

OPTICAL CHARACTER RECOGNITION



HELLO!

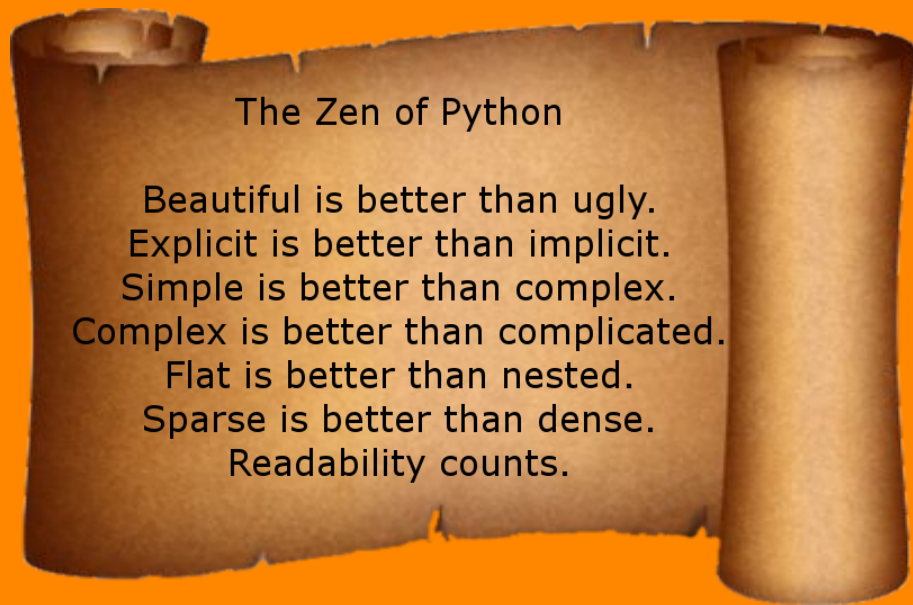
I am Toni Simu.

I am here because I want
to discuss about OCR.

(The guitar is for show-off)

What is OCR?





The Zen of Python

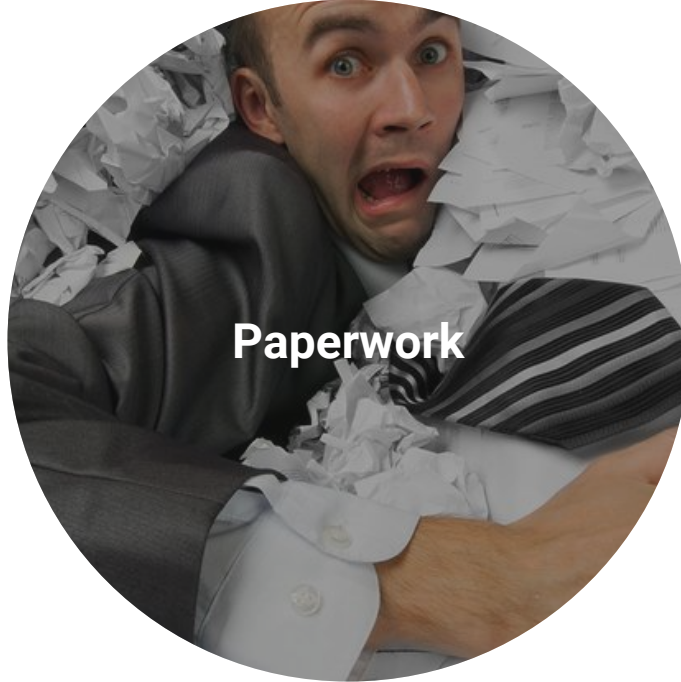
Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.
Complex is better than complicated.
Flat is better than nested.
Sparse is better than dense.
Readability counts.

Extracts text from the images

OCR

OCR Applications

OCR Applications - Paperwork



OCR Applications - Health



OCR Applications - Translation



How does OCR work?

How does OCR work?



