

---

# **w3lib Documentation**

*Release 1.18.0*

**w3lib developers**

**Aug 03, 2017**



---

# Contents

---

<b>1</b>	<b>Overview</b>	<b>1</b>
<b>2</b>	<b>Modules</b>	<b>3</b>
2.1	w3lib Package . . . . .	3
2.1.1	encoding Module . . . . .	3
2.1.2	html Module . . . . .	5
2.1.3	http Module . . . . .	8
2.1.4	url Module . . . . .	9
<b>3</b>	<b>Requirements</b>	<b>13</b>
<b>4</b>	<b>Install</b>	<b>15</b>
<b>5</b>	<b>Tests</b>	<b>17</b>
<b>6</b>	<b>Changelog</b>	<b>19</b>
6.1	1.18.0 (2017-08-03) . . . . .	19
6.2	1.17.0 (2017-02-08) . . . . .	19
6.3	1.16.0 (2016-11-10) . . . . .	19
6.4	1.15.0 (2016-07-29) . . . . .	20
6.5	1.14.3 (2016-07-14) . . . . .	20
6.6	1.14.2 (2016-04-11) . . . . .	20
6.7	1.14.1 (2016-04-07) . . . . .	20
6.8	1.14.0 (2016-04-06) . . . . .	20
6.9	1.13.0 (2015-11-05) . . . . .	20
6.10	1.12.0 (2015-06-29) . . . . .	21
6.11	1.11.0 (2015-01-13) . . . . .	21
6.12	1.10.0 (2014-08-20) . . . . .	21
6.13	1.9.0 (2014-08-16) . . . . .	21
6.14	1.8.1 (2014-08-14) . . . . .	21
6.15	1.8.0 (2014-07-31) . . . . .	21
6.16	1.7.1 (2014-07-26) . . . . .	21
6.17	1.6 (2014-06-03) . . . . .	21
6.18	1.5 (2013-11-09) . . . . .	22
6.19	1.4 (2013-10-18) . . . . .	22
6.20	1.3 (2012-05-13) . . . . .	22
6.21	1.2 (2012-05-02) . . . . .	22

6.22	1.1 (2012-04-18)	22
6.23	1.0 (2011-04-17)	22
6.24	History	22
<b>7</b>	<b>Indices and tables</b>	<b>23</b>
	<b>Python Module Index</b>	<b>25</b>

This is a Python library of web-related functions, such as:

- remove comments, or tags from HTML snippets
- extract base url from HTML snippets
- translate entites on HTML strings
- convert raw HTTP headers to dicts and vice-versa
- construct HTTP auth header
- converting HTML pages to unicode
- sanitize urls (like browsers do)
- extract arguments from urls

The w3lib library is licensed under the BSD license.



## w3lib Package

### encoding Module

Functions for handling encoding of web pages

w3lib.encoding.**html\_body\_declared\_encoding** (*html\_body\_str*)

Return the encoding specified in meta tags in the html body, or None if no suitable encoding was found

```
>>> import w3lib.encoding
>>> w3lib.encoding.html_body_declared_encoding(
...     """<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN"
...         "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
...     <html xmlns="http://www.w3.org/1999/xhtml" xml:lang="en" lang="en">
...     <head>
...         <title>Some title</title>
...         <meta http-equiv="content-type" content="text/html; charset=utf-8" />
...     </head>
...     <body>
...     ...
...     </body>
... </html>"""")
'utf-8'
>>>
```

w3lib.encoding.**html\_to\_unicode** (*content\_type\_header*, *html\_body\_str*, *default\_encoding='utf8'*,  
*auto\_detect\_fun=None*)

Convert raw html bytes to unicode

This attempts to make a reasonable guess at the content encoding of the html body, following a similar process to a web browser.

It will try in order:

- http content type header

- BOM (byte-order mark)
- meta or xml tag declarations
- auto-detection, if the `auto_detect_fun` keyword argument is not `None`
- default encoding in keyword arg (which defaults to utf8)

If an encoding other than the auto-detected or default encoding is used, overrides will be applied, converting some character encodings to more suitable alternatives.

If a BOM is found matching the encoding, it will be stripped.

