

## 21.8. `urllib.parse` — Parse URLs into components

Source code: [Lib/urllib/parse.py](https://lib/urllib/parse.py)

This module defines a standard interface to break Uniform Resource Locator (URL) strings up in components (addressing scheme, network location, path etc.), to combine the components back into a URL string, and to convert a “relative URL” to an absolute URL given a “base URL.”

The module has been designed to match the Internet RFC on Relative Uniform Resource Locators. It supports the following URL schemes: `file`, `ftp`, `gopher`, `hdl`, `http`, `https`, `imap`, `mailto`, `mms`, `news`, `nntp`, `prospero`, `rsync`, `rtsp`, `rtspu`, `sftp`, `shttp`, `sip`, `sips`, `snews`, `svn`, `svn+ssh`, `telnet`, `wais`, `ws`, `wss`.

The `urllib.parse` module defines functions that fall into two broad categories: URL parsing and URL quoting. These are covered in detail in the following sections.

### 21.8.1. URL Parsing

The URL parsing functions focus on splitting a URL string into its components, or on combining URL components into a URL string.

`urllib.parse.urlparse(urlstring, scheme="", allow_fragments=True)`

Parse a URL into six components, returning a 6-tuple. This corresponds to the general structure of a URL: `scheme://netloc/path;parameters?query#fragment`. Each tuple item is a string, possibly empty. The components are not broken up in smaller parts (for example, the network location is a single string), and % escapes are not expanded. The delimiters as shown above are not part of the result, except for a leading slash in the *path* component, which is retained if present. For example:

```
>>> from urllib.parse import urlparse
>>> o = urlparse('http://www.cwi.nl:80/%7Eguido/Python.html')
>>> o
ParseResult(scheme='http', netloc='www.cwi.nl:80', path='/%7Eguido',
              params='', query='', fragment='')
>>> o.scheme
'http'
>>> o.port
80
>>> o.geturl()
'http://www.cwi.nl:80/%7Eguido/Python.html'
```

Following the syntax specifications in [RFC 1808](#), `urlparse` recognizes a netloc only if it is properly introduced by `'/'`. Otherwise the input is presumed to be a relative URL and thus to start with a path component.

```
>>> from urllib.parse import urlparse
>>> urlparse('//www.cwi.nl:80/%7Eguido/Python.html')
ParseResult(scheme='', netloc='www.cwi.nl:80', path='/%7Eguido/Pyt
           params='', query='', fragment='')
>>> urlparse('www.cwi.nl/%7Eguido/Python.html')
ParseResult(scheme='', netloc='', path='www.cwi.nl/%7Eguido/Python
           params='', query='', fragment='')
>>> urlparse('help/Python.html')
ParseResult(scheme='', netloc='', path='help/Python.html', params=
           query='', fragment='')
```

The *scheme* argument gives the default addressing scheme, to be used only if the URL does not specify one. It should be the same type (text or bytes) as *url-string*, except that the default value `''` is always allowed, and is automatically converted to `b''` if appropriate.

If the *allow\_fragments* argument is false, fragment identifiers are not recognized. Instead, they are parsed as part of the path, parameters or query component, and *fragment* is set to the empty string in the return value.

The return value is actually an instance of a subclass of `tuple`. This class has the following additional read-only convenience attributes:

Attribute	Index	Value	Value if not present
<code>scheme</code>	0	URL scheme specifier	<i>scheme</i> parameter
<code>netloc</code>	1	Network location part	empty string
<code>path</code>	2	Hierarchical path	empty string
<code>params</code>	3	Parameters for last path element	empty string
<code>query</code>	4	Query component	empty string
<code>fragment</code>	5	Fragment identifier	empty string
<code>username</code>		User name	<code>None</code>
<code>password</code>		Password	<code>None</code>
<code>hostname</code>		Host name (lower case)	<code>None</code>
<code>port</code>			<code>None</code>

Attribute	Index	Value	Value if not present
		Port number as integer, if present	

Reading the port attribute will raise a [ValueError](#) if an invalid port is specified in the URL. See section [Structured Parse Results](#) for more information on the result object.

Unmatched square brackets in the netloc attribute will raise a [ValueError](#).

