

---

# **textract Documentation**

*Release 1.0.0*

**Dean Malmgren**

August 26, 2014



<b>1</b>	<b>Currently supporting</b>	<b>3</b>
<b>2</b>	<b>Related projects</b>	<b>5</b>
2.1	Command line interface . . . . .	5
2.2	Python package . . . . .	6
2.3	Installation . . . . .	10
2.4	Contributing . . . . .	11
2.5	Change Log . . . . .	13
<b>3</b>	<b>Indices and tables</b>	<b>15</b>
	<b>Python Module Index</b>	<b>17</b>



As undesirable as it might be, more often than not there is extremely useful information embedded in Word documents, PowerPoint presentations, PDFs, etc—so-called “dark data”—that would be valuable for further textual analysis and visualization. While *several packages* exist for extracting content from each of these formats on their own, this package provides a single interface for extracting content from any type of file, without any irrelevant markup.

This package provides two primary facilities for doing this, the *command line interface*

```
textract path/to/file.extension
```

or the *python package*

```
# some python file
import textract
text = textract.process("path/to/file.extension")
```



---

## Currently supporting

---

- `.doc` via `antiword`
- `.docx` via `python-docx`
- `.eml` via `python` builtins
- `.epub` via `ebooklib`
- `.gif` via `tesseract-ocr`
- `.jpg` and `.jpeg` via `tesseract-ocr`
- `.json` via `python` builtins
- `.html` via `beautifulsoup4`
- `.odt` via `python` builtins
- `.pdf` via `pdftotext` (default) or `pdfminer`
- `.png` via `tesseract-ocr`
- `.pptx` via `python-pptx`
- `.ps` via `ps2text`
- `.txt` via `python` builtins
- `.xls` via `xldr`
- `.xlsx` via `xldr`

Please recommend other file types by either mentioning them on the [issue tracker](#) or by *contributing*





---

## Related projects

---

Of course, `textract` isn't the first project with the aim to provide a simple interface for extracting text from any document. But this is, to the best of my knowledge, the only project that is written in python (a language commonly chosen by the natural language processing community) and is method agnostic about how content is extracted (more on this [here](#)). Here is a small sample of similar projects (feel free to add to the list):

- [Apache Tika](#) has *very similar, if not identical, aims as `textract`*. It has impressive coverage of a wide range of file formats and is written in java.
- `textract (node.js)` has similar aims as this `textract` package (including an identical name! great minds...). It is written in node.js.
- `pandoc` is intended to be a document conversion (a much more difficult task!), but it does have *the ability to convert to plain text*. It is written in Haskell.

Contents:

## 2.1 Command line interface

### 2.1.1 `textract`

Command line tool for extracting text from any document.

```
usage: textract [-h]
               [-e {aliases,ascii,base64_codec,big5,big5hkscs,bz2_codec,charmap,cp037,cp1006,cp1026,cp1140,cp1250,cp1251,
               cp1252,cp1253,cp1254,cp1255,cp1256,cp1257,cp1258,cp424,
               cp437,cp500,cp720,cp737,cp775,cp850,cp852,cp855,cp856,
               cp857,cp858,cp860,cp861,cp862,cp863,cp864,cp865,cp866,
               cp869,cp874,cp875,cp932,cp949,cp950,euc_jis_2004,euc_jisx0213,
               euc_jp,euc_kr,gb18030,gb2312,gbk,hex_codec,hp_roman8,hz,
               ...}]
               [-m METHOD] [-o OUTPUT] [-v]
               filename
```

#### Positional arguments:

**filename**                      Filename to extract text.

#### Options:

**-e=utf\_8, --encoding=utf\_8** Specify the encoding of the output.

Possible choices: `aliases`, `ascii`, `base64_codec`, `big5`, `big5hkscs`, `bz2_codec`, `charmap`, `cp037`, `cp1006`, `cp1026`, `cp1140`, `cp1250`, `cp1251`, `cp1252`, `cp1253`, `cp1254`, `cp1255`, `cp1256`, `cp1257`, `cp1258`, `cp424`, `cp437`, `cp500`, `cp720`, `cp737`, `cp775`, `cp850`, `cp852`, `cp855`, `cp856`, `cp857`, `cp858`, `cp860`, `cp861`, `cp862`, `cp863`, `cp864`, `cp865`, `cp866`, `cp869`, `cp874`, `cp875`, `cp932`, `cp949`, `cp950`, `euc_jis_2004`, `euc_jisx0213`, `euc_jp`, `euc_kr`, `gb18030`, `gb2312`, `gbk`, `hex_codec`, `hp_roman8`, `hz`,