**Pattern
recognition**

**Feature
recognition**

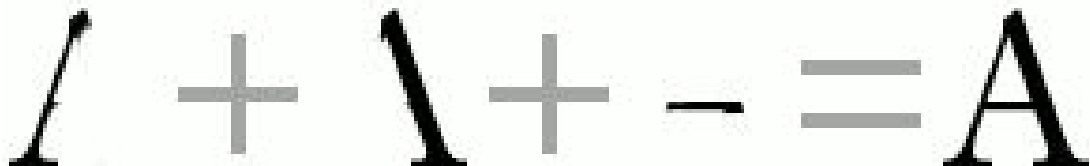
How does OCR work? – Pattern recognition

A B C D E F G
H I J K L M N O
P Q R S T U V
W X Y Z

1 2 3 4 5 6 7 8 9 0



How does OCR work? - Feature recognition



About tesseract

Beautiful is better than ugly.
Explicit is better than implicit.
Simple is better than complex.



Beautiful|is|better|than|ugly.



B|e|a|u|t|i|f|u|l

Preparing images for OCR

Image processing

Rescaling

Binarization

Noise
Removal

Rotation /
Deskewing

Border
Removal

Image processing - Rescaling



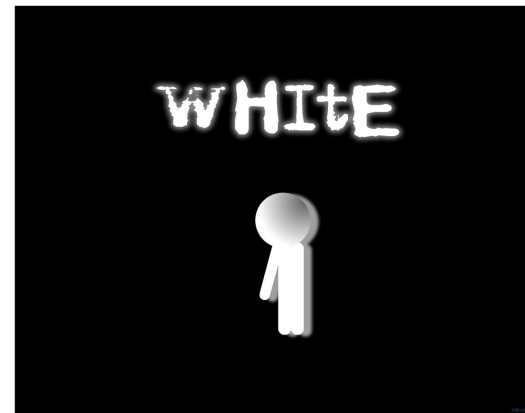
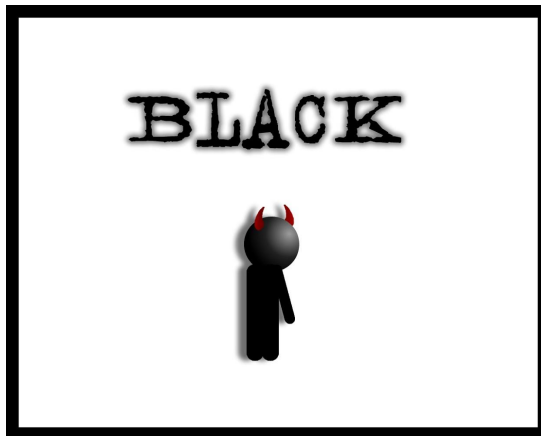
300 PPI



200 PPI

The image should be rescaled to have at least 300 ppi for a better accuracy.