The `auto_detect_fun` argument can be used to pass a function that will sniff the encoding of the text. This function must take the raw text as an argument and return the name of an encoding that python can process, or `None`. To use `chardet`, for example, you can define the function as:

```
auto_detect_fun=lambda x: chardet.detect(x).get('encoding')
```

or to use `UnicodeDammit` (shipped with the `BeautifulSoup` library):

```
auto_detect_fun=lambda x: UnicodeDammit(x).originalEncoding
```

If the locale of the website or user language preference is known, then a better default encoding can be supplied.

If `content_type_header` is not present, `None` can be passed signifying that the header was not present.

This method will not fail, if characters cannot be converted to unicode, `\uufffd` (the unicode replacement character) will be inserted instead.

Returns a tuple of (<encoding used>, <unicode\_string>)

Examples:

```
>>> import w3lib.encoding
>>> w3lib.encoding.html_to_unicode(None,
... b"<!DOCTYPE html>
... <head>
... <meta charset='UTF-8' />
... <meta name='viewport' content='width=device-width' />
... <title>Creative Commons France</title>
... <link rel='canonical' href='http://creativecommons.fr/' />
... <body>
... <p>Creative Commons est une organisation \xc3\xa0 but non lucratif
... qui a pour dessein de faciliter la diffusion et le partage des oeuvres
... tout en accompagnant les nouvelles pratiques de cr\xc3\xa9ation \xc3\xa0
↳l\xe2\x80\x99\xc3\xa8re numerique.</p>
... </body>
... </html>")
('utf-8', u'<!DOCTYPE html>\n<head>\n<meta charset="UTF-8" />\n<meta name=
↳"viewport" content="width=device-width" />\n<title>Creative Commons France</
↳title>\n<link rel=\'canonical\' href=\'http://creativecommons.fr/\' />\n<body>\n
↳<p>Creative Commons est une organisation \xe0 but non lucratif\nqui a pour
↳dessein de faciliter la diffusion et le partage des oeuvres\ntout en
↳accompagnant les nouvelles pratiques de cr\xe9ation \xe0 l\u2019\xe8re
↳numerique.</p>\n</body>\n</html>')
>>>
```

`w3lib.encoding.http_content_type_encoding` (`content_type`)

Extract the encoding in the content-type header



```
>>> import w3lib.encoding
>>> w3lib.encoding.http_content_type_encoding("Content-Type: text/html;_
↳ charset=ISO-8859-4")
'iso8859-4'
```

`w3lib.encoding.read_bom(data)`

Read the byte order mark in the text, if present, and return the encoding represented by the BOM and the BOM.

If no BOM can be detected, (None, None) is returned.

```
>>> import w3lib.encoding
>>> w3lib.encoding.read_bom(b'\xfe\xff\x6c\x34')
('utf-16-be', '\xfe\xff')
>>> w3lib.encoding.read_bom(b'\xff\xfe\x34\x6c')
('utf-16-le', '\xff\xfe')
>>> w3lib.encoding.read_bom(b'\x00\x00\xfe\xff\x00\x00\x6c\x34')
('utf-32-be', '\x00\x00\xfe\xff')
>>> w3lib.encoding.read_bom(b'\xff\xfe\x00\x00\x34\x6c\x00\x00')
('utf-32-le', '\xff\xfe\x00\x00')
>>> w3lib.encoding.read_bom(b'\x01\x02\x03\x04')
(None, None)
>>>
```

`w3lib.encoding.resolve_encoding(encoding_alias)`

Return the encoding that *encoding\_alias* maps to, or None if the encoding cannot be interpreted

```
>>> import w3lib.encoding
>>> w3lib.encoding.resolve_encoding('latin1')
'cp1252'
>>> w3lib.encoding.resolve_encoding('gb_2312-80')
'gb18030'
>>>
```

`w3lib.encoding.to_unicode(data_str, encoding)`

Convert a str object to unicode using the encoding given

Characters that cannot be converted will be converted to `\ufffd` (the unicode replacement character).

## html Module

Functions for dealing with markup text

`w3lib.html.get_base_url(text, baseurl='', encoding='utf-8')`

Return the base url if declared in the given HTML *text*, relative to the given base url.

If no base url is found, the given *baseurl* is returned.

`w3lib.html.get_meta_refresh(text, baseurl='', encoding='utf-8', ignore_tags=('script', 'no-script'))`

Return the http-equiv parameter of the HTML meta element from the given HTML text and return a tuple (interval, url) where interval is an integer containing the delay in seconds (or zero if not present) and url is a string with the absolute url to redirect.