idna, iso2022\_jp, iso2022\_jp\_1, iso2022\_jp\_2, iso2022\_jp\_2004, iso2022\_jp\_3, iso2022\_jp\_ext, iso2022\_kr, iso8859\_1, iso8859\_10, iso8859\_11, iso8859\_13, iso8859\_14, iso8859\_15, iso8859\_16, iso8859\_2, iso8859\_3, iso8859\_4, iso8859\_5, iso8859\_6, iso8859\_7, iso8859\_8, iso8859\_9, johab, koi8\_r, koi8\_u, latin\_1, mac\_arabic, mac\_centeuro, mac\_croatian, mac\_cyrillic, mac\_farsi, mac\_greek, mac\_iceland, mac\_latin2, mac\_roman, mac\_romanian, mac\_turkish, mbcs, palmos, ptcp154, punycode, quopri\_codec, raw\_unicode\_escape, rot\_13, shift\_jis, shift\_jis\_2004, shift\_jisx0213, string\_escape, tactis, tis\_620, undefined, unicode\_escape, unicode\_internal, utf\_16, utf\_16\_be, utf\_16\_le, utf\_32, utf\_32\_be, utf\_32\_le, utf\_7, utf\_8, utf\_8\_sig, uu\_codec, zlib\_codec

**-m=, --method=** specify a method of extraction for formats that support it  
**-o=, --output=-** output raw text in this file  
**-v, --version** show program's version number and exit

---

**Note:** To make the command line interface as usable as possible, autocompletion of available options with `textract` is enabled by @kislyuk's amazing `argcomplete` package. Follow instructions to [enable global autocomplete](#) and you should be all set. As an example, this is also configured in the [virtual machine provisioning for this project](#).

---

## 2.2 Python package

This package is organized to make it as easy as possible to add new extensions and support the continued growth and coverage of `textract`. For almost all applications, you will just have to do something like this:

```
import textract
text = textract.process('path/to/file.extension')
```

to obtain text from a document.

For completeness, we also include here the documentation for specific file extension parsers as well as a few other essential bits in the `textract.exceptions` and `textract.shell` module.

### 2.2.1 `textract.parsers.doc_parser` module

```
class textract.parsers.doc_parser.Parser
    Bases: textract.parsers.utils.ShellParser
    Extract text from doc files using antiword.
    extract (filename, **kwargs)
```

### 2.2.2 `textract.parsers.docx_parser` module

```
class textract.parsers.docx_parser.Parser
    Bases: textract.parsers.utils.BaseParser
    Extract text from docx file using python-docx.
    extract (filename, **kwargs)
```

### 2.2.3 textextract.parsers.eml\_parser module

**class** `textextract.parsers.eml_parser.Parser`  
 Bases: `textextract.parsers.utils.BaseParser`

Extract text from email messages in .eml format. This gets the subject and all text from the contents.

**extract** (*filename*, *\*\*kwargs*)

### 2.2.4 textextract.parsers.epub\_parser module

**class** `textextract.parsers.epub_parser.Parser`  
 Bases: `textextract.parsers.utils.BaseParser`

Extract text from epub using python epub library

**extract** (*filename*, *\*\*kwargs*)

### 2.2.5 textextract.parsers.gif\_parser module

### 2.2.6 textextract.parsers.html\_parser module

**class** `textextract.parsers.html_parser.Parser`  
 Bases: `textextract.parsers.utils.BaseParser`

Extract text from html file using beautifulsoup4. Filter text to only show the visible parts of the page. Inspiration from [here](#).

**extract** (*filename*, *\*\*kwargs*)

### 2.2.7 textextract.parsers.jpg\_parser module

### 2.2.8 textextract.parsers.json\_parser module

**class** `textextract.parsers.json_parser.Parser`  
 Bases: `textextract.parsers.utils.BaseParser`

Extract all of the string values of a json file (no keys as those are, in some sense, markup). This is useful for parsing content from mongodb dumps, for example.

**extract** (*filename*, *\*\*kwargs*)

**get\_text** (*deserialized\_json*)

Recursively get text from subcomponents of a deserialized json. To enforce the same order on the documents, make sure to read keys of `deserialized_json` in a consistent (alphabetical) order.

