

Stefan Daschek / @noniq

Encodings

When "ü" != "ü" and other oddities

Past

1963

"American Standard Code for Information Interchange"

In the Beginning, there was ASCII

	_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_	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3_	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4_	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5_	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6_	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7_	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

In the Beginning, there was ASCII

	_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_	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/
3_	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4_	@	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
5_	P	Q	R	S	T	U	V	W	X	Y	Z	[\]	^	_
6_	`	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o
7_	p	q	r	s	t	u	v	w	x	y	z	{		}	~	DEL

0x41 = A

In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX										FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"										,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

0b 10 00001 = A

In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX										FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"										,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

0b 10 00001 = A

0b 11 00001 = a

In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

control characters

In the Beginning, there was ASCII

	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111	1000	1001	1010	1011	1100	1101	1110	1111	1000	1001	1010	1011	1100	1101	1110	1111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

BS
backspace

In the Beginning, there was ASCII

	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111	1000	1001	1010	1011	1100	1101	1110	1111	1000	1001	1010	1011	1100	1101	1110	1111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

BS
backspace

▶ ^H = backspace

In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111			
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

BEL
bell

- ▶ ^H = backspace
- ▶ ^G = bell

In the Beginning, there was ASCII

	0000	0001	0010	0011	00100	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111						
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

EOT
end of text

- ▶ ^H = backspace
- ▶ ^G = bell
- ▶ ^D = end of text

In the Beginning, there was ASCII

	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111	1000	1001	1010	1011	1100	1101	1110	1111	1100	1101	1110	1111	1100	1101	1110	1111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

ESC

- ▶ ^H = backspace
- ▶ ^G = bell
- ▶ ^D = end of text (file)
- ▶ ^[= escape

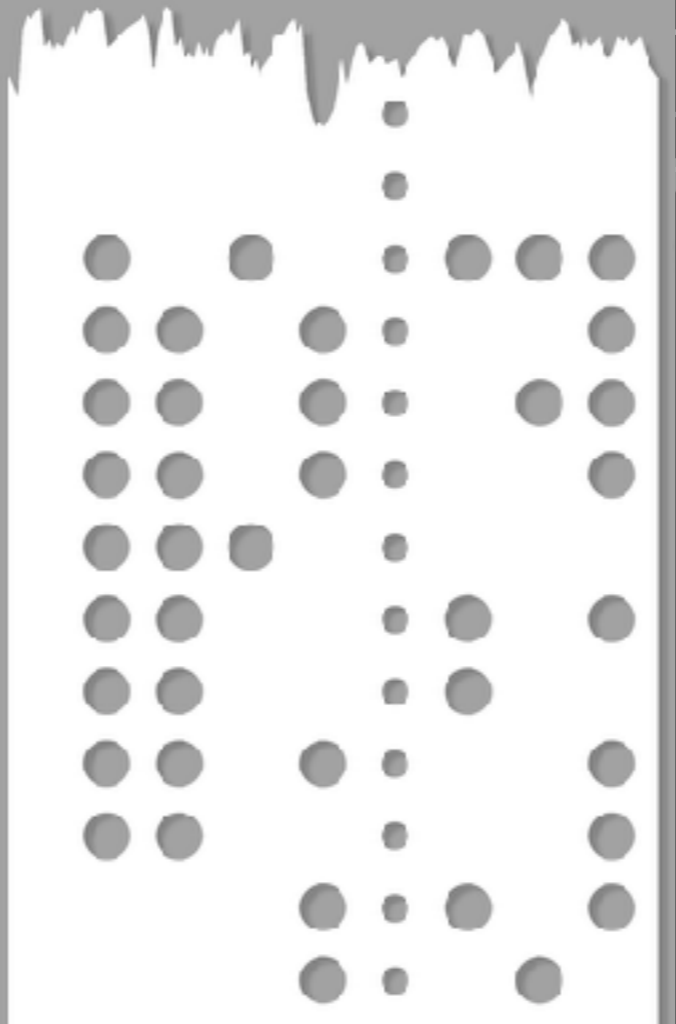
In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	DEL (11111111)				?			
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	DEL							?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

