

Steven R. Bagley

#### Introduction

- How do we get a computer to understand text?
- Unusual question to ask...
- After all, been using string and char for ages
- But a computer has no notion of text...

# Binary Encoding

- Computer only understands binary
- For it to process text, we must encode the characters in binary
- Doesn't really matter how you do it
- Providing your consistent
- And don't need to exchange it with other computers/software

# History

- Interesting to look back at the history of characters
- Helps us understand why some things are done in certain ways
- Also see how thinking about your encoding can make it much simpler to write software
- Particularly about how bits are used...

e.g. line endings...

## Teleprinter

- Computers originally output to teleprinters
- Teleprinter is basically an electromechanical typewriter
- Where the key presses could be encoded and sent over a wire to another teletype(s)
- Which would also type the message
- Or stored on paper tape...



A British Creed & Company Teleprinter No. 7



Teletype Model 33 ASR

# Teleprinter

- A computer could generate the same encoding
- Causing the teleprinter to print out...
- Of course, the computer also had to follow the same rules as a teleprinter when generating its output

Either over a wire or on paper tape Carriage returns/Line feeds etc...

#### CR/LF

- Send codes to move to the next line
- Like a typewriter, this was done by sending a Carriage Return (CR)
- This returned the carriage to the beginning of the line
- Followed by a line feed (LF) to move to the next line

# Line ending

- Actually, it had to send several CR
- Mechanical process slower than data transmission
- Hence, CR CR CR LF would be sent
- Eventually, reduced to CRLF
- To this day, on DOS/Windows systems...

# Line Endings

- If CRLF is a line ending, what does a CR or LF on its own mean?
- Multics reduced this to just an LF to avoid confusion
- Device driver converted it as necessary
- Other systems used just CR for similar reasons

#### Baudot Code

- Early Teleprinters used the Baudot Code
- Developed by Emile Baudot in 1874
- 5-bit binary code



#### Baudot Code

- 5-bits, gives us 2<sup>5</sup> or 32 codes
- Enough for all 26 letters
- Code 0 is left unused
- Five spaces left over...

00		08		10	E	18	A
01	Т	09	L	11	Z	19	W
02		0A	R	12	D	1A	J
03	0	0B	G	13	В	1B	
04		0C	I	14	S	1C	U
05	H	0 D	Р	15	Y	1D	Q
06	N	0E	С	16	F	1E	K
07	Μ	0F	V	17	Х	1F	

Position zero unused Remeber only 5 bits used!

#### Baudot Code

- Some of the unused codes used for control codes...
- Carriage Return, Line Feed, Spaces
- But what about numbers or punctuation?

#### Shift Code

- Makes use of a 'shift-code'...
- Certain values shifts the understanding of the other values
- Figures shift switches the code to a different mode
- Letters shift switches it back to letters

Position zero unused

#### LETTER CASE

00		08	LINE FEED	10	E	18	A
01	Т	09	L	11	Z	19	W
02	CARRIAGE RETURN	0A	R	12	D	1A	J
03	0	0B	G	13	В	1B	FIGURE SHIFT
04	SPACE	0C	I	14	S	1C	U
05	H	0D	Р	15	Y	1D	Q
06	N	0E	С	16	F	1E	K
07	М	OF	V	17	X	1F	LETTERS SHIFT

00		08	LINE FEED		10	3	18	
01	5	09	)		11	+	19	2
02	CARRIAGE RETURN	0A	4	2	12	WHO ARE YOU	1A	BELL
03	9	0B	&		13	?	1B	FIGURE SHIFT
04	SPACE	0C	8		14	•	1C	7
05	#	0 D	0		15	6	1D	1
06	,	0E	:		16	\$	1E	(
07	•	0F	=		17	/	1F	LETTERS SHIFT

#### FIGURE CASE

Position zero unused Explain how we can write a message in Baudot code Use whiteboard

#### Observations

- Need to know what has gone before to understand the next character
- No lower case...
- Stunt characters (CR,LF, SHIFTs) never change position
- No discernible order to the code...
- Makes processing it hard

at least what mode we are in (show how this can go wrong by decoding the message we've written twice).

e.g. convert character to integer...