### 2.2.9 textextract.parsers.odt\_parser module

**class** `textextract.parsers.odt_parser.Parser`  
 Bases: `textextract.parsers.utils.BaseParser`

Extract text from open document files.

**extract** (*filename*, *\*\*kwargs*)

**text\_to\_string** (*element*)

`to_string()`

Converts the document to a string.

## 2.2.10 `textract.parsers.pdf_parser` module

**class** `textract.parsers.pdf_parser.Parser`

Bases: `textract.parsers.utils.ShellParser`

Extract text from pdf files using either the `pdftotext` method (default) or the `pdfminer` method.

**extract** (*filename*, *method=''*, *\*\*kwargs*)

**extract\_pdfminer** (*filename*)

Extract text from pdfs using pdfminer.

**extract\_pdftotext** (*filename*)

Extract text from pdfs using the pdftotext command line utility.

## 2.2.11 `textract.parsers.png_parser` module

## 2.2.12 `textract.parsers.pptx_parser` module

**class** `textract.parsers.pptx_parser.Parser`

Bases: `textract.parsers.utils.BaseParser`

Extract text from pptx file using python-pptx

**extract** (*filename*, *\*\*kwargs*)

## 2.2.13 `textract.parsers.ps_parser` module

**class** `textract.parsers.ps_parser.Parser`

Bases: `textract.parsers.utils.ShellParser`

Extract text from postscript files using pstotext command.

**extract** (*filename*, *\*\*kwargs*)

## 2.2.14 `textract.parsers.tesseract` module

Process an image file using tesseract.

**class** `textract.parsers.tesseract.Parser`

Bases: `textract.parsers.utils.ShellParser`

Extract text from various image file formats using tesseract-ocr

**extract** (*filename*, *\*\*kwargs*)

## 2.2.15 `textract.parsers.txt_parser` module

**class** `textract.parsers.txt_parser.Parser`

Bases: `textract.parsers.utils.BaseParser`

Parse `.txt` files

**extract** (*filename*, *\*\*kwargs*)

## 2.2.16 `textract.parsers.utils` module

This module includes a bunch of convenient base classes that are reused in many of the other parser modules.

**class** `textract.parsers.utils.BaseParser`

Bases: `object`

The `BaseParser` abstracts out some common functionality that is used across all document formats. Specifically, it owns the responsibility of handling all unicode and byte-encoding problems.

Inspiration from <http://nedbatchelder.com/text/unipain.html>

**decode** (*text*)

Decode text using the `chardet` package

**encode** (*text*, *encoding*)

Encode the `text` in `encoding` byte-encoding. This ignores code points that can't be encoded in byte-strings.

**extract** (*filename*, *\*\*kwargs*)

**process** (*filename*, *encoding*, *\*\*kwargs*)

Process `filename` and encode byte-string with `encoding`.

**class** `textract.parsers.utils.ShellParser`

Bases: `textract.parsers.utils.BaseParser`

The `ShellParser` extends the `BaseParser` to make it easy to run external programs from the command line with Fabric-like behavior.

**run** (*command*)

Run `command` and return the subsequent `stdout` and `stderr`.

**temp\_filename** ()

Return a unique tempfile name.

## 2.2.17 `textract.parsers.xls_parser` module

## 2.2.18 `textract.parsers.xlsx_parser` module

**class** `textract.parsers.xlsx_parser.Parser`

Bases: `textract.parsers.utils.BaseParser`

Extract text from Excel files (`.xls/xlsx`).

**extract** (*filename*, *\*\*kwargs*)

## 2.2.19 `textract.cli` module

Use `argparse` to handle command-line arguments.

`textract.cli.get_parser` ()

Initialize the parser for the command line interface and bind the autocompletion functionality

## 2.2.20 textract.exceptions module

**exception** `textract.exceptions.CommandLineError`

Bases: `exceptions.Exception`

The traceback of all `CommandLineError`'s is suppressed when the errors occur on the command line to provide a useful command line interface.

**render** (*msg*)

**exception** `textract.exceptions.ExtensionNotSupported` (*ext*)

Bases: `textract.exceptions.CommandLineError`

This error is raised with unsupported extensions

**exception** `textract.exceptions.MissingFileError` (*filename*)

Bases: `textract.exceptions.CommandLineError`

This error is raised when the file can not be located at the specified path.