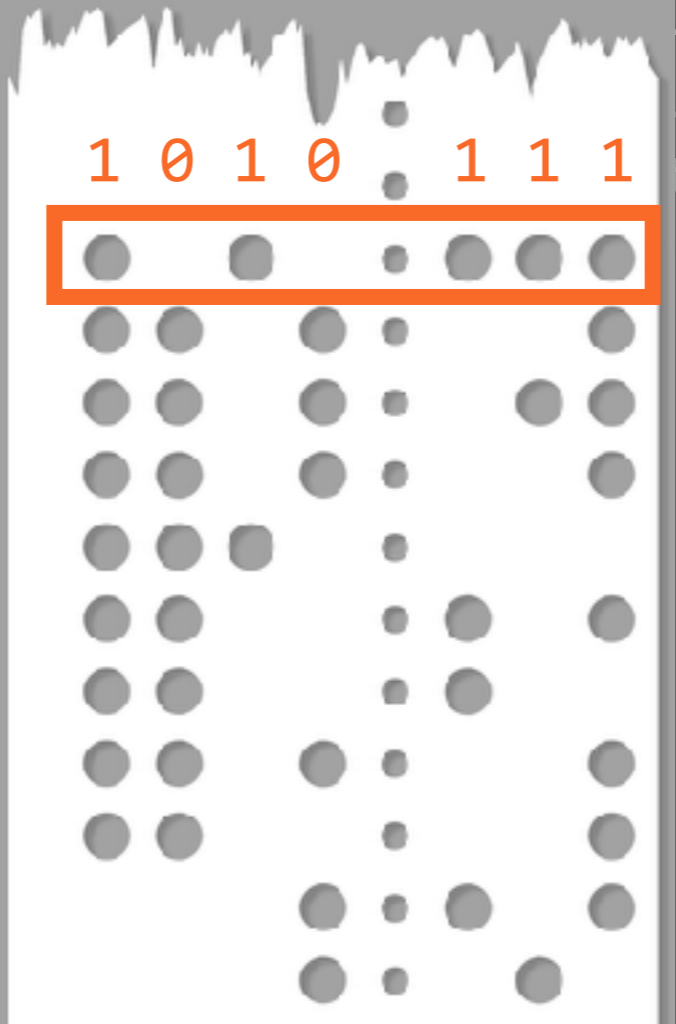
DEL
(1111111)



In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111			
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US			
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7										?	
10	@	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	[\]	^	_			
11	`	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	{		}	~				DEL

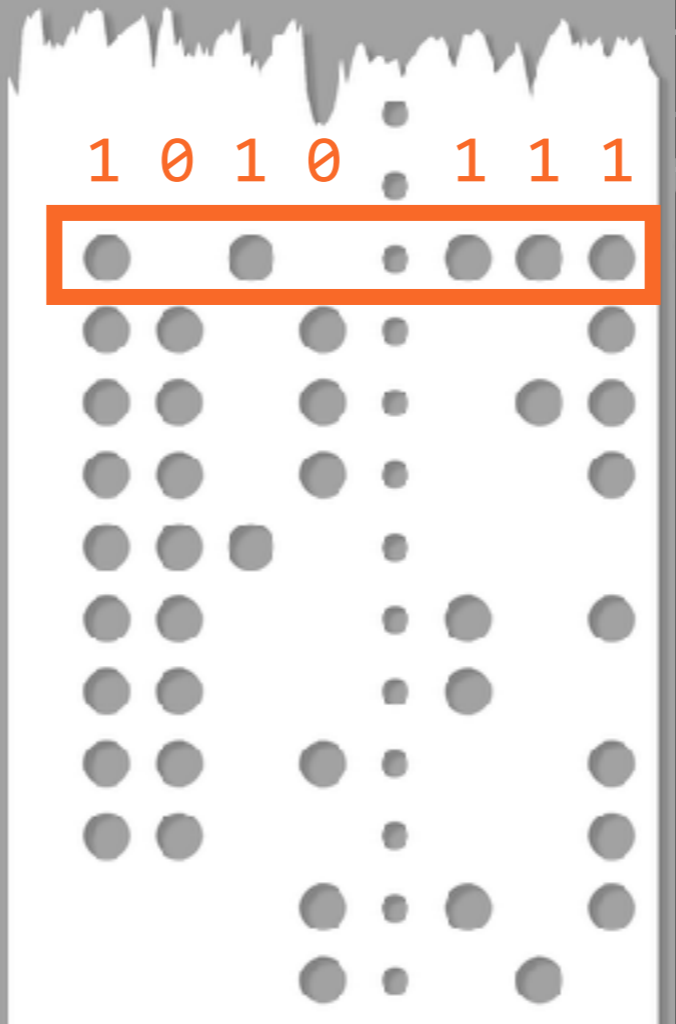
DEL
(1111111)



In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7								
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

