

Sem7-Project.docx

Date: 2017-09-30 12:07 UTC

* All sources 40 | Internet sources 9

✓ [35]	https://misteroleg.wordpress.com/2012/12...pre-processing-task/	3.9% 12 matches
✓ [37]	https://stackoverflow.com/questions/37378052/how-can-i-ocr-a-file-in-farsi	3.8% 7 matches
✓ [39]	https://www.quora.com/Where-can-I-get-th...e-with-documentation	2.8% 7 matches
✓ [40]	https://www.scribd.com/document/258117998/REport	2.8% 4 matches
✓ [41]	https://medium.com/@hdinhofer/optical-ch...-vision-76887e1d6ab0	2.4% 3 matches
✓ [42]	https://www.quora.com/How-does-the-Tesseract-API-for-OCR-work	2.4% 2 matches
✓ [43]	https://www.coursehero.com/file/pgpmjj/N...Tessesract-We-first/	2.1% 4 matches
✓ [45]	enlighten-ing.com/wp/	0.4% 1 matches
✓ [46]	ufology.wikia.com/wiki/Charles_Brown	0.3% 1 matches

17 pages, 1803 words

⚠ A very light text-color was detected that might conceal letters used to merge words.

PlagLevel: selected / overall

133 matches from 47 sources, of which 45 are online sources.

Settings

Data policy: *Compare with web sources, Check against my documents*

Sensitivity: *Medium*

Bibliography: *Consider text*

Citation detection: *Reduce PlagLevel*

Whitelist: --

Introduction

What is a Personal Info Assistant (PIA)?

Following the idea of a “Personal Digital Assistant” (PDA), a PIA can be defined as a hardware which runs applications that provide quick reference to lists and processed data through proper links.

Why PIA?

We live in the age of information.

In the entire length of time between waking up to the sound of an alarm set on our branded smartphones and setting the same alarm before going to bed at night, we encounter a wide variety of tasks every day.

Common among these activities is the fact that each of these activities expects us to be informed. Using a washing machine needs us to know how to operate the buttons. Using an air-conditioner needs us to know what buttons to press on the remote in order to get the right setting.

Well, these are simple, aren't they?

Yes, because manufacturers make their products easy to use by hiding their inner features.

What's your response when your washing machine wouldn't run no matter how many buttons you press or your AC won't cool at the right temperature?

We all get irritated, don't we? The best response we have is to call customer care or a repairman.

Right at that moment something crosses my mind. Given our dependency on these, how little do we know about the things around us!

How much do we know about our 32” LCD Plasma TV; about the 5-speed automatic washing machine; about the 4-star rated AC; about the RO Water Purifier – all of which we operate on a daily basis at our homes?

“I'm not a techie guy”, you may say, and you're correct. It's not necessary to know everything about them, but a working knowledge wouldn't hurt.

This issue is serious for students & professionals.

Not getting the desired output on the CRO no matter how well you have connected the circuits? Does it bother you why it happened?

Been there. Done that.

^[43] **The more we go out, the more we learn how little we know.** But that is no excuse to have no idea about not having any idea about things we use on a daily basis.

The only problem is – How would it be done?

No one has the time or desire to leaf through the pages of a user guide; or to search for authentic documentation on the device.

Well, what if I told you that you do not need to do any of these; the information will reveal itself to you!

What would your response be?

Well, this project may do just that. ^[37]▶

Summary of Tesseract OCR Engine

Tesseract was originally developed at Hewlett-Packard Laboratories Bristol and at Hewlett-Packard Co, Greeley Colorado between 1985 and 1994, with some more changes made in 1996 to port to Windows, and some C++-izing in 1998. In 2005 Tesseract was open sourced by HP. Since 2006 it is developed by Google.

- The 'tesseract-ocr' package contains an OCR engine - libtesseract and a command line program - tesseract.
- The lead developer is Ray Smith. The maintainer is Zdenko Podobny.
- Tesseract has unicode (UTF-8) support, and can recognize more than 100 languages "out of the box".
- Tesseract supports various output formats: plain-text, hocr(html), pdf, tsv, invisible-text-only pdf.
- You should note that in many cases, in order to get better OCR results, you'll need to improve the quality of the image you are giving Tesseract.
- This project does not include a GUI application. If you need one, please see the 3rdParty wiki page.
- Tesseract can be trained to recognize other languages. See Tesseract Training for more information.

Improving the quality of the output:

There are a variety of reasons you might not get good quality output from Tesseract. It's important to note that unless you're using a very unusual font or a new language retraining Tesseract is unlikely to help.