Image processing - Binarization



- ▶ Turn the image to grayscale
- ▶ Apply Adaptive or Global Thresholding

Image processing - Binarization - Grayscale

YELLOW

RED

GREEN

BLACK

ORANGE

BLUE

RED

GREEN

BLACK

BLUE

ORANGE

YELLOW



YELLOW

RED

GREEN

BLACK

ORANGE

BLUE

RED

GREEN

BLACK

BLUE

ORANGE

YELLOW

Image processing - Binarization - Thresholding

Less problematic

YELLOW	RED	GREEN
BLACK	ORANGE	BLUE
RED	GREEN	BLACK
BLUE	ORANGE	YELLOW



	RED	GREEN
BLACK	ORANGE	BLUE
RED	GREEN	BLACK
BLUE	ORANGE	

$t=180$

More problematic

YELLOW	RED	GREEN
BLACK	ORANGE	BLUE
RED	GREEN	BLACK
BLUE	ORANGE	YELLOW



	RED	GREEN
	ORANGE	BLUE
	GREEN	BLACK
	ORANGE	

Image processing - Binarization - Adaptive Thresholding

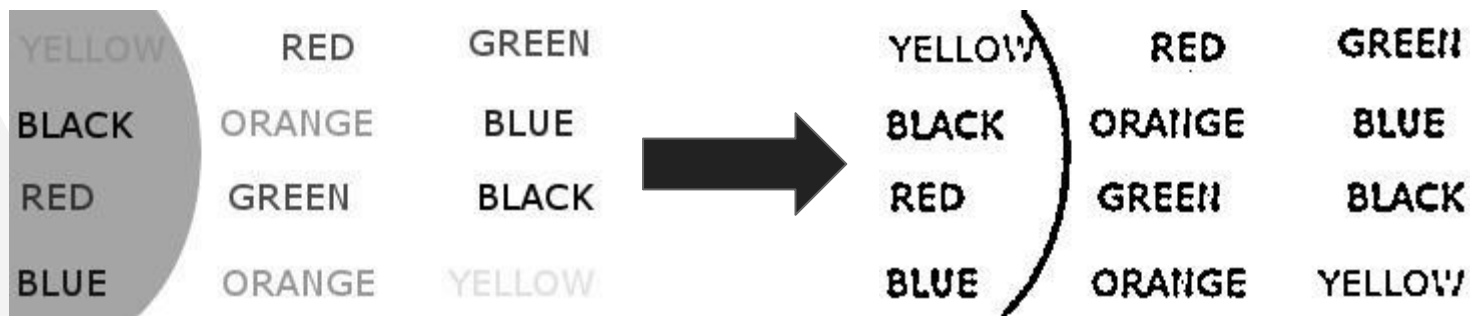


Image processing - Noise Removal

Clarification

In the last edition of the "Catholic New World," an article on the annual conference hosted by the Respect Life Office did not fully reflect the context of Cardinal George's remarks.

During an informal question-and-answer session with the archdiocese's parish Respect Life Coordinators, the cardinal emphasized that the participation by any person in the promotion of abortion, including through his or her political position, is a grave matter. While the issue of withholding Communion from some of these individuals can be complex, Cardinal George said that when any person presents him or herself to receive the Eucharist, they "take their salvation into their own hands." For a more complete explanation of this matter, reference the cardinal's column "Catholic participation in political life, revisited" (CNW, Oct. 10, 2004) online at www.catholicnewworld.com/cnw/issue/2004/cardinal_101004.html.



Clarification

In the last edition of the "Catholic New World," an article on the annual conference hosted by the Respect Life Office did not fully reflect the context of Cardinal George's remarks.

During an informal question-and-answer session with the archdiocese's parish Respect Life Coordinators, the cardinal emphasized that the participation by any person in the promotion of abortion, including through his or her political position, is a grave matter. While the issue of withholding Communion from some of these individuals can be complex, Cardinal George said that when any person presents him or herself to receive the Eucharist, they "take their salvation into their own hands." For a more complete explanation of this matter, reference the cardinal's column "Catholic participation in political life, revisited" (CNW, Oct. 10, 2004) online at www.catholicnewworld.com/cnw/issue/2004/cardinal_101004.html.

Image processing - Rotation / Deskewing

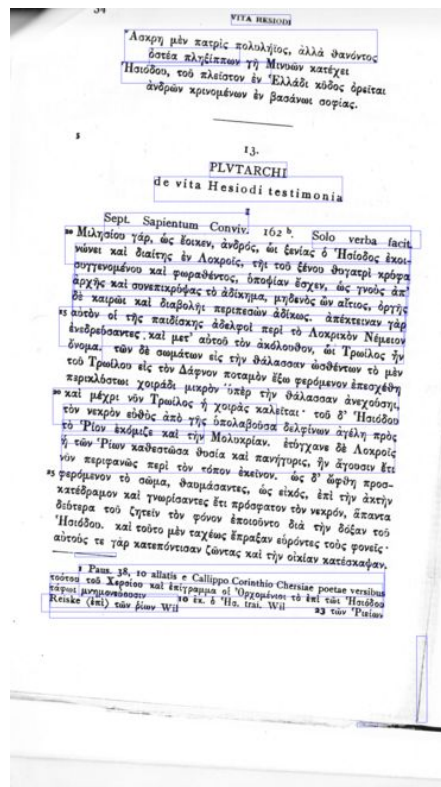
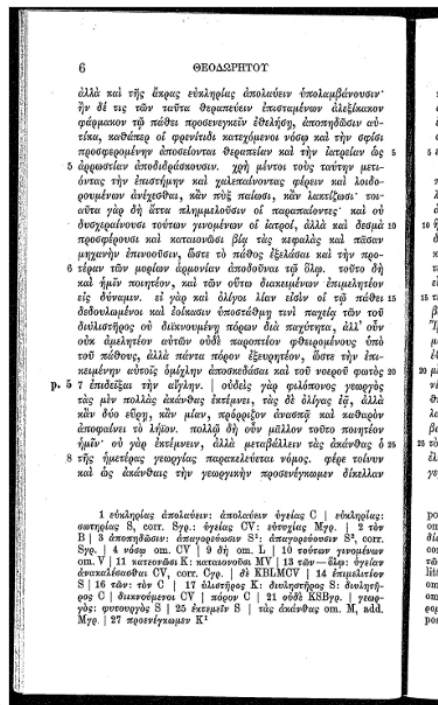


