

# XML en Unicode

# HTML - een voorbeeld

```
<HTML>
<HEAD>
<TITLE>Marc de Graauw</TITLE>
</HEAD>
<BODY>
<H1>Marc de Graauw</H1>
<P>Geslacht: Man</P>
<H2>Opleidingen</H2>
<OL>
<LI>VWO</LI>
<LI>kandidaats Biologie</LI>
<LI>doctoraal filosofie</LI>
</OL>
</BODY>
</HTML>
```

# XML - een voorbeeld

```
<?xml version="1.0" encoding="UTF-8"?>

<persoon>
  <persoonsgegevens geslacht="Man">
    <achternaam>Graauw</achternaam>
    <voorvoegsel>de</voorvoegsel>
    <voornaam>Marc</voornaam>
  </persoonsgegevens>
  <opleiding>VWO</opleiding>
  <opleiding>kandidaats Biologie</opleiding>
  <opleiding>doctoraal filosofie</opleiding>
</persoon>
```

# Vóór Unicode

bits	0110 1101	0110 0001	0111 0010	0110 0011
bytes dec	109	97	114	99
bytes hex	6d	61	72	63
ascii	m	a	r	c

- 7 bits, 128 tekens, 95 afdrukbaar
- Engels: cijfers, letters, leestekens
- niet: Frans: ê ç Duits: ß ä Grieks: ε Ω Nederlands: ë ï
- Spaans, Arabisch, Fins, Russisch, Chinees, Thais, etc. etc.

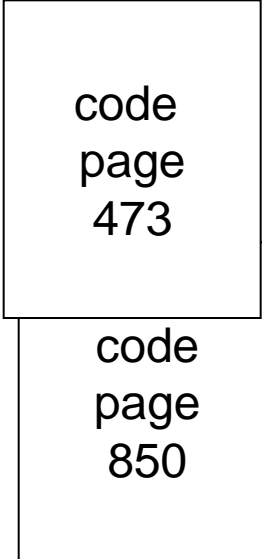
# Vóór Unicode

bytes 0 - 127



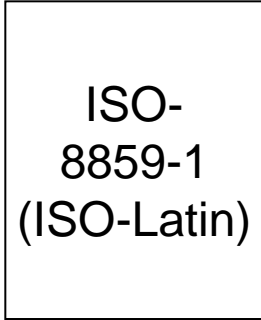
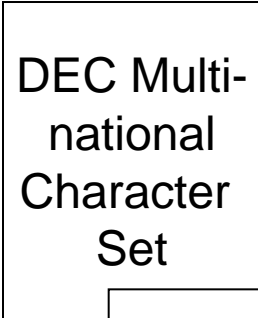
```
!"#$%&'()*+,-./
0123456789:;<=>?
@ABCDEFGHIJKLMNO
PQRSTUVWXYZ[\]^_
`abcdefghijklmno
pqrstuvwxyz{|}~
```

