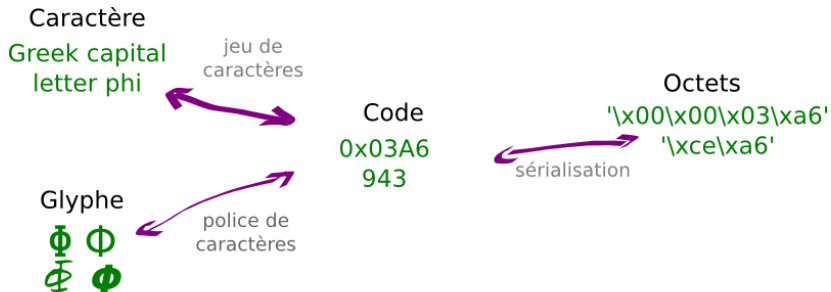
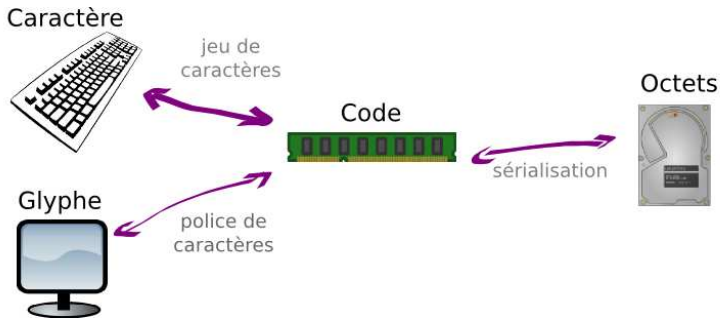


Python et Unicode

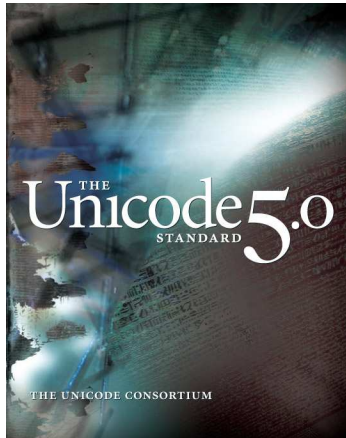
Victor Stinner

Pycon FR, Paris, mai 2009

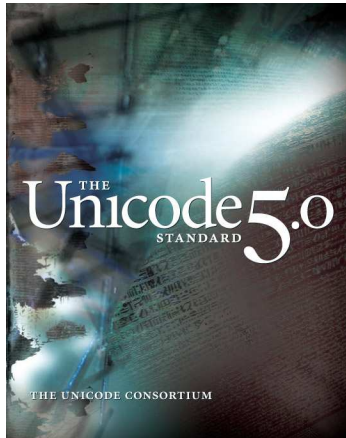




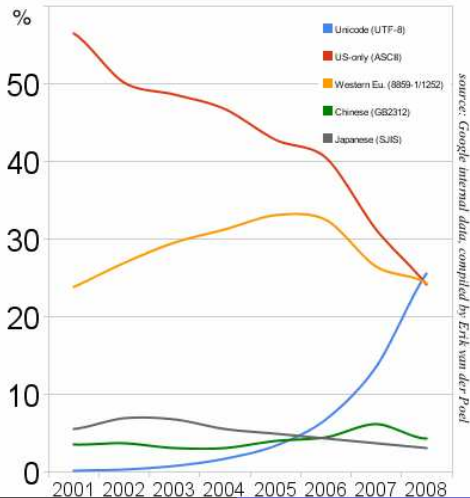
- 1961 : ASCII, 128 codes, Américain
- 1987 : ISO-8859-1, 256 codes, Europe Ouest
- 1990 : Unicode, millions de codes, Universel



- Surensemble de tous les jeux de caractères existants
- Indépendant de la langue et du matériel
- 1.114.112 codes (0x00..0x10FFFF)



Growth of Unicode on the Web





- Types str (octets) et unicode (caractères)
- "abc" : chaîne d'**octets**
- u"def" : chaîne de **caractères**



- Conversion **implicite** dans les deux sens
- "abc" + u"def" → "abcdef"
- "ab\xe9" + u"def" → UnicodeDecodeError
- Maux de tête douloureux



- Types bytes (octets) et str (caractères)
- "abc" : chaîne de **caractères**
- b"def" : chaîne d'**octets**
- b"abc" + "def" → TypeError
- Nouvelle bibliothèque io



- Python 2.x : ASCII par défaut
- Python 3.x : UTF-8 par défaut
- PEP 263 : "#coding:utf8"

```
print u"E accent aigu : é" (PEP 263)
# 0x00..0xff
print u"E accent aigu : \xe9"
# 0x00..0xffff
print u"Phi : \u03a6"
# 0x00..0x10ffff
print u"A (cunéiforme) : \U00120000"
print u"Phi : \N{GREEK CAPITAL LETTER PHI}"
```

```
>>> unicodedata.name(u"é")  
'LATIN SMALL LETTER E WITH ACUTE'  
>>> unicodedata.decomposition(u"é")  
'0065 0301'  
>>> unicodedata.normalize("NFC", u"u\u0308")  
u'ü'
```

- `codecs.open("fichier", "r",
encoding="utf-8")`
- `texte =
octets.decode("utf-8")`



- Fiable pour ASCII et UTF-8
- Heuristique pour les autres
- chardet (Mozilla)



- "caractères" →
u"caractères"
- from __future__ import
unicode_literals
- str(var) → unicode(var)
- def __str__(self) → def
__unicode__(self)



