


```
>>> "naïve" == "naïve"  
True
```

```
>>> "naïve" == "naïve"
```

```
True
```

```
>>> "naïve" == "naïve"
```

```
False
```

```
>>> "naïve" == "naïve"
```

```
True
```

```
>>> "naïve" == "naïve"
```

```
False
```

```
>>> "A" == "A"
```

```
False
```

```
>>> for c in "HELLO":  
...     print(c)
```

O
L
L
E
H

```
>>> "Piñata" [:3]
```

```
'Pin'
```

```
>>> "Piñata" [3:]
```

```
'ñata'
```

```
>>> int("8")  
4
```

Why `len(" ") == 4` and other weird things you should know about strings in Python

Marie Roald & Yngve Mardal Moe

Slides:
github.com/MarieRoald/PyConUS25





Encoding strings



Comparing strings



Slicing strings

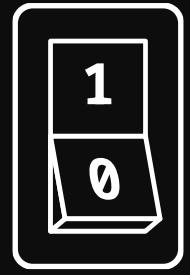
Slides:
github.com/MarieRoald/PyConUS25



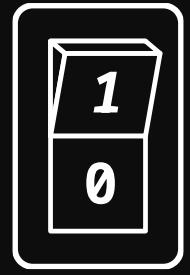
PART I

Encoding strings

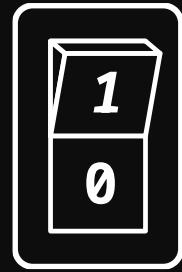
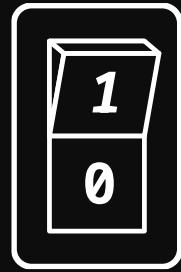
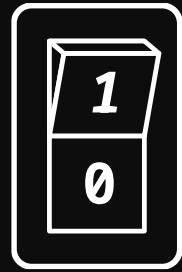
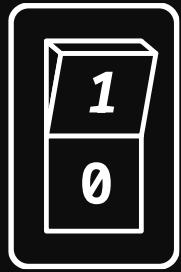
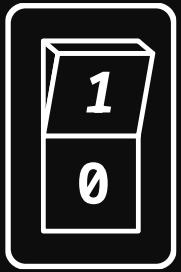
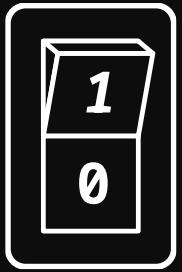
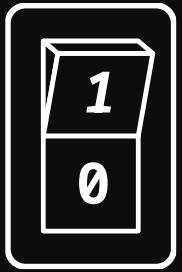
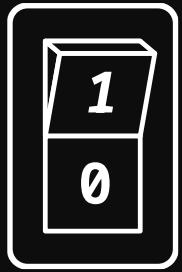




1



0



0

0

0

0

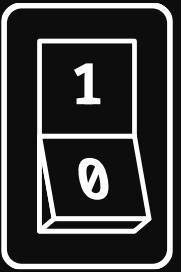
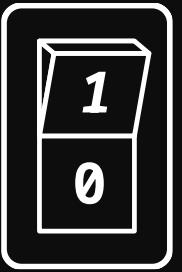
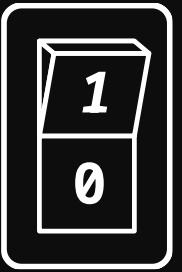
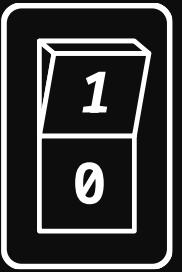
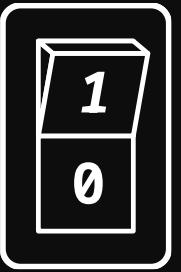
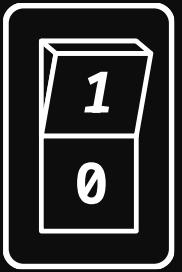
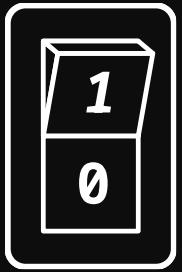
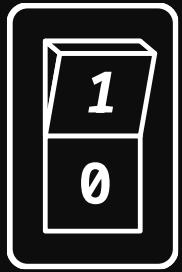
0

0

0

0

0



0

0

0

0

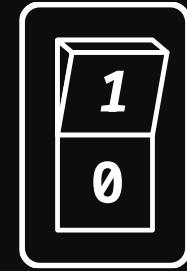
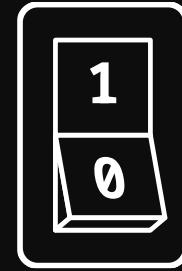
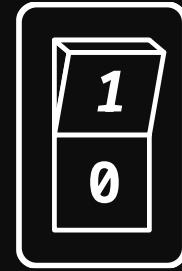
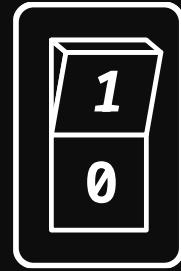
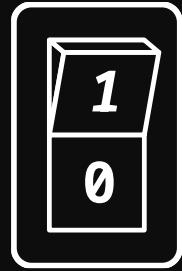
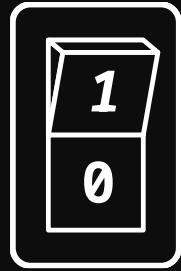
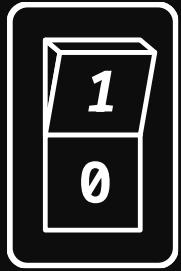
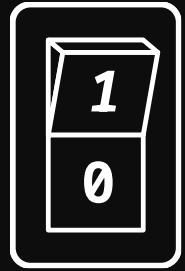
0

0

0

1

1



0

0

0

0

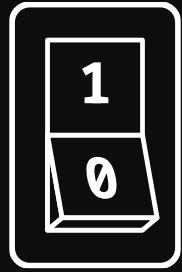
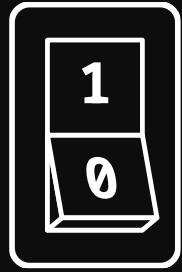
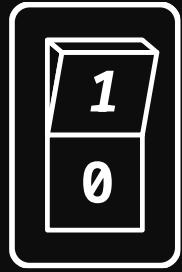
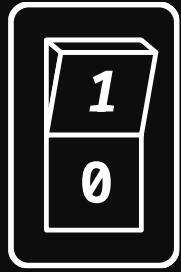
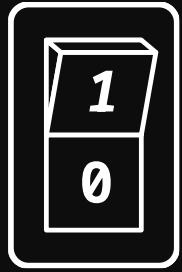
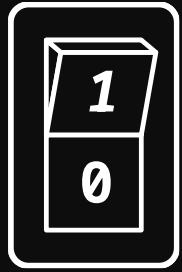
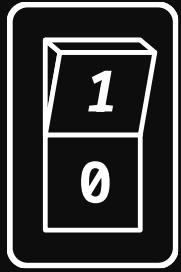
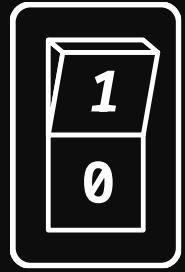
0

0

1

0

2



0

0

0

0

0

0

1

1

3

By enumerating letters, we can store them as integers

A	1
B	2
C	3
D	4
...	...
Z	26

By enumerating letters, we can store them as integers

	Code point
A	1
B	2
C	3
D	4
...	...
Z	26



Dette er en
morsom vits
~
~
(END)

ørttw rw rn
mÃwsÃm vrts
~
~
(END)

