

2. Lexical analysis

A Python program is read by a *parser*. Input to the parser is a stream of *tokens*, generated by the *lexical analyzer*. This chapter describes how the lexical analyzer breaks a file into tokens.

Python reads program text as Unicode code points; the encoding of a source file can be given by an encoding declaration and defaults to UTF-8, see [PEP 3120](#) for details. If the source file cannot be decoded, a [SyntaxError](#) is raised.

2.1. Line structure

A Python program is divided into a number of *logical lines*.

2.1.1. Logical lines

The end of a logical line is represented by the token NEWLINE. Statements cannot cross logical line boundaries except where NEWLINE is allowed by the syntax (e.g., between statements in compound statements). A logical line is constructed from one or more *physical lines* by following the explicit or implicit *line joining* rules.

2.1.2. Physical lines

A physical line is a sequence of characters terminated by an end-of-line sequence. In source files and strings, any of the standard platform line termination sequences can be used - the Unix form using ASCII LF (linefeed), the Windows form using the ASCII sequence CR LF (return followed by linefeed), or the old Macintosh form using the ASCII CR (return) character. All of these forms can be used equally, regardless of platform. The end of input also serves as an implicit terminator for the final physical line.

When embedding Python, source code strings should be passed to Python APIs using the standard C conventions for newline characters (the `\n` character, representing ASCII LF, is the line terminator).

2.1.3. Comments

A comment starts with a hash character (`#`) that is not part of a string literal, and ends at the end of the physical line. A comment signifies the end of the logical line unless the implicit line joining rules are invoked. Comments are ignored by the syntax; they are not tokens.

2.1.4. Encoding declarations

If a comment in the first or second line of the Python script matches the regular expression `coding[=:]s*([-\\w.]+)`, this comment is processed as an encoding declaration; the first group of this expression names the encoding of the source code file. The encoding declaration must appear on a line of its own. If it is the second line, the first line must also be a comment-only line. The recommended forms of an encoding expression are

```
# -*- coding: <encoding-name> -*-
```

which is recognized also by GNU Emacs, and

```
# vim:fileencoding=<encoding-name>
```

which is recognized by Bram Moolenaar's VIM.

If no encoding declaration is found, the default encoding is UTF-8. In addition, if the first bytes of the file are the UTF-8 byte-order mark (`b'\xef\xbb\xbf'`), the declared file encoding is UTF-8 (this is supported, among others, by Microsoft's **notepad**).

If an encoding is declared, the encoding name must be recognized by Python. The encoding is used for all lexical analysis, including string literals, comments and identifiers.

2.1.5. Explicit line joining

Two or more physical lines may be joined into logical lines using backslash characters (`\`), as follows: when a physical line ends in a backslash that is not part of a string literal or comment, it is joined with the following forming a single logical line, deleting the backslash and the following end-of-line character. For example:

```
if 1900 < year < 2100 and 1 <= month <= 12 \
    and 1 <= day <= 31 and 0 <= hour < 24 \
    and 0 <= minute < 60 and 0 <= second < 60:    # Looks like a valid d
    return 1
```

A line ending in a backslash cannot carry a comment. A backslash does not continue a comment. A backslash does not continue a token except for string literals (i.e., tokens other than string literals cannot be split across physical lines using a backslash). A backslash is illegal elsewhere on a line outside a string literal.

2.1.6. Implicit line joining

Expressions in parentheses, square brackets or curly braces can be split over more than one physical line without using backslashes. For example:

```
month_names = ['Januari', 'Februari', 'Maart',      # These are the
               'April',   'Mei',      'Juni',      # Dutch names
               'Juli',    'Augustus', 'September', # for the months
               'Oktober', 'November', 'December']  # of the year
```

Implicitly continued lines can carry comments. The indentation of the continuation lines is not important. Blank continuation lines are allowed. There is no NEWLINE token between implicit continuation lines. Implicitly continued lines can also occur within triple-quoted strings (see below); in that case they cannot carry comments.

2.1.7. Blank lines

A logical line that contains only spaces, tabs, formfeeds and possibly a comment, is ignored (i.e., no NEWLINE token is generated). During interactive input of statements, handling of a blank line may differ depending on the implementation of the read-eval-print loop. In the standard interactive interpreter, an entirely blank logical line (i.e. one containing not even whitespace or a comment) terminates a multi-line statement.