## Other Character Sets

- IBM used EBCDIC until 1981 and the PC
- 8-bit code, based on Binary Coded Decimal
- Lots of others about as well
- All supported a different set of characters
- Needed a standardized set
- Enter ASCII in the 1960s

Although took till the eighties to really catch on...

#### ASCII

- 7-bit code 127 characters
- Decided against using a shift character and 6-bits
- 7-bit saved space on eight-bits
- Lack of shift key made transmission more reliable

#### ASCII

- Took over two years to decide what characters to encode...
- A lot more characters encoded including lower case
- A very discernible order to the layout as well...

**USASCII** code chart

		_					6	· ' ı	0	0	10	1
1	D 3	P 5	Ь , +	Row	0	I	2	3	4	5	6	7
0	0	0	0	0	NUL .	DLE	SP	0	0	Ρ	ì	P
0	0	0	1	1	SOH	DC1	!	1	Α.	Q	a	q
0	0	1	0	2	STX	DC2	"	2	В	R	b	r
0	0	1	Ī	3	ETX	DC 3	#	3	C	S	C	5
0	1	0	0	4	EOT	DC4	\$	4	D	т	d	1
0	1	0	1	5	ENQ	NAK	%	5	E	υ	e	U
0	1	1	0	6	ACK	SYN	8	6	F	V	f	V
0	1	1	1	7	BEL	ETB	•	7	G	W	9	w
L	0	0	0	8	BS	CAN	(	8	н	×	h	X
L	0	0	1	9	нт	EM	)	9	1	Y	i	у
Ĺ	0	1	0	10	LF	SUB	*	:	J	Z	j	z
1	0	T	1		VT	ESC	+	;	к	C	k.	{
1	1	0	0	12	FF	FS	•	<	L	N	1	1
1	1	0	1	13	CR	GS	-	¥	м	3	m	}
1	1	1	0	4	SO	RS		>	N	^	n	$\sim$
1	1	1	1	15	<b>S1</b>	US	1	?	0		0	DEL
		0       0         0       0         0       0         0       0         0       1         0       1         0       1         0       1         0       1         0       1         0       1         1       0         1       0         1       1         1       1         1       1         1       1         1       1	0       0       0         0       0       0         0       0       1         0       0       1         0       1       0         0       1       0         0       1       1         0       1       1         0       1       1         0       1       1         1       0       1         1       0       1         1       0       1         1       0       1         1       1       0         1       1       1         1       1       1	0       0       0       0         0       0       0       1         0       0       1       0         0       0       1       1         0       0       1       0         0       1       0       0         0       1       0       0         0       1       1       0         0       1       1       1         0       1       1       1         0       1       1       1         1       0       0       1         1       0       1       1         1       0       1       1         1       0       1       1         1       1       0       1         1       1       0       1         1       1       1       0         1       1       1       1	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	0       0       0       0       0       NUL         0       0       0       1       1       SOH         0       0       1       0       2       STX         0       0       1       0       2       STX         0       0       1       1       3       ETX         0       0       1       1       3       ETX         0       1       0       0       4       EOT         0       1       0       0       4       EOT         0       1       0       1       5       ENQ         0       1       0       6       ACK         0       1       1       7       BEL         1       0       0       1       9       HT         1       0       1       9       HT         1       0       1       1       VT         1       0       1       1       VT         1       1       0       14       S0         1       1       1       15       S1	0       0       0       0       0       NUL       DLE         0       0       0       1       1       SOH       DC1         0       0       1       0       2       STX       DC2         0       0       1       0       2       STX       DC2         0       0       1       1       3       ETX       DC3         0       1       0       0       4       EOT       DC4         0       1       0       0       4       EOT       DC4         0       1       0       1       5       ENQ       NAK         0       1       1       0       6       ACK       SYN         0       1       1       1       7       BEL       ETB         1       0       0       8       BS       CAN         1       0       0       8       BS       CAN         1       0       1       9       HT       EM         1       0       1       1       VT       ESC         1       0       1       13       CR       GS     <	D       O       O       O       O       NUL       DLE       SP         D       O       I       I       SOH       DC1       !         D       O       I       O       2       STX       DC2       "         D       O       I       I       3       ETX       DC2       "         D       O       I       I       3       ETX       DC2       "         D       I       O       2       STX       DC2       "         D       I       O       2       STX       DC2       "         D       I       O       2       STX       DC2       "         D       I       O       A       EOT       DC4       #         D       I       O       I       SENQ       NAK       %         D       I       I       O       6       ACK       SYN       B         D       I       I       I       7       BEL       ETB       '         I       O       I       9       HT       EM       )       I         I       O       I <td< th=""><th>D       O       O       O       NUL       DLE       SP       O         D       O       O       I       I       SOH       DC1       !       1         D       O       I       O       2       STX       DC2       "       2         D       O       I       I       3       ETX       DC2       "       2         D       O       I       I       3       ETX       DC3       #       3         D       I       O       Q       A       EOT       DC4       4       4         D       I       O       I       SENQ       NAK       %       5         D       I       O       I       SENQ       NAK       %       5         O       I       I       O       6       ACK       SYN       8       6         D       I       I       I       O       6       ACK       SYN       8       6         D       I       I       I       T       T       E       I       7         I       O       O       I       B       BS       CAN</th><th>D       O       O       O       O       NUL       DLE       SP       O       Ø         D       O       O       I       I       SOH       DC1       I       I       A         D       O       I       O       2       STX       DC2       "       2       B         D       O       I       O       2       STX       DC2       "       2       B         D       O       I       I       3       ETX       DC3       #       3       C         D       I       O       Q       4       EOT       DC4       1       4       D         D       I       O       I       5       ENQ       NAK       %       5       E         D       I       O       6       ACK       SYN       B       6       F         D       I       I       7       BEL       ETB       '       7       G         I       I       O       I       9       HT       EM       )       9       I         I       O       I       I       H       VT       ESC       <t< th=""><th>D       O       O       O       O       NUL       DLE       SP       O       Ø       P         D       O       O       I       I       SOH       DC1       !       