The ASCII standard from 1963 from contains 128 code points

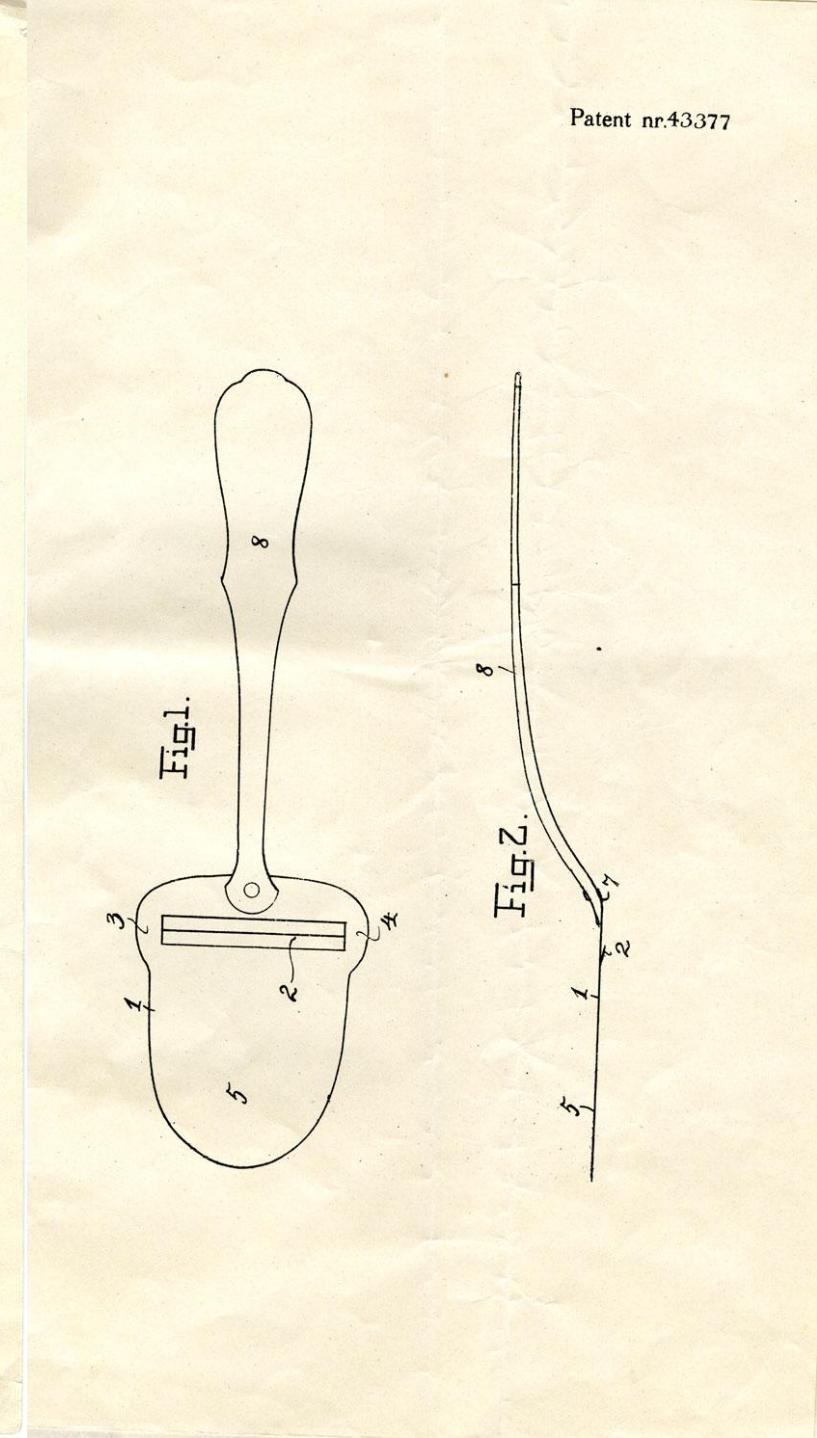
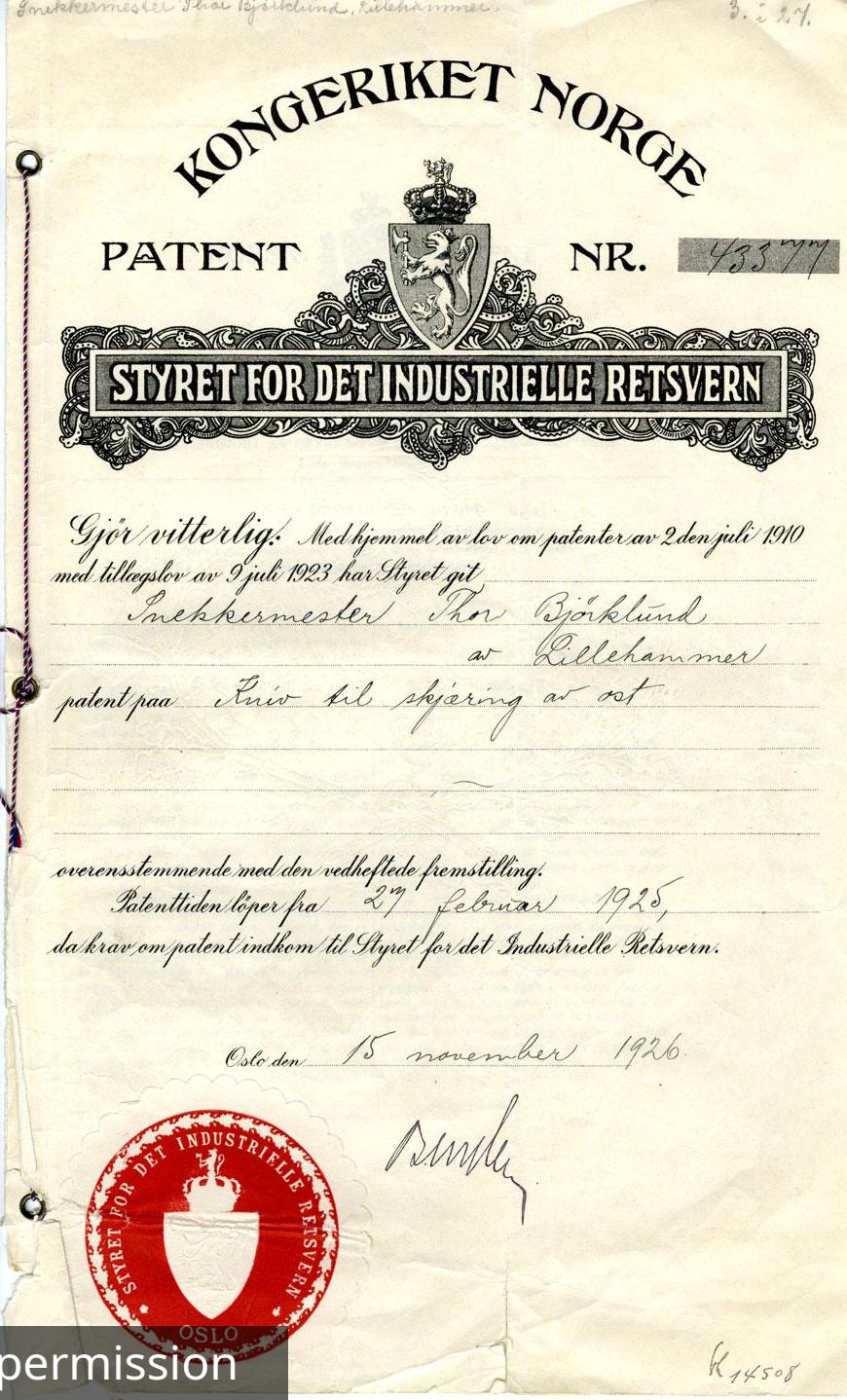
Row	Column	0	1	2	3	4	5	6	7
Bit Pattern	b7 b6 b5 b4 b3 b2 b1	0 0 0 0	0 0 1 1	0 1 0 0	1 1 0 0	1 0 1 1	1 1 1 1	1 1 1 1	1 1 1 1
		NUL	DLE	SP	0	@	P	~	P
0	0 0 0 0								
1	0 0 0 1	SOH	DC1	!	1	A	Q	a	q
2	0 0 1 0	STX	DC2	"	2	B	R	b	r
3	0 0 1 1	ETX	DC3	#	3	C	S	c	s
4	0 1 0 0	EOT	DC4	\$	4	D	T	d	t
5	0 1 0 1	ENQ	NAK	%	5	E	U	e	u
6	0 1 1 0	ACK	SYN	&	6	F	V	f	v
7	0 1 1 1	BEL	ETB	'	7	G	W	g	w
8	1 0 0 0	BS	CAN	(8	H	X	h	x
9	1 0 0 1	HT	EM)	9	I	Y	i	y
10	1 0 1 0	LF	SUB	*	:	J	Z	j	z
11	1 0 1 1	VT	ESC	+	;	K	[k	{
12	1 1 0 0	FF	FS	,	<	L	\	l	
13	1 1 0 1	CR	GS	-	=	M]	m	}
14	1 1 1 0	SO	RS	.	>	N	^	n	~
15	1 1 1 1	SI	US	/	?	O	-	o	DEL

