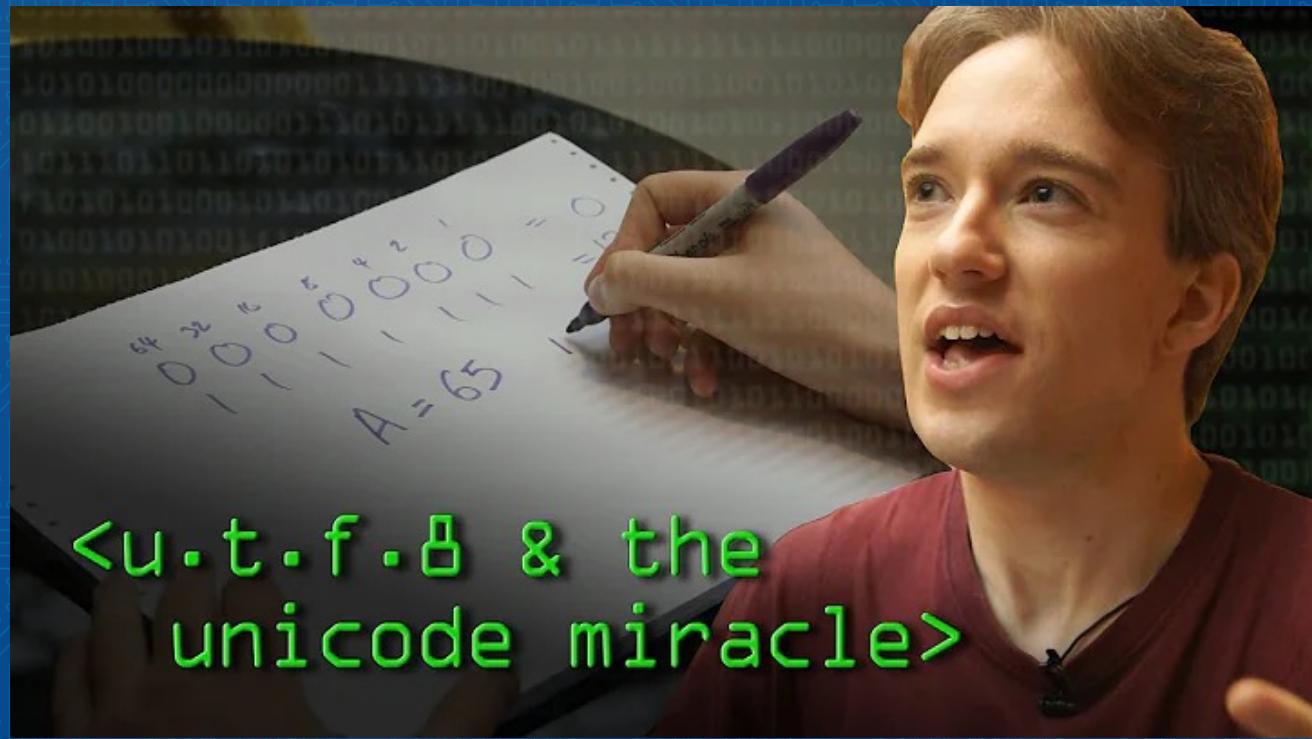


A Practical Guide to UTF-8 and Character Encodings

Utah Stack JS meetup
Ken Snyder



Tom Scott: “The most beautiful hack”



<u.t.f.8 & the
unicode miracle>

What we'll cover

- History
- Design
- Landmines
- Web technologies
- Tools
- Emoji

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	 	Space	64	40	100	@	Ø	96	60	140	`	'
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A	97	61	141	a	a
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B	98	62	142	b	b
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C	99	63	143	c	c
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D	100	64	144	d	d
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E	101	65	145	e	e
6	6	006	ACK (acknowledge)	38	26	046	&	&	70	46	106	F	F	102	66	146	f	f
7	7	007	BEL (bell)	39	27	047	'	'	71	47	107	G	G	103	67	147	g	g
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H	104	68	150	h	h
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I	105	69	151	i	i
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J	106	6A	152	j	j
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K	107	6B	153	k	k
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L	108	6C	154	l	l
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M	109	6D	155	m	m
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N	110	6E	156	n	n
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O	111	6F	157	o	o
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P	112	70	160	p	p
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q	113	71	161	q	q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R	114	72	162	r	r
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S	115	73	163	s	s
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T	116	74	164	t	t
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U	117	75	165	u	u
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V	118	76	166	v	v
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W	119	77	167	w	w
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X	120	78	170	x	x
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y	121	79	171	y	y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z	122	7A	172	z	z
27	1B	033	ESC (escape)	59	3B	073	;	;	91	5B	133	[[123	7B	173	{	{
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\	124	7C	174	|	
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]]	125	7D	175	}	}
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^	126	7E	176	~	~
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_	127	7F	177		DEL