Image processing – Border removal



6

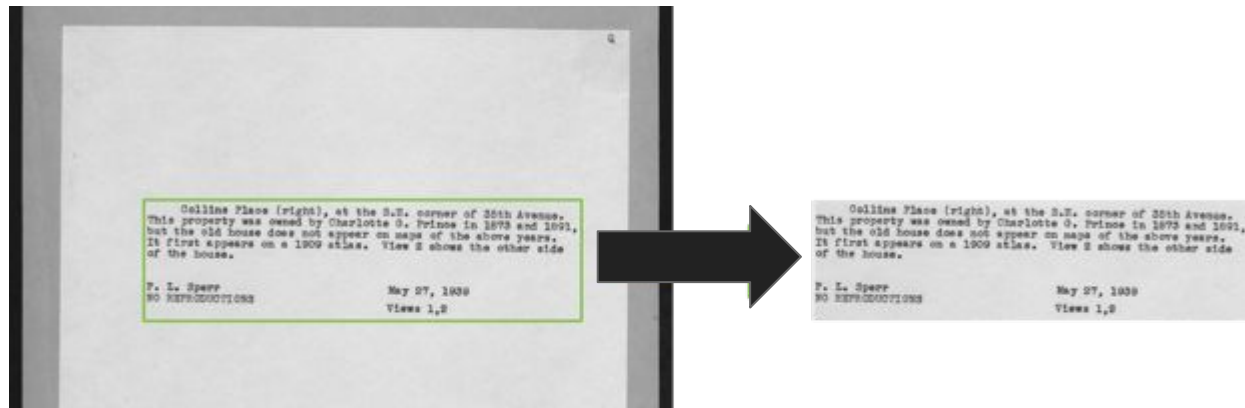
ΘΕΟΔΩΡΗΤΟΥ

ἀλλὰ καὶ τῆς ἑκφρασεως ἐκφρασεως ἐκφρασεως
 ἦν δὲ τις τῶν τούτων θεοποιῶν ἐκφρασεως ἐκφρασεως
 φράσεων τῶν πᾶσι προσεγγιζῶν ἰδιότης, ἀποπερδῶν αὐ-
 τῶν, καθάπερ οἱ φρενίτιδι κατεγόμενοι νόσῳ καὶ τὴν σφί-
 προσεγγιζομένην ἀποσιώπῃσι θεοποιῶν καὶ τὴν ἰατρῶν ὡς
 ἀρρωστῶν ἀποδιδράσκουσιν. γὰρ μὲντοι τοὺς ταύτην μι-
 ὄντας τὴν ἐκφρασὴν καὶ χαλεπαίνοντας φέρειν καὶ λαιδο-
 ρομένους ἀνέχεσθαι, κἂν πῶς παύσαι, κἂν κατίζωαι· τοι-
 αῦτα γὰρ δὲ ἔστι πλημμελοῦναι οἱ παρασιώπῃτες· καὶ οὐ
 δοχευοῦνται καὶ κατανοοῦναι βίη τῶν κεφαλῶν καὶ πᾶσαν
 μηχανὴν ἐκπορεύουσαν, ὥστε τὸ πάθος ἐξέλκεται καὶ τὴν πο-
 τῆν τῶν μερῶν ἀγωνίαν ἀποδίδουσι τῶν ἰατρῶν. τοῦτο δὲ
 καὶ ἡμῖν ποιητέον, καὶ τῶν οὕτω διακειμένων ἐπιμελητέον
 εἰς δύνανται. εἰ γὰρ καὶ ὁλόγῳ λίαν εἶναι οἱ τῶν πᾶσι
 δεδουλωμένοι καὶ ὁλοκῶν ἐπιστάθμη τινὶ παχέῃ τῶν τοῦ
 διώλεσθαι οὐ δύνανται πᾶσαν διὰ παχέτηρα, ἀλλ' οὐκ
 οὐκ ἀμειλίχῃ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ
 τῶν πᾶσι, ἀλλὰ πάντα πᾶσαν ἐκφρασεως, ὥστε τὴν ἐκ-
 φρασὴν ἐντοὺς ἑκφρασεως ἀποσιώπῃσι καὶ τοῦ νοσήτος φωνῆς
 τῶν πᾶσι τῶν πᾶσι. ὁδοὺς γὰρ φιλανθρωπίας γενομένης
 τῶν μὴ πολλὰς ἀνάσθαι ἰατρῶν, τῶν δὲ ὁλόγῳ ἔξ, ἀλλὰ
 κἂν δύο εὐχῇ, κἂν μίαν, πᾶσαν ἀνάσθαι καὶ καθαρὴν
 ἀποκαταίει τὸ λῆψιν. πολλὰ δὲ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ οὐκ
 ἡμῖν· οὐ γὰρ ἐκφρασεως, ἀλλὰ μεταβάλλει τῶν ἀνάσθαι ὡς
 τῆς ἡμετέρας γενομένης παρακαλεῖται νόμος. φέρει τούτων
 καὶ ὡς ἀνάσθαι τὴν γενομένην προσεγγιζομένην διακρίαν

1 ἐκφρασεως ἐκφρασεως: ἐκφρασεως ἐκφρασεως C | ἐκφρασεως:
 ἀποπερδῶν S, corr. Srg. | ἐκφρασεως: ἐκφρασεως Mrg. | 2 τῶν
 B | 3 ἀποπερδῶν: ἀποπερδῶν S²: ἀποπερδῶν S¹, corr.
 Srg. | 4 νόσῳ om. CV | 9 δὲ om. L | 10 τούτων γενομένων
 om. V | 11 κατενοοῦναι K: κατενοοῦναι MV | 13 τῶν—ἐκφρ:
 ἀποκατενοοῦναι CV, corr. Srg. | 14 ἐκφρασεως S: ἀποκατε-
