

---

**urlextract**

***Release 1.8.0***

**Dec 19, 2022**



---

## Contents

---

<b>1 urlextract - command line</b>	<b>3</b>
<b>2 URLExtract class</b>	<b>5</b>
<b>3 Indices and tables</b>	<b>9</b>
<b>Index</b>	<b>11</b>



urlextract is package with python class and command line script used for extraction of URLs from given text.



# CHAPTER 1

---

## urlextract - command line

---

urlextract - command line program that will print all URLs to stdout

Usage:                   \$ urlextract [-h] [-v] [-u] [-dl] [-c] [-i <ignore\_file>] [-p <permit\_file>] [-l LIMIT] [<input\_file>]

**positional arguments:** <input\_file> input text file with URLs to extract

**optional arguments:**

- h, --help               show this help message and exit
- v, --version           show program's version number and exit
- u, --unique           print out only unique URLs found in file
- dl, --disable-localhost    disable extracting "localhost" as URL
- c, --check-dns        print out only URLs for existing domain names
- i <ignore\_file>, --ignore-file <ignore\_file>   input text file with URLs to exclude from extraction
- p <permit\_file>, --permit-file <permit\_file>   input text file with URLs that can be processed
- l LIMIT, --limit LIMIT   Maximum count of URLs that can be processed. Set 0 to disable the limit. Default: 10000



# CHAPTER 2

---

## URLExtract class

---

```
class urlextract.URLExtract(extract_email=False, cache_dns=True, extract_localhost=True, limit=10000, allow_mixed_case_hostname=True, **kwargs)
```

Class for finding and extracting URLs from given string.

### Examples:

```
from urlextract import URLExtract

extractor = URLExtract()
urls = extractor.find_urls("Let's have URL example.com example.")
print(urls) # prints: ['example.com']

# Another way is to get a generator over found URLs in text:
for url in extractor.gen_urls(example_text):
    print(url) # prints: ['example.com']

# Or if you want to just check if there is at least one URL in text:
if extractor.has_urls(example_text):
    print("Given text contains some URL")
```

### add\_enclosure(left\_char: str, right\_char: str)

Add new enclosure pair of characters. That should be removed when their presence is detected at beginning and end of found URL

#### Parameters

- **left\_char** (str) – left character of enclosure pair - e.g. “(”
- **right\_char** (str) – right character of enclosure pair - e.g. “)”

### allow\_mixed\_case\_hostname

If set to True host should contain mixed case letters (upper-case and lower-case)

#### Return type

bool

### extract\_email

If set to True email will be extracted from text

**Return type** bool

**extract\_localhost**

If set to True ‘localhost’ will be extracted as URL from text

**Return type** bool

**find\_urls** (*text: str, only\_unique=False, check\_dns=False, get\_indices=False,*

*with\_schema\_only=False*) → List[Union[str, Tuple[str, Tuple[int, int]]]]

Find all URLs in given text.

**Parameters**

- **text** (*str*) – text where we want to find URLs
- **only\_unique** (*bool*) – return only unique URLs
- **check\_dns** (*bool*) – filter results to valid domains
- **get\_indices** (*bool*) – whether to return beginning and ending indices as (<url>, (idx\_begin, idx\_end))
- **with\_schema\_only** (*bool*) – get domains with schema only (e.g. <https://janlipovsky.cz> but not example.com)

**Returns** list of URLs found in text

**Return type** list

**Raises** **URLExtractError** – Raised when count of found URLs reaches given limit. Processed URLs are returned in *data* argument.

**gen\_urls** (*text: str, check\_dns=False, get\_indices=False, with\_schema\_only=False*) → Generator[Union[str, Tuple[str, Tuple[int, int]]], None, None]

Creates generator over found URLs in given text.

**Parameters**

- **text** (*str*) – text where we want to find URLs
- **check\_dns** (*bool*) – filter results to valid domains
- **get\_indices** (*bool*) – whether to return beginning and ending indices as (<url>, (idx\_begin, idx\_end))
- **with\_schema\_only** (*bool*) – get domains with schema only

**Yields** URL or URL with indices found in text or empty string if nothing was found

**Return type** str|tuple(str, tuple(int, int))

**get\_after\_tld\_chars()** → List[str]

Returns list of chars that are allowed after TLD

**Returns** list of chars that are allowed after TLD

**Return type** list

**get\_enclosures()** → Set[Tuple[str, str]]

Returns set of enclosure pairs that might be used to enclosure URL. For example brackets (example.com), [example.com], {example.com}

**Returns** set of tuple of enclosure characters

**Return type** set(tuple(str,str))

**get\_stop\_chars\_left()** → Set[str]

Returns set of stop chars for text on left from TLD.

**Returns** set of stop chars

**Return type** set

**get\_stop\_chars\_left\_from\_scheme()** → Set[str]

Returns set of stop chars for text on left from scheme.