Source: www.LookupTables.com

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

Low Ascii											
000:	013: J	026: →	039: '	052: 4	065: À	078: Μ	091: [104: h	117: u		
001: Θ	014: Π	027: ←	040: (053: 5	066: Β	079: Ο	092: \`	105: i	118: v		
002: Ω	015: *	028: ↶	041:)	054: 6	067: Ϲ	080: Ρ	093:]	106: j	119: w		
003: ♥	016: ►	029: ↔	042: *	055: 7	068: Ⓓ	081: Ϙ	094: ^	107: k	120: x		
004: ♦	017: ◀	030: ▲	043: +	056: 8	069: Ⓔ	082: 𝑹	095: _	108: l	121: y		
005: ◆	018: ‡	031: ▼	044: ,	057: 9	070: Ⓕ	083: Ϻ	096: ˇ	109: m	122: z		
006: ◆	019: !!	032: :	045: -	058: :	071: Ⓖ	084: Ⓣ	097: a	110: n	123: {		
007: •	020: ¶	033: !	046: .	059: :	072: Ⓗ	085: 𝑼	098: b	111: o	124: :		
008: •	021: §	034: "	047: /	060: <	073: I	086: 𝑽	099: c	112: p	125: }		
009: ◦	022: ▬	035: #	048: 0	061: =	074: J	087: ₩	100: d	113: q	126: ~		
010: ◦	023: ‡	036: \$	049: 1	062: >	075: K	088: X	101: e	114: r	127: △		
011: δ	024: ↑	037: ٪	050: 2	063: ?	076: L	089: Ƴ	102: f	115: s			
012: ♀	025: ↓	038: ؑ	051: ۳	064: ۰	077: M	090: Z	103: g	116: t			
High Ascii											
128: Ҫ	141: ି	154: Ӯ	167: ܂	180: ܃	193: ܄	206: ܅	219: ܆	232: ܇	245: ܈		
129: ӻ	142: Ӯ	155: ܉	168: ܊	181: ܋	194: ܌	207: ܍	220: ܎	233: ܏	246: ܐ		
130: ܀	143: ܁	156: ܂	169: ܃	182: ܄	195: ܅	208: ܆	221: ܇	234: ܈	247: ܔ		
131: ܁	144: ܁	157: ܂	170: ܃	183: ܄	196: ܅	209: ܆	222: ܇	235: ܈	248: ܔ		
132: ܁	145: ܁	158: ܁	171: ܁	184: ܁	197: ܁	210: ܁	223: ܁	236: ܁	249: ܁		
133: ܁	146: ܁	159: ܁	172: ܁	185: ܁	198: ܁	211: ܁	224: ܁	237: ܁	250: ܁		
134: ܁	147: ܁	160: ܁	173: ܁	186: ܁	199: ܁	212: ܁	225: ܁	238: ܁	251: ܁		
135: ܁	148: ܁	161: ܁	174: ܁	187: ܁	200: ܁	213: ܁	226: ܁	239: ܁	252: ܁		
136: ܁	149: ܁	162: ܁	175: ܁	188: ܁	201: ܁	214: ܁	227: ܁	240: ܁	253: ܁		
137: ܁	150: ܁	163: ܁	176: ܁	189: ܁	202: ܁	215: ܁	228: ܁	241: ܁	254: ܁		
138: ܁	151: ܁	164: ܁	177: ܁	190: ܁	203: ܁	216: ܁	229: ܁	242: ܁	255: ܁		
139: ܁	152: ܁	165: ܁	178: ܁	191: ܁	204: ܁	217: ܁	230: ܁	243: ܁			
140: ܁	153: ܁	166: ܁	179: ܁	192: ܁	205: =	218: ܁	231: ܁	244: ܁			