Image processing

- Rescaling
- Binariesation
- Noise Removal
- Rotation / Deskewing
- Border Removal
- Tools / Libraries

Page segmentation method

Dictionaries, word lists, and patterns

Image processing

Tesseract does various image processing operations internally (using the Leptonica library) before doing the actual OCR. It generally does a very good job of this, but there will inevitably be cases where it isn't good enough, which can result in a significant reduction in accuracy.

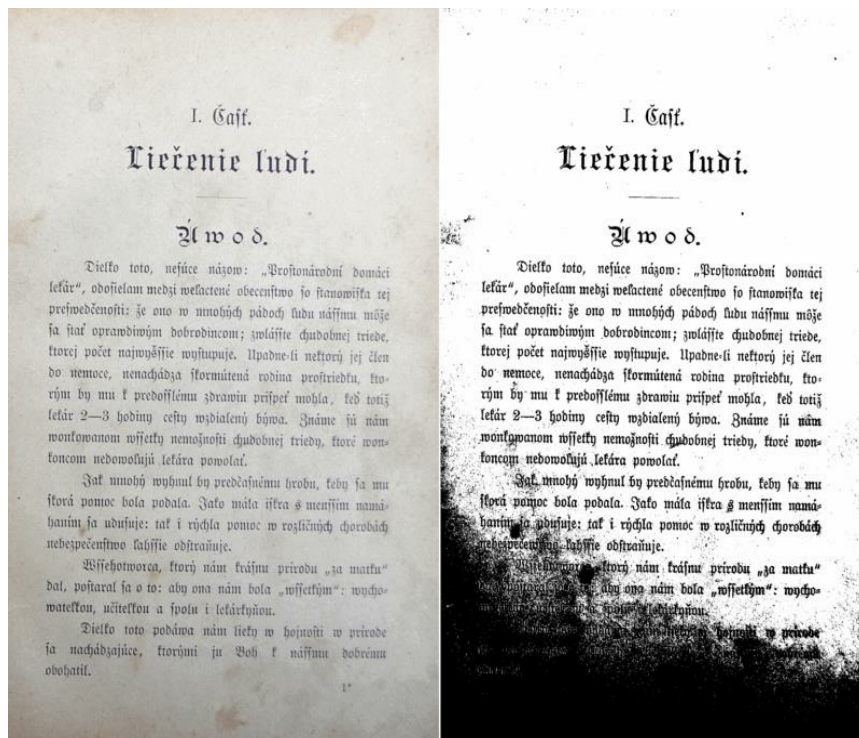
You can see how Tesseract has processed the image by using the configuration variable `tessedit_write_images` to `true` when running Tesseract. If the resulting `tessinput.tif` file looks problematic, try some of these image processing operations before passing the image to Tesseract.

Rescaling

Tesseract works best on images which have a DPI of at least 300 dpi, so it may be beneficial to resize images.

Binarisation

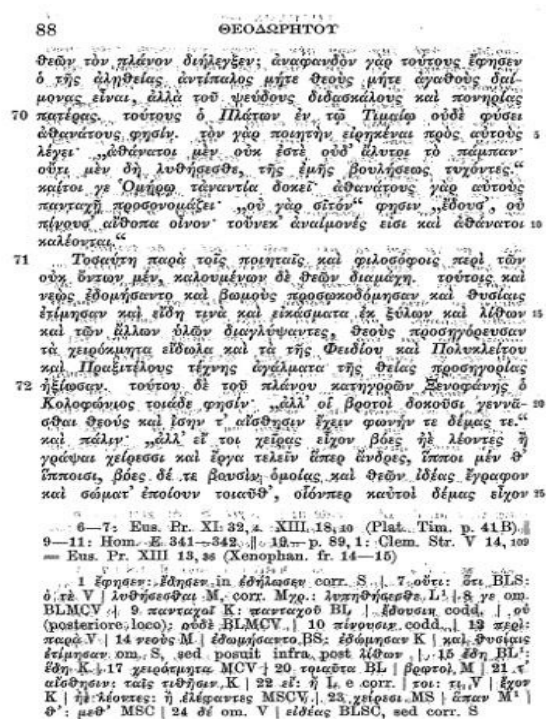
This is converting an image to black and white.^[40] Tesseract does this internally, but the result can be suboptimal, particularly if the page background is of uneven darkness.