If no meta redirect is found, (None, None) is returned.

`w3lib.html.remove_comments(text, encoding=None)`

Remove HTML Comments.

```
>>> import w3lib.html
>>> w3lib.html.remove_comments(b"test <!--textcoment--> whatever")
u'test whatever'
>>>
```

w3lib.html.**remove\_entities** (*text*, *keep*=(), *remove\_illegal*=True, *encoding*='utf-8')

**Warning:** This function is deprecated and will be removed in future. Please use `replace_entities()` instead.

w3lib.html.**remove\_tags** (*text*, *which\_ones*=(), *keep*=(), *encoding*=None)

Remove HTML Tags only.

*which\_ones* and *keep* are both tuples, there are four cases:

which_ones	keep	what it does
<b>not empty</b>	empty	remove all tags in which_ones
empty	<b>not empty</b>	remove all tags except the ones in keep
empty	empty	remove all tags
<b>not empty</b>	<b>not empty</b>	not allowed

Remove all tags:

```
>>> import w3lib.html
>>> doc = '<div><p><b>This is a link:</b> <a href="http://www.example.com">example
↳</a></p></div>'
>>> w3lib.html.remove_tags(doc)
u'This is a link: example'
>>>
```

Keep only some tags:

```
>>> w3lib.html.remove_tags(doc, keep=('div',))
u'<div>This is a link: example</div>'
>>>
```

Remove only specific tags:

```
>>> w3lib.html.remove_tags(doc, which_ones=('a', 'b'))
u'<div><p>This is a link: example</p></div>'
>>>
```

You can't remove some and keep some:

```
>>> w3lib.html.remove_tags(doc, which_ones=('a',), keep=('p',))
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/usr/local/lib/python2.7/dist-packages/w3lib/html.py", line 101, in _
↳remove_tags
    assert not (which_ones and keep), 'which_ones and keep can not be given at _
↳the same time'
AssertionError: which_ones and keep can not be given at the same time
>>>
```

w3lib.html.**remove\_tags\_with\_content** (*text*, *which\_ones*=(), *encoding*=None)

Remove tags and their content.

*which\_ones* is a tuple of which tags to remove including their content. If is empty, returns the string unmodified.

```
>>> import w3lib.html
>>> doc = '<div><p><b>This is a link:</b> <a href="http://www.example.com">example
↳</a></p></div>'
>>> w3lib.html.remove_tags_with_content(doc, which_ones=('b',))
u'<div><p> <a href="http://www.example.com">example</a></p></div>'
>>>
```

`w3lib.html.replace_entities` (*text*, *keep*=(), *remove\_illegal*=True, *encoding*='utf-8')

Remove entities from the given *text* by converting them to their corresponding unicode character.

*text* can be a unicode string or a byte string encoded in the given *encoding* (which defaults to 'utf-8').

If *keep* is passed (with a list of entity names) those entities will be kept (they won't be removed).

It supports both numeric entities (&#nnnn; and &#hhhh;) and named entities (such as &nbsp; or &gt;).

If *remove\_illegal* is True, entities that can't be converted are removed. If *remove\_illegal* is False, entities that can't be converted are kept "as is". For more information see the tests.

Always returns a unicode string (with the entities removed).

```
>>> import w3lib.html
>>> w3lib.html.replace_entities(b'Price: &pound;100')
u'Price: \xa3100'
>>> print(w3lib.html.replace_entities(b'Price: &pound;100'))
Price: £100
>>>
```

`w3lib.html.replace_escape_chars` (*text*, *which\_ones*=('n', 't', 'r'), *replace\_by*=u'', *encoding*=None)

Remove escape characters.

*which\_ones* is a tuple of which escape characters we want to remove. By default removes \n, \t, \r.

*replace\_by* is the string to replace the escape characters by. It defaults to ' ', meaning the escape characters are removed.

`w3lib.html.replace_tags` (*text*, *token*='', *encoding*=None)

Replace all markup tags found in the given *text* by the given token. By default *token* is an empty string so it just removes all tags.

*text* can be a unicode string or a regular string encoded as *encoding* (or 'utf-8' if *encoding* is not given.)

Always returns a unicode string.