1       A       Q         D       O       I       I       SOH       DC1       !       1       A       Q         D       O       I       O       2       STX       DC2       "       2       B       R         D       O       I       I       3       ETX       DC3       #       3       C       S         D       I       O       0       4       EOT       DC4       4       D       T         D       I       0       1       5       ENQ       NAK       %       5       E       U         D       I       0       6       ACK       SYN       8       6       F       V         D       I       I       7       BEL       ETB       '       7       G       W         I       0       0       8       BS       CAN       (       8       H       X</th><th>D       O       O       O       NUL       DLE       SP       O       @       P       ``         D       O       O       I       I       SOH       DC1       I       I       A       Q       a         D       O       O       I       I       SOH       DC1       I       I       A       Q       a         D       O       I       O       2       STX       DC2       "       2       B       R       b         D       O       I       I       3       ETX       DC3       #       3       C       S       c         D       I       O       0       4       EOT       DC4       4       D       T       d         D       I       O       6       ACK       SYN       8       6       F       V       f         D       I       I       O       6       ACK       SYN       8       6       F       V       f         D       I       I       7       BEL       ETB       '       7       G       W       g         I       0       0</th></t<></th></td<>	D       O       O       O       NUL       DLE       SP       O         D       O       O       I       I       SOH       DC1       !       1         D       O       I       O       2       STX       DC2       "       2         D       O       I       I       3       ETX       DC2       "       2         D       O       I       I       3       ETX       DC3       #       3         D       I       O       Q       A       EOT       DC4       4       4         D       I       O       I       SENQ       NAK       %       5         D       I       O       I       SENQ       NAK       %       5         O       I       I       O       6       ACK       SYN       8       6         D       I       I       I       O       6       ACK       SYN       8       6         D       I       I       I       T       T       E       I       7         I       O       O       I       B       BS       CAN	D       O       O       O       O       NUL       DLE       SP       O       Ø         D       O       O       I       I       SOH       DC1       I       I       A         D       O       I       O       2       STX       DC2       "       2       B         D       O       I       O       2       STX       DC2       "       2       B         D       O       I       I       3       ETX       DC3       #       3       C         D       I       O       Q       4       EOT       DC4       1       4       D         D       I       O       I       5       ENQ       NAK       %       5       E         D       I       O       6       ACK       SYN       B       6       F         D       I       I       7       BEL       ETB       '       7       G         I       I       O       I       9       HT       EM       )       9       I         I       O       I       I       H       VT       ESC <t< th=""><th>D       O       O       O       O       NUL       DLE       SP       O       Ø       P         D       O       O       I       I       SOH       DC1       !       1       A       Q         D       O       I       I       SOH       DC1       !       1       A       Q         D       O       I       O       2       STX       DC2       "       2       B       R         D       O       I       I       3       ETX       DC3       #       3       C       S         D       I       O       0       4       EOT       DC4       4       D       T         D       I       0       1       5       ENQ       NAK       %       5       E       U         D       I       0       6       ACK       SYN       8       6       F       V         D       I       I       7       BEL       ETB       '       7       G       W         I       0       0       8       BS       CAN       (       8       H       X</th><th>D       O       O       O       NUL       DLE       SP       O       @       P       ``         D       O       O       I       I       SOH       DC1       I       I       A       Q       a         D       O       O       I       I       SOH       DC1       I       I       A       Q       a         D       O       I       O       2       STX       DC2       "       2       B       R       b         D       O       I       I       3       ETX       DC3       #       3       C       S       c         D       I       O       0       4       EOT       DC4       4       D       T       d         D       I       O       6       ACK       SYN       8       6       F       V       f         D       I       I       O       6       ACK       SYN       8       6       F       V       f         D       I       I       7       BEL       ETB       '       7       G       W       g         I       0       0</th></t<>	D       O       O       O       O       NUL       DLE       SP       O       Ø       P         D       O       O       I       I       SOH       DC1       !       1       A       Q         D       O       I       I       SOH       DC1       !       1       A       Q         D       O       I       O       2       STX       DC2       "       2       B       R         D       O       I       I       3       ETX       DC3       #       3       C       S         D       I       O       0       4       EOT       DC4       4       D       T         D       I       0       1       5       ENQ       NAK       %       5       E       U         D       I       0       6       ACK       SYN       8       6       F       V         D       I       I       7       BEL       ETB       '       7       G       W         I       0       0       8       BS       CAN       (       8       H       X	D       O       O       O       NUL       DLE       SP       O       @       P       ``         D       O       O       I       I       SOH       DC1       I       I       A       Q       a         D       O       O       I       I       SOH       DC1       I       I       A       Q       a         D       O       I       O       2       STX       DC2       "       2       B       R       b         D       O       I       I       3       ETX       DC3       #       3       C       S       c         D       I       O       0       4       EOT       DC4       4       D       T       d         D       I       O       6       ACK       SYN       8       6       F       V       f         D       I       I       O       6       ACK       SYN       8       6       F       V       f         D       I       I       7       BEL       ETB       '       7       G       W       g         I       0       0