 νοοῦναι CV | 15 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 16 τῶν: τῶν C | 17 ἀποκατενοοῦναι K: ἀποκατε-
 νοοῦναι CV | 18 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 19 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 20 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 21 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 22 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 23 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 24 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 25 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 26 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV | 27 ἀποκατενοοῦναι K: ἀποκατενοοῦναι S: ἀποκατε-
 νοοῦναι CV

Tips and tricks

Tips and tricks – Crop the text



Try to crop the text part from the image

Tips and tricks – Tesseract configuration

You have a parameter `-psm`(page segmentation method) that can take the next values:

- 0 Orientation and script detection (OSD) only.
- 1 Automatic page segmentation with OSD.
- 2 Automatic page segmentation, but no OSD, or OCR.
- 3 Fully automatic page segmentation, but no OSD. (Default)
- 4 Assume a single column of text of variable sizes.
- 5 Assume a single uniform block of vertically aligned text.
- 6 Assume a single uniform block of text.
- 7 Treat the image as a single text line.
- 8 Treat the image as a single word.
- 9 Treat the image as a single word in a circle.
- 10 Treat the image as a single character.

Demo Time





Conclusions / what are you going to do with it?

Resources

- ▶ <https://github.com/tesseract-ocr/tesseract/wiki/ImproveQuality>
- ▶ https://github.com/danvk/oldnyc/blob/master/ocr/tess/crop_morphology.py
- ▶ <http://static.googleusercontent.com/media/research.google.com/ro/pubs/archive/33418.pdf>

You can find the code from demos on:

- ▶ https://github.com/simutoni/ocr_examples

THANKS!

Any questions?

You can find me at simutoni@gmail.com