Examples:

```
>>> import w3lib.html
>>> w3lib.html.replace_tags(u'This text contains <a>some tag</a>')
u'This text contains some tag'
>>> w3lib.html.replace_tags('<p>Je ne parle pas <b>fran\xe7ais</b></p>', ' -- ',
↳'latin-1')
u' -- Je ne parle pas -- fran\xe7ais -- -- '
>>>
```

`w3lib.html.strip_html5_whitespace` (*text*)

Strip all leading and trailing space characters (as defined in <https://www.w3.org/TR/html5/infrastructure.html#space-character>).

Such stripping is useful e.g. for processing HTML element attributes which contain URLs, like `href`, `src` or form `action` - HTML5 standard defines them as “valid URL potentially surrounded by spaces” or “valid non-empty URL potentially surrounded by spaces”.

```
>>> strip_html5_whitespace(' hello\n')
'hello'
```

`w3lib.html.unquote_markup` (*text*, *keep=()*, *remove\_illegal=True*, *encoding=None*)

This function receives markup as a text (always a unicode string or a UTF-8 encoded string) and does the following:

- 1.removes entities (except the ones in *keep*) from any part of it that is not inside a CDATA
- 2.searches for CDATAs and extracts their text (if any) without modifying it.
- 3.removes the found CDATAs

## http Module

`w3lib.http.basic_auth_header` (*username*, *password*)

Return an *Authorization* header field value for HTTP Basic Access Authentication (RFC 2617)

```
>>> import w3lib.http
>>> w3lib.http.basic_auth_header('someuser', 'somepass')
'Basic c29tZXVzZXI6c29tZXBhc3M='
```

`w3lib.http.headers_dict_to_raw` (*headers\_dict*)

Returns a raw HTTP headers representation of headers

For example:

```
>>> import w3lib.http
>>> w3lib.http.headers_dict_to_raw({b'Content-type': b'text/html', b'Accept': b
↳ 'gzip'})
'Content-type: text/html\r\nAccept: gzip'
>>>
```

Note that keys and values must be bytes.

Argument is None (returns None):

```
>>> w3lib.http.headers_dict_to_raw(None)
>>>
```

`w3lib.http.headers_raw_to_dict` (*headers\_raw*)

Convert raw headers (single multi-line bytestring) to a dictionary.

For example:

```
>>> import w3lib.http
>>> w3lib.http.headers_raw_to_dict(b"Content-type: text/html\n\rAccept: gzip\n\n")
{'Content-type': ['text/html'], 'Accept': ['gzip']}
```

Incorrect input:

```
>>> w3lib.http.headers_raw_to_dict(b"Content-typt gzip\n\n")
{}
>>>
```

Argument is None (return None):

```
>>> w3lib.http.headers_raw_to_dict (None)
>>>
```

## url Module

This module contains general purpose URL functions not found in the standard library.

`w3lib.url.add_or_replace_parameter` (*url*, *name*, *new\_value*)

Add or remove a parameter to a given url

```
>>> import w3lib.url
>>> w3lib.url.add_or_replace_parameter('http://www.example.com/index.php', 'arg',
↳ 'v')
'http://www.example.com/index.php?arg=v'
>>> w3lib.url.add_or_replace_parameter('http://www.example.com/index.php?arg1=v1&
↳ arg2=v2&arg3=v3', 'arg4', 'v4')
'http://www.example.com/index.php?arg1=v1&arg2=v2&arg3=v3&arg4=v4'
>>> w3lib.url.add_or_replace_parameter('http://www.example.com/index.php?arg1=v1&
↳ arg2=v2&arg3=v3', 'arg3', 'v3new')
'http://www.example.com/index.php?arg1=v1&arg2=v2&arg3=v3new'
>>>
```

`w3lib.url.any_to_uri` (*uri\_or\_path*)

If given a path name, return its File URI, otherwise return it unmodified

`w3lib.url.canonicalize_url` (*url*, *keep\_blank\_values=True*, *keep\_fragments=False*, *encoding=None*)

Canonicalize the given url by applying the following procedures:

- sort query arguments, first by key, then by value
- percent encode paths ; non-ASCII characters are percent-encoded using UTF-8 (RFC-3986)
- percent encode query arguments ; non-ASCII characters are percent-encoded using passed *encoding* (UTF-8 by default)
- normalize all spaces (in query arguments) '+' (plus symbol)
- normalize percent encodings case (%2f -> %2F)
- remove query arguments with blank values (unless *keep\_blank\_values* is True)
- remove fragments (unless *keep\_fragments* is True)