*Changed in version 3.2:* Added IPv6 URL parsing capabilities.

*Changed in version 3.3:* The fragment is now parsed for all URL schemes (unless *allow\_fragment* is false), in accordance with [RFC 3986](#). Previously, a whitelist of schemes that support fragments existed.

*Changed in version 3.6:* Out-of-range port numbers now raise [ValueError](#), instead of returning [None](#).

`urllib.parse.parse_qs(qs, keep_blank_values=False, strict_parsing=False, encoding='utf-8', errors='replace')`

Parse a query string given as a string argument (data of type *application/x-www-form-urlencoded*). Data are returned as a dictionary. The dictionary keys are the unique query variable names and the values are lists of values for each name.

The optional argument *keep\_blank\_values* is a flag indicating whether blank values in percent-encoded queries should be treated as blank strings. A true value indicates that blanks should be retained as blank strings. The default false value indicates that blank values are to be ignored and treated as if they were not included.

The optional argument *strict\_parsing* is a flag indicating what to do with parsing errors. If false (the default), errors are silently ignored. If true, errors raise a [ValueError](#) exception.

The optional *encoding* and *errors* parameters specify how to decode percent-encoded sequences into Unicode characters, as accepted by the [bytes.decode\(\)](#) method.

Use the [urllib.parse.urlencode\(\)](#) function (with the *doseq* parameter set to True) to convert such dictionaries into query strings.

*Changed in version 3.2:* Add *encoding* and *errors* parameters.

`urllib.parse.parse_qs1(qs, keep_blank_values=False, strict_parsing=False, encoding='utf-8', errors='replace')`

Parse a query string given as a string argument (data of type *application/x-www-form-urlencoded*). Data are returned as a list of name, value pairs.

The optional argument *keep\_blank\_values* is a flag indicating whether blank values in percent-encoded queries should be treated as blank strings. A true value indicates that blanks should be retained as blank strings. The default false value indicates that blank values are to be ignored and treated as if they were not included.

The optional argument *strict\_parsing* is a flag indicating what to do with parsing errors. If false (the default), errors are silently ignored. If true, errors raise a [ValueError](#) exception.

The optional *encoding* and *errors* parameters specify how to decode percent-encoded sequences into Unicode characters, as accepted by the [bytes.decode\(\)](#) method.

Use the [urllib.parse.urlencode\(\)](#) function to convert such lists of pairs into query strings.

*Changed in version 3.2:* Add *encoding* and *errors* parameters.

`urllib.parse.urlunparse(parts)`

Construct a URL from a tuple as returned by `urlparse()`. The *parts* argument can be any six-item iterable. This may result in a slightly different, but equivalent URL, if the URL that was parsed originally had unnecessary delimiters (for example, a ? with an empty query; the RFC states that these are equivalent).

`urllib.parse.urlsplit(urlstring, scheme="", allow_fragments=True)`

This is similar to [urlparse\(\)](#), but does not split the params from the URL. This should generally be used instead of [urlparse\(\)](#) if the more recent URL syntax allowing parameters to be applied to each segment of the *path* portion of the URL (see [RFC 2396](#)) is wanted. A separate function is needed to separate the path segments and parameters. This function returns a 5-tuple: (addressing scheme, network location, path, query, fragment identifier).

The return value is actually an instance of a subclass of [tuple](#). This class has the following additional read-only convenience attributes:

Attribute	Index	Value	Value if not present
scheme	0	URL scheme specifier	<i>scheme</i> parameter
netloc	1	Network location part	empty string

Attribute	Index	Value	Value if not present
path	2	Hierarchical path	empty string
query	3	Query component	empty string
fragment	4	Fragment identifier	empty string
username		User name	<a href="#">None</a>
password		Password	<a href="#">None</a>
hostname		Host name (lower case)	<a href="#">None</a>
port		Port number as integer, if present	<a href="#">None</a>

Reading the port attribute will raise a [ValueError](#) if an invalid port is specified in the URL. See section [Structured Parse Results](#) for more information on the result object.

Unmatched square brackets in the netloc attribute will raise a [ValueError](#).