2.1.8. Indentation

Leading whitespace (spaces and tabs) at the beginning of a logical line is used to compute the indentation level of the line, which in turn is used to determine the grouping of statements.

Tabs are replaced (from left to right) by one to eight spaces such that the total number of characters up to and including the replacement is a multiple of eight (this is intended to be the same rule as used by Unix). The total number of spaces preceding the first non-blank character then determines the line's indentation. Indentation cannot be split over multiple physical lines using backslashes; the whitespace up to the first backslash determines the indentation.

Indentation is rejected as inconsistent if a source file mixes tabs and spaces in a way that makes the meaning dependent on the worth of a tab in spaces; a `TabError` is raised in that case.

Cross-platform compatibility note: because of the nature of text editors on non-UNIX platforms, it is unwise to use a mixture of spaces and tabs for the indentation

in a single source file. It should also be noted that different platforms may explicitly limit the maximum indentation level.

A formfeed character may be present at the start of the line; it will be ignored for the indentation calculations above. Formfeed characters occurring elsewhere in the leading whitespace have an undefined effect (for instance, they may reset the space count to zero).

The indentation levels of consecutive lines are used to generate INDENT and DEDENT tokens, using a stack, as follows.

Before the first line of the file is read, a single zero is pushed on the stack; this will never be popped off again. The numbers pushed on the stack will always be strictly increasing from bottom to top. At the beginning of each logical line, the line's indentation level is compared to the top of the stack. If it is equal, nothing happens. If it is larger, it is pushed on the stack, and one INDENT token is generated. If it is smaller, it *must* be one of the numbers occurring on the stack; all numbers on the stack that are larger are popped off, and for each number popped off a DEDENT token is generated. At the end of the file, a DEDENT token is generated for each number remaining on the stack that is larger than zero.

Here is an example of a correctly (though confusingly) indented piece of Python code:

```
def perm(l):
    # Compute the list of all permutations of l
    if len(l) <= 1:
        return [l]
    r = []
    for i in range(len(l)):
        s = l[:i] + l[i+1:]
        p = perm(s)
        for x in p:
            r.append(l[i:i+1] + x)
    return r
```

The following example shows various indentation errors:

```
def perm(l):
    for i in range(len(l)):
        s = l[:i] + l[i+1:]
        p = perm(l[:i] + l[i+1:])
        for x in p:
            r.append(l[i:i+1] + x)
    return r
```

(Actually, the first three errors are detected by the parser; only the last error is found by the lexical analyzer — the indentation of `return r` does not match a level popped off the stack.)

2.1.9. Whitespace between tokens

Except at the beginning of a logical line or in string literals, the whitespace characters space, tab and formfeed can be used interchangeably to separate tokens. Whitespace is needed between two tokens only if their concatenation could otherwise be interpreted as a different token (e.g., `ab` is one token, but `a b` is two tokens).

2.2. Other tokens

Besides `NEWLINE`, `INDENT` and `DEDENT`, the following categories of tokens exist: *identifiers*, *keywords*, *literals*, *operators*, and *delimiters*. Whitespace characters (other than line terminators, discussed earlier) are not tokens, but serve to delimit tokens. Where ambiguity exists, a token comprises the longest possible string that forms a legal token, when read from left to right.

2.3. Identifiers and keywords

Identifiers (also referred to as *names*) are described by the following lexical definitions.

The syntax of identifiers in Python is based on the Unicode standard annex UAX-31, with elaboration and changes as defined below; see also [PEP 3131](#) for further details.

Within the ASCII range (U+0001..U+007F), the valid characters for identifiers are the same as in Python 2.x: the uppercase and lowercase letters A through Z, the underscore `_` and, except for the first character, the digits 0 through 9.

Python 3.0 introduces additional characters from outside the ASCII range (see [PEP 3131](#)). For these characters, the classification uses the version of the Unicode Character Database as included in the `unicodedata` module.