[40]

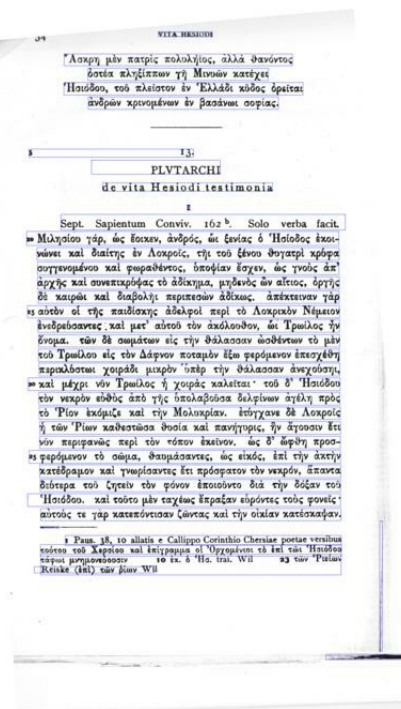
Noise Removal

Noise is random variation of brightness or colour in an image, that can make the text of the image more difficult to read.^[40] Certain types of noise cannot be removed by Tesseract in the binarisation step, which can cause accuracy rates to drop.

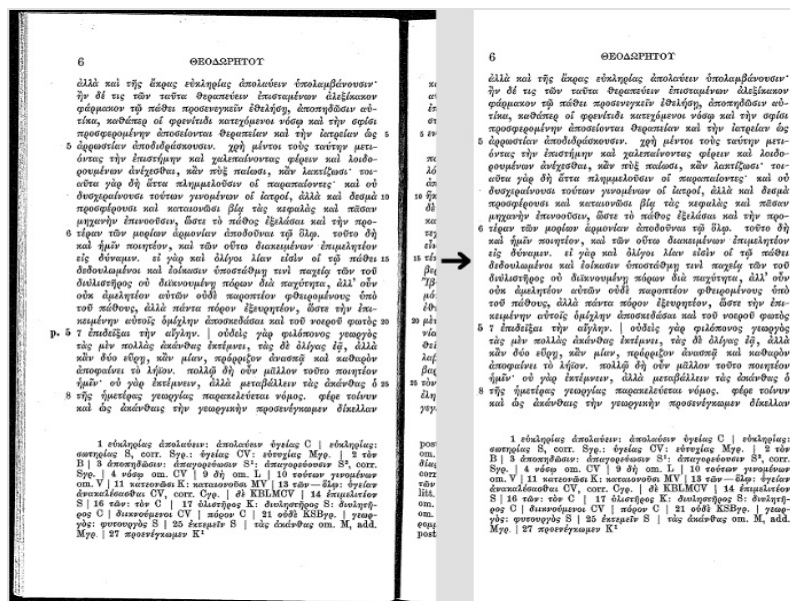


Rotation / Deskewing

A skewed image is when an page has been scanned when not straight. The quality of Tesseract's line segmentation reduces significantly if a page is too skewed, which severely impacts the quality of the OCR. To address this rotating the page image so that the text lines are horizontal.



Scanned pages often have dark borders around them. These can be erroneously picked up as extra characters, especially if they vary in shape and gradation.



Tools / Libraries

[Leptonica](#)
[OpenCV](#)
[Scan Tailor](#)
[ImageMagick](#)
[unpaper](#)
[ImageJ](#)
[Gimp](#)

Page segmentation method

By default Tesseract expects a page of text when it segments an image. If you're just seeking to OCR a small region try a different segmentation mode, using the `-psm` argument. Note that adding a white border to text which is too tightly cropped may also help, see issue 398.

To see a complete list of supported page segmentation modes, use `tesseract -h`. Here's the list as of 3.21:

- ^[35] ▶ 0 Orientation and script detection (OSD) only. ^[35] ▶
- 1 Automatic page segmentation with OSD. ^[35] ▶
- 2 Automatic page segmentation, but no OSD, or OCR. ^[35] ▶
- 3 Fully automatic page segmentation, but no OSD. (Default) ^[35] ▶
- 4 Assume a single column of text of variable sizes. ^[35] ▶
- 5 Assume a single uniform block of vertically aligned text. ^[35] ▶
- 6 Assume a single uniform block of text. ^[35] ▶
- 7 Treat the image as a single text line. ^[35] ▶
- 8 Treat the image as a single word. ^[35] ▶
- 9 Treat the image as a single word in a circle. ^[35] ▶
- 10 Treat the image as a single character.
- 11 Sparse text. Find as much text as possible in no particular order.
- 12 Sparse text with OSD.
- 13 Raw line. Treat the image as a single text line, bypassing hacks that are Tesseract-specific.

Dictionaries, word lists, and patterns

By default Tesseract is optimized to recognize sentences of words. If you're trying to recognize something else, like receipts, price lists, or codes, there are a few things you can do to improve the accuracy of your results, as well as double-checking that the appropriate segmentation method is selected.

