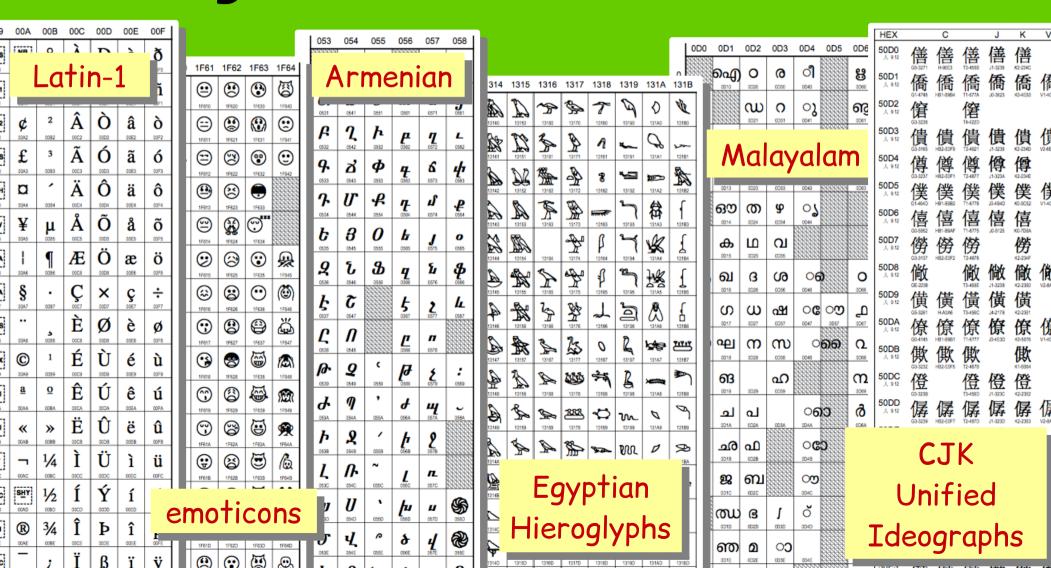
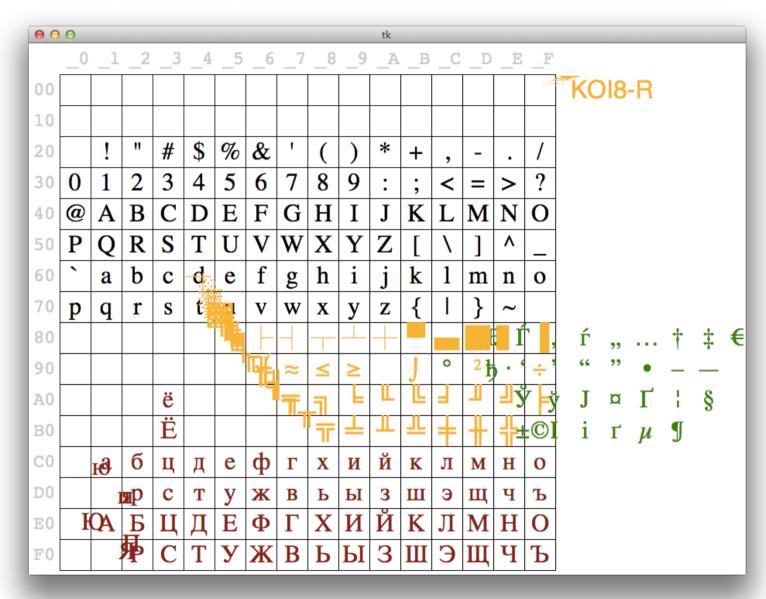
Unicode solutions in Python 2 and 3



Unicode solutions in Python 2 and 3



The single-byte codepage ballet



Video: https://www.youtube.com/watch?v=J4qioAacrYo

Source code: http://bit.ly/10qt0MZ



Why Unicode

- Too many incompatible byte encodings
- Separate concepts:
 - character identity: oneode point or each abstract character
 - U+0041 → LATIN CAPITAL LETTER A
 - U+096C → DEVANAGARI DIGIT SIX
 - binary representation: multipelacodings

```
• U+0041 \rightarrow 0x41
```

0x41 0x00

• U+096C \rightarrow 0xE0 0xA5 0xA¢ 0x6C 0x09



UTF-8

UTF-16LE

A sample of encodings

char.	code point	ascii	latin1	cp1252	ср437	gb2312	utf-8	utf-16le
Α	U+0041	41	41	41	41	41	41	41 00
Ś	U+00BF	*	BF	BF	A8	*	C2 BF	BF 00
Ã	U+00C3	*	C3	C3	*	*	C3 83	C3 00
á	U+00E1	*	E1	E1	A0	A8 A2	C3 A1	E1 00
Ω	U+03A9	*	*	*	EA	A6 B8	CE A9	A9 03
È	U+06BF	*	*	*	*	*	DA BF	BF 06
66	U+201C	*	*	93	*	A1 B0	E2 80 9C	1C 20
€	U+20AC	*	*	80	*	*	E2 82 AC	AC 20
Г	U+250C	*	*	*	DA	A9 B0	E2 94 8C	0C 25
气	U+6C14	*	*	*	*	C6 F8	E6 B0 94	14 6C
氣	U+6C23	*	*	*	*	*	E6 B0 A3	23 6C
&	U+1D11E	*	*	*	*	*	F0 9D 84 9E	34 D8 1E DD