Fig. 24.2 ISO 7-Bit Code



Ostehøvel





Patent nr.43377

```
>>> "ostehøvel".encode("ascii")
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
UnicodeEncodeError: 'ascii' codec
can't encode character '\xf8' in
position 5: ordinal not in range(128)
```

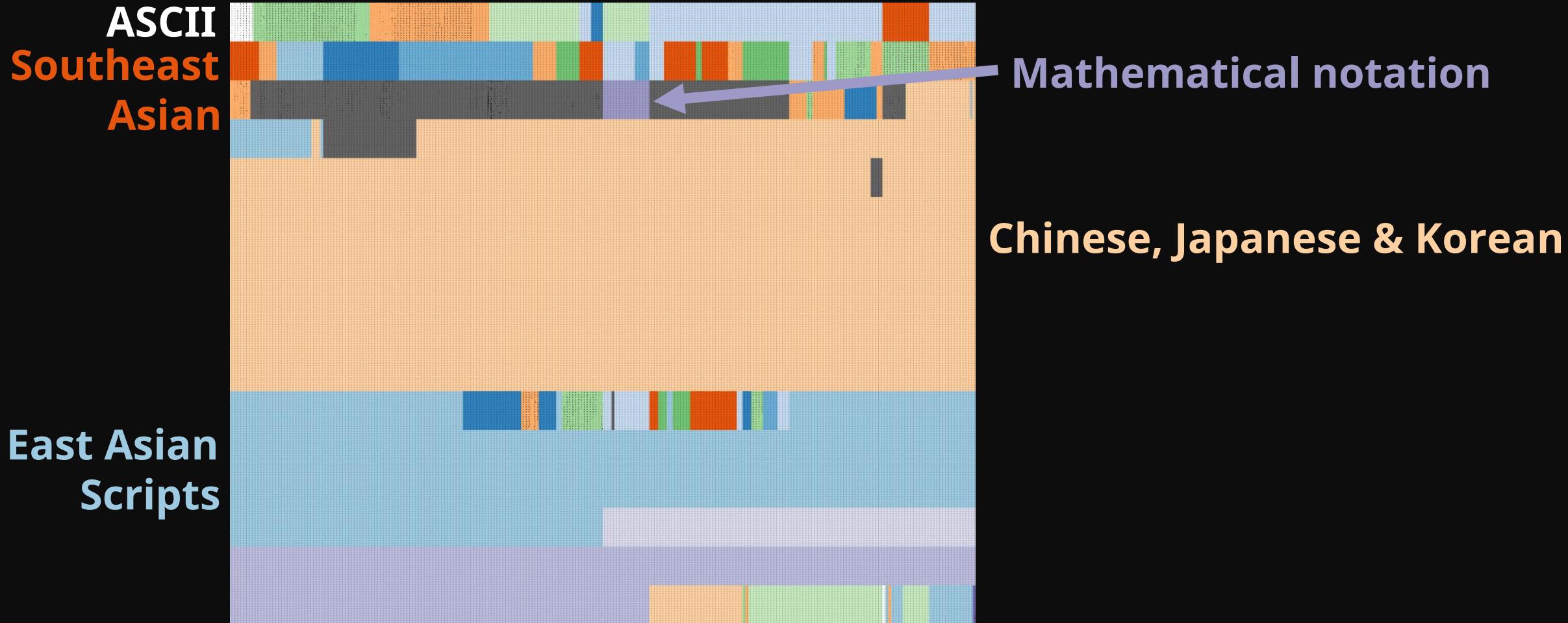
Since 1963, we have seen many language-specific standards for encoding text

Standard Encodings		
Codec	Aliases	Languages
ascii	646, us-ascii	English
big5	big5-tw, csbig5	Traditional Chinese
big5hkscs	big5-hkscs, hkscs	Traditional Chinese
cp037	IBM037, IBM039	English
cp273	273, IBM273, csIBM273	German <small>Added in version 3.4.</small>
cp424	EBCDIC-CP-HE, IBM424	Hebrew
cp437	437, IBM437	English
cp500	EBCDIC-CP-BE, EBCDIC-CP-CH, IBM500	Western Europe
cp720		Arabic
cp737		Greek

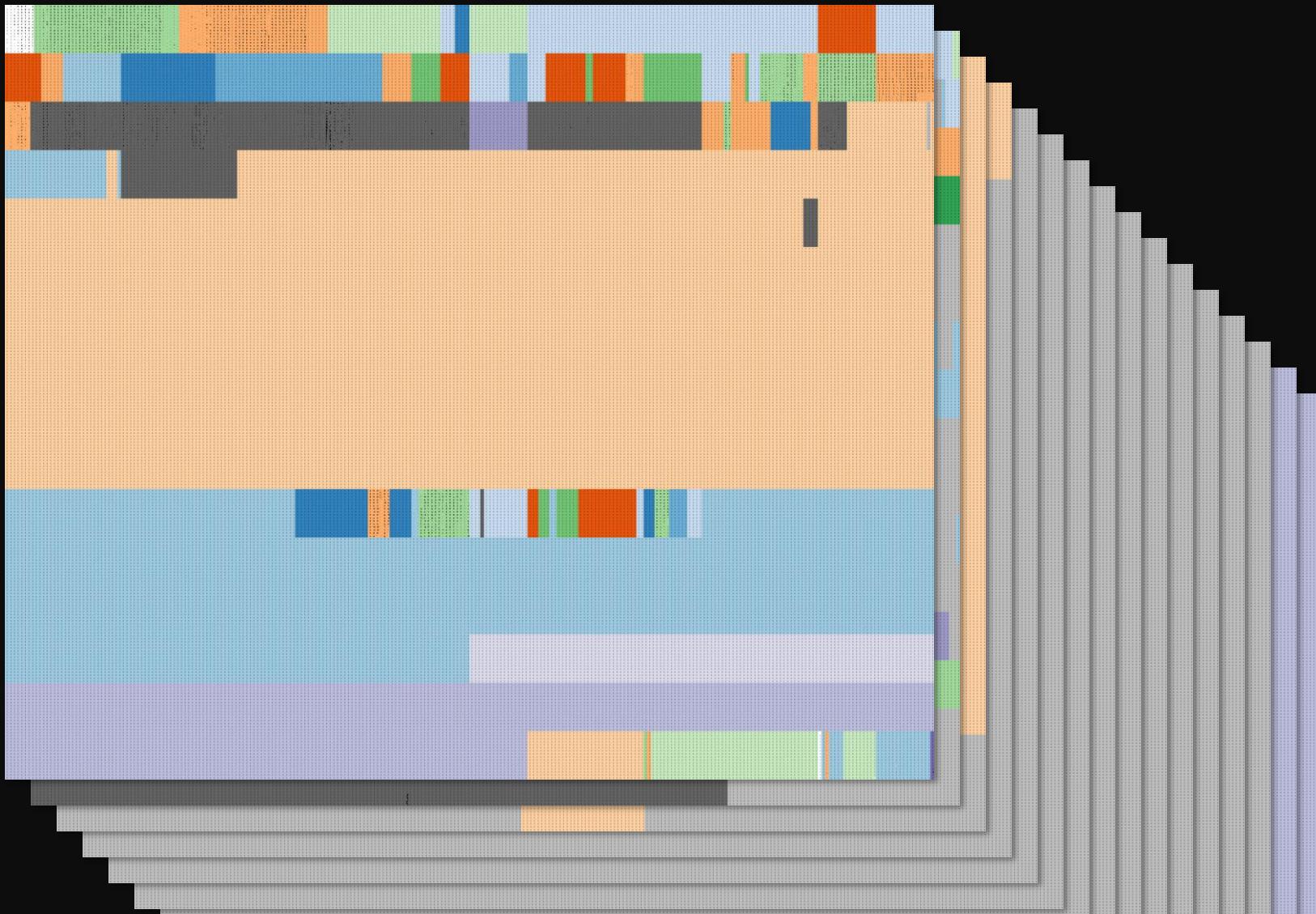
[...]

