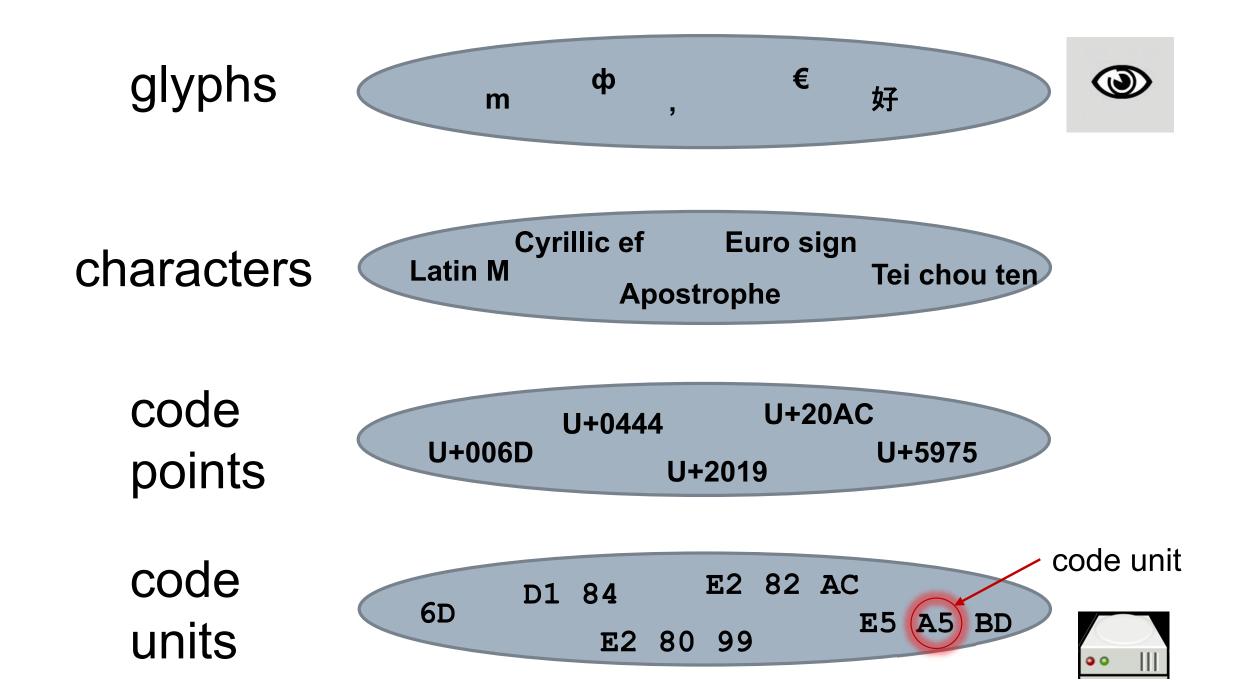
### Unicode and UTF-8

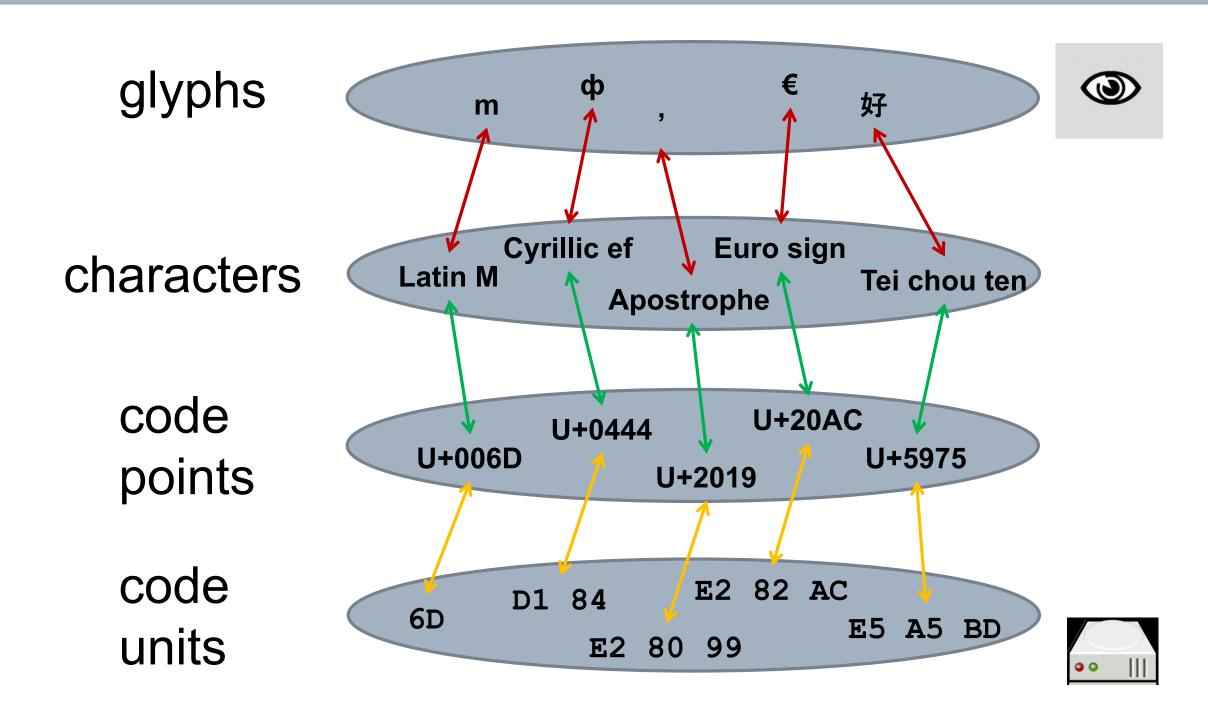
Computer Science and Engineering ■ College of Engineering ■ The Ohio State University