.encode() vs .decode()

- "Humans use text. Computers speak bytes."
 - Esther Nam and Travis Fischer in Character encoding and Unicode in Python (Pycon US 2014)
- Use .encode() to convertumantext tobytes
- Use .decode() to convertes to humantext

2.7 gotcha:
the methods
.encode() and .decode()
exist in **str** and **unicode**



Data types for text and bytes

	Python 2.7	Python 3.4
Human text	unicode u'café', u'caf\xe9'	str 'café', u'café'
Bytes (imutável)	str 'café', 'caf\xe9', b'café'	bytes b'caf\xc3\xa9'
Bytes (mutável)	<pre>bytearray bytearray(b'caf\xc3\xa9')</pre>	<pre>bytearray bytearray(b'caf\xc3\xa9')</pre>



Best practice

The Unicode sandwich

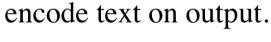


bytes→str 100% str

str→bytes encode

Decode bytes on input,

process text only,





Coping with Unicode Errors

- SyntaxError
 - A .py file is loaded with contents in an unexpected encoding
- UnicodeDecodeError
 - A binary sequence is contains bytes that are not valid in the expected encoding
- UnicodeEncodeError
 - A Unicode string contains codepoints that have no representation in the desired encoding



Coping with SyntaxError

- A .py file is loaded with contents in an unexpected encoding
 - The source file encoding is not the default, and no # encoding comment was found.
 - The source file encoding is not the one declared in the # encoding comment

encoding2 7

Characters: 143 · Words: 21

Defaults:

Unicode database

\$ python3	3 nume	rics_demo	.py		5. bash	
U+0031	1	re_dig	isdig	isnum	1.00	DIGIT ONE
U+00bc	1 ⁄ ₄	-	-	isnum	0.25	VULGAR FRACTION ONE QUARTER
U+00b2	2	-	isdig	isnum	2.00	SUPERSCRIPT TWO
U+0969	રૂ	re_dig	isdig	isnum	3.00	DEVANAGARI DIGIT THREE
U+136b	<u>c</u>	-	isdig	isnum	3.00	ETHIOPIC DIGIT THREE
U+216b	XII	_	-	isnum	12.00	ROMAN NUMERAL TWELVE
U+2466	7	_	isdig	isnum	7.00	CIRCLED DIGIT SEVEN
U+2480	(13)	_	-	isnum	13.00	PARENTHESIZED NUMBER THIRTEEN
U+3285	\bigcirc	_	-	isnum	6.00	CIRCLED IDEOGRAPH SIX
\$						



Unicode database

```
$ python3 numerics_demo.py
U + 0031
                  re_dig isdig
                                                1.00
                                                        DIGIT ONE
                                     isnum
U+00bc
                                     isnum
                                                0.25
                                                        VULGAR FRACTION ONE OUARTER
U+00b2
                                                        SUPERSCRIPT TWO
                            isdia
                                               2.00
                                     isnum
U+0969
                            isdia
                                                3.00
                                                        DEVANAGART DIGIT THREE
                  re_dia
                                     isnum
U+136b
                            isdia
                                                3.00
                                                        ETHIOPIC DIGIT THREE
                                     isnum
U+216b
                                              12.00
                                                        ROMAN NUMERAL TWELVE
                                     isnum
           XII
                                     ichum
                                                7 00
U+2466
           (7)
                            isdia
                                                        CTRCLED DIGIT SEVEN
                             \Theta \Theta \Theta
                                                      numerics demo.py — Edited
U+2480
                              U+3285
                              1 import unicodedata
                              2 import re
                               re_digit = re.compile(r'\d')
                              6 sample = '1\xbc\xb2\u0969\u136b\u216b\u2466\u2480\u3285'
                               for char in sample:
                                   print('U+%04x' % ord(char),
                              9
                                                                                        # <A>
                                          char.center(6),
                                                                                        # <B>
                             10
                                          're dig' if re digit.match(char) else '-',
                                                                                        # <C>
                             11
                                          'isdig' if char.isdigit() else '-',
                             12
                                                                                        # <D>
                                          'isnum' if char.isnumeric() else '-',
                             13
                                                                                       # <F>
                                          format(unicodedata.numeric(char), '5.2f'),
                                                                                        # <F>
                             14
                             15
                                          unicodedata.name(char).
                                                                                        # <G>
                                          sep='\t')
                             16
                             17
                                                        Characters: 578 · Words: 57
```