87 text encoding schemes

By using two bytes per code point, we can store
65 536 different characters



By using four bytes per code-point Unicode supports more than one million characters



UTF-8 strings take approximately four times more space than ASCII strings

```
>>> len("pycon".encode("ascii"))  
5
```

```
>>> len("pycon".encode("utf-32"))  
24
```

By using the first bit(s) of each code-point, we can encode characters compactly

00000000 = 0

By using the first bit(s) of each code-point, we can encode characters compactly

00000001 = 1

By using the first bit(s) of each code-point, we can encode characters compactly

00000010 = 2

By using the first bit(s) of each code-point, we can encode characters compactly

01111111 = 127

By using the first bit(s) of each code-point, we can encode characters compactly

11000010 10000000 = 128

By using the first bit(s) of each code-point, we can encode characters compactly

11000010 10000001 = 129

By using the first bit(s) of each code-point, we can encode characters compactly

11000010 10000010 = 130

The UTF-8 encoding can store ASCII-strings compactly

```
>>> len("pycon".encode("ascii"))  
5
```

```
>>> len("pycon".encode("utf-32"))  
24
```

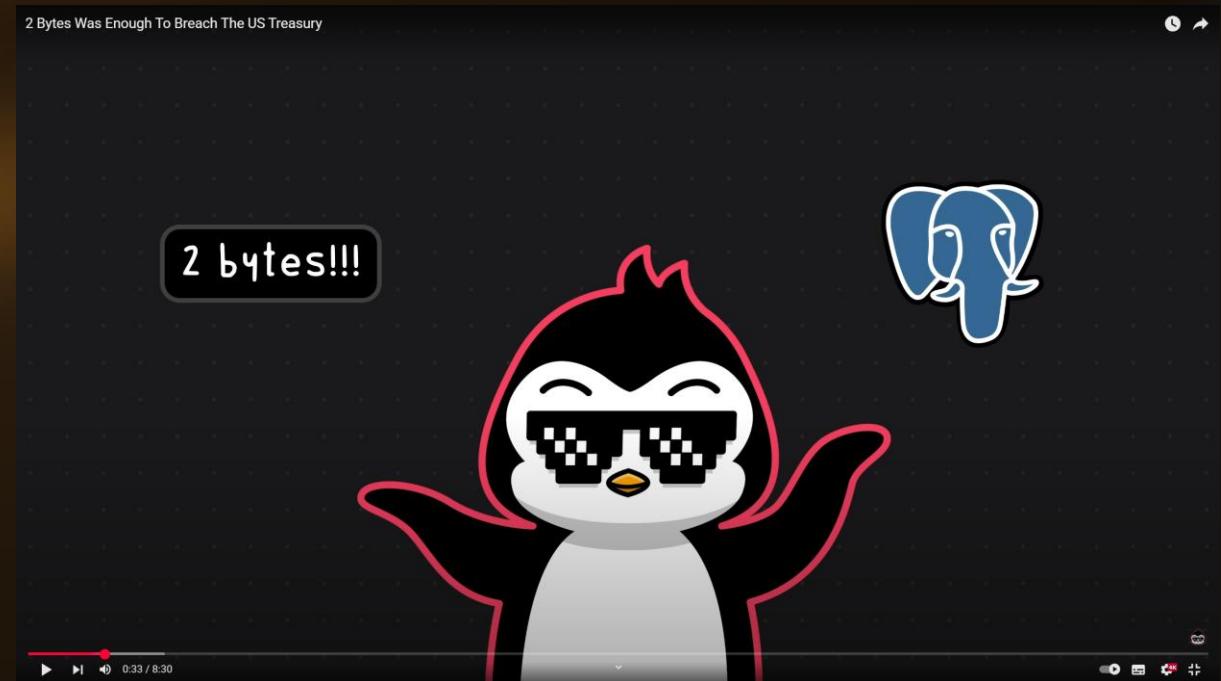
```
>>> len("pycon".encode("utf-8"))  
5
```

The variable-length encoding makes UTF-8 much more difficult to parse correctly

The screenshot shows a news article from The Register. The headline reads: "PostgreSQL bug lets bad guys break into US Treasury". The main text discusses how PostgreSQL's string escaping routines handle invalid UTF-8 characters, which can be exploited through SQL injection. It also mentions that running meta-commands like ! can execute shell commands on the operating system.

Fewer said: "Because of how PostgreSQL string escaping routines handle invalid UTF-8 characters, in combination with how invalid byte sequences within the invalid UTF-8 characters are processed by psql, an attacker can leverage CVE-2025-1094 to generate a SQL injection."

Running meta-commands can extend psql's functionality, and it's through these that an attacker can feasibly achieve ACE by using the exclamation mark meta-command to execute a shell command on the operating system. Attackers can also use the vulnerability to execute SQL statements of their choosing.



https://www.theregister.com/2025/02/14/postgresql_bug_treasury/
<https://www.youtube.com/watch?v=rgsIkZkfIMw> (PwnFunction)

Python has three types of unicode string encodings

Unicode Type

These are the basic Unicode object types used for the Unicode implementation in Python:

```
type Py_UCS4  
type Py_UCS2  
type Py_UCS1
```

Part of the [Stable ABI](#).

These types are typedefs for unsigned integer types wide enough to contain characters of 32 bits, 16 bits and 8 bits, respectively. When dealing with single Unicode characters, use [Py_UCS4](#).

Added in version 3.3.

If all characters in a string are in ASCII, then each code point is stored with one byte

```
>>> sys.getsizeof("a"*4097 + "\u0010")  
1042
```

If just one code point needs two bytes, then all code points are stored using two bytes

```
>>> sys.getsizeof("a"*4097 + "\u0010")  
1042  
>>> sys.getsizeof("a"*4097 + "\u0100")  
2060
```

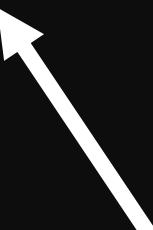
If just one code point needs more than two bytes,
then all are stored using four bytes

```
>>> sys.getsizeof("a"*4097+\u0010)  
1042  
>>> sys.getsizeof("a"*4097+\u0100)  
2060  
>>> sys.getsizeof("a"*4097+\u00040000)\ud83d\udc02  
4064
```

Each Unicode character is part of a category

```
>>> from unicodedata import category  
>>> category("8")  
'Nd'
```

Each Unicode character is part of a category

```
>>> from unicodedata import category  
>>> category("8")  
'Nd'  
  
Decimal number
```

```
>>> int("8")  
4
```

We can get more information about characters by looking at their names

```
>>> from unicodedata import name  
>>> name("8")  
'BENGALI DIGIT FOUR'
```

Unicode code-points can be used to do more than just encode characters

```
>>> from unicodedata import category  
>>> category("HELLO"[0])  
'Cf'  
  
Format
```

Unicode code-points can be used to do more than just encode characters

```
>>> from unicodedata import name  
>>> name("HELLO"[0])  
'LEFT-TO-RIGHT EMBEDDING'
```

Unicode code-points can be used to do more than just encode characters

```
>>> from unicodedata import name  
>>> name("HELLO"[0])  
'LEFT-TO-RIGHT EMBEDDING'  
>>> name("HELLO"[1])  
'RIGHT-TO-LEFT OVERRIDE'
```

Text-formatting characters can sometimes lead to surprising results

```
>>> print("HELLO"[2:])  
OLLEH
```

Part II

Comparing strings



Do you approach the dragon? [y/n]



```
>>> answer = input()  
>>> answer is "y"
```

```
>>> "y" is "y"
```

```
>>> "y" is "y"
01_is_compare.py:1: SyntaxWarning:
"is" with 'str' literal. Did you
mean "=="?
```

```
>>> "y" is "y"
01_is_compare.py:1: SyntaxWarning:
"is" with 'str' literal. Did you
mean "=="?
True
```

To save time and memory, Python stores some strings just once in memory

String interning

Interned strings are conceptually part of an interpreter-global set of interned strings, meaning that:

- no two interned strings have the same content (across an interpreter);
- two interned strings can be safely compared using pointer equality (Python `is`).