bytes 128 - 255

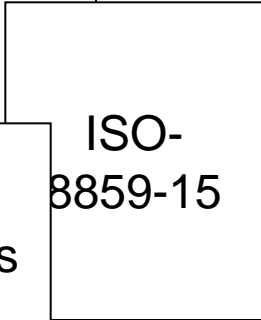


Engels,  
meeste Frans, Duits  
IBM PC

West  
Europees  
+ Á ß



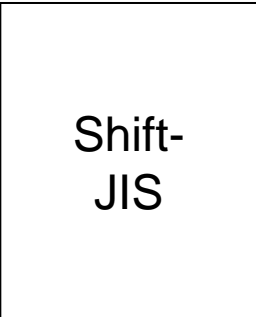
West-  
Europees



ISO-Latin  
+ 'IJ'sland  
+ 'œ'uf



ISO-Latin  
+ œ, €



Japans



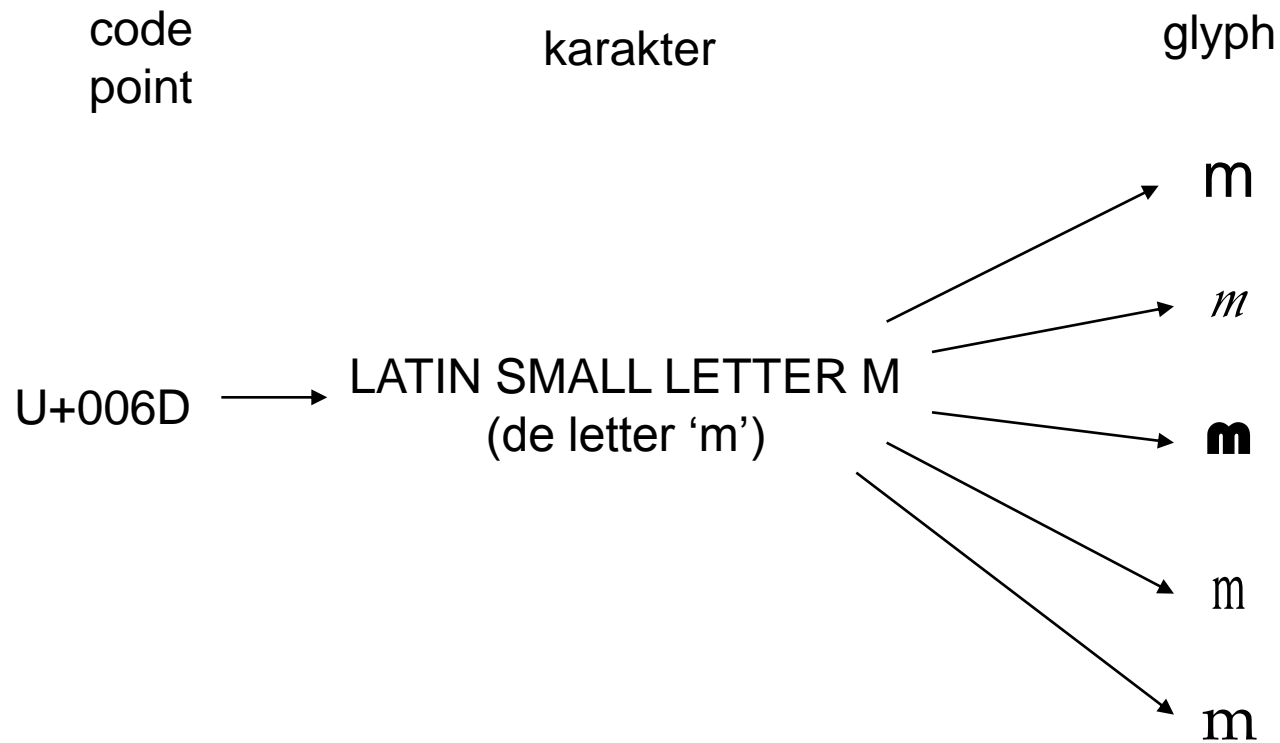
IBM  
mainframe

etc. etc. etc....

# Vóór Unicode

	CP 437	CP 850	Windows 1252	ISO- Latin-1	ISO- Latin-15	Unicode
a	61	61	61	61	61	0061
á	a0	a0	e1	e1	e1	00e1
ä	84	84	e4	e4	e4	00e4
€	-	-	80	-	a4	20ac
õ	-	e4	f5	f5	f5	00f5
Ç	80	80	c7	c7	c7	00c7
œ	-	-	9c	-	bd	0153

# Unicode



0000

### C0 Controls and Basic Latin

007F

	000	001	002	003	004	005	006	007
0	NUL 0000	DLE 0010	SP 0020	0 0030	@ 0040	P 0050	` 0060	p 0070
1	SOH 0001	DC1 0011	! 0021	1 0031	A 0041	Q 0051	a 0061	q 0071
2	STX 0002	DC2 0012	" 0022	2 0032	B 0042	R 0052	b 0062	r 0072
3	ETX 0003	DC3 0013	# 0023	3 0033	C 0043	S 0053	c 0063	s 0073
4	EOT 0004	DC4 0014	\$ 0024	4 0034	D 0044	T 0054	d 0064	t 0074
5	ENQ 0005	NAK 0015	% 0025	5 0035	E 0045	U 0055	e 0065	u 0075
6	ACK 0006	SYN 0016	& 0026	6 0036	F 0046	V 0056	f 0066	v 0076
7	BEL 0007	ETB 0017	' 0027	7 0037	G 0047	W 0057	g 0067	w 0077