Go through some of the ordered parts

#### ASCII order

- Contiguous control codes, letters and numbers...
- Control codes grouped at the low end
- In a sortable order, so separator characters before letters/numbers
- So can sort text purely on the ASCII code numbers

#### ASCII order

- Punctuation matches the numbers in the next column (similar to that found on typewriter keyboards)
- Upper-case and lower-case are precisely one bit different
- Start position of alphabet chosen to match a British standard...

Makes it easy to change case (and 1/7th of bit errors during transmission won't make the message unintelligble).

## ASCII limitations

- Only encodes 127 characters...
- US centric, e.g. no encoding for a £ symbol
- Often used in an extended 8-bit
- But these extensions aren't standard
- ISO 8859 defines several different ones for different areas...

e.g. left right double quotes are 210,211 in Mac Roman encoding, but this O grave or O acute in WinAnsi

- Need a character set with a fixed position for every character
- Including international character sets other than the Latin Alphabet
- Unicode was designed as such a Universal Character System
- Dates back to the late-1980s

- Aimed to be able to encode the union of all characters used in 1988
- More characters than could fit in 8-bits
- Originally 16-bit, now 32-bit characters
- But will only go up U+10FFFF
- Often called wide characters

That's character with hex code 0x10FFFF Modern OSes and Systems tend to be unicode NT et al are 16bitUnicode internally Others tend to be UTF8

- Defined in terms of I7 Code Planes (0-16)
- Not all used...
- Plane 0 is the Basic Multilingual Plane
- Other planes define CJK ideographs, historic scripts etc.