Problem: No interoperability

Common character encodings [edit]

- ISO 646
 - ASCII
- EBCDIC
- ISO 8859:
 - ISO 8859-1 Western Europe
 - ISO 8859-2 Western and Central Europe
 - ISO 8859-3 Western Europe and South European (Turkish, Maltese plus Esperanto)
 - ISO 8859-4 Western Europe and Baltic countries (Lithuania, Estonia, Latvia and Lapp)
 - ISO 8859-5 Cyrillic alphabet
 - ISO 8859-6 Arabic
 - ISO 8859-7 Greek
 - ISO 8859-8 Hebrew
 - ISO 8859-9 Western Europe with amended Turkish character set
 - ISO 8859-10 Western Europe with rationalised character set for Nordic languages, including complete Icelandic set
 - ISO 8859-11 Thai
 - ISO 8859-13 Baltic languages plus Polish
 - ISO 8859-14 Celtic languages (Irish Gaelic, Scottish, Welsh)
 - ISO 8859-15 Added the Euro sign and other rationalisations to ISO 8859-1
 - ISO 8859-16 Central, Eastern and Southern European languages (Albanian, Bosnian, Croatian, Hungarian, Polish, Romanian, Serbian and Slovenian, but also French, German, Italian and Irish Gaelic)
 - CP437, CP720, CP737, CP850, CP852, CP855, CP857, CP858, CP860, CP861,
- Mac OS Roman
- KOI8-R, KOI8-U, KOI7
- MIK
- ISCII
- TSCII
- VISCII
- JIS X 0208 is a widely deployed standard for Japanese character encoding that has several encoding forms.
 - Shift JIS (Microsoft [Code page 932](#) is a dialect of Shift_JIS)
 - EUC-JP
 - ISO-2022-JP
- JIS X 0213 is an extended version of JIS X 0208.
 - Shift_JIS-2004
 - EUC-JIS-2004
 - ISO-2022-JP-2004
- Chinese Guobiao
 - GB 2312
 - GBK (Microsoft [Code page 936](#))
 - GB 18030
- Taiwan Big5 (a more famous variant is Microsoft [Code page 950](#))
 - Hong Kong HKSCS
- Korean
 - KS X 1001 is a Korean double-byte character encoding standard
 - EUC-KR
 - ISO-2022-KR

UTF-8 vs utf-8 vs utf8

- UTF-8 is correct
- utf-8 is okay in case-insensitive situations
- utf8 works in most places
- Including Node, HTTP and HTML

Requirements

1. English text must be one byte per character
2. Can't have 8 zeros in a row (NULL terminator)
3. Backwards compatible with ASCII
4. Easily traversable/navigable

First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4	Usable Bits	Max Possible Characters
U+0000	U+007F	0xxxxxxxx				7	128
U+0080	U+07FF	110xxxxx	10xxxxxxxx			11	1,920
U+0800	U+FFFF	1110xxxx	10xxxxxxxx	10xxxxxxxx		16	61,440
U+10000	U+1FFFFFF	11110xxx	10xxxxxxxx	10xxxxxxxx	10xxxxxxxx	21	1,048,576

Will theoretically work to 6 bytes

First code point	Last code point	Byte 1	Byte 2	Byte 3	Byte 4	Byte 5	Byte 6	Usable Bits	Max Possible Characters
U+0000	U+007F	0xxxxxxxx						7	128
U+0080	U+07FF	110xxxxx	10xxxxxx					11	1,920
U+0800	U+FFFF	1110xxxx	10xxxxxx	10xxxxxx				16	61,440
U+10000	U+1FFFFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx			21	1,048,576
U+200000	U+3FFFFFF	111110xx	10xxxxxx	10xxxxxx	10xxxxxx	10xxxxxx		26	67,108,864
U+4000000	U+7FFFFFFF	1111110x	10xxxxxx	10xxxxxx	10xxxxxx	10xxxxxx	10xxxxxx	31	2,147,483,648

Navigation in a string

Starts with a zero: 1 byte ASCII character = a U+0061 or 97

Starts with 3 ones: 3 bytes = ✈️ U+2708 or 14,851,208

Starts with 4 ones: 4 bytes = 🚶 U+1F680 or 4,036,991,616

01100001	11100010	10011100	10001000	11101111	10111000
10001111	11110000	10011111	10011010	10000000	11110000
10011111	10010000	10111000	11000011	10100110	00001010