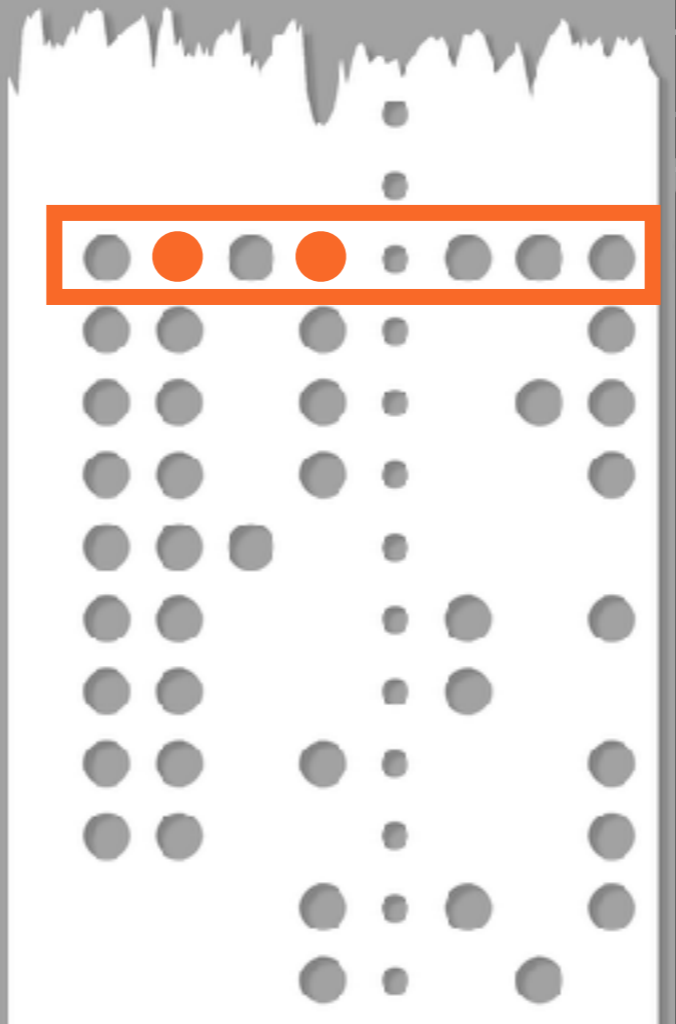
DEL
(1111111)



In the Beginning, there was ASCII

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7								
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

DEL
(11111111)



But ASCII wasn't good enough

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

But ASCII wasn't good enough

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

Mötley Crüe

Motörhead

But ASCII wasn't good enough

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
01	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

Mötley Crüe

Motörhead

äÜßØáçÆ§£«»©

aka "Latin-1"

So we got ISO-8859-1 ...

	0000	0001	0010	0011	0100	0101	0110	0111	0100	0101	0110	0111	0110	0111	0111	1000	1001	1010	1011	1010	1011	1011	1011	1100	1101	1101	1101	1110	1111	1111	1111	
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
010	@	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	[\]	^	_
011	`	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	{		}	~	
100	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA	HTS	HTJ	VTS	PLD	PLU	RI	SS2	SS3	DCS	PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
101	NBSP	ı	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	-	®	¯	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
110	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
111	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

So we got ISO-8859-1 ...

	0000	0001	0010	0011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
010	@	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	[\]	^	_
011	`	a	b															q	r	s	t	u	v	w	x	y	z	{		}	~	
100	PAD	HOP	BPH															PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
101	NBSP	ı	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	-	®	¯	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
110	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
111	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

currency sign
 0b 101 00100 = 0xA4 = 164

So we got ISO-8859-1, but ...

	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
010	@	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	[\]	^	_
011	`	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	{		}	~	
100	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA	HTS	HTJ	VTS	PLD	PLU	RI	SS2	SS3	DCS	PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
101	NBSP	ı	¢	£	¤	¥	¦	§	¨	©	ª	«	¬	-	®	¯	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
110	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
111	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

So we got ISO-8859-1, but ...

	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111	1000	1001	1010	1011	1100	1101	1110	1111	1100	1101	1110	1111				
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
010	@	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	[\]	^	_
011	`	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	{		}	~	
100	PAD	HOP	BPH	NBH	IND	NEL	SSA	ESA																	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
101	NBSP	ı	ç	£	¤	¥	¦	§																	¸	¹	º	»	¼	½	¾	¿
110	À	Á	Â	Ã	Ä	Å	Æ	Ç																	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
111	à	á	â	ã	ä	å	æ	ç																	ø	ù	ú	û	ü	ý	þ	ÿ

No €!