#### Lecture 41

A standard for the discrete representation of written text

### The Big Picture: Sets



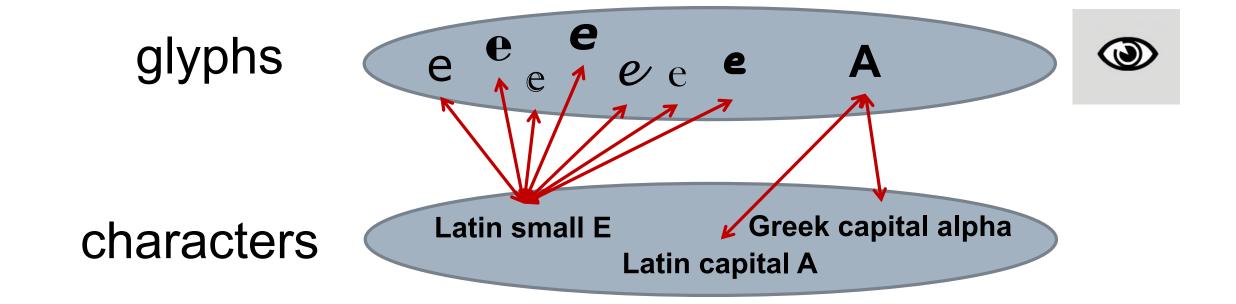




Ласкаво просимо! добро пожаловать Willkommen 霍/ Benvenuti ようこそ 한 영합니다 Bienvenue ยินดีต้อนรับ WELKOM Soo dhawow

THE OHIO STATE
UNIVERSITY
UNIVERSITY LIBRARIES

### Glyphs vs Characters



- ☐ Glyph: "An individual mark on a written medium that contributes to the meaning of what is written."
  - See foyer floor in main library
- One character can have different glyphs
  - Example: Latin Small E could be e, e, e, e, e...
- □ One *glyph* can be different *characters* 
  - 0 is both Digit Zero and (capital) Latin 0
  - A is both (capital) Latin A and Greek Alpha
- One unit of text can consist of multiple glyphs
  - An accented letter (é) is two glyphs
  - The ligature of f+i (fi) is two glyphs