Disabling the dictionaries Tesseract uses should increase recognition if most of your text isn't dictionary words. They can be disabled by setting the both of the configuration variables `load_system_dawg` and `load_freq_dawg` to `false`.

It is also possible to add words to the word list Tesseract uses to help recognition, or to add common character patterns, which can further help to improve accuracy if you have a good idea of the sort of input you expect. This is explained in more detail in the Tesseract manual.

If you know you will only encounter a subset of the characters available in the language, such as only digits, you can use the `tessedit_char_whitelist` configuration variable. See the FAQ for an example.

Implementation

1. Block Diagram
2. Progress Remarks
3. Softwares
4. Source Code Listing
 - pytesseract.py
 - ocr_pia_1.py
 - ocr_pia_2.py
 - thresholding_adaptive_naive.py
 - thresholding_global_naive.py
 - thresholding_global_naive2.py
 - thresholding_global_otsu.py
 - thresholding_global_otsu2.py

Progress Remarks

- OCR Pre-processing

The first step towards a mobile-platform implementation for our project was to identify the key areas to be resolved to improve text recognition.

We identified these areas to be:

1. Choice of lossless compression image format (eg. PNG)
2. Adaptive Binarization of image
3. Fragmentation
4. Filter Application

- Source Code

The source-code will be made available (alongwith necessary datasets) on my Github page:

<http://www.github.com/CRT13/Projects/>

Results

Basic Dataset

Our basic dataset consisted of 5 images, captured using a Lenovo K6 Power smartphone. The testing system was a low-end laptop running Windows, with the following specs:

imple Thresholding using OpenCV

OpenCV offers 5 options for “Simple Thresholding”. These are listed below:

- **THRESH_BINARY**

$$\text{dst}(x, y) = \begin{cases} \text{maxval} & \text{if } \text{src}(x, y) > \text{thresh} \\ 0 & \text{otherwise} \end{cases}$$

- **THRESH_BINARY_INV**

$$\text{dst}(x, y) = \begin{cases} 0 & \text{if } \text{src}(x, y) > \text{thresh} \\ \text{maxval} & \text{otherwise} \end{cases}$$

- **THRESH_TRUNC**

$$\text{dst}(x, y) = \begin{cases} \text{threshold} & \text{if } \text{src}(x, y) > \text{thresh} \\ \text{src}(x, y) & \text{otherwise} \end{cases}$$

- **THRESH_TOZERO**

$$\text{dst}(x, y) = \begin{cases} \text{src}(x, y) & \text{if } \text{src}(x, y) > \text{thresh} \\ 0 & \text{otherwise} \end{cases}$$

- **THRESH_TOZERO_INV**

$$\text{dst}(x, y) = \begin{cases} 0 & \text{if } \text{src}(x, y) > \text{thresh} \\ \text{src}(x, y) & \text{otherwise} \end{cases}$$

The details can be referred from pg. 294 of “The OpenCV Reference Manual”.
So, we generated sample plots for each of these.