0E00

Thai

0E7F

	0E0	0E1	0E2	0E3	0E4	0E5	0E6	0E7
0		๐ 0E10	๑ 0E20	๒ 0E30	๓ 0E40	๔ 0E50		
1	๕ 0E01	๖ 0E11	๗ 0E21	๘ 0E31	๙ 0E41	๐ 0E51		
2	๑ 0E02	๒ 0E12	๓ 0E22	๔ 0E32	๕ 0E42	๖ 0E52		
3	๗ 0E03	๘ 0E13	๙ 0E23	๐ 0E33	๑ 0E43	๒ 0E53		
4	๓ 0E04	๔ 0E14	๕ 0E24	๖ 0E34	๗ 0E44	๘ 0E54		
5	๙ 0E05	๐ 0E15	๑ 0E25	๒ 0E35	๓ 0E45	๔ 0E55		
6	๑ 0E06	๒ 0E16	๓ 0E26	๔ 0E36	๕ 0E46	๖ 0E56		
7	๗ 0E07	๘ 0E17	๙ 0E27	๐ 0E37	๑ 0E47	๒ 0E57		

**4E00** **CJK Unified Ideographs** **4EFF**














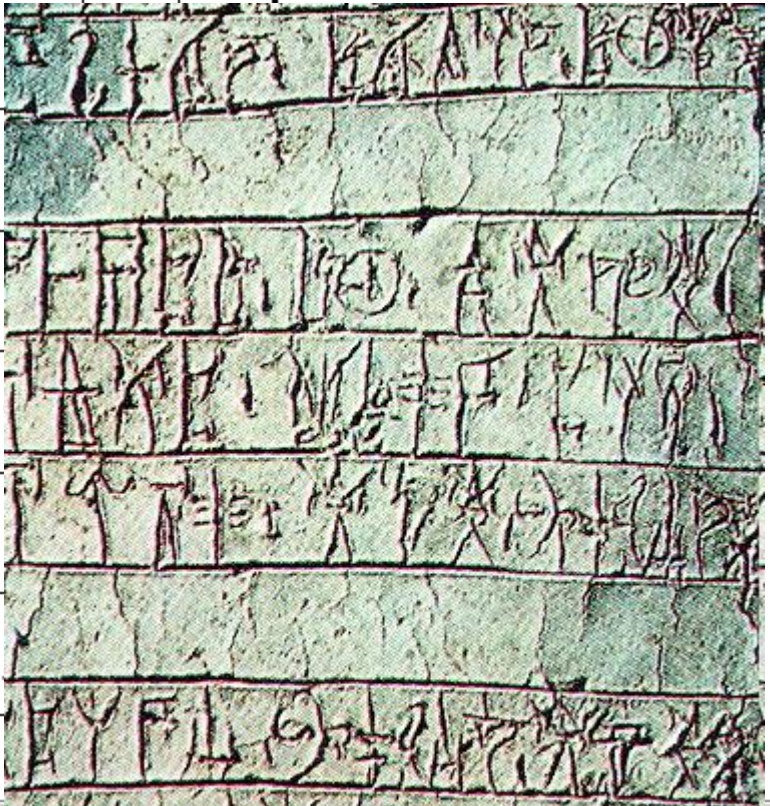





























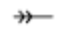
	4E00	4E01	4E02	4E03	4E04	4E05	4E06	4E07	4E08	4E09	4EA0	4EB0	4EC0	4ED0	4EE0	4EF0
0	一	丂	北	丰	乂	乐	习	买	龜	亏	宀	京	什	伞	仟	仰