This is used to optimize dict and attribute lookups, among other things.

Python uses two different mechanisms to intern strings: singletons and dynamic interning.

To save time and memory, Python stores some strings just once in memory

String interning

Interned strings are conceptually part of an interpreter-global set of interned strings, meaning that:

- no two interned strings have the same content (across an interpreter);
- two interned strings can be safely compared using pointer equality (Python `is`).

This is used to optimize dict and attribute lookups, among other things.

Python uses two different mechanisms to intern strings: singletons and dynamic interning.

Not all strings are interned by Python

```
>>> "a"*4097 is "a"*4097  
False
```

The equal operator checks if the contents of two strings are equal

```
>>> "y" == "y"
```

```
True
```

Sometimes, we want to compare strings without caring about case

```
>>> "Y" == "y"
```

```
False
```

Sometimes, we want to compare strings without caring about case

```
>>> "Y".lower() == "y"  
True
```



+43471320300

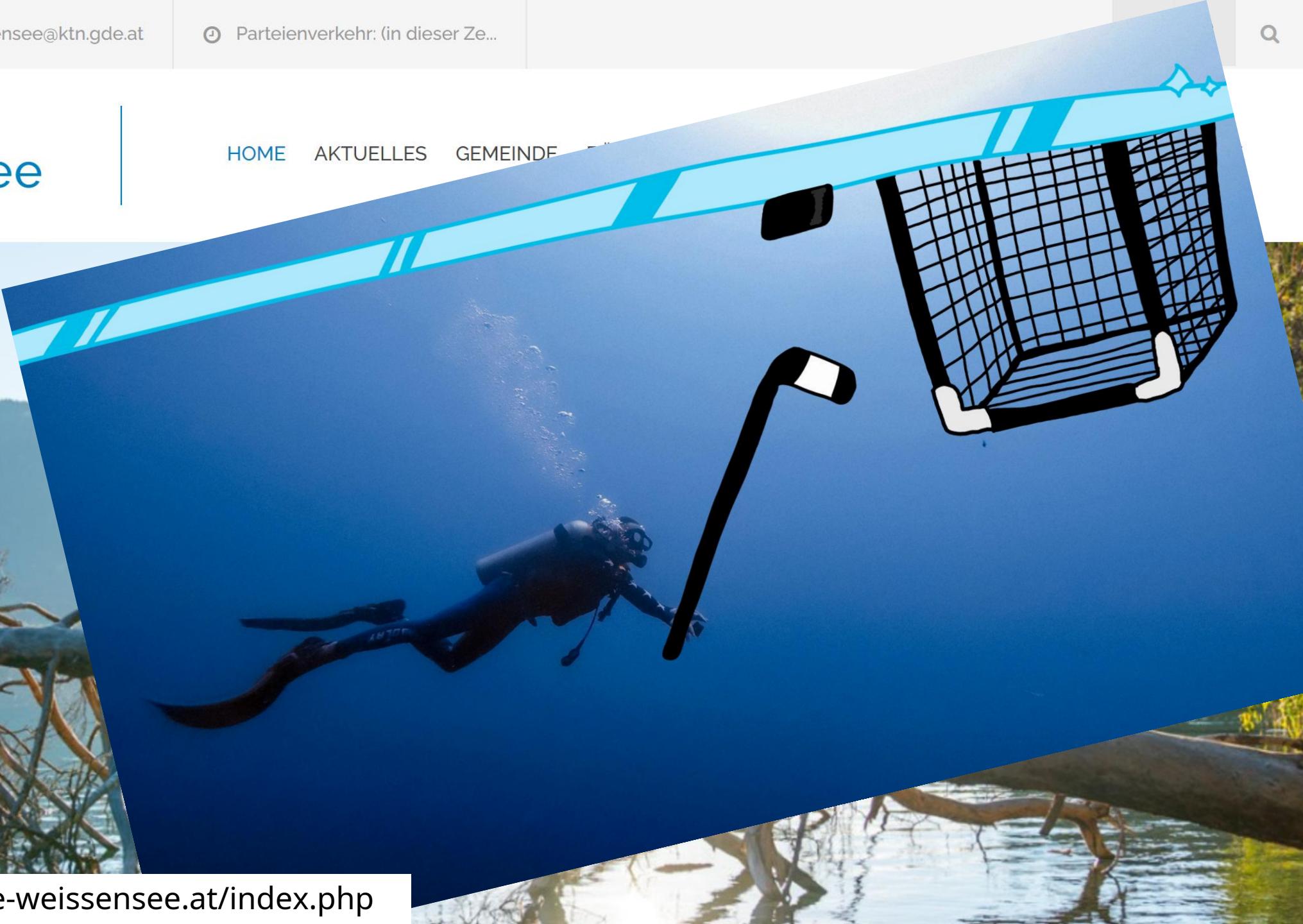
weissensee@ktn.gde.at

Parteienverkehr: (in dieser Ze...)