Growth and evolution over time

Year	Version	Total Assigned	Emoji Count	Max	Bytes
1991	v1.0	12,795	0	3	
1993	v1.1	40,635	0	3	
1996	v2.0	178,500	0	4	
1998	V2.1	178,502	0	4	
1999	v3.0	188,809	0	4	
2001	v3.1	233,787	0	4	
2003	v4.0	236,029	0	4	
2006	v5.0	238,671	0	4	
2010	v6.0	249,031	0	4	
2014	v7.0	252,603	0	4	
2015	v8.0	260,319	0	4	
2016	v9.0	267,819	1,788	4	
2017	v10.0	276,337	2,666	4	
2018	v11.0	277,021	2,789	4	
2019	v12.0	277,575	3,019	4	
2020	v13.0	283,506	3,304	4	
2021	v14.0	284,344	3,633	4	
2022	v15.0	288,833	3,664	4	

00	01	02	03	04	05	06	07	08	09	0A	0B	0C	0D	0E	0F
10	11	12	13	14	15	16	17	18	19	1A	1B	1C	1D	1E	1F
20	21	22	23	24	25	26	27	28	29	2A	2B	2C	2D	2E	2F
30	31	32	33	34	35	36	37	38	39	3A	3B	3C	3D	3E	3F
40	41	42	43	44	45	46	47	48	49	4A	4B	4C	4D	4E	4F
50	51	52	53	54	55	56	57	58	59	5A	5B	5C	5D	5E	5F
60	61	62	63	64	65	66	67	68	69	6A	6B	6C	6D	6E	6F
70	71	72	73	74	75	76	77	78	79	7A	7B	7C	7D	7E	7F
80	81	82	83	84	85	86	87	88	89	8A	8B	8C	8D	8E	8F
90	91	92	93	94	95	96	97	98	99	9A	9B	9C	9D	9E	9F
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

- Latin script
- Non-Latin European scripts
- African scripts
- Middle Eastern and Southwest Asian scripts
- South and Central Asian scripts
- Southeast Asian scripts
- East Asian scripts
- CJK characters
- Indonesian and Oceanic scripts
- American scripts
- Notational systems
- Symbols
- Private use
- UTF-16 surrogates

As of Unicode 15.0

Mojibake



ウィキペディア
フリー百科事典

āñjā, ñäf³äfšäf¹ä.,
 ä, äf³ÿäf¹äf²äf³äf¹, ïäf²
 äf¹ äf¹ä, ïäf²«
 æœéè, 'ä @äf²æ ¥ä²°,
 æ-ä -ä „äfšäf¹ä.,
 æœéè, 'ä @æ'æ-°
 ä Sä %ä <ä> ðéïçä°
 ç'ç'ç'ç'äfšäf¹ä.,
 ä, cäffäf -äfäf¹äf%
 (ä, ä, ä, -
 äfjäf²äf³äf¹, ä, cäf²äf¹
 f²ä, °)

$\tilde{a}f\tilde{a}f \llcorner \tilde{a}f -$

$\tilde{a}f\tilde{a}f\llcorner\tilde{a}f-$

ä°•æ^ç«-

a ScYYa, %

af a,a ®a

a „a”
 $\approx 1 \tilde{a}$ $\Omega \tilde{a}$

āfšāf†ā, £ā

$\tilde{a} \rightarrow$

ãf..ãf¼ãf«

ãf^aãf³ã -å f

æ-‡ å—åŒ-ã

å†°å...: ãf•ãf^aãf^{1/4}C^{TM3/4}ç§‘ä°å...ãEžã,ïã,£ã,ãfšãf†ã,£ã,ç†/4Wikipediaï/4%å€

W ă, ā; řă, āf'săf ſă, řă, că ſă ⑧æ-ťă - åŒ-ă 'ă «ă oă „ă ;ă -ā€ Help:c%1æ@Sæ-ťă -ă, 'ă "ě; ſă ă ă •ă „ă€