Kept the first 127 chars identical to ASCII
The next 128 are roughly similar to common usage

# Multibyte

- Multibyte characters are subject to endian issues like any other multibyte word
- Unicode provides a way to identify the order using a special character
- The Byte-Order Marker (вом) U+FEFF
- The inverse (U+FFFE) isn't a legal character
- Used as the first character in the stream

XML spec shows how to recover the byte order from the first few bytes whethere it is there or not (due to a known opening sequence).

- Multibyte nature means it takes up more space than ASCII
- Yet for latin text, the code points are identical lots of zero bytes
- Wouldn't it be great if we only had to use multibyte characters when necessary

## UTF

- UCS Transformation Format aims to do just that
- Variable length byte sequences
- Aims to be on average shorter than 16-bit chars
- But for some characters will be longer

- Uses a modulo 190 based system
- Characters U+0000 U+009F encoded as is
- Characters from U+00A0 processed through a set of rules

#### Code Point

#### Encoding

Code Point	Encoding
x < U+009F	X

Code Point	Encoding
x < U+009F	X
U+00A0 <= x <= U+00FF	A0 X

Code Point	Encoding				
x < U+009F	X				
U+00A0 <= x <= U+00FF	A0 x				
U+0100 <= x <= U+4015 y = (x - 0x100)	A1 + y / 0xBE T(y % 0xBE)				

Code Point	Encoding			
x < U+009F	X			
U+00A0 <= x <= U+00FF	A0 x			
U+0100 <= x <= U+4015 y = (x - 0x100)	A1 + y / 0xBE T(y % 0xBE)			
U+4016 <= x < U+38E2D y = x - 0x4016	F6 + y / 0xBE2 T(y / BE % BE) T(y % BE)			

Code Point	Encoding			
x < U+009F	X			
U+00A0 <= x <= U+00FF	A0 x			
$\begin{array}{llllllllllllllllllllllllllllllllllll$	A1 + y / 0xBE T(y % 0xBE)			
U+4016 <= x < U+38E2D y = x - 0x4016	F6 + y / 0xBE2 T(y / BE % BE) T(y % BE)			
•••	•••			

#### T(z) defined:

Z	T(z)
0x00-0x5D	z + 0x21
0x5E—0xBD	z + 0x42
0xBE-0xDE	z – OxBE
0xBF-0xFF	z - 0x60

- Problem with UTF-1 is that it is not selfsyncing
- You need to know where you are in the stream to understand a byte
- Because it uses all byte values as both single characters and parts of multibyte characters

- Designed by Ken Thompson and Rob Pike over dinner on the back of an envelope
- Aimed to solve the problems of UTF-I
- Self-syncing
- Type of each byte uniquely decodable
- Varies in length from 1-4 bytes

Yes really see -- <u>http://www.cl.cam.ac.uk/~mgk25/ucs/utf-8-history.txt</u>Or Not quite as compact as UTF-1

Bits	Last Code Point	Byte I	Byte 2	Byte 3	Byte 4
------	--------------------	--------	--------	--------	--------

Bits	Last Code Point	Byte I	Byte 2	Byte 3	Byte 4
7	U+007F	0xxxxxxx			

Bits	Last Code Point	Byte I	Byte 2	Byte 3	Byte 4
7	U+007F	0xxxxxxx			
11	U+07FF	110xxxxx	10xxxxxx		

Bits	Last Code Point	Byte I	Byte 2	Byte 3	Byte 4
7	U+007F	0xxxxxxx			
11	U+07FF	110xxxxx	10xxxxxx		
16	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	

Bits	Last Code Point	Byte I	Byte 2	Byte 3	Byte 4
7	U+007F	0xxxxxxx			
11	U+07FF	110xxxxx	10xxxxxx		
16	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	
21	U+1FFFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

# UTF-8 self syncing

- Top two bits tell us the type of byte
  - Bit 7 is 0, a single byte character
  - Bit 6-7 is 11, start of a MBCS, a leading byte
  - Bit 6-7 is 10, part of a MBCS, a continuing byte
- Can easily find our place in the stream for decoding...

#### Conclusion

- Text requires us to define a mapping to store it
- But how we define that mapping affects how the software must work
- By thinking carefully about the mapping, we can make it much easier to write software