Gemeinde
Weißensee

HOME AKTUELLES GEMEINDE



<https://www.gemeinde-weissensee.at/index.php>

+43471320300

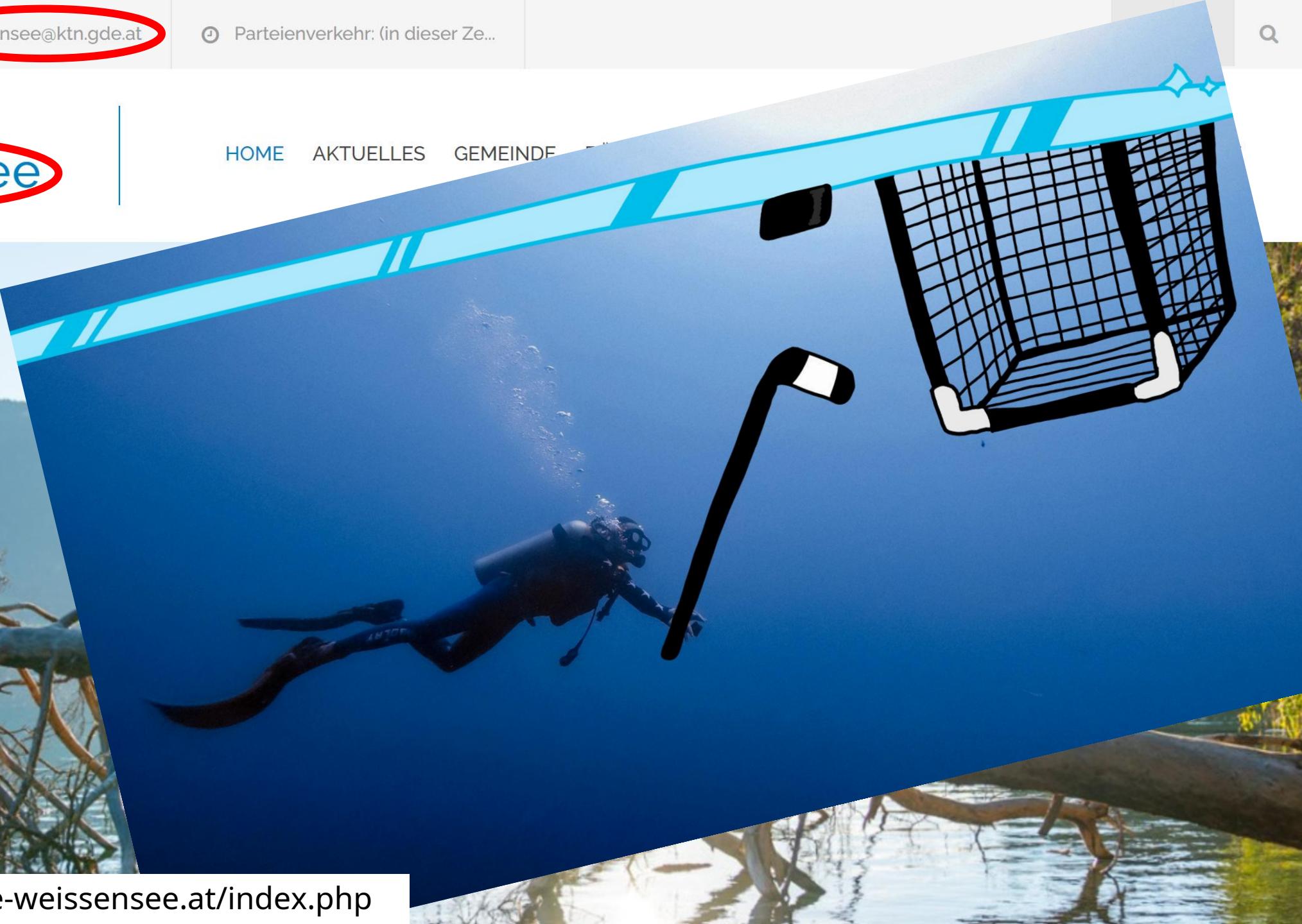
✉ weissensee@ktn.gde.at

ⓘ Parteienverkehr: (in dieser Ze...



Gemeinde
Weißensee

HOME AKTUELLES GEMEINDE



Using case transformations to compare strings is not robust

```
>>> "Weissensee".upper() ==  
    "Weißensee".upper()  
True
```

Using case transformations to compare strings is not robust

```
>>> "Weissensee".lower() ==  
    "Weißensee".lower()  
False
```

We should casifold for caseless string comparison

3.13.3 Default Case Folding

Case folding is related to case conversion. However, the main purpose of case folding is to contribute to caseless matching of strings, whereas the main purpose of case conversion is to put strings into a particular cased form.

Default Case Folding does not preserve normalization forms. A string in a particular Unicode normalization form may not be in that normalization form after it has been casfolded.

Python has builtin support for casefolding

```
>>> "Weissensee".casfold() ==  
"Weißensee".casfold()  
True
```

D R T ካ O i S ቅ F Y A J E ቁ ቃ A T ገ ቃ W
ዶ የ G M ዓ ቀ O H ፕ የ ቅ ቃ G ለ h Z ዓ O I ስ
ጥ ማ ም ሰ ሰ 4 የ ቅ ስ R ለ W S ሰ J J V S ሰ
፩ ሰ L C ቅ ም P C ስ h K d C G ዘ ሰ ቃ 6 ስ
፪ ካ R G B G ዓ 'A l 'A l 'A l 'A l 'A l 'A l 'H l 'H l 'H l
'H l 'H l 'O l 'O l 'O l 'O l 'O l 'O l A l H l O p B 'n ገ ዕ .



Pålegg





Pålegg

Pålegg

Pålegg

Pålegg

Some strings are equivalent, but different

```
>>> "pålegg" == "pålegg"  
False
```

Some strings are equivalent, but different

\u00e5

```
>>> "p\u00e5legg" == "p\u00e5legg"  
False
```

Some strings are equivalent, but different

```
\u00E5 \u0061\u030A  
[ ] [ ]  
>>> "pålegg" == "p\u00e5legg"  
False
```

Piñata != **Piñata**

FJÄLLBO != **FJÄLLBO**

Beyoncé != **Beyoncé**

Normalization makes equivalent strings equal

5 Normalization and Case

This section discusses issues that must be taken into account when considering normalization and case folding of identifiers in programming languages or scripting languages. Using normalization avoids many problems where apparently identical identifiers are not treated equivalently. Such problems can appear both during compilation and during linking—in particular across different programming languages. To avoid such problems, programming languages can normalize identifiers before storing or comparing them. Generally if the programming language has case-sensitive identifiers, then Normalization Form C is appropriate; whereas, if the programming language has case-insensitive identifiers, then Normalization Form KC is more appropriate.

Implementations that take normalization and case into account have two choices: to treat variants as equivalent, or to disallow variants.

There are four Unicode normalization forms

	Composed	Decomposed
Canonical	NFC	NFD
Compatibility	NFKC	NFKD

There are four Unicode normalization forms

	Composed	Decomposed
Canonical	NFC	NFD
Compatibility	NFKC	NFKD

The canonical composed form tries to represent the text with as few code-points as possible

```
>>> from unicodedata import  
... normalize  
>>> len(normalize("NFC", "pålegg"))  
6
```

The canonical decomposed form tries to split the “base characters” from their accents

```
>>> from unicodedata import  
... normalize  
>>> len(normalize("NFD", "pålegg"))
```

	Composed	Decomposed
Canonical	NFC	NFD
Compatibility	NFKC	NFKD

Some characters are stylistic subsets of other characters

Symbols:

½ ₣ % «»

Superscript:

² ⁵ ₜ ₛ

Fancy letters:

ℒ a.m. ™ ⁷

**And over
3800 more!**

The compatibility normalizations transform these characters into their unrestricted forms

```
>>> from unicodedata import  
... normalize  
>>> normalize("NFKC", "½")  
'1/2'
```

Python uses compatibility normalisation for parsing identifiers

```
>>> message = "Hello world"  
>>> print(message)  
Hello world
```

README



Fancifypy

Are you tired of boring Python code? Or do you want to make a real impression with your next PR? Look no further! **Fancify_{py}** will transform any plain and boring Python code into a feast for the eyes.

```
pipx run fancifypy code "print(range(5))" --seed 42
```



```
print ( range ( 5 ) )
```



Alternatively, you can fancify files directly

```
pipx run fancifypy path script.py
```



and defancified files (if you for some reason want boring files again)

```
pipx run fancifypy boring script.py
```



Note that **Fancify_{py}** leaves somewhat odd whitespaces, so we recommend using [Ruff](#) or [Black](#) afterwards to format the code.

Is this safe?

“Confusables” are characters that look alike but are different

```
>>> "A" == "A"
```

```
False
```

Unicode® Technical Standard #39

UNICODE SECURITY MECHANISMS

Version	16.0.0
Editors	Mark Davis (markdavis@google.com), Michel Suignard (michel@suignard.com)
Date	2024-09-03
This Version	https://www.unicode.org/reports/tr39/tr39-30.html
Previous Version	https://www.unicode.org/reports/tr39/tr39-28.html
Latest Version	https://www.unicode.org/reports/tr39/
Latest Proposed Update	https://www.unicode.org/reports/tr39/proposed.html
Revision	30

Summary

Because Unicode contains such a large number of characters and incorporates the varied writing systems of the world, incorrect usage can expose programs or systems to possible security attacks. This document specifies mechanisms that can be used to detect possible security problems.

Status

This document has been reviewed by Unicode members and other interested parties, and has been approved for publication by the Unicode Consortium. This is a stable document and may be used as reference material or cited as a normative reference by other specifications.

A **Unicode Technical Standard (UTS)** is an independent specification. Conformance to the Unicode Standard does not imply conformance to any UTS.

Please submit corrigenda and other comments with the online reporting form [[Feedback](#)]. Related information that is useful in understanding this document is found in the [References](#). For the latest version of the Unicode Standard, see [[Unicode](#)]. For a list of current Unicode Technical Reports, see [[Reports](#)]. For more information about versions of the Unicode Standard, see [[Versions](#)].