(because it was created in 1987)

aka "Latin-9"

... thus ISO-8859-15 was created

	0000	0001	0010	0011	0100	0101	0110	0111	0100	0101	0110	0111	0110	0111	0111	1000	1001	1010	1011	1010	1011	1011	1011	1100	1101	1101	1101	1110	1111	1111	1111	
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
010	@	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	[\]	^	_
011	`	a	b															q	r	s	t	u	v	w	x	y	z	{		}	~	
100	PAD	HOP	BPH															PU1	PU2	STS	CCH	MW	SPA	EPA	SOS	SGCI	SCI	CSI	ST	OSC	PM	APC
101	NBSP	ı	ç	£	€	¥	Š	§	š	©	ª	«	¬	-	®	-	°	±	²	³	Ž	μ	¶	·	ž	¹	º	»	Œ	œ	ÿ	ı
110	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
111	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

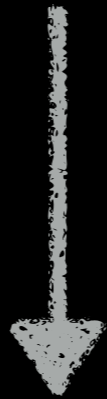
0b101 00100 = 0xA4 = 164
in ISO-8859-1: ¤

Metaphor: Translating Text

"GIFT"

Metaphor: Translating Text

"GIFT"



"CADEAU"

Metaphor: Translating Text

"GIFT"



"CADEAU"

"POISON"

How to determine the right encoding

- ▶ Metadata (e.g. charset header)
- ▶ Standards (e.g. UTF-8 for XML, Latin-1 for HTTP headers)
- ▶ Context
- ▶ Best guess `_(ツ)_/`

More encodings: Windows-1252 ...

	0000	0001	0010	0011	0100	0101	0110	0111	0100	0101	0110	0111	0110	0111	0111	1000	1001	1010	1011	1010	1011	1011	1011	1100	1101	1101	1110	1111	1111	1111		
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
010	@	^	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	€	
011	`	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	~	
100	€	,	f	"	...	†	‡	^	%	Š	‹	Œ	Ž																			
101	NBSP	ı	ç	£	¤	¥	¦	§	¨	©	ª	«	¬	-	®	¯	°	±	²	³	´	µ	¶	·	¸	¹	º	»	¼	½	¾	¿
110	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	Í	Î	Ï	Ð	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
111	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï	ð	ñ	ò	ó	ô	õ	ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

0b100 00000 = 0x80 = 128

... and Mac-Roman

	0000	0001	0010	0011	0100	0101	0110	0111	0100	0101	0110	0111	0110	0111	0111	0111	1000	1001	1010	1011	1010	1011	1011	1011	1100	1101	1101	1101	1110	1111	1110	1111	1111	1111
000	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US		
001	SP	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?		
010	@	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	[\]	^	_		
011	`	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	{		}	~			
100	Ä	Å	Ç	É	Ñ	Ö	Ü	á	à	â	ä	ã	å	ç	é	è	ê	ë	í	ì	î	ï	ñ	ó	ò	ô	ö	õ	ú	ù	û	ü		
101	†	°	¢	£	§	•	¶	β	®	©	™	'	¨	≠	Æ	Ø	∞	±	≤	≥	¥	μ	∂	Σ	Π	π	∫	ª	º	Ω	æ	ø		
110	¿	¡	¬	√	f	≈	Δ	«	»	...		À	Ã	Õ	Œ	œ	-	-	"	"			÷	◊	ÿ	ÿ	/	œ	<	>	fi	fl		
111	‡	·		„	%	Â	Ê	Á	Ë	È	Í	Î	Ï	Ì	Ó	Ô			Ò	Ú	Û	Ü	ı	ˆ	˜	-	˘	·	°	˙	˚	˛	ˇ	

Present

We need more characters!

We need more characters!

- ▶ **ASCII:** 128 characters

We need more characters!

- ▶ **ASCII:** 128 characters
- ▶ **ISO-8859-1 & Co:** 256 characters

We need more characters!