æ-‡å-åŒ-ã “!/4~ã,ã ~ã °ã “!1/4%ã “ã -æf ã,³ãf³ãf”ãf¥ãf¹/4ã,ã §æ-‡å-ã,’è|çøºã T™ã,·éš·ã «ã€ æfã —ã è|çøºã ·ã,Œã “ã „ç %é±†ã ®ã “ã “ã€,

„**CE**æ-**þ**-**CE**-**a** ‘**a** „**a** †**e**’**e**’%**a** -**a** **a**,**ð****f****ð****a**’**f****ð****y****ð****f****1**/**4****a**,**ð****c****°****ð****c****f****a** **þ****ð****Y****ð****%****þ****a** -**a** **i****ð****f****ð****a**’**f****ð****a** **þ****a**’**æ**-**þ**-**a**,**ð****1**/**2****ð****c****”****a** -**a** **æ**-**ð****c****±****ð****c****%****a** **®****ð****a**’**f****ð****o****ð****a**’**f****ð****a**,**ð****c****f****“****ð****a**’**f****ð****a**,**ð****f****TM****ð****a**’**f****ð****a**’**æ****1**/**2****ð****c****”****ð****e**’**e**’**a****ž****a** «**a** **Š****a** „**a** **i****ð****o****2****ð****1**/**2****a**’**TM****a**,**ð****c****”****ð****e**’**a****ž****a** **CE****ð****a**’**æ** -**a** **æ** „**a** **ð****f****a** **Þ****Y****a** „**a** **ð****a**’**æ****1**/**2****ð****c****”****ð****Y****a** „**a** **ð****f****a**’**æ** -**a** **®**

ç>®æ-ŋ [é žèj̪çəŋ]

1 ä, »ä «åŽÝå»

- æ-‡å-ã, ³ãf¹äf‰ã @åø§ã ã .
 - æ-‡å-ã, ³ãf¹äf‰ã @é •ã „
 - ã, äf³ã, ³ãf¹äf‡ã, £ãf³ã, °
 - æ-‡å-äf•ã, ©ãf³ãfã @é •ã „
 - ç‰¹å@šã @æ©Ýèf¹æŒ‡å@š
 - æ-‡å-å‡ã, œ

Strange to undecipherable

Swedish example:		Smörgås (open sandwich)
File encoding	Setting in browser	Result
MS-DOS 437	ISO 8859-1	Sm"rḡs
ISO 8859-1	Mac Roman	Sm^rgÅs
UTF-8	ISO 8859-1	SmÃ¶rgÃ¥s
UTF-8	Mac Roman	Smvðrgv·s
UTF-8	ibm/cp037 (EBCDIC)	�_C��EÅCv�

Character encoding vs character set

- UTF-8 is the character encoding
- Unicode is the character set
- Unicode Consortium curates the character set
- UTF-* stands for Unicode Transform Protocol *

Character encoding vs character set

character	encoding	bits
A	UTF-8	01000001
A	UTF-16	00000000 01000001
A	UTF-32	00000000 00000000 00000000 01000001
あ	UTF-8	11100011 10000001 10000010
あ	UTF-16	00110000 01000010
あ	UTF-32	00000000 00000000 00110000 01000010

Character encoding vs character set

Character	Binary code point	Binary UTF-8	Hex UTF-8
\$	U+0024	010 0100 00100100	24
£	U+00A3	000 1010 0011 11000010 10100011	C2 A3
€	U+0939	0000 1001 0011 1001 11100000 10100100 10111001	E0 A4 B9
₩	U+20AC	0010 0000 1010 1100 11100010 10000010 10101100	E2 82 AC
한	U+D55C	1101 0101 0101 1100 11101101 10010101 10011100	ED 95 9C
⌚	U+10348	0 0001 0000 0011 0100 1000 11110000 10010000 10001101 10001000	F0 90 8D 88

Emoji modifiers

U+	1F385	1F474	1F44A	1F44D
base	🤔	RGBO	✊	👍
base + FITZ-1-2	🎅	RGBO	✊	👍
base + FITZ-3	🎅	RGBO	✊	👍
base + FITZ-4	🎅	RGBO	✊	👍
base + FITZ-5	🎅	RGBO	✊	👍
base + FITZ-6	🎅	RGBO	✊	👍