The url passed can be bytes or unicode, while the url returned is always a native str (bytes in Python 2, unicode in Python 3).

```
>>> import w3lib.url
>>>
>>> # sorting query arguments
>>> w3lib.url.canonicalize_url('http://www.example.com/do?c=3&b=5&b=2&a=50')
'http://www.example.com/do?a=50&b=2&b=5&c=3'
>>>
>>> # UTF-8 conversion + percent-encoding of non-ASCII characters
>>> w3lib.url.canonicalize_url(u'http://www.example.com/r\u00e9sum\u00e9')
'http://www.example.com/r%C3%A9sum%C3%A9'
>>>
```

For more examples, see the tests in `tests/test_url.py`.

`w3lib.url.file_uri_to_path(uri)`

Convert File URI to local filesystem path according to: [http://en.wikipedia.org/wiki/File\\_URI\\_scheme](http://en.wikipedia.org/wiki/File_URI_scheme)

`w3lib.url.parse_data_uri(uri)`

Parse a data: URI, returning a 3-tuple of media type, dictionary of media type parameters, and data.

`w3lib.url.path_to_file_uri(path)`

Convert local filesystem path to legal File URIs as described in: [http://en.wikipedia.org/wiki/File\\_URI\\_scheme](http://en.wikipedia.org/wiki/File_URI_scheme)

`w3lib.url.safe_download_url(url)`

Make a url for download. This will call `safe_url_string` and then strip the fragment, if one exists. The path will be normalised.

If the path is outside the document root, it will be changed to be within the document root.

`w3lib.url.safe_url_string(url, encoding='utf8', path_encoding='utf8')`

Convert the given URL into a legal URL by escaping unsafe characters according to RFC-3986.

If a bytes URL is given, it is first converted to `str` using the given encoding (which defaults to 'utf-8'). 'utf-8' encoding is used for URL path component (unless overridden by `path_encoding`), and given encoding is used for query string or form data. When passing an encoding, you should use the encoding of the original page (the page from which the URL was extracted from).

Calling this function on an already "safe" URL will return the URL unmodified.

Always returns a native `str` (bytes in Python2, unicode in Python3).

`w3lib.url.url_query_cleaner(url, parameterlist=(), sep='&', kvsep='=', remove=False, unique=True, keep_fragments=False)`

Clean URL arguments leaving only those passed in the `parameterlist` keeping order

```
>>> import w3lib.url
>>> w3lib.url.url_query_cleaner("product.html?id=200&foo=bar&name=wired", ('id',))
'product.html?id=200'
>>> w3lib.url.url_query_cleaner("product.html?id=200&foo=bar&name=wired", ['id',
↳ 'name'])
'product.html?id=200&name=wired'
>>>
```

If `unique` is `False`, do not remove duplicated keys

```
>>> w3lib.url.url_query_cleaner("product.html?d=1&e=b&d=2&d=3&other=other", ['d'],
↳ unique=False)
'product.html?d=1&d=2&d=3'
>>>
```

If `remove` is `True`, leave only those **not in** `parameterlist`.

```
>>> w3lib.url.url_query_cleaner("product.html?id=200&foo=bar&name=wired", ['id'],
↳ remove=True)
'product.html?foo=bar&name=wired'
>>> w3lib.url.url_query_cleaner("product.html?id=2&foo=bar&name=wired", ['id',
↳ 'foo'], remove=True)
'product.html?name=wired'
>>>
```

By default, URL fragments are removed. If you need to preserve fragments, pass the `keep_fragments` argument as `True`.

```
>>> w3lib.url.url_query_cleaner('http://domain.tld/?bla=123#123123', ['bla'],
↳ remove=True, keep_fragments=True)
'http://domain.tld/#123123'
```

`w3lib.url.url_query_parameter(url, parameter, default=None, keep_blank_values=0)`

Return the value of a url parameter, given the url and parameter name

General case:

```
>>> import w3lib.url
>>> w3lib.url.url_query_parameter("product.html?id=200&foo=bar", "id")
'200'
>>>
```

Return a default value if the parameter is not found:

```
>>> w3lib.url.url_query_parameter("product.html?id=200&foo=bar", "notthere",
↳ "mydefault")
'mydefault'
>>>
```

Returns None if `keep_blank_values` not set or 0 (default):

```
>>> w3lib.url.url_query_parameter("product.html?id=", "id")
>>>
```

Returns an empty string if `keep_blank_values` set to 1:

```
>>> w3lib.url.url_query_parameter("product.html?id=", "id", keep_blank_values=1)
''
>>>
```

`w3lib.url.urljoin_rfc(base, ref, encoding='utf-8')`

**Warning:** This function is deprecated and will be removed in future. It is not supported with Python 3. Please use `urlparse.urljoin` instead.