*Changed in version 3.6:* Out-of-range port numbers now raise [ValueError](#), instead of returning [None](#).

`urllib.parse.urlunsplit(parts)`

Combine the elements of a tuple as returned by [urlsplit\(\)](#) into a complete URL as a string. The *parts* argument can be any five-item iterable. This may result in a slightly different, but equivalent URL, if the URL that was parsed originally had unnecessary delimiters (for example, a ? with an empty query; the RFC states that these are equivalent).

`urllib.parse.urljoin(base, url, allow_fragments=True)`

Construct a full (“absolute”) URL by combining a “base URL” (*base*) with another URL (*url*). Informally, this uses components of the base URL, in particular the addressing scheme, the network location and (part of) the path, to provide missing components in the relative URL. For example:

```
>>> from urllib.parse import urljoin
>>> urljoin('http://www.cwi.nl/%7Eguido/Python.html', 'FAQ.html')
'http://www.cwi.nl/%7Eguido/FAQ.html'
```

The *allow\_fragments* argument has the same meaning and default as for [urlparse\(\)](#).

**Note:** If *url* is an absolute URL (that is, starting with `//` or `scheme://`), the *url*’s host name and/or scheme will be present in the result. For example:

```
>>> urljoin('http://www.cwi.nl/%7Eguido/Python.html',
...         '//www.python.org/%7Eguido')
'http://www.python.org/%7Eguido'
```

If you do not want that behavior, preprocess the *url* with `urlsplit()` and `urlunsplit()`, removing possible *scheme* and *netloc* parts.

*Changed in version 3.5:* Behaviour updated to match the semantics defined in [RFC 3986](#).

`urllib.parse.urldefrag(url)`

If *url* contains a fragment identifier, return a modified version of *url* with no fragment identifier, and the fragment identifier as a separate string. If there is no fragment identifier in *url*, return *url* unmodified and an empty string.

The return value is actually an instance of a subclass of `tuple`. This class has the following additional read-only convenience attributes:

Attribute	Index	Value	Value if not present
<code>url</code>	0	URL with no fragment	empty string
<code>fragment</code>	1	Fragment identifier	empty string

See section [Structured Parse Results](#) for more information on the result object.

*Changed in version 3.2:* Result is a structured object rather than a simple 2-tuple.

## 21.8.2. Parsing ASCII Encoded Bytes

The URL parsing functions were originally designed to operate on character strings only. In practice, it is useful to be able to manipulate properly quoted and encoded URLs as sequences of ASCII bytes. Accordingly, the URL parsing functions in this module all operate on `bytes` and `bytearray` objects in addition to `str` objects.

If `str` data is passed in, the result will also contain only `str` data. If `bytes` or `bytearray` data is passed in, the result will contain only `bytes` data.

Attempting to mix `str` data with `bytes` or `bytearray` in a single function call will result in a `TypeError` being raised, while attempting to pass in non-ASCII byte values will trigger `UnicodeDecodeError`.

To support easier conversion of result objects between `str` and `bytes`, all return values from URL parsing functions provide either an `encode()` method (when the result contains `str` data) or a `decode()` method (when the result contains `bytes`

data). The signatures of these methods match those of the corresponding `str` and `bytes` methods (except that the default encoding is `'ascii'` rather than `'utf-8'`). Each produces a value of a corresponding type that contains either `bytes` data (for `encode()` methods) or `str` data (for `decode()` methods).

Applications that need to operate on potentially improperly quoted URLs that may contain non-ASCII data will need to do their own decoding from bytes to characters before invoking the URL parsing methods.

The behaviour described in this section applies only to the URL parsing functions. The URL quoting functions use their own rules when producing or consuming byte sequences as detailed in the documentation of the individual URL quoting functions.

*Changed in version 3.2:* URL parsing functions now accept ASCII encoded byte sequences

## 21.8.3. Structured Parse Results

The result objects from the `urlparse()`, `urlsplit()` and `urldefrag()` functions are subclasses of the `tuple` type. These subclasses add the attributes listed in the documentation for those functions, the encoding and decoding support described in the previous section, as well as an additional method:

`urllib.parse.SplitResult.geturl()`

Return the re-combined version of the original URL as a string. This may differ from the original URL in that the scheme may be normalized to lower case and empty components may be dropped. Specifically, empty parameters, queries, and fragment identifiers will be removed.