Contents

- 1 [Introduction](#)
- 2 [Conformance](#)

<https://www.unicode.org/reports/tr39/>

<https://www.unicode.org/faq/security.html>

<https://www.unicode.org/Public/security/16.0.0/intentional.txt>

<https://www.unicode.org/Public/security/16.0.0/confusables.txt>

Some characters require locale-specific handling

```
>>> "I".lower() == "ı"  
False
```

Some letters share the same code point

\u014a

N

Noto sans

Some letters share the same code point

\u014a

N

Noto sans

\u014a

n

Arial

Part III

Slicing strings

```
>>> len(" ")
```

```
>>> len(" ")
```

4

```
>>> for c in "█":  
...     print(c, hex(ord(c)))
```

```
>>> for c in "FLAG":  
...     print(c, hex(ord(c)))  
FLAG: 1f3f3
```

```
>>> for c in "FLAG":  
...     print(c, hex(ord(c)))  
FLAG: 1f3f3  
FLAG: fe0f
```

```
>>> for c in "FLAG":  
...     print(c, hex(ord(c)))  
FLAG: 1f3f3  
L: fe0f  
G: 200d
```

```
>>> for c in "🏳️‍🌈":
...     print(c, hex(ord(c)))
🏳️: 1f3f3
⚧: fe0f
⋮: 200d
🌈: 1f308
```

```
>>> print("🌈)[:2]
```



```
>>> for c in "🏳️‍🌈":
...     print(c, hex(ord(c)))
🏳️: 1f3f3
⚧: fe0f
⋮: 200d
🌈: 1f308
```

The image shows two smartphones side-by-side, both displaying the PyCon US 2025 mobile application. The phones are set against a dark wooden background.

Left Phone Screen:

- Header:** Shows the URL us.pycon.org/2025/schedule, the time 12:45, and battery level 57%.
- Event Title:** "Why `len('😊') == 4` and other weird things you should know about strings in Python" (with a yellow emoji face icon).
- Details:** Saturday, May 17th, 2025 3:15 p.m.–3:45 p.m. in Ballroom BC.
- Presented by:** Marie Roald, Yngve Mardal Moe.
- Experience Level:** Some experience.
- Description:** A brief text explaining string behavior in Python, followed by a snippet of code: `len('😊') == 4, 'ñ' != 'ñ', 'world hello'.split('t') == ['olleh'] How is this`.

Right Phone Screen:

- Header:** Shows the URL us.pycon.org/2025/schedule, the time 12:45, and battery level 60%.
- Event Title:** "Why `len('😊') == 4` and other weird things you should know about strings in Python" (with a yellow emoji face icon).
- Details:** Saturday, May 17th, 2025 3:15 p.m.–3:45 p.m. in Ballroom BC.
- Presented by:** Marie Roald, Yngve Mardal Moe.
- Experience Level:** Some experience.
- Description:** A detailed description of the talk's content.

Accented letters can also form grapheme clusters

```
>>> "pålegg" [:1]  
'pa'
```

Accented letters can also form grapheme clusters

```
>>> "pålegg"[:1]  
'pa'
```

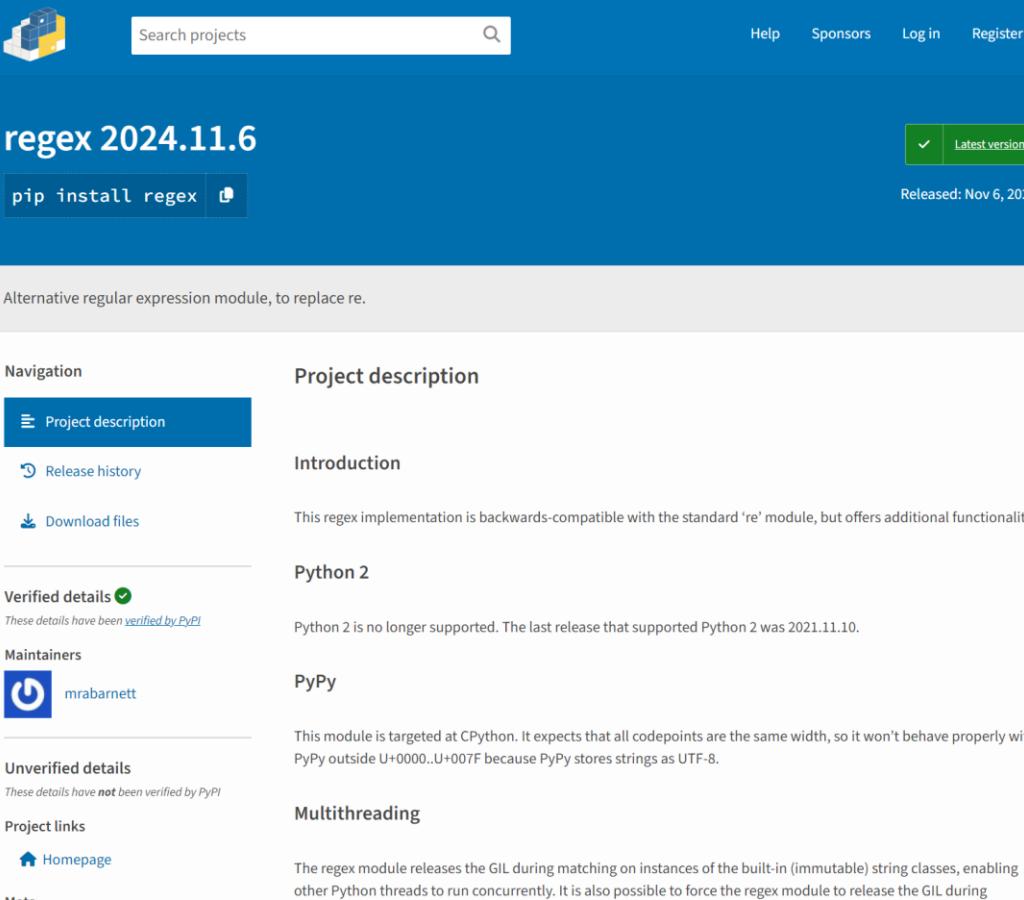
```
>>> from unicodedata import  
... normalize  
>>> len(normalize("NFC", "á"))  
2
```

Korean characters can also form grapheme clusters

>>> "ㅃ" [:2]
'ㅃ'

There are several open source tools for better Unicode handling in Python

Regex - Expands the builtin re module



The screenshot shows the PyPI project page for the `regex` package. The header includes the PyPI logo, a search bar, and navigation links for Help, Sponsors, Log in, and Register. The main title is `regex 2024.11.6`, with a green button labeled "Latest version". Below the title, there's a pip install link and a release date of Nov 6, 2024. A brief description states: "Alternative regular expression module, to replace re." The left sidebar has a "Navigation" section with "Project description" (selected), "Release history", and "Download files". The "Project description" section contains an "Introduction" paragraph about being backwards-compatible with the standard `re` module. It also mentions "Python 2" support, noting it's no longer supported. The "Maintainers" section lists "mrabarnett" with a GitHub icon. The "Unverified details" section notes that details have not been verified by PyPI. The "Project links" section includes a "Homepage" link.

regex 2024.11.6

pip install regex

Released: Nov 6, 2024

Alternative regular expression module, to replace re.