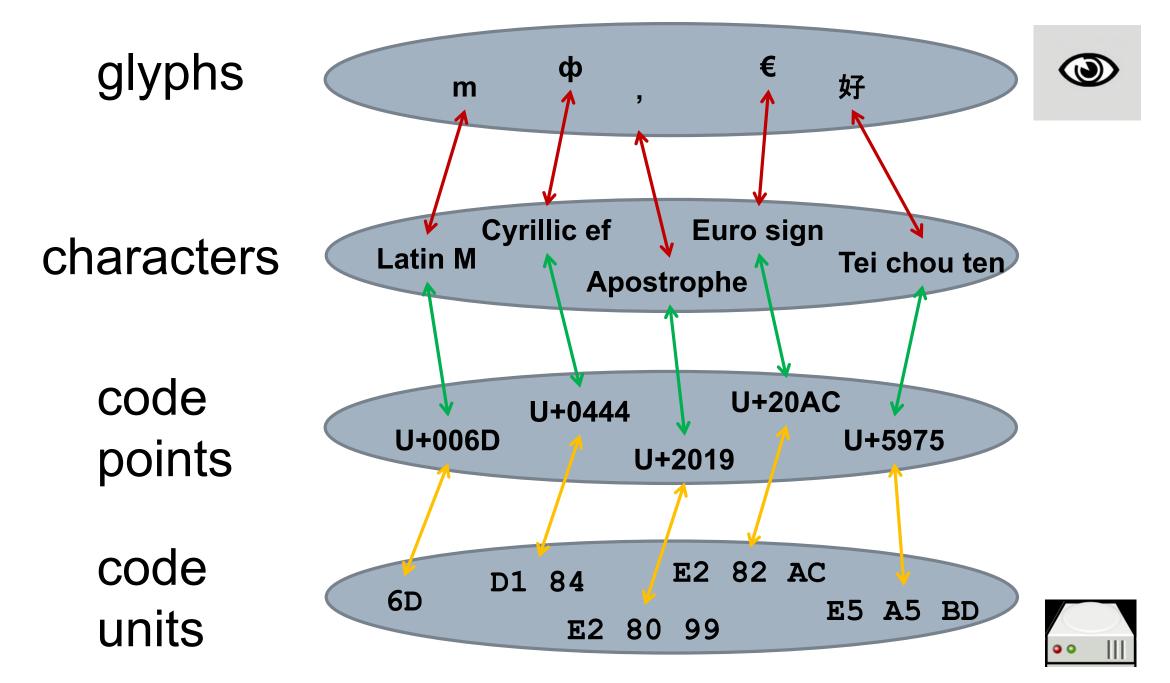
### Security Issue: Eyes Deceive

- Visual homograph: Two different characters that look the same
  - Would you click here: <u>www.paypal.com</u>?

## Security Issue: Eyes Deceive (Cont'd)

- Visual homograph: Two different characters that look the same
  - Would you click here: <u>www.paypal.com</u>?
  - Oops! The second 'a' is actually CYRILLIC SMALL LETTER
  - This site successfully registered in 2005
- Other examples: combining characters
  - $\tilde{n}$  = LATIN SMALL LETTER N WITH TILDE
  - $\tilde{n}$  = LATIN SMALL LETTER N + COMBINING TILDE
- □ "Solution"
  - Heuristics that warn users when languages are mixed and homographs are possible

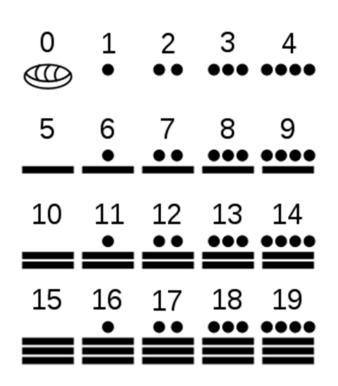
### Unicode: Characters, Code Points



- □ Each character is assigned a unique code point
- A code point is defined by an integer value, and is given a name
  - one hundred and nine (109, or 0x6d)
  - LATIN SMALL LETTER M
- □ Convention: Write code points as U+hex
  - Example: U+006D
- □ As of Sept '24, v16.0 (see unicode.org):
  - Contains 154,998 code points emoji-versions.html
  - Covers 168 scripts (and counting...) unicode.org/charts/