For `urldefrag()` results, only empty fragment identifiers will be removed. For `urlsplit()` and `urlparse()` results, all noted changes will be made to the URL returned by this method.

The result of this method remains unchanged if passed back through the original parsing function:

```
>>> from urllib.parse import urlsplit
>>> url = 'HTTP://www.Python.org/doc/#'
>>> r1 = urlsplit(url)
>>> r1.geturl()
'http://www.Python.org/doc/'
>>> r2 = urlsplit(r1.geturl())
>>> r2.geturl()
'http://www.Python.org/doc/'
```

```
>>>
```

The following classes provide the implementations of the structured parse results when operating on `str` objects:

`class urllib.parse.DefragResult(url, fragment)`

Concrete class for `urldefrag()` results containing `str` data. The `encode()` method returns a `DefragResultBytes` instance.

*New in version 3.2.*

`class urllib.parse.ParseResult(scheme, netloc, path, params, query, fragment)`

Concrete class for `urlparse()` results containing `str` data. The `encode()` method returns a `ParseResultBytes` instance.

`class urllib.parse.SplitResult(scheme, netloc, path, query, fragment)`

Concrete class for `urlsplit()` results containing `str` data. The `encode()` method returns a `SplitResultBytes` instance.

The following classes provide the implementations of the parse results when operating on `bytes` or `bytearray` objects:

`class urllib.parse.DefragResultBytes(url, fragment)`

Concrete class for `urldefrag()` results containing `bytes` data. The `decode()` method returns a `DefragResult` instance.

*New in version 3.2.*

`class urllib.parse.ParseResultBytes(scheme, netloc, path, params, query, fragment)`

Concrete class for `urlparse()` results containing `bytes` data. The `decode()` method returns a `ParseResult` instance.

*New in version 3.2.*

`class urllib.parse.SplitResultBytes(scheme, netloc, path, query, fragment)`

Concrete class for `urlsplit()` results containing `bytes` data. The `decode()` method returns a `SplitResult` instance.

*New in version 3.2.*

## 21.8.4. URL Quoting

The URL quoting functions focus on taking program data and making it safe for use as URL components by quoting special characters and appropriately encoding non-



ASCII text. They also support reversing these operations to recreate the original data from the contents of a URL component if that task isn't already covered by the URL parsing functions above.

`urllib.parse.quote(string, safe='/', encoding=None, errors=None)`

Replace special characters in *string* using the %xx escape. Letters, digits, and the characters `'_.'` are never quoted. By default, this function is intended for quoting the path section of URL. The optional *safe* parameter specifies additional ASCII characters that should not be quoted — its default value is `'/'`.

*string* may be either a [str](#) or a [bytes](#).

The optional *encoding* and *errors* parameters specify how to deal with non-ASCII characters, as accepted by the [str.encode\(\)](#) method. *encoding* defaults to `'utf-8'`. *errors* defaults to `'strict'`, meaning unsupported characters raise a [UnicodeEncodeError](#). *encoding* and *errors* must not be supplied if *string* is a [bytes](#), or a [TypeError](#) is raised.

Note that `quote(string, safe, encoding, errors)` is equivalent to `quote_from_bytes(string.encode(encoding, errors), safe)`.

Example: `quote('/El Niño/')` yields `'/El%20Ni%C3%B1o/'`.

`urllib.parse.quote_plus(string, safe='', encoding=None, errors=None)`

Like [quote\(\)](#), but also replace spaces by plus signs, as required for quoting HTML form values when building up a query string to go into a URL. Plus signs in the original string are escaped unless they are included in *safe*. It also does not have *safe* default to `'/'`.

Example: `quote_plus('/El Niño/')` yields `'%2FE1+Ni%C3%B1o%2F'`.

`urllib.parse.quote_from_bytes(bytes, safe='/')`

Like [quote\(\)](#), but accepts a [bytes](#) object rather than a [str](#), and does not perform string-to-bytes encoding.

Example: `quote_from_bytes(b'a&\xef')` yields `'a%26%EF'`.