1	丁	丑	兩	卯	乂	采	乡	乱	乾	云	亡	宥	仁	仑	仵	伶
2	丂	刃	丢	串	乂	兵	屮	姿	亂	互	亢	亲	仉	令	仉	仲
3	七	专	卵	弗	乃	兵	纟	乳	粼	亅	亅	毫	仃	仓	代	仉
4	上	且	两	临	乂	乔	丂	哲	𡗗	五	交	寔	仄	仔	令	仉
5	丁	丕	严	莘	久	帛	丂	乳	丂	井	亥	衰	仅	仕	以	仉
6	丂	世	並	丂	乂	乖	书	丂	了	三	亦	亶	仆	他	仉	仉
7	万	卅	丧	丂	乂	乘	芝	丂	尔	𡗗	产	廉	仇	仗	夫	价

Page Display

10080

Linear B Ideograms

100FF

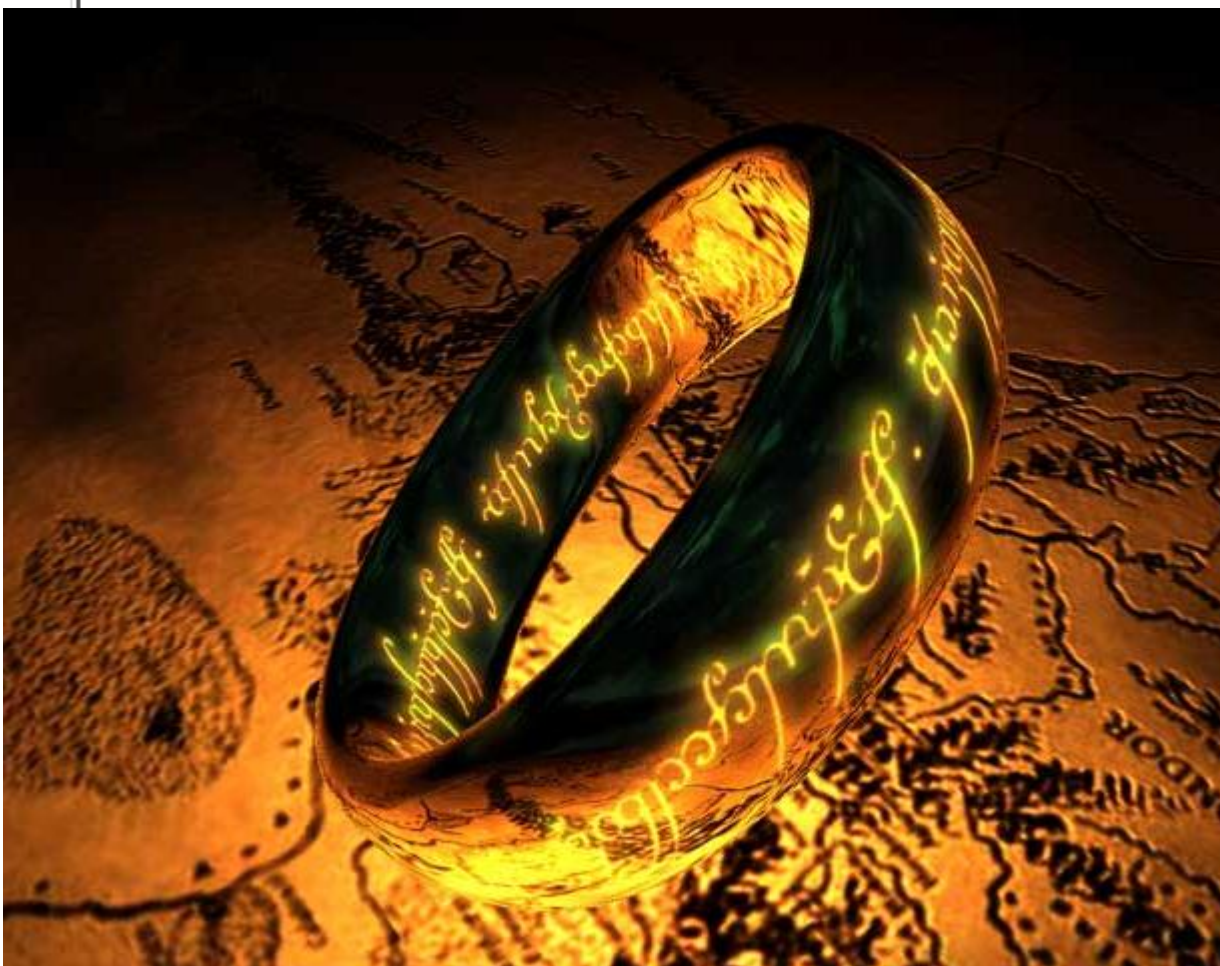
	1008	1009	100A	100B	100C	100D	100E	100F
0	 10080	 10090	 100A0	 100B0	 100C0	 100D0	 100E0	 100F0
1	 10081	 10091	 100A1	 100B1	 100C1			
2	 10082	 10092	 100A2	 100B2	 100C2			
3	 10083	 10093	 100A3	 100B3	 100C3			
4	 10084	 10094	 100A4	 100B4	 100C4			
5	 10085	 10095	 100A5	 100B5	 100C5			
6	 10086	 10096	 100A6	 100B6	 100C6			
7	 10087	 10097	 100A7	 100B7	 100C7			