Navigation

Project description

Release history

Download files

Verified details ✓

These details have been [verified by PyPI](#)

Maintainers

PyPy

Unverified details

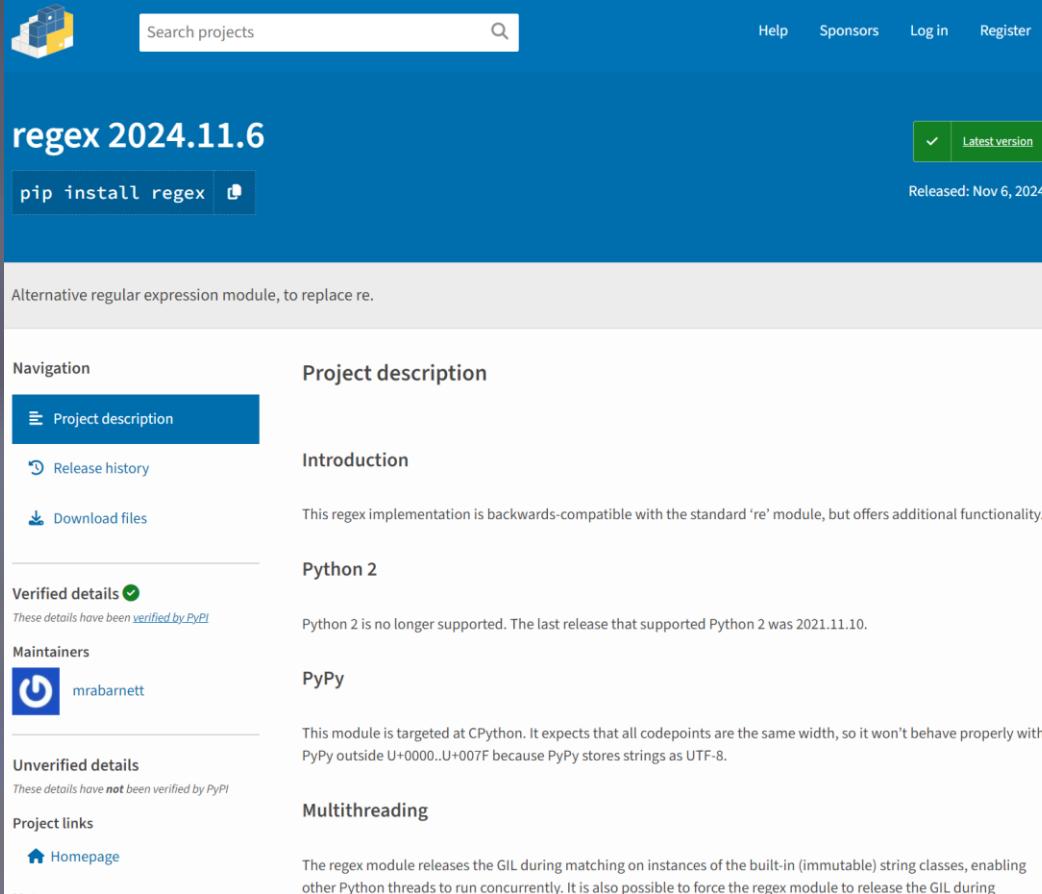
These details have [not](#) been verified by PyPI

Project links

Homepage

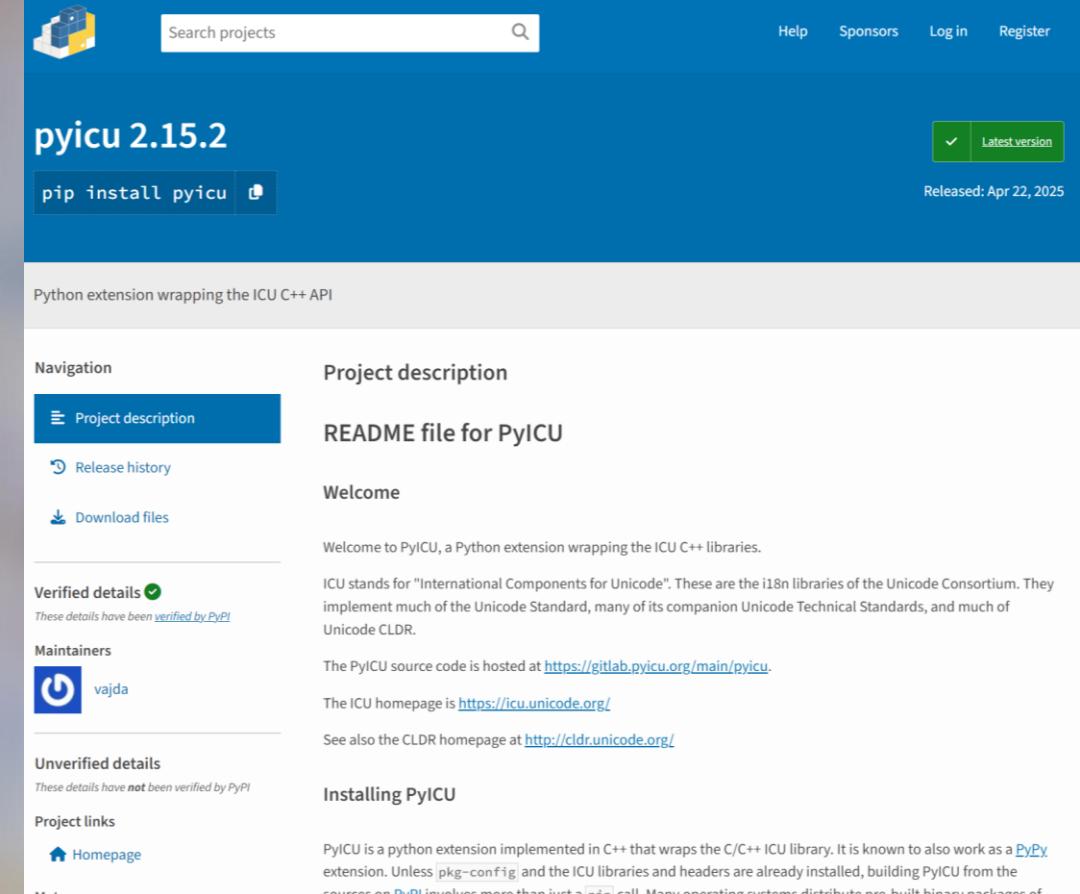
There are several open source tools for better Unicode handling in Python

Regex - Expands the builtin `re` module



The screenshot shows the PyPI project page for the `regex` package. The header includes the Python logo, a search bar, and navigation links for Help, Sponsors, Log in, and Register. The main title is `regex 2024.11.6`, with a green "Latest version" button. Below the title, there's a "pip install regex" button and a release date of Nov 6, 2024. A brief description states: "Alternative regular expression module, to replace re.". The left sidebar has a "Navigation" section with "Project description" (selected), "Release history", and "Download files". The "Project description" section contains sections for "Introduction", "Python 2" (noted as unsupported), "PyPy" (targeted at CPython), and "Multithreading".

PyICU - Wraps the C/C++ ICU library



The screenshot shows the PyPI project page for the `pyicu` package. The header includes the Python logo, a search bar, and navigation links for Help, Sponsors, Log in, and Register. The main title is `pyicu 2.15.2`, with a green "Latest version" button. Below the title, there's a "pip install pyicu" button and a release date of Apr 22, 2025. A brief description states: "Python extension wrapping the ICU C++ API". The left sidebar has a "Navigation" section with "Project description" (selected), "Release history", and "Download files". The "Project description" section contains sections for "README file for PyICU", "Welcome", and "Installing PyICU". It also mentions the ICU library and its homepage.

You can't just count grapheme clusters when you measure string length

weird text

 Use UTF-8

open(f, encoding="utf-8")

 Use UTF-8

```
open(f, encoding="utf-8")
```

 Casefold for caseless matching

```
"Weissensee".casefold()
```

 Use UTF-8

```
open(f, encoding="utf-8")
```

 Casefold for caseless matching

```
"Weissensee".casefold()
```

 Normalise your strings

```
normalize("NFC", "Pålegg")
```

 Use UTF-8

`open(f, encoding="utf-8")`

 Consider security implications

Read [Unicode TR39](#)

 Casefold for caseless matching

`"Weissensee".casefold()`

 Normalise your strings

`normalize("NFC", "Pålegg")`

 Use UTF-8

```
open(f, encoding="utf-8")
```

 Casefold for caseless matching

```
"Weissensee".casefold()
```

 Normalise your strings

```
normalize("NFC", "Pålegg")
```

 Consider security implications

Read **Unicode TR39**



Characters are not code points

Use **extended grapheme clusters** for visual length

 Use UTF-8

```
open(f, encoding="utf-8")
```

 Casefold for caseless matching

```
"Weissensee".casefold()
```

 Normalise your strings

```
normalize("NFC", "Pålegg")
```

 Consider security implications

Read **Unicode TR39**



Characters are not code points

Use **extended grapheme clusters** for visual length



Unicode is imperfect and not all aspects of Unicode is builtin in Python!