**exception** `textract.exceptions.ShellError` (*command, exit\_code, stdout, stderr*)

Bases: `textract.exceptions.CommandLineError`

This error is raised when a shell.run returns a non-zero exit code (meaning the command failed).

**failed\_message** ()

**is\_uninstalled** ()

**uninstalled\_message** ()

**exception** `textract.exceptions.UnknownMethod` (*method*)

Bases: `textract.exceptions.CommandLineError`

This error is raised when the specified `-method` on the command line is unknown.

## 2.3 Installation

One of the main goals of textract is to make it as easy as possible to start using textract (meaning that installation should be as quick and painless as possible). This package is built on top of several python packages and other source libraries. Assuming you are using `pip` or `easy_install` to install textract, the `python packages` are all installed by default with textract. The source libraries are a separate matter though and largely depend on your operating system.

### 2.3.1 Ubuntu / Debian

There are two steps required to run this package on Ubuntu/Debian. First you must install some system packages using the `apt-get` package manager before installing textract from `pypi`.

```
apt-get install python-dev libxml2-dev libxslt1-dev antiword poppler-utils pstotext tesseract-ocr
pip install textract
```

---

**Note:** It may also be necessary to install `zlib1g-dev` on Docker instances of Ubuntu. See [issue #19](#) for details

---

### 2.3.2 OSX

These steps rely on you having `homebrew` installed as well as the `cask` plugin (`brew install caskroom/cask/brew-cask`). The basic idea is to first install `XQuartz` before installing a bunch of system packages before installing `textextract` from `pypi`.

```
brew cask install xquartz
brew install poppler antiword tesseract
pip install textextract
```

---

**Note:** `pstotext` is not currently a part of `homebrew` so `.ps` extraction must be enabled by manually installing from source.

---

**Note:** Depending on how you have python configured on your system with `homebrew`, you may also need to install the python development header files for `textextract` to properly install.

---

### 2.3.3 Don't see your operating system installation instructions here?

My apologies! Installing system packages is a bit of a drag and its hard to anticipate all of the different environments that need to be accomodated (wouldn't it be awesome if there were a system-agnostic package manager or, better yet, if python could install these system dependencies for you?!?!). If you're operating system doesn't have documentation about how to install the `textextract` dependencies, please *contribute a pull request* with:

1. A new section in here with the appropriate details about how to install things. In particular, please give instructions for how to install the following libraries before running `pip install textextract`:
  - `libxml2 2.6.21 or later` is required by the `.docx` parser which uses `lxml` via `python-docx`.
  - `libxslt 1.1.15 or later` is required by the `.docx` parser which users `lxml` via `python-docx`.
  - python header files are required for building `lxml`.
  - `antiword` is required by the `.doc` parser.
  - `pdftotext` is *optionally* required by the `.pdf` parser (there is a pure python fallback that works if `pdftotext` isn't installed).
  - `pstotext` is required by the `.ps` parser.
  - `tesseract-ocr` is required by the `.jpg`, `.png` and `.gif` parser.
2. Add a requirements file to the `requirements` directory of the project with the lower-cased name of your operating system (e.g. `requirements/windows`) so we can try to keep these things up to date in the future.

## 2.4 Contributing

The overarching goal of this project is to make it as easy as possible to extract raw text from any document for the purposes of most natural language processing tasks. In practice, this means that this project should preferentially provide tools that correctly produce output that has words in the correct order but that whitespace between words, formatting, etc is totally irrelevant.

Importantly, this project is committed to being as agnostic about how the content is extracted as it is about the means in which the text is analyzed downstream. This means that `textextract` should support multiple modes of extracting text from any document and provide reasonably good defaults (defaulting to tools that tend to produce the correct word sequence).

Another important aspect of this project is that we want to have extremely good documentation. If you notice a type-o, error, confusing statement etc, please fix it!

### 2.4.1 Quick start

1. Fork and clone the project:

```
git clone https://github.com/YOUR-USERNAME/textract.git
```

2. Contribute! There are several [open issues](#) that provide good places to dig in. Check out the [contribution guidelines](#) and send pull requests; your help is greatly appreciated!