JavaScript strings: an Array of code points

	A	𠂇	好	不
Code point	U+0041	U+05D0	U+597D	U+233B4
UTF-8	41	D7 90	E5 A5 BD	F0 A3 8E B4
UTF-16	00 41	05 D0	59 7D	D8 4C DF B4
UTF-32	00 00 00 41	00 00 05 D0	00 00 59 7D	00 02 33 B4

Parts of the stack

- Content-Type header to serve document
- Browser chooses which encoding to use
- JavaScript strings treated correctly
- Content-Type from fetch() to server (defaults to document)
- Back-end server & language
- Connection to database
- Database column type

JSON encoding differences - / /\n

- Browser - "/ /\n"
- Node - "/ /\n"
- Go - "/ /\n"
- Python - "\ud83d\ude80\n"
- PHP - "\ud83d\ude80\n" OR "\v \v\n"
- Java, C++ - depends

Databases

- MySQL 5.6 & AWS Aurora require specifying “utf8mb4”
- Everything else handles UTF-8 by default

String literals: encoding code points

```
const airplane = "\u2708";  
const airplane = "\u{2708}";  
const airplane = "✈";
```

```
const rocket = "\u{1F680}";
```

HTML entities code points

- Airplane: ✈
- Rocket: 🚀
- `parselnt("1F680", 16) => 128640 => 🚀`

String methods are problematic

```
> "🚀".length;  
< 2  
  
> "🚀-rocket".charCodeAt(1);  
< 56960  
  
> "🚀-rocket".codePointAt(1);  
< 56960  
  
> "🚀-rocket"[1];  
< '\uDE80'
```

Slicing strings

```
> "🚀-rocket".slice(1);  
< '\uDE80-rocket'
```

```
> [..."🚀-rocket"].slice(1).join("");  
< '-rocket'
```

```
> Array.from("🚀-rocket").slice(1).join("");  
< '-rocket'
```

BOM (Byte order mark)

- UTF-8/UTF-16 BOM is \uFEFF
- It has no meaning and should be ignored
- You'll want to strip it

```
> '\uFEFFabc'.length
< 4
> '\uFEFFabc'.slice(0, 1);
< ''
```

Hello npm?

```
1 function length(str: string) : number {
2     return [...str].length;
3 }
4
5 function slice(str: string, from: number = 0, to: number = undefined) : string {
6     return [...str].slice(from, to).join('');
7 }
8
9 function codePointAt(str: string, at: number = 0) : number {
10    if (at < 0) {
11        return undefined;
12    }
13    for (const char of str) {
14        if (at-- === 0) {
15            return char.codePointAt(0);
16        }
17    }
18    return undefined;
19 }
20
21 function charAt(str: string, at: number = 0): string {
22    if (at < 0) {
23        return '';
24    }
25    for (const char of str) {
26        if (at-- === 0) {
27            return char;
28        }
29    }
30    return '';
31 }
32
33 module.exports = function (string, start, end) {
34     var realStart = toNumber(start, 0);
35     var realEnd = toNumber(end, string.length);
36     if (realEnd === realStart) {
37         return '';
38     } else if (realEnd > realStart) {
39         return slice(string, realStart, realEnd);
40     } else {
41         return slice(string, realEnd, realStart);
42     }
43 }
```

```
1 function length(str: string) : number {
2     return [...str].length;
3 }
4
5 function slice(str: string, from: number = 0, to: number = undefined) : string {
6     return [...str].slice(from, to).join('');
7 }
8
9 function codePointAt(str: string, at: number = 0) : number {
10    if (at < 0) {
11        return undefined;
12    }
13    for (const char of str) {
14        if (at-- === 0) {
15            return char.codePointAt(0);
16        }
17    }
18    return undefined;
19 }
20
21 function charAt(str: string, at: number = 0): string {
22    if (at < 0) {
23        return '';
24    }
25    for (const char of str) {
26        if (at-- === 0) {
27            return char;
28        }
29    }
30    return '';
31 }
32
33 module.exports = function (string, start, end) {
34     var realStart = toNumber(start, 0);
35     var realEnd = toNumber(end, string.length);
36     if (realEnd === realStart) {
37         return '';
38     } else if (realEnd > realStart) {
39         return slice(string, realStart, realEnd);
40     } else {
41         return slice(string, realEnd, realStart);
42     }
43 }
```