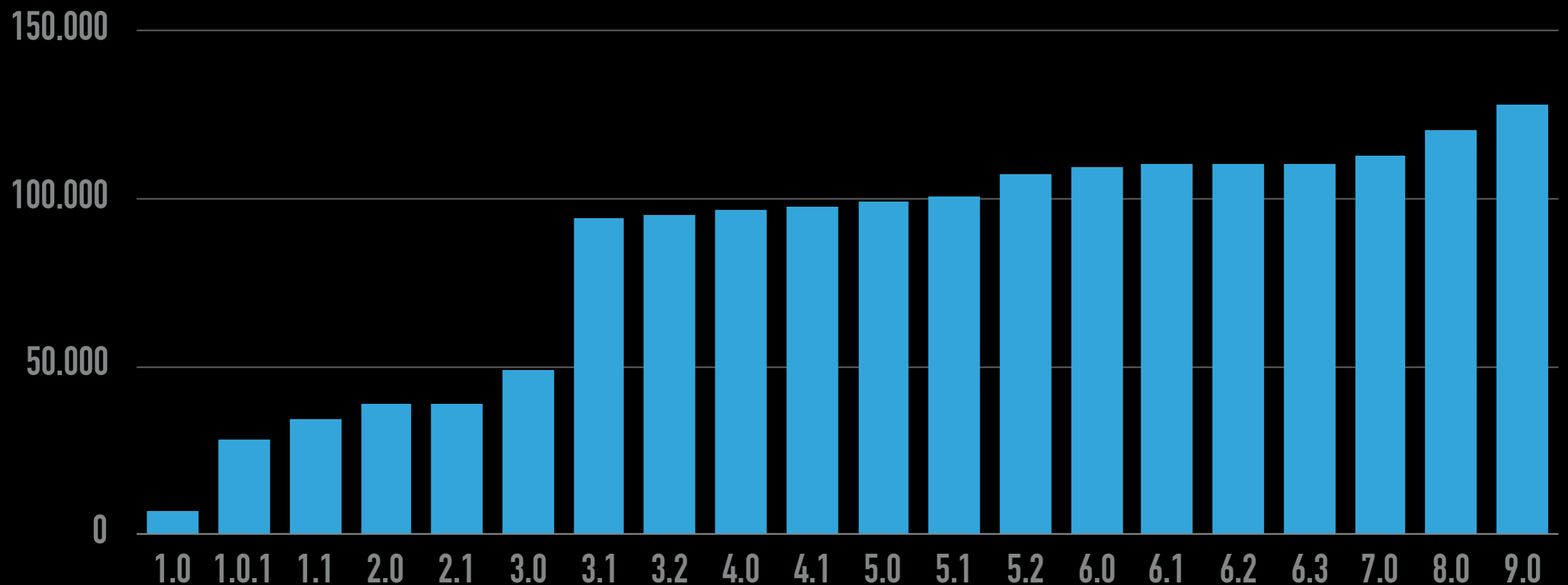
- ▶ **ASCII:** 128 characters
- ▶ **ISO-8859-1 & Co:** 256 characters
- ▶ **Unicode:** ???

Unicode

- ▶ **128,237** characters

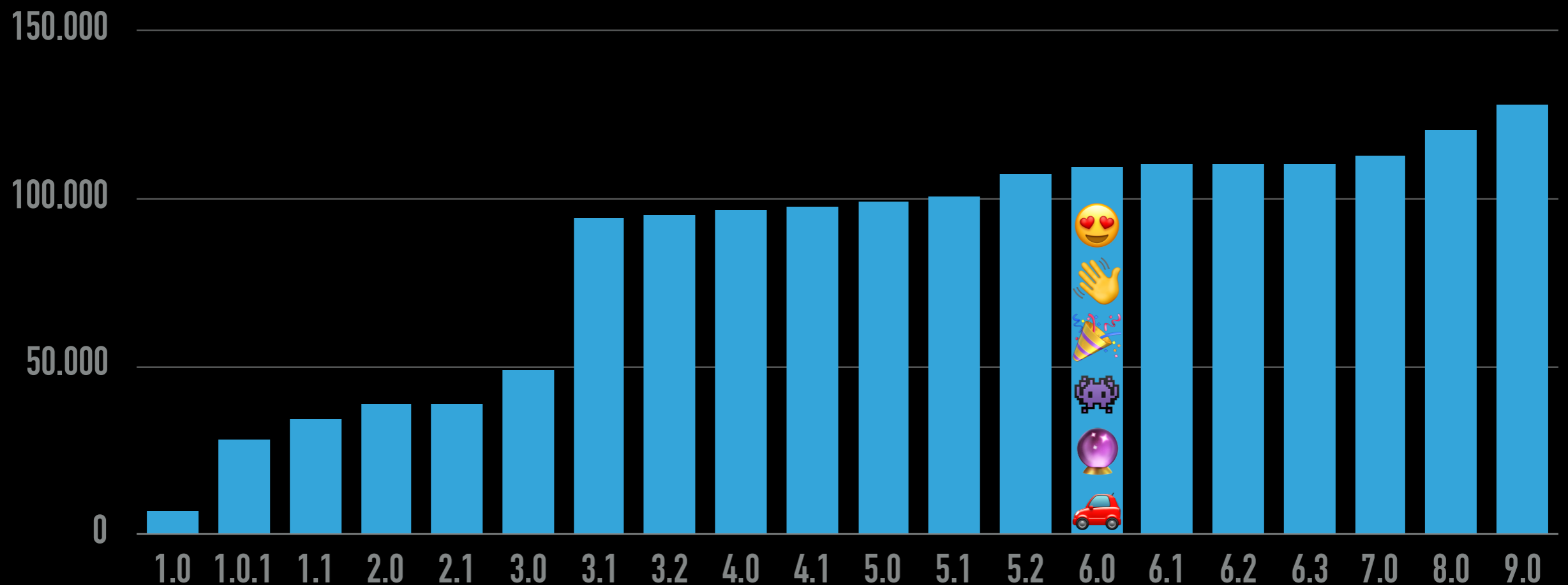
Unicode: Proudly adding characters since 1991