Identifiers are unlimited in length. Case is significant.

```
identifier ::= xid_start xid_continue*
id_start   ::= <all characters in general categories Lu, Ll, Lt, Lr,
id_continue ::= <all characters in id_start, plus characters in the
xid_start   ::= <all characters in id_start whose NFKC normalization
xid_continue ::= <all characters in id_continue whose NFKC normalization
```

The Unicode category codes mentioned above stand for:

- *Lu* - uppercase letters
- *Li* - lowercase letters
- *Lt* - titlecase letters
- *Lm* - modifier letters
- *Lo* - other letters
- *Nl* - letter numbers
- *Mn* - nonspacing marks
- *Mc* - spacing combining marks
- *Nd* - decimal numbers
- *Pc* - connector punctuations
- *Other_ID_Start* - explicit list of characters in [PropList.txt](#) to support backwards compatibility
- *Other_ID_Continue* - likewise

All identifiers are converted into the normal form NFKC while parsing; comparison of identifiers is based on NFKC.

A non-normative HTML file listing all valid identifier characters for Unicode 4.1 can be found at <https://www.dcl.hpi.uni-potsdam.de/home/loewis/table-3131.html>.

2.3.1. Keywords

The following identifiers are used as reserved words, or *keywords* of the language, and cannot be used as ordinary identifiers. They must be spelled exactly as written here:

False	class	finally	is	return
None	continue	for	lambda	try
True	def	from	nonlocal	while
and	del	global	not	with
as	elif	if	or	yield
assert	else	import	pass	
break	except	in	raise	

2.3.2. Reserved classes of identifiers

Certain classes of identifiers (besides keywords) have special meanings. These classes are identified by the patterns of leading and trailing underscore characters:

`_*`

Not imported by `from module import *`. The special identifier `_` is used in the interactive interpreter to store the result of the last evaluation; it is stored in the

`builtins` module. When not in interactive mode, `_` has no special meaning and is not defined. See section [The import statement](#).

Note: The name `_` is often used in conjunction with internationalization; refer to the documentation for the `gettext` module for more information on this convention.

— *

System-defined names. These names are defined by the interpreter and its implementation (including the standard library). Current system names are discussed in the [Special method names](#) section and elsewhere. More will likely be defined in future versions of Python. Any use of `__*__` names, in any context, that does not follow explicitly documented use, is subject to breakage without warning.

— *

Class-private names. Names in this category, when used within the context of a class definition, are re-written to use a mangled form to help avoid name clashes between “private” attributes of base and derived classes. See section [Identifiers \(Names\)](#).

2.4. Literals

Literals are notations for constant values of some built-in types.

2.4.1. String and Bytes literals

String literals are described by the following lexical definitions:

```
stringliteral ::= [stringprefix](shortstring | longstring)
stringprefix  ::= "r" | "u" | "R" | "U" | "f" | "F"
                | "fr" | "Fr" | "fR" | "FR" | "rf" | "rF" | "Rf"
shortstring   ::= ''' shortstringitem* ''' | ''' shortstringitem* '
longstring    ::= ''' longstringitem* ''' | ''' longstringite
shortstringitem ::= shortstringchar | stringescapeseq
longstringitem  ::= longstringchar | stringescapeseq
shortstringchar ::= <any source character except "\" or newline or th
longstringchar  ::= <any source character except "\">
stringescapeseq ::= "\" <any source character>
```

```
bytesliteral  ::= bytesprefix(shortbytes | longbytes)
bytesprefix   ::= "b" | "B" | "br" | "Br" | "bR" | "BR" | "rb" | "rB"
shortbytes    ::= ''' shortbytesitem* ''' | ''' shortbytesitem* '
longbytes     ::= ''' longbytesitem* ''' | ''' longbytesitem* '
shortbytesitem ::= shortbyteschar | bytesescapeseq
```

```
longbytesitem ::= longbyteschar | bytesescapeseq
shortbyteschar ::= <any ASCII character except "\" or newline or the
longbyteschar  ::= <any ASCII character except "\">
bytesescapeseq ::= "\" <any ASCII character>
```

One syntactic restriction not indicated by these productions is that whitespace is not allowed between the [stringprefix](#) or [bytesprefix](#) and the rest of the literal. The source character set is defined by the encoding declaration; it is UTF-8 if no encoding declaration is given in the source file; see section [Encoding declarations](#).

In plain English: Both types of literals can be enclosed in matching single quotes (') or double quotes ("). They can also be enclosed in matching groups of three single or double quotes (these are generally referred to as *triple-quoted strings*). The backslash (\) character is used to escape characters that otherwise have a special meaning, such as newline, backslash itself, or the quote character.