12000 Cuneiform 120FF

	1200	1201	1202	1203	1204	1205	1206	1207	1208	1209	120A	120B	120C	120D	120E	120F
0																
1																
2																
3																
4																
5																
6																
7																

E00	E01	E02	E03	E04	E05	E06	E07
p	ᄁ	λ	ᄂ	ᄃ	·		
p	ᄂ	ᄃ	ᄄ	ᄅ	:		
q	ᄃ	λ	ᄅ	ᄆ	ᄇ	ᄈ	
q	ᄄ	ᄅ	ᄆ	ᄇ	ᄉ	ᄊ	
ᄉ	ᄋ			ᄌ	ᄍ	ᄎ	
ᄏ	ᄐ			ᄑ	ᄒ	ᄓ	
ᄔ	ᄕ			ᄌ	ᄍ	ᄎ	
ᄏ	ᄐ	ᄑ		ᄒ	ᄓ	ᄔ	
ᄕ	ᄌ	ᄍ		ᄎ	ᄏ	ᄐ	
ᄑ	ᄒ	ᄓ		ᄔ	ᄕ	ᄌ	
ᄍ	ᄎ	ᄏ		ᄐ	ᄑ	ᄒ	
ᄓ	ᄔ	ᄕ		ᄌ	ᄍ	ᄎ	
ᄏ	ᄐ	ᄑ		ᄒ	ᄓ	ᄔ	

Tengwar – Tolkien - niet officieel (private use range)