▶ **128,237** characters



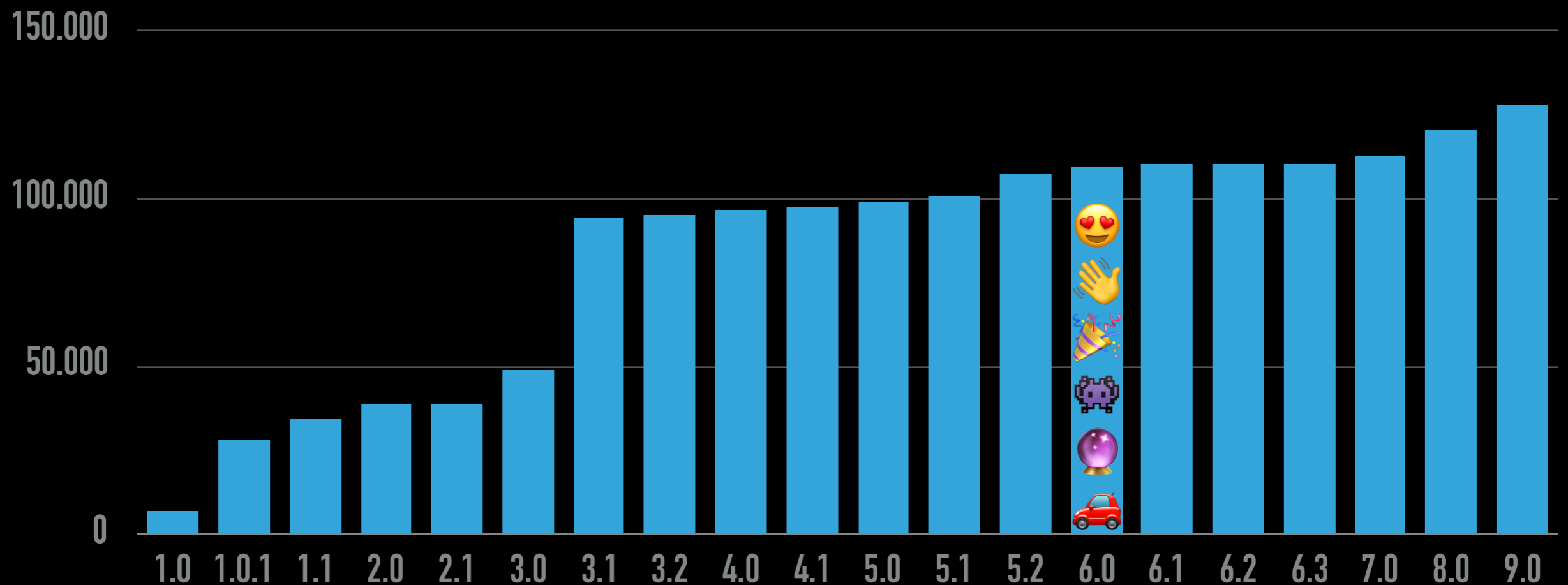
Unicode: Proudly adding characters since 1991

▶ **128,237** characters



Unicode: Proudly adding characters since 1991

▶ **128,237 characters**



▶ **1,114,112 code points: U+0000 to U+10FFFF**

Unicode Encodings

U+0000 to U+10FFFF (21 bits)

Unicode Encodings

U+0000 to U+10FFFF (21 bits)

- ▶ **UTF-32:** 4 bytes (32 bits) per codepoint

Unicode Encodings

U+0000 to U+10FFFF (21 bits)

- ▶ **UTF-32:** 4 bytes (32 bits) per codepoint
- ▶ **UTF-16:** 2 or 4 bytes per codepoint

Unicode Encodings

U+0000 to U+10FFFF (21 bits)

- ▶ **UTF-32:** 4 bytes (32 bits) per codepoint
- ▶ **UTF-16:** 2 or 4 bytes per codepoint
- ▶ **UTF-8:** 1-4 bytes per codepoint