Same as `urlparse.urljoin` but supports unicode values in base and ref parameters (in which case they will be converted to str using the given encoding).

Always returns a str.

```
>>> import w3lib.url
>>> w3lib.url.urljoin_rfc('http://www.example.com/path/index.html', u'/otherpath/
↳ index2.html')
'http://www.example.com/otherpath/index2.html'
>>>
```

```
>>> # Note: the following does not work in Python 3
>>> w3lib.url.urljoin_rfc(b'http://www.example.com/path/index.html', u
↳ 'fran\u00e7ais/d\u00e9part.htm')
'http://www.example.com/path/fran\xc3\xa7ais/d\xc3\xa9part.htm'
>>>
```





## CHAPTER 3

---

### Requirements

---

Python 2.7 or Python 3.3+



## CHAPTER 4

---

### Install

---

```
pip install w3lib
```



## CHAPTER 5

---

### Tests

---

`nose` is the preferred way to run tests. Just run: `nosetests` from the root directory to execute tests using the default Python interpreter.

`tox` could be used to run tests for all supported Python versions. Install it (using ‘`pip install tox`’) and then run `tox` from the root directory - tests will be executed for all available Python interpreters.



### 1.18.0 (2017-08-03)

- Include additional assets used for distribution packages in the source tarball
- Consider [ and ] as safe characters in path and query components of URLs, i.e. they are not escaped anymore
- Disable codecov project coverage check

### 1.17.0 (2017-02-08)

- Add Python 3.5 and 3.6 support
- Add `w3lib.url.parse_data_uri` helper for parsing “data:” URIs
- Add `w3lib.html.strip_html5_whitespace` function to strip leading and trailing whitespace as per W3C recommendations, e.g. for cleaning “href” attribute values
- Fix `w3lib.http.headers_raw_to_dict` for multiple headers with same name
- Do not distribute tests/test\_\*.pyc artifacts

### 1.16.0 (2016-11-10)

- `canonicalize_url()` and `safe_url_string()`: strip “:” when no port is specified (as per RFC 3986; see also <https://github.com/scrapy/scrapy/issues/2377>)
- `url_query_cleaner()`: support new `keep_fragments` argument (defaulting to `False`)

## 1.15.0 (2016-07-29)

- Add `canonicalize_url()` to `w3lib.url`

## 1.14.3 (2016-07-14)

Bugfix release:

- Handle IDNA encoding failures in `safe_url_string()` (issue #62)

## 1.14.2 (2016-04-11)

Bugfix release:

- fix function import for (deprecated) `urljoin_rfc` (issue #51)
- only expose wanted functions from `w3lib.url`, via `__all__` (see issue #54, <https://github.com/scrapy/scrapy/issues/1917>)

## 1.14.1 (2016-04-07)

Bugfix release:

- For bytes URLs, when supplied encoding (or default UTF8) is wrong, `safe_url_string` falls back to percent-encoding offending bytes.

## 1.14.0 (2016-04-06)

Changes to `safe_url_string`:

- proper handling of non-ASCII characters in Python2 and Python3
- support IDNs
- new `path_encoding` to override default UTF-8 when serializing non-ASCII characters before percent-encoding

`html_body_declared_encoding` also detects encoding when not sole attribute in `<meta>`.

Package is now properly marked as `zip_safe`.

## 1.13.0 (2015-11-05)

- `remove_tags` removes uppercase tags as well;
- ignore meta-redirects inside script or noscript tags by default, but add an option to not ignore them;
- `replace_entities` now handles entities without trailing semicolon;
- fixed uncaught `UnicodeDecodeError` when decoding entities.



## 1.12.0 (2015-06-29)

- meta\_refresh regex now handles leading newlines and whitespaces in the url;
- include tests folder in source distribution.