Bytes literals are always prefixed with 'b' or 'B'; they produce an instance of the [bytes](#) type instead of the [str](#) type. They may only contain ASCII characters; bytes with a numeric value of 128 or greater must be expressed with escapes.

Both string and bytes literals may optionally be prefixed with a letter 'r' or 'R'; such strings are called *raw strings* and treat backslashes as literal characters. As a result, in string literals, '\u' and '\u' escapes in raw strings are not treated specially. Given that Python 2.x's raw unicode literals behave differently than Python 3.x's the 'ur' syntax is not supported.

New in version 3.3: The 'rb' prefix of raw bytes literals has been added as a synonym of 'br'.

New in version 3.3: Support for the unicode legacy literal (u'value') was reintroduced to simplify the maintenance of dual Python 2.x and 3.x codebases. See [PEP 414](#) for more information.

A string literal with 'f' or 'F' in its prefix is a *formatted string literal*; see [Formatted string literals](#). The 'f' may be combined with 'r', but not with 'b' or 'u', therefore raw formatted strings are possible, but formatted bytes literals are not.

In triple-quoted literals, unescaped newlines and quotes are allowed (and are retained), except that three unescaped quotes in a row terminate the literal. (A "quote" is the character used to open the literal, i.e. either ' or ".)

Unless an 'r' or 'R' prefix is present, escape sequences in string and bytes literals are interpreted according to rules similar to those used by Standard C. The recognized escape sequences are:

Escape Sequence	Meaning	Notes
\newline	Backslash and newline ignored	
\\	Backslash (\)	
\'	Single quote (')	
\"	Double quote (")	
\a	ASCII Bell (BEL)	
\b	ASCII Backspace (BS)	
\f	ASCII Formfeed (FF)	
\n	ASCII Linefeed (LF)	
\r	ASCII Carriage Return (CR)	
\t	ASCII Horizontal Tab (TAB)	
\v	ASCII Vertical Tab (VT)	
\ooo	Character with octal value <i>ooo</i>	(1,3)
\xhh	Character with hex value <i>hh</i>	(2,3)

Escape sequences only recognized in string literals are:

Escape Sequence	Meaning	Notes
\N{name}	Character named <i>name</i> in the Unicode database	(4)
\uxxxx	Character with 16-bit hex value <i>xxxx</i>	(5)
\Uxxxxxxxx	Character with 32-bit hex value <i>xxxxxxxx</i>	(6)

Notes:

1. As in Standard C, up to three octal digits are accepted.
2. Unlike in Standard C, exactly two hex digits are required.
3. In a bytes literal, hexadecimal and octal escapes denote the byte with the given value. In a string literal, these escapes denote a Unicode character with the given value.
4. *Changed in version 3.3:* Support for name aliases [\[1\]](#) has been added.
5. Exactly four hex digits are required.
6. Any Unicode character can be encoded this way. Exactly eight hex digits are required.

Unlike Standard C, all unrecognized escape sequences are left in the string unchanged, i.e., *the backslash is left in the result*. (This behavior is useful when debugging: if an escape sequence is mistyped, the resulting output is more easily recognized as broken.) It is also important to note that the escape sequences only recognized in string literals fall into the category of unrecognized escapes for bytes literals.

Changed in version 3.6: Unrecognized escape sequences produce a `DeprecationWarning`. In some future version of Python they will be a `SyntaxError`.

Even in a raw literal, quotes can be escaped with a backslash, but the backslash remains in the result; for example, `r"\"` is a valid string literal consisting of two characters: a backslash and a double quote; `r"\` is not a valid string literal (even a raw string cannot end in an odd number of backslashes). Specifically, *a raw literal cannot end in a single backslash* (since the backslash would escape the following quote character). Note also that a single backslash followed by a newline is interpreted as those two characters as part of the literal, *not* as a line continuation.