Unicode Encodings

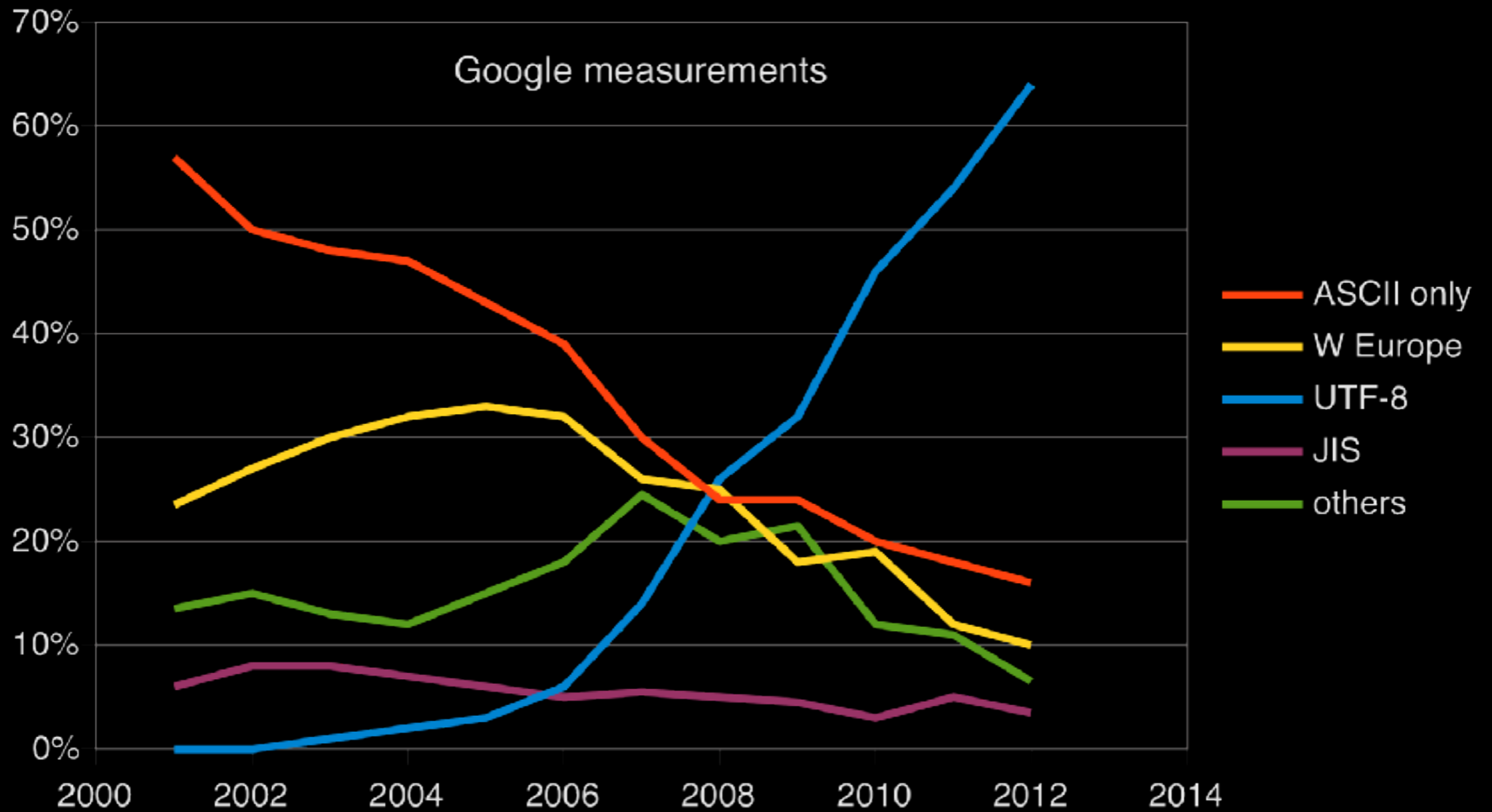
U+0000 to U+10FFFF (21 bits)

- ▶ **UTF-32:** 4 bytes (32 bits) per codepoint
- ▶ **UTF-16:** 2 or 4 bytes per codepoint
- ▶ **UTF-8:** 1-4 bytes per codepoint

	Bytes / Codepoint	Fixed Width	File size compared to ASCII
UTF-32	4	✓	4 times larger
UTF-16	2 or 4	-	double
UTF-8	1-4	-	almost identical

UTF-8 rules the ~~World~~ Web

Share of web pages with different encodings



UTF-8: How does it even work?

UTF-8: 1-4 bytes per codepoint



UTF-8: How does it even work?

UTF-8: 1-4 bytes per codepoint

Codepoint	Byte 1	Byte 2	Byte 3	Byte 4
U+0000 to U+007F	0xxx xxxx	–	–	–

UTF-8: How does it even work?

UTF-8: 1-4 bytes per codepoint

Codepoint	Byte 1	Byte 2	Byte 3	Byte 4
U+0000 to U+007F	0xxx xxxx	–	–	–
U+0080 to U+07FF	110x xxxx	10xx xxxx	–	–

UTF-8: How does it even work?

