

Navigation

- [index](#)
- [modules](#) |
- [next](#) |
- [previous](#) |
- [Python v2.6.4 documentation](#) »
- [The Python Standard Library](#) »
- [8. String Services](#) »

8.9. unicodedata — Unicode Database¶

This module provides access to the Unicode Character Database which defines character properties for all Unicode characters. The data in this database is based on the `UnicodeData.txt` file version 5.1.0 which is publicly available from <ftp://ftp.unicode.org/>.

The module uses the same names and symbols as defined by the UnicodeData File Format 5.1.0 (see <http://www.unicode.org/Public/5.1.0/ucd/UCD.html>). It defines the following functions:

`unicodedata.lookup(name)`¶

Look up character by name. If a character with the given name is found, return the corresponding Unicode character. If not found, `KeyError` is raised.

`unicodedata.name(unichr, default)`¶

Returns the name assigned to the Unicode character `unichr` as a string. If no name is defined, `default` is returned, or, if not given, `ValueError` is raised.

`unicodedata.decimal(unichr, default)`¶

Returns the decimal value assigned to the Unicode character `unichr` as integer. If no such value is defined, `default` is returned, or, if not given, `ValueError` is raised.

`unicodedata.digit(unichr, default)`¶

Returns the digit value assigned to the Unicode character `unichr` as integer. If no such value is defined, `default` is returned, or, if not given, `ValueError` is raised.

`unicodedata.numeric(unichr, default)`¶

Returns the numeric value assigned to the Unicode character `unichr` as float. If no such value is defined, `default` is returned, or, if not given, `ValueError` is raised.

`unicodedata.category(unichr)`¶

Returns the general category assigned to the Unicode character `unichr` as string.

`unicodedata.bidirectional(unichr)`¶

Returns the bidirectional category assigned to the Unicode character `unichr` as string. If no such value is defined, an empty string is returned.

`unicodedata.combining(unichr)`¶

Returns the canonical combining class assigned to the Unicode character `unichr` as integer. Returns 0 if no combining class is defined.

`unicodedata.east_asian_width(unichr)`¶

Returns the east asian width assigned to the Unicode character `unichr` as string.

New in version 2.4.

`unicodedata.mirrored(unichr)`¶

Returns the mirrored property assigned to the Unicode character `unichr` as integer. Returns 1 if the character has been identified as a “mirrored” character in bidirectional text, 0 otherwise.

`unicodedata.decomposition(unichr)`¶

Returns the character decomposition mapping assigned to the Unicode character `unichr` as string. An empty string is returned in case no such mapping is defined.

`unicodedata.normalize(form, unistr)`¶

Return the normal form `form` for the Unicode string `unistr`. Valid values for `form` are ‘NFC’, ‘NFKC’, ‘NFD’, and ‘NFKD’.

The Unicode standard defines various normalization forms of a Unicode string, based on the definition of canonical equivalence and compatibility equivalence. In Unicode, several characters can be expressed in various way. For example, the character U+00C7 (LATIN CAPITAL LETTER C WITH CEDILLA) can also be expressed as the sequence U+0327 (COMBINING CEDILLA) U+0043 (LATIN CAPITAL LETTER C).

For each character, there are two normal forms: normal form C and normal form D. Normal form D (NFD) is also known as canonical decomposition, and translates each character into its decomposed form. Normal form C (NFC) first applies a canonical decomposition, then composes pre-combined characters again.

In addition to these two forms, there are two additional normal forms based on compatibility equivalence. In Unicode, certain characters are supported which normally would be unified with other characters. For example, U+2160 (ROMAN NUMERAL ONE) is really the same thing as U+0049 (LATIN CAPITAL LETTER I). However, it is supported in Unicode for compatibility with existing character sets (e.g. gb2312).

The normal form KD (NFKD) will apply the compatibility decomposition, i.e. replace all compatibility characters with their equivalents. The normal form KC (NFKC) first applies the compatibility decomposition, followed by the canonical composition.

Even if two unicode strings are normalized and look the same to a human reader, if one has combining characters and the other doesn't, they may not compare equal.

New in version 2.3.

In addition, the module exposes the following constant:

```
unicodedata.unidata_version
```

The version of the Unicode database used in this module.

New in version 2.3.

```
unicodedata.ucd_3_2_0
```

This is an object that has the same methods as the entire module, but uses the Unicode database version 3.2 instead, for applications that require this specific version of the Unicode database (such as IDNA).

New in version 2.5.

Examples:

```
>>> import unicodedata
>>> unicodedata.lookup('LEFT CURLY BRACKET')
u'{'
>>> unicodedata.name(u'/')
'SOLIDUS'
>>> unicodedata.decimal(u'9')
9
>>> unicodedata.decimal(u'a')
Traceback (most recent call last):
  File "<stdin>", line 1, in ?
ValueError: not a decimal
>>> unicodedata.category(u'A') # 'L'etter, 'u'ppercase
'Lu'
>>> unicodedata.bidirectional(u'\u0660') # 'A'rabic, 'N'umber
'AN'
```

Previous topic

[8.8. codecs — Codec registry and base classes](#)

Next topic

[8.10. stringprep — Internet String Preparation](#)

This Page

- [Show Source](#)

Navigation

- [index](#)
- [modules](#) |
- [next](#) |
- [previous](#) |
- [Python v2.6.4 documentation](#) »
- [The Python Standard Library](#) »
- [8. String Services](#) »

© [Copyright](#) 1990-2010, Python Software Foundation.

The Python Software Foundation is a non-profit corporation. [Please donate.](#)

Last updated on Feb 26, 2010. Created using [Sphinx](#) 0.6.3.