005	<u>ACK</u> 0006	<u>BEL</u> 0007	<u>BS</u> 0008	<u>HT</u> 0009	<u>LF</u> 000A	<u>VT</u> 000B	<u>FF</u> 000C	<u>CR</u> 000D	<u>SO</u> 000E	<u>SI</u> 000F
10	<u>DLE</u> 0010	<u>DC1</u> 0011	<u>DC2</u> 0012	<u>DC3</u> 0013	<u>DC4</u> 0014	<u>NAK</u> 0015	<u>SYN</u> 0016	<u>ETB</u> 0017	<u>CAN</u> 0018	<u>EM</u> 0019	<u>SUB</u> 001A	<u>ESC</u> 001B	<u>FS</u> 001C	<u>GS</u> 001D	<u>RS</u> 001E	<u>US</u> 001F
20	<u>SP</u> 0020	<u>!</u> 0021	<u>"</u> 0022	<u>#</u> 0023	<u>\$</u> 0024	<u>%</u> 0025	<u>&amp;</u> 0026	<u>'</u> 0027	<u>(</u> 0028	<u>)</u> 0029	<u>*</u> 002A	<u>+</u> 002B	<u>,</u> 002C	<u>-</u> 002D	<u>.</u> 002E	<u>/</u> 002F
30	<u>0</u> 0030	<u>1</u> 0031	<u>2</u> 0032	<u>3</u> 0033	<u>4</u> 0034	<u>5</u> 0035	<u>6</u> 0036	<u>7</u> 0037	<u>8</u> 0038	<u>9</u> 0039	<u>:</u> 003A	<u>;</u> 003B	<u>&lt;</u> 003C	<u>=</u> 003D	<u>&gt;</u> 003E	<u>?</u> 003F
40	<u>@</u> 0040	<u>A</u> 0041	<u>B</u> 0042	<u>C</u> 0043	<u>D</u> 0044	<u>E</u> 0045	<u>F</u> 0046	<u>G</u> 0047	<u>H</u> 0048	<u>I</u> 0049	<u>J</u> 004A	<u>K</u> 004B	<u>L</u> 004C	<u>M</u> 004D	<u>N</u> 004E	<u>O</u> 004F
50	<u>P</u> 0050	<u>Q</u> 0051	<u>R</u> 0052	<u>S</u> 0053	<u>T</u> 0054	<u>U</u> 0055	<u>V</u> 0056	<u>W</u> 0057	<u>X</u> 0058	<u>Y</u> 0059	<u>Z</u> 005A	<u>[</u> 005B	<u>\</u> 005C	<u>]</u> 005D	<u>^</u> 005E	<u>_</u> 005F
60	<u>`</u> 0060	<u>a</u> 0061	<u>b</u> 0062	<u>c</u> 0063	<u>d</u> 0064	<u>e</u> 0065	<u>f</u> 0066	<u>g</u> 0067	<u>h</u> 0068	<u>i</u> 0069	<u>j</u> 006A	<u>k</u> 006B	<u>l</u> 006C	<u>m</u> 006D	<u>n</u> 006E	<u>o</u> 006F
70	<u>p</u> 0070	<u>q</u> 0071	<u>r</u> 0072	<u>s</u> 0073	<u>t</u> 0074	<u>u</u> 0075	<u>v</u> 0076	<u>w</u> 0077	<u>x</u> 0078	<u>y</u> 0079	<u>z</u> 007A	<u>{</u> 007B	<u> </u> 007C	<u>}</u> 007D	<u>~</u> 007E	<u>DEL</u> 007F
80																
90																
A0	<u>NEBSP</u> 00A0	<u>†</u> 00A1	<u>‡</u> 00A2	<u>£</u> 00A3	<u>¤</u> 00A4	<u>¥</u> 00A5	<u>¦</u> 00A6	<u>§</u> 00A7	<u>¨</u> 00A8	<u>©</u> 00A9	<u>ª</u> 00AA	<u>«</u> 00AB	<u>¬</u> 00AC	<u>­</u> 00AD	<u>®</u> 00AE	<u>¯</u> 00AF
B0	<u>°</u> 00B0	<u>±</u> 00B1	<u>²</u> 00B2	<u>³</u> 00B3	<u>´</u> 00B4	<u>µ</u> 00B5	<u>¶</u> 00B6	<u>·</u> 00B7	<u>¸</u> 00B8	<u>¹</u> 00B9	<u>º</u> 00BA	<u>»</u> 00BB	<u>¼</u> 00BC	<u>½</u> 00BD	<u>¾</u> 00BE	<u>¿</u> 00BF
C0	<u>À</u> 00C0	<u>Á</u> 00C1	<u>Â</u> 00C2	<u>Ã</u> 00C3	<u>Ä</u> 00C4	<u>Å</u> 00C5	<u>Æ</u> 00C6	<u>Ç</u> 00C7	<u>È</u> 00C8	<u>É</u> 00C9	<u>Ê</u> 00CA	<u>Ë</u> 00CB	<u>Ì</u> 00CC	<u>Í</u> 00CD	<u>Î</u> 00CE	<u>Ï</u> 00CF
D0	<u>Ð</u> 00D0	<u>Ñ</u> 00D1	<u>Ò</u> 00D2	<u>Ó</u> 00D3	<u>Ô</u> 00D4	<u>Õ</u> 00D5	<u>Ö</u> 00D6	<u>×</u> 00D7	<u>Ø</u> 00D8	<u>Ù</u> 00D9	<u>Ú</u> 00DA	<u>Û</u> 00DB	<u>Ü</u> 00DC	<u>Ý</u> 00DD	<u>Þ</u> 00DE	<u>ß</u> 00DF
E0	<u>à</u> 00E0	<u>á</u> 00E1	<u>â</u> 00E2	<u>ã</u> 00E3	<u>ä</u> 00E4	<u>å</u> 00E5	<u>æ</u> 00E6	<u>ç</u> 00E7	<u>è</u> 00E8	<u>é</u> 00E9	<u>ê</u> 00EA	<u>ë</u> 00EB	<u>ì</u> 00EC	<u>í</u> 00ED	<u>î</u> 00EE	<u>ï</u> 00EF
F0	<u>ø</u> 00F0	<u>ñ</u> 00F1	<u>ò</u> 00F2	<u>ó</u> 00F3	<u>ô</u> 00F4	<u>õ</u> 00F5	<u>÷</u> 00F6	<u>ø</u> 00F7	<u>ù</u> 00F8	<u>ú</u> 00F9	<u>û</u> 00FA	<u>ü</u> 00FB	<u>ý</u> 00FC	<u>ÿ</u> 00FD	<u>þ</u> 00FE	<u>ÿ</u> 00FF