`urllib.parse.unquote(string, encoding='utf-8', errors='replace')`

Replace %xx escapes by their single-character equivalent. The optional *encoding* and *errors* parameters specify how to decode percent-encoded sequences into Unicode characters, as accepted by the [bytes.decode\(\)](#) method.

*string* must be a [str](#).

*encoding* defaults to 'utf-8'. *errors* defaults to 'replace', meaning invalid sequences are replaced by a placeholder character.

Example: `unquote('/E1%20Ni%C3%B1o/')` yields `'/E1 Niño/'`.

`urllib.parse.unquote_plus(string, encoding='utf-8', errors='replace')`

Like `unquote()`, but also replace plus signs by spaces, as required for unquoting HTML form values.

*string* must be a `str`.

Example: `unquote_plus('/E1+Ni%C3%B1o/')` yields `'/E1 Niño/'`.

`urllib.parse.unquote_to_bytes(string)`

Replace `%xx` escapes by their single-octet equivalent, and return a `bytes` object.

*string* may be either a `str` or a `bytes`.

If it is a `str`, unescaped non-ASCII characters in *string* are encoded into UTF-8 bytes.

Example: `unquote_to_bytes('a%26%EF')` yields `b'a&\xef'`.

`urllib.parse.urlencode(query, doseq=False, safe="", encoding=None, errors=None, quote_via=quote_plus)`

Convert a mapping object or a sequence of two-element tuples, which may contain `str` or `bytes` objects, to a percent-encoded ASCII text string. If the resultant string is to be used as a *data* for POST operation with the `urlopen()` function, then it should be encoded to bytes, otherwise it would result in a `TypeError`.

The resulting string is a series of key=value pairs separated by '&' characters, where both *key* and *value* are quoted using the *quote\_via* function. By default, `quote_plus()` is used to quote the values, which means spaces are quoted as a '+' character and '/' characters are encoded as %2F, which follows the standard for GET requests (application/x-www-form-urlencoded). An alternate function that can be passed as *quote\_via* is `quote()`, which will encode spaces as %20 and not encode '/' characters. For maximum control of what is quoted, use `quote` and specify a value for *safe*.

When a sequence of two-element tuples is used as the *query* argument, the first element of each tuple is a key and the second is a value. The value element in itself can be a sequence and in that case, if the optional parameter *doseq* is evaluates to True, individual key=value pairs separated by '&' are generated

for each element of the value sequence for the key. The order of parameters in the encoded string will match the order of parameter tuples in the sequence.

The *safe*, *encoding*, and *errors* parameters are passed down to *quote\_via* (the *encoding* and *errors* parameters are only passed when a query element is a [str](#)).

To reverse this encoding process, [parse\\_qs\(\)](#) and [parse\\_qsl\(\)](#) are provided in this module to parse query strings into Python data structures.

Refer to [urlib examples](#) to find out how urlencode method can be used for generating query string for a URL or data for POST.

*Changed in version 3.2:* Query parameter supports bytes and string objects.

*New in version 3.5:* *quote\_via* parameter.

**See also:**

**[RFC 3986](#) - Uniform Resource Identifiers**

This is the current standard (STD66). Any changes to `urllib.parse` module should conform to this. Certain deviations could be observed, which are mostly for backward compatibility purposes and for certain de-facto parsing requirements as commonly observed in major browsers.

**[RFC 2732](#) - Format for Literal IPv6 Addresses in URL's.**

This specifies the parsing requirements of IPv6 URLs.

**[RFC 2396](#) - Uniform Resource Identifiers (URI): Generic Syntax**

Document describing the generic syntactic requirements for both Uniform Resource Names (URNs) and Uniform Resource Locators (URLs).

**[RFC 2368](#) - The mailto URL scheme.**

Parsing requirements for mailto URL schemes.

**[RFC 1808](#) - Relative Uniform Resource Locators**

This Request For Comments includes the rules for joining an absolute and a relative URL, including a fair number of “Abnormal Examples” which govern the treatment of border cases.

**[RFC 1738](#) - Uniform Resource Locators (URL)**

This specifies the formal syntax and semantics of absolute URLs.