

# Optical Character Recognition

All Versions

Optical Character Recognition is converting images of text into actual text. In these examples find ways of using OCR in python.

## Examples

### PyOCR

Another module of some use is PyOCR , source code of which is [here](#) .

Also simple to use and has more features than PyTesseract .

To initialize:

```
from PIL import Image
import sys

import pyocr
import pyocr.builders

tools = pyocr.get_available_tools()
# The tools are returned in the recommended order of usage
tool = tools[0]

langs = tool.get_available_languages()
lang = langs[0]
# Note that languages are NOT sorted in any way. Please refer
# to the system locale settings for the default language
# to use.
```

And some examples of usage:

```
txt = tool.image_to_string(
    Image.open('test.png'),
    lang=lang,
    builder=pyocr.builders.TextBuilder()
)
# txt is a Python string

word_boxes = tool.image_to_string(
    Image.open('test.png'),
    lang="eng",
    builder=pyocr.builders.WordBoxBuilder()
)
# list of box objects. For each box object:
# box.content is the word in the box
# box.position is its position on the page (in pixels)
#
# Beware that some OCR tools (Tesseract for instance)
# may return empty boxes

line_and_word_boxes = tool.image_to_string(
    Image.open('test.png'), lang="fra",
    builder=pyocr.builders.LineBoxBuilder()
)
# list of line objects. For each line object:
# line.word_boxes is a list of word boxes (the individual words in the line)
# line.content is the whole text of the line
# line.position is the position of the whole line on the page (in pixels)
#
# Beware that some OCR tools (Tesseract for instance)
# may return empty boxes

# Digits - Only Tesseract (not 'libtesseract' yet !)
digits = tool.image_to_string(
    Image.open('test-digits.png'),
    lang=lang,
    builder=pyocr.tesseract.DigitBuilder()
```

### PyTesseract

PyTesseract is an in-development python package for OCR.

Using PyTesseract is pretty easy:

```
try:
    import Image
except ImportError:
    from PIL import Image
```

```
from PIL import image

import pytesseract

#Basic OCR
print(pytesseract.image_to_string(Image.open('test.png')))

#In French
print(pytesseract.image_to_string(Image.open('test-european.jpg'), lang='fra'))
```

PyTesseract is open source and can be found [here](#) .

Syntax

Parameters

Remarks