Figure 13: ISO-8859-1 code page

# Windows Code Pages

- **Overview:**

- Multiple code pages for different regions and languages.
- Used in Microsoft Windows from the 1980s and 1990s.

- **Examples:**

- CP1252 (Western Europe), CP932 (Japan)

- **Issues:**

- Inconsistencies across different systems.

## Windows Encoding Transition to Unicode

### ■ **UCS-2 (Unicode Character Set - 2 bytes):**

- **Introduction:** Windows NT 3.1 (1993)
- **Details:** 16-bit fixed-width encoding for the first 65,536 Unicode characters, used internally for Windows APIs.

### ■ **UTF-16:**

- **Adoption:** Windows 2000 (2000)
- **Details:** An extension of UCS-2, accommodating all Unicode characters by using surrogate pairs for characters beyond the Basic Multilingual Plane (BMP).
  - **Surrogate Pair Example:** The emoji 😊 (Unicode U+1F60A) would be represented as: U+D83D (High Surrogate) + U+DE0A (Low Surrogate).

### ■ **UTF-8:**

- **Support Added:** Windows 10 version 1803 (April 2018 Update)
- **Details:** Variable-width encoding, backward compatible with ASCII. Became more prominently supported for developers with the introduction of the `ActiveCodePage` property in Windows 10 1903 (May 2019 Update).

# Unicode

- Universal character set covering all scripts, supporting over 143,000 characters.
- It assigns a unique number (called a “code point”) to each character, regardless of platform, program, or language.
- 1.112.064 valid code points within the codespace.
- As of Unicode 16.0, released in September 2024, 299,056 (27%) of these code points are allocated, 155,063 (14%) have been assigned characters, 137,468 (12%) are reserved for private use, 2,048 are used to enable the mechanism of surrogates, and 66 are designated as noncharacters, leaving the remaining 815,056 (73%) unallocated.
- Unicode has different encoding forms: UTF-8, UTF-16, and UTF-32.

# UTF

- **UTF-8:**
  - Variable-length encoding, backward compatible with ASCII, byte-order independent.
- **UTF-16:**
  - Variable-length encoding (2 or 4 bytes per character).
  - Latin and most commonly used CJK<sup>1</sup> characters are encoded in 2 bytes.
- **UTF-32:**
  - Fixed-length encoding (4 bytes per character).

---

<sup>1</sup>Chinese, Japanese, and Korean.