Depending on your development preferences, there are lots of ways to get started developing with textract:

#### Developing in a native Ubuntu environment

3. Install all the necessary system packages:

```
./provision/travis-mock.sh
./provision/debian.sh

# optionally run some of the steps in these scripts, but you
# may want to be selective about what you do as they alter global
# environment states
./provision/python.sh
./provision/development.sh
```

4. On the virtual machine, make sure everything is working by running the suite of functional tests:

```
nosetests
```

These functional tests are designed to be run on an Ubuntu 12.04 LTS server, just like the virtual machine and the server that runs the travis-ci test suite. There are some other tests that have been added along the way in the [Travis configuration](#). For your convenience, you can run all of these tests with:

```
./tests/run.py
```

Current build status:

#### Developing with Vagrant virtual machine

3. Install [Vagrant](#) and [Virtualbox](#) and launch the development virtual machine:

```
vagrant plugin install iniparse
vagrant up && vagrant provision
```

On vagrant sshing to the virtual machine, note that the `PYTHONPATH` and `PATH` [environment variables](#) have been altered in this virtual machine so that any changes you make to textract in development are automatically incorporated into the command.

4. See [step 4](#) in the Ubuntu development environment. Current build status:

#### Developing with Docker container

3. Go to the [Docker documentation](#) and follow the instructions under “If you’d like to try the latest version of Docker” to install Docker.



4. Just run `tests/run_docker_tests.sh` to run the full test suite. Current build status:

## 2.5 Change Log

This project uses [semantic versioning](#) to track version numbers, where backwards incompatible changes (highlighted in **bold**) bump the major version of the package.

### 2.5.1 latest changes in development for next release

#### 2.5.2 1.0.0

- **standardized encoding of output with “-e/-encoding“ option (#39)**
- support for `.xls` and `.xlsx` files (#42 and #55 by @levivm)
- support for `.epub` files (#40 by @kokxx)
- several bug fixes, including:
  - removing tesseract version info from output of image parsers (#48)
  - problems with spaces in filenames (#53)
  - concurrency problems with tesseract (#44 by @ShawnMilo, #41 by @christomitov)
- several internal improvements, including:
  - switching to using class-based parsers to abstract away the common functionality between different parser classes (#39)
  - switching to using a python-based test suite and added standardized text tests to make sure output is consistent across file types (#49)
  - including support for Docker-based testing (#46 by @ShawnMilo)

#### 2.5.3 0.5.1

- several bug fixes, including:
  - documentation fixes
  - shell commands hanging on large files (#33)

#### 2.5.4 0.5.0

- support for `.json` files (#13 by @anthonygarvan)
- support for `.odt` files (#29 by @christomitov)
- support for `.ps` files (#25)
- support for `.gif`, `.jpg`, `.jpeg`, and `.png` files (#30 by @christomitov)
- several bug fixes, including:
  - improved fallback handling in `.pdf` parser if the `pdftotext` command line utility isn't installed (#26)
  - improved documentation for installation instructions on non-Ubuntu operating systems (#21, #26)

- several internal improvements, including:
  - cleaned up implementation of extension parsers to avoid magic

### **2.5.5 0.4.0**

- support for `.html` files (#7)
- support for `.eml` files (#4)
- automated the documentation for the python package using sphinx-apidoc in docs/Makefile (#9)

### **2.5.6 0.3.0**

- support for `.txt` files, haha (#8)
- fixed installation bug with not properly including requirements files in the manifest

### **2.5.7 0.2.0**

- support for `.doc` files (#2)
- support for `.pdf` files (#3)
- several bug fixes, including:
  - fixing tab complete bug no file paths (#6)
  - fixing tests to make sure the work properly on travis-ci

### **2.5.8 0.1.0**

- Initial release, support for `.docx` and `.pptx`

---

## Indices and tables

---

- *genindex*
- *modindex*
- *search*



## t

`textract.cli`, 9  
`textract.exceptions`, 10  
`textract.parsers.doc_parser`, 6  
`textract.parsers.docx_parser`, 6  
`textract.parsers.eml_parser`, 7  
`textract.parsers.epub_parser`, 7  
`textract.parsers.gif_parser`, 7  
`textract.parsers.html_parser`, 7  
`textract.parsers.jpg_parser`, 7  
`textract.parsers.json_parser`, 7  
`textract.parsers.odt_parser`, 7  
`textract.parsers.pdf_parser`, 8  
`textract.parsers.png_parser`, 8  
`textract.parsers.pptx_parser`, 8  
`textract.parsers.ps_parser`, 8  
`textract.parsers.tesseract`, 8  
`textract.parsers.txt_parser`, 8  
`textract.parsers.utils`, 9  
`textract.parsers.xls_parser`, 9  
`textract.parsers.xlsx_parser`, 9