## 1.11.0 (2015-01-13)

- url\_query\_cleaner now supports str or list parameters;
- add support for resolving base URLs in <base> tags with attributes before href.

## 1.10.0 (2014-08-20)

- reverted all 1.9.0 changes.

## 1.9.0 (2014-08-16)

- all url-related functions accept bytes and unicode and now return bytes.

## 1.8.1 (2014-08-14)

- w3lib.http.basic\_auth\_header now returns bytes

## 1.8.0 (2014-07-31)

- add support for big5-hkscs encoding.

## 1.7.1 (2014-07-26)

- PY3 fixed headers\_raw\_to\_dict and headers\_dict\_to\_raw;
- documentation improvements;
- provide wheels.

## 1.6 (2014-06-03)

- w3lib.form.encode\_multipart is deprecated;
- docstrings and docs are improved;
- w3lib.url.add\_or\_replace\_parameter is re-implemented on top of stdlib functions;
- remove\_entities is renamed to replace\_entities.

## 1.5 (2013-11-09)

- Python 2.6 support is dropped.

## 1.4 (2013-10-18)

- Python 3 support;
- `get_meta_refresh` encoding handling is fixed;
- check for '?' in `add_or_replace_parameter`;
- ISO-8859-1 is used for HTTP Basic Auth;
- fixed unicode handling in `replace_escape_chars`;

## 1.3 (2012-05-13)

- support non-standard `gb_2312_80` encoding;
- drop Python 2.5 support.

## 1.2 (2012-05-02)

- Detect encoding for content attr before `http-equiv` in meta tag.

## 1.1 (2012-04-18)

- `w3lib.html.remove_comments` handles multiline comments;
- Added `w3lib.encoding` module, containing functions for working with character encoding, like encoding autodetection from HTML pages.
- `w3lib.url.urljoin_rfc` is deprecated.

## 1.0 (2011-04-17)

First release of w3lib.

## History

The code of w3lib was originally part of the [Scrapy framework](#) but was later stripped out of Scrapy, with the aim of make it more reusable and to provide a useful library of web functions without depending on Scrapy.

## CHAPTER 7

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



**W**

w3lib.encoding, 3

w3lib.html, 5

w3lib.http, 8

w3lib.url, 9



**A**

add\_or\_replace\_parameter() (in module w3lib.url), 9  
any\_to\_uri() (in module w3lib.url), 9

**B**

basic\_auth\_header() (in module w3lib.http), 8

**C**

canonicalize\_url() (in module w3lib.url), 9

**F**

file\_uri\_to\_path() (in module w3lib.url), 10

**G**

get\_base\_url() (in module w3lib.html), 5  
get\_meta\_refresh() (in module w3lib.html), 5

**H**

headers\_dict\_to\_raw() (in module w3lib.http), 8  
headers\_raw\_to\_dict() (in module w3lib.http), 8  
html\_body\_declared\_encoding() (in module w3lib.encoding), 3  
html\_to\_unicode() (in module w3lib.encoding), 3  
http\_content\_type\_encoding() (in module w3lib.encoding), 4

**P**

parse\_data\_uri() (in module w3lib.url), 10  
path\_to\_file\_uri() (in module w3lib.url), 10

**R**

read\_bom() (in module w3lib.encoding), 5  
remove\_comments() (in module w3lib.html), 5  
remove\_entities() (in module w3lib.html), 6  
remove\_tags() (in module w3lib.html), 6  
remove\_tags\_with\_content() (in module w3lib.html), 6  
replace\_entities() (in module w3lib.html), 7  
replace\_escape\_chars() (in module w3lib.html), 7

replace\_tags() (in module w3lib.html), 7  
resolve\_encoding() (in module w3lib.encoding), 5

**S**

safe\_download\_url() (in module w3lib.url), 10  
safe\_url\_string() (in module w3lib.url), 10  
strip\_html5\_whitespace() (in module w3lib.html), 7

**T**

to\_unicode() (in module w3lib.encoding), 5

**U**

unquote\_markup() (in module w3lib.html), 8  
url\_query\_cleaner() (in module w3lib.url), 10  
url\_query\_parameter() (in module w3lib.url), 11  
urljoin\_rfc() (in module w3lib.url), 11

**W**

w3lib.encoding (module), 3  
w3lib.html (module), 5  
w3lib.http (module), 8  
w3lib.url (module), 9