Number of bytes	First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4	Byte 5	Byte 6
1	U+0000	U+007F	0xxxxxxx					
2	U+0080	U+07FF	110xxxxx	10xxxxxx				
3	U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx			
4	U+10000	U+1FFFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx		
5	U+200000	U+3FFFFFFF	111110xx	10xxxxxx	10xxxxxx	10xxxxxx	10xxxxxx	
6	U+4000000	U+7FFFFFFF	1111110x	10xxxxxx	10xxxxxx	10xxxxxx	10xxxxxx	10xxxxxx

Figure 14: UTF-8 Structure

## UTF-8 takes over

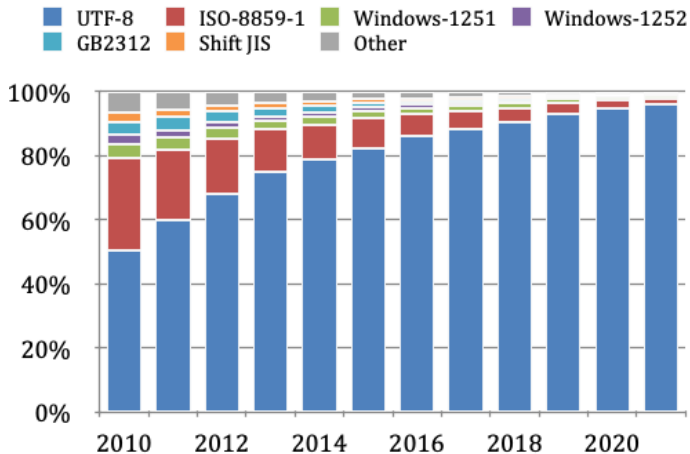


Figure 15: Declared character set for the 10 million most popular websites since 2010



# Endianness

- **Big Endian vs. Little Endian:**
  - Byte order in memory representation.
  - Impacts how multi-byte characters are read.

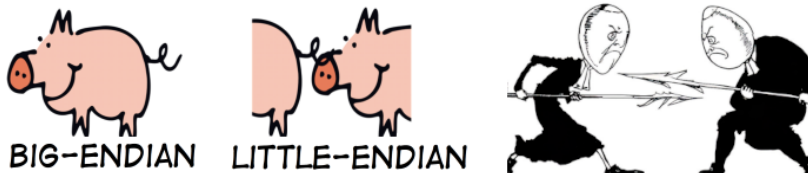
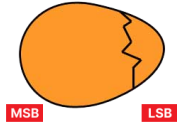


Figure 16: Endianness

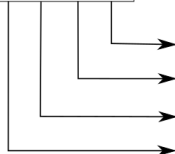
LITTLE ENDIAN - The  
Lilliputians break their eggs  
at the smaller end



Little-endian

32-bit integer

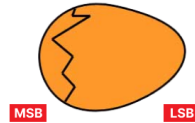
0A0B0C0D



Memory

⋮		⋮
0D	$a$	0A
0C	$a+1$	0B
0B	$a+2$	0C
0A	$a+3$	0D
⋮		⋮

BIG ENDIAN - The  
Blefuscutians break their  
eggs at the big end.



Big-endian

32-bit integer

0A0B0C0D

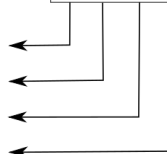


Figure 17: Big-Endian and Little-Endian

## ■ Example:

- UTF-16 and UTF-32 can be big or little endian.
- Byte Order Mark (BOM) to indicate the endianness being used.
- The BOM is the code point U+FEFF (BOM, ZWNBSP<sup>2</sup>).
  - Big-endian (UTF-16BE): FE FF
  - Little-endian (UTF-16LE): FF FE
  - Big-endian (UTF-32BE): 00 00 FE FF
  - Little-endian (UTF-32LE): FF FE 00 00
  - BOM in UTF-8: EF BB BF, serves more as a signature to indicate that the file is encoded in UTF-8 rather than specifying byte order.

---

<sup>2</sup>zero width no-break space

## Text File Formats

- **.txt:** Usually ASCII or UTF-8.
- **.csv:** Can use various encodings; important for data exchange.
- **.json:** JavaScript Object Notation, for data interchange.
- **.yaml:** YAML Ain't Markup Language, for data serialization.
- **.log:** Log files for recording events, errors, and system activities.
- **.ini:** Initialization files for configuration settings.
- **.conf:** Configuration files, similar to .ini, used by many applications.

# Markup Files

- **Markdown:**

- Lightweight markup language for formatting text.
- Markdown itself doesn't have a built-in mechanism for declaring encoding in the file header.