So I did it

- <https://www.npmjs.com/package/multibyte>

The screenshot shows the npmjs.com package page for 'multibyte'. The page has a header with the package name 'multibyte' and its version '1.0.0-beta.1'. It includes tabs for 'Readme' (selected), 'Code' (Beta), 'Dependencies' (0), 'Dependents' (0), and 'Versions' (1). Below the tabs, there's a brief description: 'multibyte provides common string functions that respect multibyte Unicode characters.' A command line input field shows the command 'npm i multibyte'. To the right, there's a sidebar with repository information ('github.com/kensnyder/multibyte'), a 'Homepage' link ('github.com/kensnyder/multibyte#readme'), and various statistics: Version 1.0.0-beta.1, License ISC, Unpacked Size 30.8 kB, Total Files 55, Issues 0, Pull Requests 0, and a 'Last publish' time of 2 minutes ago. At the bottom, there's a 'Collaborators' section showing a profile picture and a 'Try on RunKit' button.

multibyte ts

1.0.0-beta.1 • Public • Published 2 minutes ago

Readme Code Beta Dependencies 0 Dependents 0 Versions 1

multibyte

npm v1.0.0-beta.1 build passing codecov 100% license package not found

multibyte provides common string functions that respect multibyte Unicode characters.

```
npm install multibyte
```

The problem and the solution

On one hand, JavaScript strings use UTF-16 encoding, and on the other hand, JavaScript strings behave like an Array of code points. Unicode characters that take more than 2 bytes (like newer emoji) get split into 2 code points in many situations.

If you display Unicode text from a UTF-8 source, you need these multibyte functions that take advantage of the fact that `Array.from()` is Unicode safe.

```
import {  
  charAt,  
  codePointAt,  
  length,  
  slice,  
  split,  
  truncateBytes,  
} from 'multibyte';  
  
// JavaScript String.prototype.charAt() is not Unicode aware  
'𠮷𠮷'.charAt(1); // "\ud83d" ✗  
charAt('𠮷𠮷', 1); // "𠮷" ✓
```

Install

```
> npm i multibyte
```

Repository

github.com/kensnyder/multibyte

Homepage

github.com/kensnyder/multibyte#readme

Version	License
1.0.0-beta.1	ISC

Unpacked Size	Total Files
30.8 kB	55

Issues	Pull Requests
0	0

Last publish

2 minutes ago

Collaborators

Try on RunKit

Intl.Collator

JavaScript Demo: Intl.Collator

```
1 console.log(['z', 'a', 'z', 'ä'].sort(new Intl.Collator('de').compare));
2 // Expected output: Array ["a", "ä", "z", "Z"]
3
4 console.log(['z', 'a', 'z', 'ä'].sort(new Intl.Collator('sv').compare));
5 // Expected output: Array ["a", "z", "Z", "ä"]
6
7 console.log(['z', 'a', 'z', 'ä'].sort(
8   new Intl.Collator('de', { caseFirst: 'upper' } ).compare
9 ));
10 // Expected output: Array ["a", "ä", "Z", "z"]
11 |
```

Intl.Segmenter

JavaScript Demo: Intl.Segmenter

```
1 const segmenterFr = new Intl.Segmenter('fr', { granularity: 'word' });
2 const string1 = 'Que ma joie demeure';
3
4 const iterator1 = segmenterFr.segment(string1)[Symbol.iterator]();
5
6 console.log(iterator1.next().value.segment);
7 // Expected output: 'Que'
8
9 console.log(iterator1.next().value.segment);
10 // Expected output: ' '
11
```

Emoji Implementation



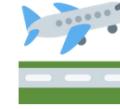
apple



google



facebook



twitter



mozilla



microsoft



samsung



lg

Emoji combinations with ZWJ

$$\text{🏳️} + \text{🌈} = \text{🏳️‍🌈}$$

$$\text{👁} + \text{🗨} = \text{👁🗨}$$

$$\text{👯} + \text{♂} = \text{👯♂}$$

$$\text{👱} + \text{🟤} + \text{🏭} = \text{👷}$$

$$\text{👨} + \text{👨} + \text{👱} = \text{👨👨👱}$$

$$\text{👱} + \text{❤️} + \text{💋} + \text{👱} = \text{💏}$$

Obligatory Zalgo Text