UTF-8: 1-4 bytes per codepoint

Codepoint	Byte 1	Byte 2	Byte 3	Byte 4
U+0000 to U+007F	0xxx xxxx	–	–	–
U+0080 to U+07FF	110x xxxx	10xx xxxx	–	–
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	–

UTF-8: How does it even work?

UTF-8: 1-4 bytes per codepoint

Codepoint	Byte 1	Byte 2	Byte 3	Byte 4
U+0000 to U+007F	0xxx xxxx	–	–	–
U+0080 to U+07FF	110x xxxx	10xx xxxx	–	–
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	–
U+10000 to U+10FFFF	1111 0xxx	10xx xxxx	10xx xxxx	10xx xxxx

UTF-8: How does it even work?

UTF-8: 1-4 bytes per codepoint

Codepoint	Byte 1	Byte 2	Byte 3	Byte 4
U+0000 to U+007F	0xxx xxxx	-	-	-
U+0080 to U+07FF	110xxxxx	10xxxxxx	-	-
U+0800 to U+FFFF	1110xxxx	10xxxxxx	10xxxxxx	-
U+10000 to U+10FFFF	11110xxx	10xxxxxx	10xxxxxx	10xxxxxx

	00000	00001	00010	00011	00100	00101	00110	00111	01000	01001	01010	01011	01100	01101	01110	01111	10000	10001	10010	10011	10100	10101	10110	10111	11000	11001	11010	11011	11100	11101	11110	11111
00	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣	␣		
01	!	"	#	\$	%	&	'	()	*	+	,	-	.	/	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?	
10	@	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	[\]	^	_
11	`	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	{		}	~	DEL

UTF-8 basic characters = ASCII

UTF-8 done manually



U+20AC

0b 0010 0000 1010 1100

	Byte 1	Byte 2	Byte 3	Byte 4
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	

UTF-8 done manually



U+20AC

0b 0010 0000 1010 1100

	Byte 1	Byte 2	Byte 3	Byte 4
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	
	1100 0010			

UTF-8 done manually



U+20AC

0b 0010 0000 1010 1100

	Byte 1	Byte 2	Byte 3	Byte 4
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	
	1100 0010	1000 0010		

UTF-8 done manually



U+20AC

0b 0010 0000 1010 1100

	Byte 1	Byte 2	Byte 3	Byte 4
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	
	1100 0010	1000 0010	1010 1100	

UTF-8 done manually



U+20AC

0b 0010 0000 1010 1100

	Byte 1	Byte 2	Byte 3	Byte 4
U+0800 to U+FFFF	1110 xxxx	10xx xxxx	10xx xxxx	
	1100 0010	1000 0010	1010 1100	
	0xE2	0x82	0xAC	

Gotchas

~~Problems~~

Müller vs. Müller

Müller vs. Müller: Precomposed vs. combining characters

Character	M	ü		l	l	e	r
Code Point	U+004D	U+00FC		U+006C	U+006C	U+0065	U+0072
Code Point	U+004D	U+0075	U+0308	U+006C	U+006C	U+0065	U+0072
Character	M	u	ö	l	l	e	r

Combining
diaeresis

Müller vs. Müller: Precomposed vs. combining characters

Character	M	ü		l	l	e	r
Code Point	U+004D	U+00FC		U+006C	U+006C	U+0065	U+0072
Code Point	U+004D	U+0075	U+0308	U+006C	U+006C	U+0065	U+0072
Character	M	u	ö	l	l	e	r

Combining
diaeresis

Müller

Unicode Normalization

	Canonical Equivalence	Compatibility Equivalence
Composed	NFC	NFKC
Decomposed	NFD	NFKD

foo vs. foo

foo vs. foo: Invisible characters in Unicode

- ▶ U+200B = zero width space
- ▶ U+200C = zero width non-joiner
- ▶ U+200D = zero width joiner
- ▶ U+2062 = invisible times
- ▶ U+2063 = invisible separator
- ▶ U+2064 = invisible plus
- ▶ U+FEFF = zero width no-break space
- ▶ ...

foo vs. foo: Invisible characters in Unicode

- ▶ U+200B = zero width space
- ▶ U+200C = zero width non-joiner
- ▶ U+200D = zero width joiner
- ▶ U+2062 = invisible times
- ▶ U+2063 = invisible separator
- ▶ U+2064 = invisible plus
- ▶ U+FEFF = zero width no-break space
- ▶ ...

also BOM
(Byte Order Mark)

Family Splitting



Thank you!

Stefan Daschek / @noniq