- **TeX:**

- Typesetting language for high-quality typography.
- Encoding: Often UTF-8, but can be sensitive to non-ASCII characters without proper preamble setup.
- `\usepackage[utf8]{inputenc}`

- **XML (eXtensible Markup Language):**

- Used for structured data storage and transmission.
- Encoding: Declared in XML declaration, typically UTF-8 or UTF-16. Encoding declaration is crucial for correct parsing.
- `<?xml version="1.0" encoding="UTF-8"?>`

- **HTML (HyperText Markup Language):**

- Standard markup language for documents designed to be displayed in a web browser.
- Encoding: Default is often UTF-8, but can be specified with the charset attribute in the `<meta>` tag. Incorrect encoding can lead to garbled text.
- `<head><meta charset="UTF-8"></head>`
- HTTP Content-Type header: `Content-Type: text/html; charset=UTF-8`

## Linux Tools for Text Encoding

- **iconv:** Converts text from one encoding to another.
  - Example: `iconv -f ISO-8859-1 -t UTF-8 input.txt > output.txt`
- **file:** Identifies file types and encodings.
  - Example: `file -i example.txt`
- **uconv (from ICU):** More advanced conversion with Unicode support.
  - Example: `uconv -f UTF-8 -t UTF-16 input.txt -o output.txt`
- **dos2unix / unix2dos:** Converts between Windows and Unix line endings.
  - Example: `dos2unix file.txt` (converts CRLF to LF)
  - Example: `unix2dos file.txt` (converts LF to CRLF)

## ■ **base64:**

- Example: Used for encoding e-mail attachemets.
- Usage: `echo "test" | base64` to encode, `echo "dGVzdA==" | base64 -d` to decode.

## ■ **base58:**

- Example: Useful for encoding addresses in cryptocurrencies (e.g., Bitcoin).
- Usage: `echo "test" | base58` to encode, `echo "E8f4pE5" | base58 -d` to decode.

## ■ **base32:**

- Example: Used to encode addresses in Bitcoin's Segregated Witness (SegWit) protocol.
- Usage: `echo "test" | base32` for encoding, `echo "ORSXG5A=" | base32 -d` for decoding.



## ■ **recode:**

- Function: Similar to iconv but with additional capabilities.
- Usage: `recode latin1..utf-8 file.txt`

## ■ **xxd:**

- Function: Create a hex dump of a binary file, useful for understanding byte-level data.
- Usage: `xxd -p file.bin` for plain hex, `xxd -r -p hex.txt` to revert.

## ■ **Use Cases:**

- Data migration, cleaning, and internationalization.

## File Extensions

- **Just a Name:** Extensions don't define file content.
- **Content Matters:** True type determined by data inside.
- **Beware:** Misleading extensions can be risky.
- **Purpose:** Created to indicate file type for ease of use.

# How file Detects Encoding

## Magic Numbers / File Signatures

- `file` looks for specific byte sequences at the start of files that uniquely identify file formats or types.
  - For text in UTF-8, looks for the BOM marker `EF BB BF`.
  - For JPEG images, it looks for `FF D8 FF`.
  - A PDF file starts with `%PDF`.
  - PNG images start with the bytes `89 50 4E 47 0D 0A 1A 0A`.
  - WAV files start with `52 49 46 46`.
  - MP3 files might begin with `49 44 33`.

## Text Files - Encoding Detection

- **Heuristics:** When magic numbers aren't conclusive, `file` uses heuristics.
  - **Character Analysis:** Examines byte sequences for patterns typical of specific encodings.
  - **Frequency Analysis:** Looks at the frequency and distribution of characters to guess language and thus encoding.
  - **Control Characters:** Presence or absence of certain control characters can hint at encoding.

## MIME Database

- Maps file content to MIME types and encodings.
- Location: `/usr/share/file/magic.mgc` or similar (compiled database).
  - The `magic.mgc` database is generated from a set of “magic” text files (e.g., `/etc/magic`) which define the rules for recognizing various file formats.
  - The rules consist of:
    - Byte offsets
    - Byte patterns
    - Regular expressions
    - Human-readable descriptions
  - Example:

0	string	<code>\x89PNG\r\n\x1a\n</code>	PNG image data
0	string	<code>%PDF-</code>	PDF document

List of file signatures (Wikipedia)

## False Recognition/Confusion

- **Ambiguity:** Some files can be interpreted as multiple encodings, especially if they contain only ASCII characters.
  - Example: A file with only ASCII might be reported as `us-ascii`, but could be UTF-8 or ISO-8859-1.
- **Incomplete Information:** Short files or files with limited character sets might not provide enough data for accurate detection.
- **Encoding Overlap:** Encodings that share a subset of characters (like ASCII in UTF-8) can lead to confusion.
- **Binary in Text:** Files with embedded binary data might confuse the tool into thinking it's a binary file rather than text with encoding.
- **False Positives:** Sometimes, `file` might guess wrong due to patterns that mimic another encoding or due to an updated but not comprehensive magic database.

# Conclusion

## ■ Key Points:

- **Journey Through Encoding:** From historical codes like Morse, Baudot to modern standards like UTF-8 and Unicode.
- **Evolution:** Text encoding has moved from simple to complex systems.
- **Universal Solution:** Unicode provides a global text representation.
- **Concepts, Challenges, Solutions, Applications, and Mindfulness:** Understanding these is crucial.

## ■ Future:

- Continued evolution of encoding standards to accommodate new scripts and symbols.