- `codecs.open("fichier", "w", encoding="utf-8")`
- `octets =
texte.encode("utf-8")`





- gucharmap / KCharSelect
- Livre Unicode
- BeautifulSoup
- Django 1.0 parle unicode

Fichier Affichage Rechercher Aller à Aide

Sans Gras Italique 129

Script

Table de caractères Détails du caractère

N'ko	◌̊	◌̋	அ	ஆ	இ	ஈ	உ	ஊ	எ
Nouveau Tai lü	ᨆ	ᨇ	ᨈ	ᨉ	ᨊ	ᨋ	ᨌ	ᨍ	ᨎ
Ogam	᠋	᠋	᠋	᠋	᠋	᠋	᠋	᠋	᠋
Oriya	ୠ	ୡ	ୢ	ୣ	୤	୦	୧	୨	୩
Osmanya	ᲀ	ᲁ	ᲂ	ᲃ	ᲄ	ᲅ	ᲆ	ᲇ	ᲈ
Ougantique	ᲀ	ᲁ	ᲂ	ᲃ	ᲄ	ᲅ	ᲆ	ᲇ	ᲈ
Phags-pa	ᠪ	ᠣ	ᠤ	ᠨ	ᠠ	ᠲ	ᠦ	ᠦ	ᠦ
Phénicien	𐤀	𐤁	𐤂	𐤃	𐤄	𐤅	𐤆	𐤇	𐤈
Rejang	ᨆ	ᨇ	ᨈ	ᨉ	ᨊ	ᨋ	ᨌ	ᨍ	ᨎ
Runes	ᚠ	ᚡ	ᚢ	ᚣ	ᚤ	ᚥ	ᚦ	ᚧ	ᚨ
Santālī	ᱠ	ᱡ	ᱢ	ᱣ	ᱤ	ᱥ	ᱦ	ᱧ	ᱨ
Saurashtra	𑀀	𑀁	𑀂	𑀃	𑀄	𑀅	𑀆	𑀇	𑀈
Shavien	ᲀ	ᲁ	ᲂ	ᲃ	ᲄ	ᲅ	ᲆ	ᲇ	ᲈ
Singhalais	ආ	භ	ඈ	ඉ	ඊ	උ	ඌ	ඍ	ඎ
Soudanais	ᲀ	ᲁ	ᲂ	ᲃ	ᲄ	ᲅ	ᲆ	ᲇ	ᲈ
Syloti nāgrī	ᱠ	ᱡ	ᱢ	ᱣ	ᱤ	ᱥ	ᱦ	ᱧ	ᱨ
Syriaque	ܐ	܂	܄	܆	܈	܊	܌	܎	ܐ
Tagal	ᜀ	ᜁ	ᜂ	ᜃ	ᜄ	ᜅ	ᜆ	ᜇ	ᜈ
Tagbanoua	ᜀ	ᜁ	ᜂ	ᜃ	ᜄ	ᜅ	ᜆ	ᜇ	ᜈ
Tai-le	ᨆ	ᨇ	ᨈ	ᨉ	ᨊ	ᨋ	ᨌ	ᨍ	ᨎ
Tamoul	ஊ	஋	஌	஍	எ	ஏ	உ	ஊ	஋
Télougou	ఊ	ఋ	ఌ	఍	ఎ	఑	ఒ	ఓ	ఔ
Thai	ก	ข	ฃ	ค	ฅ	ฉ	ช	ซ	ฌ
Thâna	ᨆ	ᨇ	ᨈ	ᨉ	ᨊ	ᨋ	ᨌ	ᨍ	ᨎ
Tibétain	ཀ	ཁ	ག	ང	ཅ	ཆ	ཇ	ཉ	ཏ
Tifinaghe	ⵀ	ⵁ	ⵂ	ⵃ	ⵄ	ⵅ	ⵆ	ⵇ	ⵈ
Vai	ᲀ	ᲁ	ᲂ	ᲃ	ᲄ	ᲅ	ᲆ	ᲇ	ᲈ

Texte à copier : Copier

U+0B82 TAMIL SIGN ANUSVARA * not used in Tamil

The screenshot shows the gucharmap application window. The menu bar includes 'Fichier', 'Affichage', 'Rechercher', 'Aller à', and 'Aide'. The font settings are 'Sans', 'Gras', 'Italique', and '29'. The left sidebar lists various scripts, with 'Tamil' selected. The main area displays the character 'U+0B82 TAMIL SIGN ANUSVARA' as a circle of dots with a larger dot above it. Below the character, the text reads: 'Propriétés générales du caractère', 'Présent dans Unicode depuis : 1.1', 'Catégorie Unicode : Marque, à chasse nulle', 'Diverses représentations utiles', 'UTF-8 : 0xE0 0xAE 0x82', 'UTF-16 : 0x0B82', 'UTF-8 en C octal échappé : \340\256\202', 'Entité décimale XML : ஂ', 'Annotations et références croisées', and 'Remarques : not used in Tamil'. At the bottom, there is a 'Texte à copier' field containing 'U+0B82 TAMIL SIGN ANUSVARA not used in Tamil' and a 'Copier' button.

- http://www.haypocalc.com/wiki/Python_Unicode
- *Unicode e Python 3* par Ezio Melotti
<http://wolf.netsons.org/pycon/unicode-python3.html>

- <http://www.flickr.com/photos/thorhakonsen/2349601149/>
- <http://www.flickr.com/photos/serendipitypeace/2267477928/>
- <http://www.flickr.com/photos/ecosnake/477735690/>