**Returns** set of stop chars

**Return type** set

**get\_stop\_chars\_right()** → Set[str]

Returns set of stop chars for text on right from TLD.

**Returns** set of stop chars

**Return type** set

**static get\_version()** → str

Returns version number.

**Returns** version number

**Return type** str

**has\_urls** (*text*: str, *check\_dns*=False, *with\_schema\_only*=False) → bool

Checks if text contains any valid URL. Returns True if text contains at least one URL.

**Parameters**

- **text** – text where we want to find URLs
- **check\_dns** (bool) – filter results to valid domains
- **with\_schema\_only** (bool) – consider domains with schema only

**Returns** True if at least one URL was found, False otherwise

**Return type** bool

**ignore\_list**

Set of URLs to be ignored (not returned) while extracting from text

**Returns** Returns set of ignored URLs

**Return type** set(str)

**load\_ignore\_list** (*file\_name*)

Load URLs from file into ignore list

**Parameters** **file\_name** (str) – path to file containing URLs

**load\_permit\_list** (*file\_name*)

Load URLs from file into permit list

**Parameters** **file\_name** (str) – path to file containing URLs

**permit\_list**

Set of URLs that can be processed

**Returns** Returns set of URLs that can be processed

**Return type** set(str)

**remove\_enclosure** (*left\_char*: str, *right\_char*: str)

Remove enclosure pair from set of enclosures.

**Parameters**

- **left\_char** (*str*) – left character of enclosure pair - e.g. “(”
- **right\_char** (*str*) – right character of enclosure pair - e.g. “)”

**set\_after\_tld\_chars** (*after\_tld\_chars: Iterable[str]*)  
Set chars that are allowed after TLD.

**Parameters** **after\_tld\_chars** (*list*) – list of characters

**set\_stop\_chars\_left** (*stop\_chars: Set[str]*)  
Set stop characters for text on left from TLD. Stop characters are used when determining end of URL.

**Parameters** **stop\_chars** (*set*) – set of characters

**Raises** `TypeError`

**set\_stop\_chars\_left\_from\_scheme** (*stop\_chars: Set[str]*)  
Set stop characters for text on left from scheme. Stop characters are used when determining end of URL.

**Parameters** **stop\_chars** (*set*) – set of characters

**Raises** `TypeError`

**set\_stop\_chars\_right** (*stop\_chars: Set[str]*)  
Set stop characters for text on right from TLD. Stop characters are used when determining end of URL.

**Parameters** **stop\_chars** (*set*) – set of characters

**Raises** `TypeError`

**update()**  
Update TLD list cache file.

**Returns** True if update was successful False otherwise

**Return type** `bool`

**update\_when\_older** (*days: int*) → `bool`

Update TLD list cache file if the list is older than number of days given in parameter *days* or if it does not exist.

**Parameters** **days** (*int*) – number of days from last change

**Returns** True if update was successful, False otherwise

**Return type** `bool`

# CHAPTER 3

---

## Indices and tables

---

- genindex
- modindex
- search



---

## Index

---

### A

add\_enclosure () (*urlextract.URLEExtract method*), 5  
allow\_mixed\_case\_hostname (*urlextract.URLEExtract attribute*), 5

### E

extract\_email (*urlextract.URLEExtract attribute*), 5  
extract\_localhost (*urlextract.URLEExtract attribute*), 6

### F

find\_urls () (*urlextract.URLEExtract method*), 6

### G

gen\_urls () (*urlextract.URLEExtract method*), 6  
get\_after\_tld\_chars () (*urlextract.URLEExtract method*), 6  
get\_enclosures () (*urlextract.URLEExtract method*), 6  
get\_stop\_chars\_left () (*urlextract.URLEExtract method*), 6  
get\_stop\_chars\_left\_from\_scheme () (*urlextract.URLEExtract method*), 7  
get\_stop\_chars\_right () (*urlextract.URLEExtract method*), 7  
get\_version () (*urlextract.URLEExtract static method*), 7

### H

has\_urls () (*urlextract.URLEExtract method*), 7

### I

ignore\_list (*urlextract.URLEExtract attribute*), 7

### L

load\_ignore\_list () (*urlextract.URLEExtract method*), 7  
load\_permit\_list () (*urlextract.URLEExtract method*), 7

### P

permit\_list (*urlextract.URLEExtract attribute*), 7

### R

remove\_enclosure () (*urlextract.URLEExtract method*), 7

### S

set\_after\_tld\_chars () (*urlextract.URLEExtract method*), 8  
set\_stop\_chars\_left () (*urlextract.URLEExtract method*), 8  
set\_stop\_chars\_left\_from\_scheme () (*urlextract.URLEExtract method*), 8  
set\_stop\_chars\_right () (*urlextract.URLEExtract method*), 8

### U

update () (*urlextract.URLEExtract method*), 8  
update\_when\_older () (*urlextract.URLEExtract method*), 8  
URLEExtract (*class in urlextract*), 5