## Example Recent Addition (v11)

**Computer Science and Engineering** ■ The Ohio State University

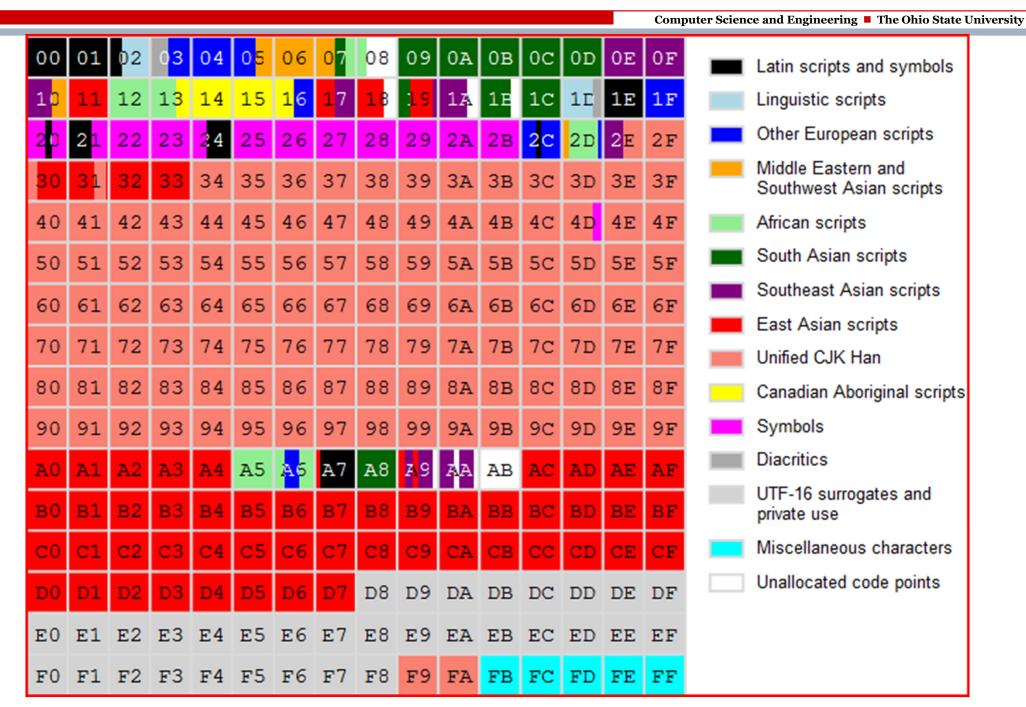


#### Mayan numerals

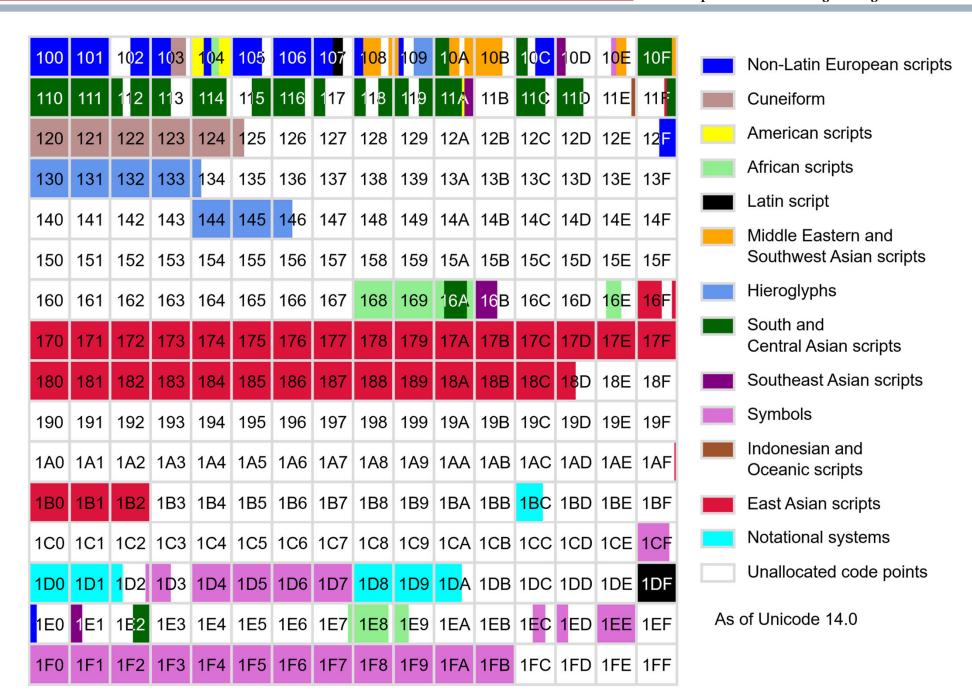
```
1D2E0 MAYAN NUMERAL ZERO
1D2E1 . MAYAN NUMERAL ONE
1D2E2 . . MAYAN NUMERAL TWO
1D2E3 ... MAYAN NUMERAL THREE
1D2E4 .... MAYAN NUMERAL FOUR
1D2E5 ___ MAYAN NUMERAL FIVE
1D2E6 _ MAYAN NUMERAL SIX
1D2E7 • • MAYAN NUMERAL SEVEN
1D2E8 ••• MAYAN NUMERAL EIGHT
1D2E9 •••• MAYAN NUMERAL NINE
1D2EA ___ MAYAN NUMERAL TEN
1D2EB 📥 MAYAN NUMERAL ELEVEN
1D2EC MAYAN NUMERAL TWELVE
1D2ED 🔐 MAYAN NUMERAL THIRTEEN
1D2EE *** MAYAN NUMERAL FOURTEEN
1D2EF 🚃 MAYAN NUMERAL FIFTEEN
1D2F0 📥 MAYAN NUMERAL SIXTEEN
1D2F1 🗮 MAYAN NUMERAL SEVENTEEN
1D2F2 👑 MAYAN NUMERAL EIGHTEEN
1D2F3 MAYAN NUMERAL NINETEEN
```