2.4.2. String literal concatenation

Multiple adjacent string or bytes literals (delimited by whitespace), possibly using different quoting conventions, are allowed, and their meaning is the same as their concatenation. Thus, `"hello" 'world'` is equivalent to `"helloworld"`. This feature can be used to reduce the number of backslashes needed, to split long strings conveniently across long lines, or even to add comments to parts of strings, for example:

```
re.compile("[A-Za-z_]"      # Letter or underscore
           "[A-Za-z0-9_]*"  # Letter, digit or underscore
           )
```

Note that this feature is defined at the syntactical level, but implemented at compile time. The `+` operator must be used to concatenate string expressions at run time. Also note that literal concatenation can use different quoting styles for each component (even mixing raw strings and triple quoted strings), and formatted string literals may be concatenated with plain string literals.

2.4.3. Formatted string literals

New in version 3.6.

A *formatted string literal* or *f-string* is a string literal that is prefixed with `'f'` or `'F'`. These strings may contain replacement fields, which are expressions delimited by

curly braces `{}`. While other string literals always have a constant value, formatted strings are really expressions evaluated at run time.

Escape sequences are decoded like in ordinary string literals (except when a literal is also marked as a raw string). After decoding, the grammar for the contents of the string is:

```
f_string      ::= (literal_char | "{" | "}") | replacement_field
replacement_field ::= "{" f_expression ["!" conversion] [":" format_spec]
f_expression   ::= (conditional_expression | "*" or_expr)
                  ("," conditional_expression | "," "*" or_expr
                  | yield_expression
conversion      ::= "s" | "r" | "a"
format_spec     ::= (literal_char | NULL | replacement_field)*
literal_char    ::= <any code point except "{", "}" or NULL>
```

The parts of the string outside curly braces are treated literally, except that any doubled curly braces `'{{' or '}}` are replaced with the corresponding single curly brace. A single opening curly bracket `'{'` marks a replacement field, which starts with a Python expression. After the expression, there may be a conversion field, introduced by an exclamation point `!''`. A format specifier may also be appended, introduced by a colon `':'`. A replacement field ends with a closing curly bracket `'}'`.

Expressions in formatted string literals are treated like regular Python expressions surrounded by parentheses, with a few exceptions. An empty expression is not allowed, and a `lambda` expression must be surrounded by explicit parentheses. Replacement expressions can contain line breaks (e.g. in triple-quoted strings), but they cannot contain comments. Each expression is evaluated in the context where the formatted string literal appears, in order from left to right.

If a conversion is specified, the result of evaluating the expression is converted before formatting. Conversion `!'s'` calls `str()` on the result, `!'r'` calls `repr()`, and `!'a'` calls `ascii()`.

The result is then formatted using the `format()` protocol. The format specifier is passed to the `__format__()` method of the expression or conversion result. An empty string is passed when the format specifier is omitted. The formatted result is then included in the final value of the whole string.

Top-level format specifiers may include nested replacement fields. These nested fields may include their own conversion fields and `format specifiers`, but may not include more deeply-nested replacement fields. The `format specifier mini-language` is the same as that used by the string `.format()` method.

Formatted string literals may be concatenated, but replacement fields cannot be split across literals.

Some examples of formatted string literals:

```
>>> name = "Fred"
>>> f"He said his name is {name!r}."
"He said his name is 'Fred'."
>>> f"He said his name is {repr(name)}." # repr() is equivalent to !r
"He said his name is 'Fred'."
>>> width = 10
>>> precision = 4
>>> value = decimal.Decimal("12.34567")
>>> f"result: {value:{width}.{precision}}" # nested fields
'result:      12.35'
>>> today = datetime(year=2017, month=1, day=27)
>>> f"{today:%B %d, %Y}" # using date format specifier
'January 27, 2017'
>>> number = 1024
>>> f"{number:#0x}" # using integer format specifier
'0x400'
```

A consequence of sharing the same syntax as regular string literals is that characters in the replacement fields must not conflict with the quoting used in the outer formatted string literal:

```
f"abc {a["x"]} def" # error: outer string literal ended prematurely
f"abc {a['x']} def" # workaround: use different quoting
```

Backslashes are not allowed in format expressions and will raise an error:

```
f"newline: {ord('\n')}}" # raises SyntaxError
```

To include a value in which a backslash escape is required, create a temporary variable.

```
>>> newline = ord('\n')
>>> f"newline: {newline}"
'newline: 10'
```

Formatted string literals cannot be used as docstrings, even if they do not include expressions.

```
>>> def foo():
...     f"Not a docstring"
... 
```

```
>>> foo.__doc__ is None
True
```

See also [PEP 498](#) for the proposal that added formatted string literals, and `str.format()`, which uses a related format string mechanism.