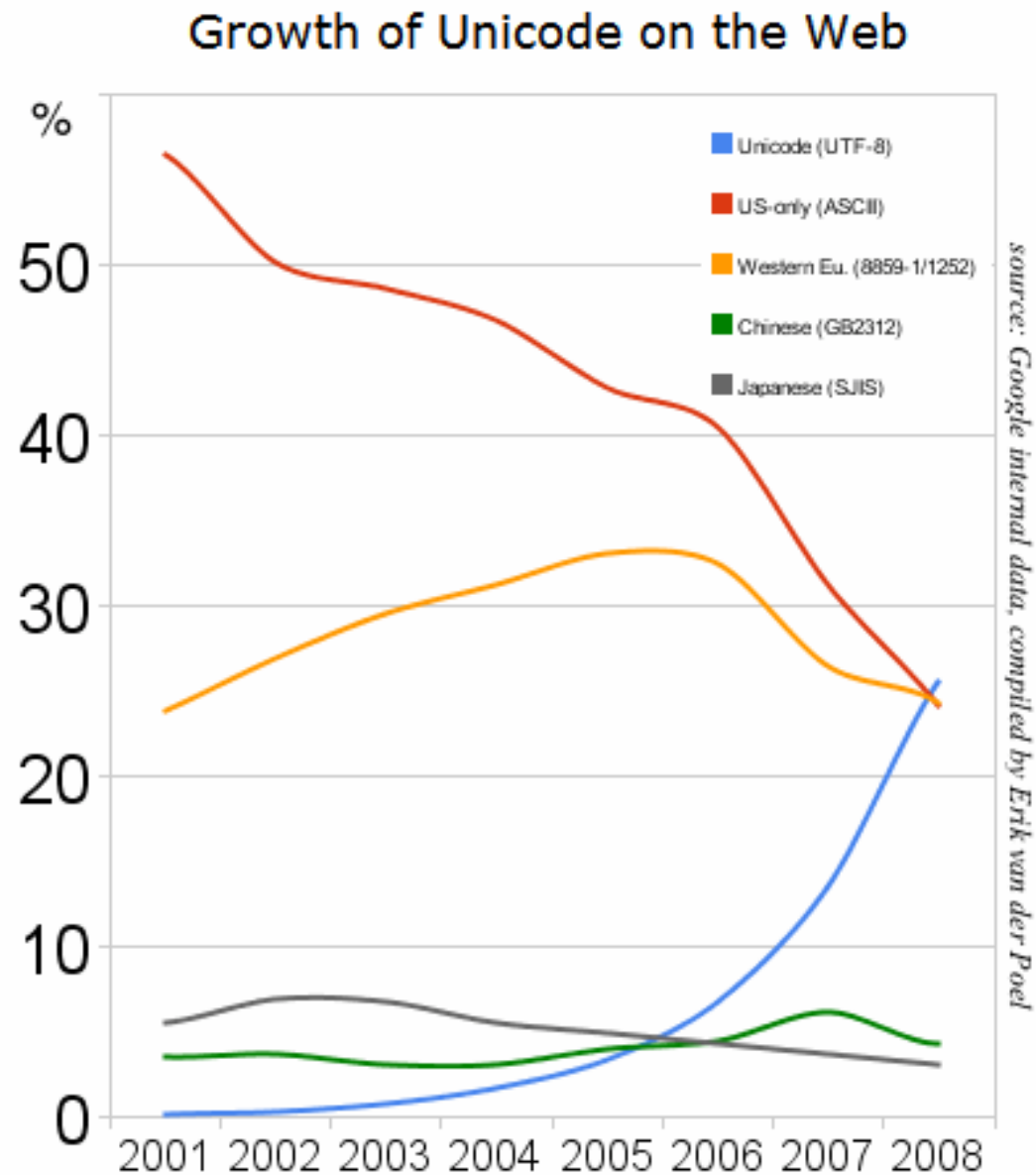
# Unicode encodings

- Unicode
  - U+006D = 'm'
  - karakter 0 – 255: gelijk aan ISO-Latin-1
  - 1.114.112 code points (0 – 10FFFF)
- UTF-16 encoding
  - 4 bytes
  - 0000 – FFFF: gelijk aan Unicode nummer
  - Byte Order Mark
    - U+FEFF (ZERO-WIDTH NO-BREAK SPACE)
    - byte-swapped = U+FFFE = geen legaal karakter
  - efficiënt voor Chinees en Japans
- UTF-8
  - 1 tot 4 bytes
  - 0 – 127: gelijk aan ASCII
  - ergo: ASCII tekst is altijd ook UTF-8 tekst
  - efficiënt voor Westerse talen

# Unicode encodings

teken	code point	UTF-8	UTF-16	ISO-Latin-1	ASCII
spatie	U+0020	20	00 20	20	20
a	U+0061	61	00 61	61	61
ä	U+00E4	C3 A4	00 E4	E4	-

- XML
- Java
- .NET



Marc de Graauw

<http://www.marcdegrauw.com/>