- Code points are grouped into categories
  - Basic Latin, Cyrillic, Arabic, Cherokee, Currency, Mathematical Operators, ...
- $\square$  Unicode allows for 17 x 2<sup>16</sup> code points
  - 0 to 1,114,111 (i.e., > 1 million)
  - $\blacksquare$  U+0000 to U+10FFFF
- $\square$  Each block of  $2^{16}$  code points is a *plane* 
  - U+nnnnnn, same green ==> same plane
  - ~64,000 code points per plane
- □ Plane 0 is the *basic multilingual plane* (BMP)
  - Has (practically) everything you could need
  - Convention: code points in BMP written U+nnnn (ie 4 digits, leading 0's if needed)
  - Others code points written without leading 0's

### Basic Multilingual Plane



### Supplemental Plane (plane 1)



### UTF-8: Code Points & Octets

glyphs **(3)** 好 m Cyrillic ef **Euro sign** characters Latin M Tei chou ten **Apostrophe** code **U+20AC** U+0444 U+006D U+5975 points U+2019 code E2 82 AC D1 84 6D E5 A5 BD units 80 99

- Encodes each code point (integer) as a sequence of bytes (octets)
- Variable length
  - Some code points require 1 octet
  - Others require 2, 3, or 4
- Consequence: Can not infer number of characters from size of file!
- □ No endian-ness: a *sequence* of octets D0 BF D1 80 D0 B8 D0 B2 D0 B5 D1 82...
- Other encodings exist!
  - Eg UTF-16 uses 2 bytes per code point (more general term: code unit)

### UTF-8 Encoding Recipe: 1 and 2-Bytes

- □ 1-byte encodings
  - First bit is 0
  - Example: 0110 1101 (encodes U+006D)
- 2-byte encodings
  - First byte starts with **110**...
  - Second byte starts with 10...
    - □ Example: **110**1 0000 **10**11 1111
    - □ Payload: 1101 0000 1011 1111 = 100 0011 1111
      - = 0x043F
    - □ Code point: U+043F i.e. п, Cyrillic small letter pe

- □ Generalization: An encoding of length k:
  - First byte starts with k 1's, then a 0
    - $\square$  Example **1110** 0110 ==> first byte of a 3-byte encoding
  - Subsequent k-1 bytes each start with 10
  - Remaining bits are the payload
- □ Example: E2 82 AC 11100010 10000010 10101100
  - Payload: 0x20AC (i.e., U+20AC, €)
- Consequence: Stream is self-synchronizing
  - Losing a byte affects only one character

# UTF-8 Encoding Summary

**Computer Science and Engineering** ■ The Ohio State University

Unicode	Byte1	Byte2	Byte3	Byte4	example
U+0000-U+007F	0xxxxxxx				'\$' U+0024 → 00100100 → 0x24
U+0080-U+07FF	110ууухх	10xxxxxx			'¢' U+00 <u>A</u> 2 → 110000 <u>10</u> , 10 <u>10</u> 0010 → 0xC2, 0xA2
U+0800-U+FFFF	1110уууу	10уууухх	10xxxxxx		'€' U+20AC → 1110 <u>0010</u> , 100000 <u>10</u> , 10 <u>10</u> 1100 → 0xE2, 0x82, 0xAC
U+10000-U+10FFFF	11110 <i>zzz</i>	10 <i>zzy</i> yyy	10уууухх	10xxxxxx	'類' U+024B62 → 11110000,10100100,10101101,10100010 → 0xF0,0xA4,0xAD,0xA2

(from wikipedia)

- For the following UTF-8 encoding, what is the corresponding code point(s)?
  - F0 A4 AD A2

- For the following Unicode code point, what is its UTF-8 encoding?
  - U+20AC

### Security Issue: Multiplicity of Encodings

- Not all octet sequences are legal encodings
  - "overlong" encodings are illegal
  - example: C0 AF
    - = **110**0 00**00 10**10 1111
    - = U+002F (encoding should be 2F)
- □ Classic security bug (IIS 2001)
  - Should reject URL requests with "../.."
    - □ Looked for 2E 2E 2F 2E 2E (in encoding)
  - Accepted "..%c0%af.." (doesn't contain x2F)
    - □ 2E 2E CO AF 2E 2E is ok to allow through
  - After accepting, server then decoded
    - □ 2E 2E CO AF 2E 2E decoded into "../.."
- ☐ Moral: String is a sequence of *code units* 
  - But we care about code points

### Other (Older) Encodings: ASCII

- ☐ In the beginning...
- Character sets were small
  - ASCII: only 128 characters (ie 2<sup>7</sup>)
  - 1 byte/character, leading bit always 0