2.4.4. Numeric literals

There are three types of numeric literals: integers, floating point numbers, and imaginary numbers. There are no complex literals (complex numbers can be formed by adding a real number and an imaginary number).

Note that numeric literals do not include a sign; a phrase like `-1` is actually an expression composed of the unary operator `'-'` and the literal `1`.

2.4.5. Integer literals

Integer literals are described by the following lexical definitions:

```
integer      ::= decinteger | bininteger | octinteger | hexinteger
decinteger   ::= nonzerodigit (["_"] digit)* | "0"+ (["_"] "0")*
bininteger   ::= "0" ("b" | "B") (["_"] bindigit)+
octinteger   ::= "0" ("o" | "O") (["_"] octdigit)+
hexinteger   ::= "0" ("x" | "X") (["_"] hexdigit)+
nonzerodigit ::= "1"..."9"
digit        ::= "0"..."9"
bindigit     ::= "0" | "1"
octdigit     ::= "0"..."7"
hexdigit     ::= digit | "a"..."f" | "A"..."F"
```

There is no limit for the length of integer literals apart from what can be stored in available memory.

Underscores are ignored for determining the numeric value of the literal. They can be used to group digits for enhanced readability. One underscore can occur between digits, and after base specifiers like `0x`.

Note that leading zeros in a non-zero decimal number are not allowed. This is for disambiguation with C-style octal literals, which Python used before version 3.0.

Some examples of integer literals:

```
7      2147483647      0o177      0b100110111
3      79228162514264337593543950336  0o377      0xdeadbeef
      100_000_000_000      0b_1110_0101
```

Changed in version 3.6: Underscores are now allowed for grouping purposes in literals.

2.4.6. Floating point literals

Floating point literals are described by the following lexical definitions:

```
floatnumber ::= pointfloat | exponentfloat
pointfloat  ::= [digitpart] fraction | digitpart "."
exponentfloat ::= (digitpart | pointfloat) exponent
digitpart   ::= digit (["_"] digit)*
fraction    ::= "." digitpart
exponent    ::= ("e" | "E") ["+" | "-"] digitpart
```

Note that the integer and exponent parts are always interpreted using radix 10. For example, 077e010 is legal, and denotes the same number as 77e10. The allowed range of floating point literals is implementation-dependent. As in integer literals, underscores are supported for digit grouping.

Some examples of floating point literals:

```
3.14    10.    .001    1e100    3.14e-10    0e0    3.14_15_93
```

Changed in version 3.6: Underscores are now allowed for grouping purposes in literals.

2.4.7. Imaginary literals

Imaginary literals are described by the following lexical definitions:

```
imagnumber ::= (floatnumber | digitpart) ("j" | "J")
```

An imaginary literal yields a complex number with a real part of 0.0. Complex numbers are represented as a pair of floating point numbers and have the same restrictions on their range. To create a complex number with a nonzero real part, add a floating point number to it, e.g., (3+4j). Some examples of imaginary literals:

```
3.14j    10.j    10j    .001j    1e100j    3.14e-10j    3.14_15_93j
```

2.5. Operators

The following tokens are operators:

+	-	*	**	/	//	%	@
<<	>>	&		^	~		
<	>	<=	>=	==	!=		

2.6. Delimiters

The following tokens serve as delimiters in the grammar:

()	[]	{	}	
,	:	.	;	@	=	->
+=	-=	*=	/=	//=	%=	@=
&=	=	=	>>=	<<=	**=	

The period can also occur in floating-point and imaginary literals. A sequence of three periods has a special meaning as an ellipsis literal. The second half of the list, the augmented assignment operators, serve lexically as delimiters, but also perform an operation.

The following printing ASCII characters have special meaning as part of other tokens or are otherwise significant to the lexical analyzer:

'	"	#	\
---	---	---	---

The following printing ASCII characters are not used in Python. Their occurrence outside string literals and comments is an unconditional error:

\$?	`
----	---	---

Footnotes

[1] <http://www.unicode.org/Public/9.0.0/ucd/NameAliases.txt>