6D = Latin small m

ASCII Code Chart																
	0	1	2	<sub>1</sub> 3	4	5	6	7	8 1	9	ΙA	В	С	D	E	L F
0	NUL	SOH	STX	ETX	EOT	ENQ	ACK	BEL	BS	HT	LF	VT	FF	CR	SO	SI
ī	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
2		!	=	#	\$	%	&	-	(	)	*	+	,	-	•	/
3	0	1	2	3	4	5	6	7	8	9	:	;	٧	=	>	?
4	@	Α	В	С	D	Ε	F	G	Н	I	J	K	L	М	N	0
5	Р	Q	R	S	T	U	٧	W	Х	Υ	Z	[	\	]	^	_
6	`	а	b	С	d	е	f	g	h	i	j	k	ι	m	n	0
7	р	q	r	s	t	u	V	W	х	у	z	{	_	<b>}</b>	~	DEL
_														T		
													V	•		

### Other (Older) Encodings: Code Pages

- ☐ In the beginning...
- Character sets were small
  - ASCII: only 128 characters (ie 2<sup>7</sup>)
  - 1 byte/character, leading bit always 0
- ☐ Globalization means more characters...
  - But 1 byte/character seems fundamental
- □ Solutions:
  - Use that leading bit!
    - □ Text data now looks just like binary data
    - □ 256 characters
  - Use more than 1 encoding!
    - Must specify data + encoding used
    - □ Each encoding gives 256 characters

## ISO-8859 family (eg -1 Latin)

**Computer Science and Engineering** ■ The Ohio State University

	-0	-1	-2	-3	-4	-5	-6	-7	-8	-9	-A	-B	-C	-D	-E	-F
0-		0001	0002	0003	0004	0005	0006	0007	0008	0009	000A	0008	000C	0000	000E	000F
1-	0010	0011	0012	0013	0014	0015	0016	0017	0018	0019	001A	001B	001C	001D	001E	001F
2-		!	"	#	\$	%	&	1	(	)	*	+	,	-		/
3-	0020	1	2	3	4	5	6	7	8	9	002A	002B	< 002C	002D =	> 002E	002F
4-	<b>@</b>	A	B	C	<b>D</b>	E	<b>F</b>	<b>G</b>	H	I 0039	J	<b>K</b>	L	M	N	003F
5-	P	Q 0041	0042 <b>R</b>	S 0043	T	U	0046 <b>V</b>	W	0048 <b>X</b>	Y	<b>Z</b>	004B	0040	004D	004E	004F
J-	0050	0051	0052	0053	0054	0055	0056	0057	0058	0059	005A	005B	005C	005D	005E	005F
6-	0060	<b>a</b>	b 0062	C 0063	<b>d</b>	e 0065	<b>f</b>	<b>g</b>	h 0068	i 0069	<b>j</b>	<b>k</b>	0060	m 006D	n 006E	O 006F
7-	<b>p</b>	<b>q</b>	r 0072	S 0073	<b>t</b>	u 0075	<b>V</b>	<b>W</b>	<b>X</b>	<b>y</b>	<b>Z</b>	<b>{</b>	007C	} 007D	~ 007E	007F
8-	0080	0081	0082	0083	0084	0085	0086	0087	0088	0089	008A	0088	0080	008D	008E	008F
9-	2722		2722		2004		2000			2222						2005
Α-	0090	0091	Ø	£	D094	¥	0096	\$	0098	© 0099	<u>a</u>	≪ ≪	0090	0090	009E	009F
	00A0	00A1	00A2	00A3	00A4	00A5	00A6	00A7	8A00	00A9	00AA	00AB	00AC	00AD	00AE	00AF
B-	0080	<u>+</u>	0082	3 00B3	0084	$\mu_{_{\scriptscriptstyle{00B5}}}$	¶ 0086	0087	د 00B8	0089	0 00BA	>> 0088	1/4 008C	1/2 00BD	3/4 00BE	Ġ OOBF
C-	À	Á 0001	Â	Ã	Ä 00C4	Å 0005	Æ 0006	Ç	È	É	Ê	Ë	Ì	Í	Î	<b>Ï</b>
D-	Ð	Ñ	Ò	Ó	Ô 00D4	Õ	Ö	X 0007	Ø	Ù	Ú	Û	Ü	Ý	Þ OODE	ß
E-	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
	00E0	00E1 ≃	00E2	00E3	00E4	00E5 ≈	00E6	00E7	00E8	00E9	00EA	00EB	00EC	00ED	00EE	00EF
F-	<b>ð</b> 00F0	<b>ñ</b> 00F1	<b>Ò</b> 00F2	<b>Ó</b> 00F3	<b>ô</b> 00F4	<b>Õ</b> 00F5	<b>Ö</b> 00F6	00F7	<b>Ø</b> 00F8	ù 00F9	Ú OOFA	û OOFB	ü oofc	<b>ý</b> 00FD	<b>p</b> OOFE	ÿ

- 0-7F match ASCII

reserved (control characters)

A0-FF differ, eg:

- -1 "Western"
- -2 "East European"
- -9 "Turkish

# Windows Family (eg 1252 Latin)

**Computer Science and Engineering** ■ The Ohio State University

Windows-1252 (CP1252)																
	х0	х1	<b>x2</b>	х3	х4	х5	х6	х7	х8	х9	хA	хВ	хС	хD	хE	хF
0x	<u>NUL</u>	<u>SOH</u>	<u>STX</u>	<u>ETX</u>	<u>EOT</u>	<u>ENQ</u>	<u>ACK</u>	<u>BEL</u>	<u>BS</u>	<u>HT</u>	<u>LF</u>	<u>VT</u>	<u>FF</u>	<u>CR</u>	<u>so</u>	<u>SI</u>
1x	DLE	DC1	DC2	DC3	DC4	<u>NAK</u>	<u>SYN</u>	<u>ETB</u>	<u>CAN</u>	<u>EM</u>	<u>SUB</u>	<u>ESC</u>	<u>FS</u>	<u>GS</u>	<u>RS</u>	<u>US</u>
2x	<u>SP</u>	ļ.	"	#	\$	%	&	•	(	)	*	+	,	-		/
3x	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
4x	@	Α	В	С	D	E	F	G	Н	I	J	K	L	М	N	0
5x	Р	Q	R	S	Т	U	V	W	X	Υ	Z	[	١	]	٨	
6x	•	а	b	С	d	е	f	g	h	i	j	k		111	n	0
7x	р	q	r	s	t	u	v	w	Х	y	Z	{		}	~	<u>DEL</u>
8x	€		,	f	"		1	‡	^	‰	Š	<	Œ		Ž	
9x			, <b>4</b>		"	•	_	_	~	тм	š	>	œ		ž	Ϋ
Ax	NBSP	i	¢	£	¤	¥	ŀ	§		©	а	«	7		®	-
Вх	0	±	2	3	•	μ	¶		د	1	0	»	1/4	1/2	3/4	ن
Сх	À	Á	Â	Ã	Ä	Å	Æ	Ç	È	É	Ê	Ë	Ì	ĺ	î	Ϊ
Dx	Đ	Ñ	Ò	Ó	Ô	Õ	Ö	×	Ø	Ù	Ú	Û	Ü	Ý	Þ	ß
Ex	à	á	â	ã	ä	å	æ	ç	è	é	ê	ë	ì	í	î	ï
Fx	ð	ñ	ò	ó	ô	õ	Ö	÷	ø	ù	ú	û	ü	ý	þ	ÿ

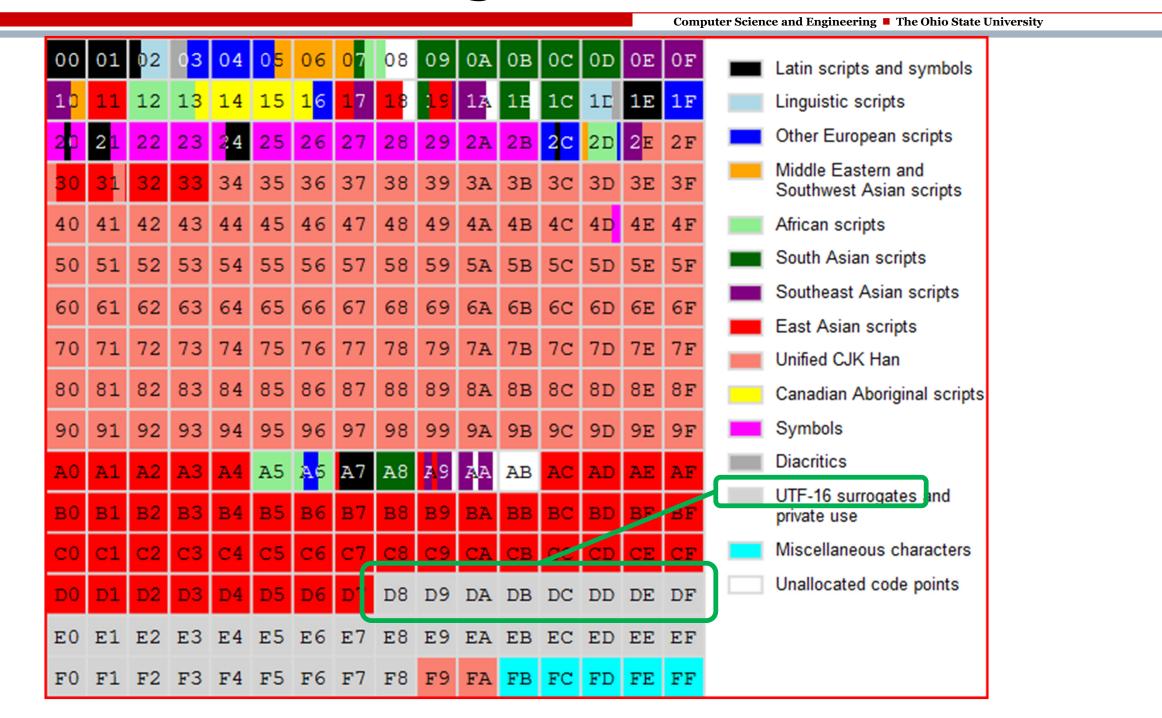
92 = apostrophe

<u>Name</u>	Labels
The Encoding	
UTF-8	"unicode-1-1-utf-8"
	"utf-8"
	"utf8"
	<b>I</b>
windows-1252	"ansi_x3.4-1968"
	"ascii"
	"cp1252"
	"cp819"
	"csisolatin1"
	"ibm819"
	"iso-8859-1"
	"iso-ir-100"
	"iso8859-1"
	"iso88591"
	"iso_8859-1"

- Unicode started as 2<sup>16</sup> code points
  - The BMP of modern Unicode
  - Bottom 256 code points match ISO-8859-1
- □ Simple 1:1 encoding (UTF-16)
  - Code point <--> 16-bit code unit (ie 2 bytes)
  - Simple, but doubles storage needed for ASCII
- □ Later, code points outside of BMP added
  - A pair of words (aka "surrogate pairs") carry 20-bit payload split, 10 bits in each word
  - First: **1101 10**xx xxxx xxxx (xD800-DBFF)
  - Second: 1101 11yy yyyy yyyy (xDC00-DFFF)
- Consequence: U+D800 to U+DFFF became reserved code points in Unicode
  - And now we are stuck with this legacy, even for UTF-8

□ JavaScript uses UTF-16 let  $x = "\{u\{1f916\}hi" // robot face + hi$ x.length //=> 4 (number of code units) x.charAt(0) //=> char from 1st code unit x.charAt(2) // surprise? [...x][2] // spread = linear time {...x} // spreads code units Ruby supports multiple encodings  $x = "\{u\{1f916\}\}"$ x.length x.bytes.map { |b| b.to s(2) } x.encoding x.encode! Encoding::UTF 16 x.bytes.map { |b| b.to s(16) }

### Recall: Basic Multilingual Plane



- A multi-byte representation must distinguish between big & little endian
  - Example: 00 25 00 25 00 25
  - "%%%" if LE, "— —" if BE
- One solution: Specify encoding in name
  - UTF-16BE or UTF-16LE
- Another solution: require byte order mark (BOM) at the start of the file
  - U+FEFF (ZERO WIDTH NO BREAK SPACE)
  - There is *no* U+FFFE code point
  - So FE FF → BigE, while FF FE → LittleE
  - Not considered part of the text

- □ Should we add a BOM to the start of UTF-8 files too?
  - UTF-8 encoding of U+FEFF is EF BB BF
- Advantages:
  - Forms magic-number for UTF-8 encoding
- Disadvantages:
  - Not backwards-compatible to ASCII
  - Existing programs may no longer work
  - *E.g.*, In Unix, shebang (#!, *i.e.* 23 21) at *start* of file is significant: file is a script
    - #! /bin/bash

- Using U+FEFF as ZWNBSP deprecated
  - Reserved for BOM uses (at start of file)
- □ Alternative: U+200D ("zwidge")
- □ Joined characters may be rendered as a single glyph
  - Co-opted for use with emojis
- □ Example: 1 character?
  - U+1F3F4 U+200D U+2620
  - let  $x = "\{u\{1F3F4\}\{u\{200D\}\{u\{2620\}\}"\}\}$
  - WAVING BLACK FLAG, ZWJ, SKULL AND CROSSBONES







- □ What is a "text" file? (vs "binary" file)
  - Given a file, how can you tell which it is?
- A JavaScript program reads in a 5MB file of ASCII text, storing it in the string £
  - How many characters are in £?
  - How much memory does **f** occupy?
- □ How many characters are in string s?

```
let s = ... // a JavaScript string
console.assert(s.length == 7) // true
```

- □ Which is better: UTF-8 or UTF-16?
- ☐ What's so scary about:

```
..%c0%af..
```

### Summary

- □ Text vs binary
  - In pre-historic times: most significant bit
  - Now: data is data
- Unicode code points
  - Integers U+0000..U+10FFFF
  - BMP: Basic Multilingual Plane
- □ UTF-8
  - A variable-length, self-synchronizing encoding of unicode code points
  - Backwards compatible with ISO 8859-1, and hence with ASCII too