Original Image



BINARY



BINARY_INV



TRUNC



TOZERO



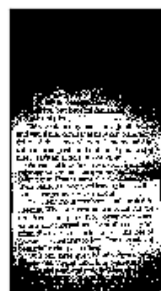
TOZERO_INV



Original Image



BINARY



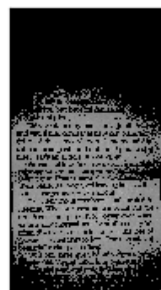
BINARY_INV



TRUNC



TOZERO



TOZERO_INV



Original Image



BINARY



BINARY_INV



TRUNC



TOZERO



TOZERO_INV



Original Image



BINARY



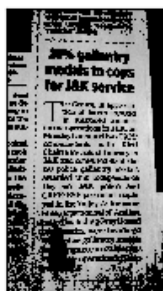
BINARY_INV



TRUNC



TOZERO



TOZERO_INV



Original Image



BINARY



BINARY_INV



TRUNC



TOZERO



TOZERO_INV



Next, we tried to observe the impact of “Simple Thresholding” on Tesseract's OCR output. We present the results for THRESH_TRUNC. First things first, there was an obvious reduction in file-size, but the runtimes were arbitrary, possibly due to inconsistencies caused by the lossy nature of JPEG images.


```
Python 3.6.1 Shell
File Edit Shell Debug Options Window Help
Python 3.6.1 (v3.6.1:69c0db5, Mar 21 2017, 18:41:36) [MSC v.1900 64 bit (AMD64)] on win32
Type "copyright", "credits" or "license()" for more information.
>>>
===== RESTART: C:\Users\CRT13\Desktop\ocrPIA\1.py =====
Image to be processed: 3a.jpg
"mm mm _ _ , " ' I t I ? 1' mt! 113'); iili;f4p;$4 g. 3' 8 3 Md ml (1.9 5 1
1119 Wnuv. in applmim king in 110111 of fowes 9113111 gal m the . , in susmled comma gecmb term
l Opel mans m I li on I . Ionda3 honoured two L RPF mked oommldants with Kim a mob S Chakm for a
cts of bravery in under J K and conferred 40 of the eIndi- 190 police gallanm' medals a the E anar
ded this Independence ham. 5 Day on J 5.1K police and Atro- CRPF1B8F personnel deplo- adm-
yed inthe allex At the same Din- S timga2pezsonnel of Andhra
>>>
===== RESTART: C:\Users\CRT13\Desktop\ocrPIA\1.py =====
Image to be processed: 3b.jpg
1",3331EMH 911111 1; 111M136 ml t3133-33 33313 3.3!; 113111 1.111111 11131111 111113.131 131119
111111111111 11 111111111111 ??Pa e 20% gallantry 11 medals to cons --- for
J K service . . Zt ~n . 1 111 131'1 ' 31111 0111111311 1111 311111111111 1 S gxy; V a
1 1 . 1.1.1 531 -1. 2101. 051011 ~ 111..11...1 p p E K 1 11 .311 1 11. ~11 11111111111 60
'1 111113- 390311 3 811.1111 1 L111 '21. 11115111 JMx 1 1 31111111111 1' 11111111111
two 1 RFF 01111111141 3 1111121311113111 x 1111 11111 . 5 111121 E L 11 1311'11'0 31115 11111
'1111'111'1111 Lmder 1, 11111 and 1 11115911111 40 0 11111 91nd 90 1131111151 g2. anu'v 1111
111 5 91' the 5 111131121911 1115 h depe 1111911119. bedu- D111 On J K police and. Am 5 C RPF
BSF personnel depk Edna- yed in the Vane At the same D110 3 maimrsoxmelo of Andhra ' police, mCl
dmgG1lex Hound xx ere honoured Q g .; Ei E g reaunbnagnaipnvcuutngmil ""i .."
```

```
Python 3.6.1 Shell
File Edit Shell Debug Options Window Help
Python 3.6.1 (v3.6.1:69c0db5, Mar 21 2017, 18:41:36) [MSC v.1900 64 bit (AMD64)] on win32
Type "copyright", "credits" or "license()" for more information.
>>>
===== RESTART: C:\Users\CRT13\Desktop\ocrPIA\1.py =====
Image to be processed: 4a.jpg
H S 3 I-Day Bareilly . I 9 I D I madrassas get govt warning Barely/Pilibhit The admi~ nistrat
ion will take legal ac- tion against madrassas that dont. organise singing of na- tional anthem and
d record pro- ceedings on Independence Day. Bareilly divisional com- missionerPV-"Jagmohan said.
We are Indians rst and our religion. caste or creed is secondary educational insti- tutes are cons
idered public places. We will enforce g0 vemm em. order in madrassas and if them is blatant viola
tion, we will take legal action, Jagmohan said. He added ac- v tion would bemkn only after going m
g evidence of violation MP 8 tvclan-tc;xaruanl""""""""""
>>>
===== RESTART: C:\Users\CRT13\Desktop\ocrPIA\1.py =====
Image to be processed: 4b.jpg
' l-Day Bareilly lts madrassas get govt warning Bareilly/Pibhit The admi- nistmtion will take l
egal 30 mm." against maclrassas that. dont cnganise singing Of na tirmal anthem and mxnd pm- (,irr
xiings (m I nderxzendence 1.7);3 Bmwilly divisional com- rwlissi ,;xl rx P " .Jamnohzm said.
We am Indizms rst. and our mligmn, caste or creed is , . qxw n lcl2lry.. mlucxnkmal insti U " lm
m arr ( rimsiderml public ' ' plums. W will enfnm g0~ - ; wrmmmm rnder m nwdrgzxssas i 8,1') .d
if ham is 1 )laturit V1012 , ticm. we w i H rake hgzg'z'sfl act ion Jagmohzm 5; id. I (73 mum.
31 t; ion wuu Id be? w kf 51 only after Wink! 11 1;me Widmm Of Jiriititmw. TWPS 20% gallantry
3 s ; "i v w wumwm mother iii? medals to caps
```

```
Python 3.6.1 Shell
File Edit Shell Debug Options Window Help
Python 3.6.1 (v3.6.1:69c0db5, Mar 21 2017, 18:41:36) [MSC v.1900 64 bit (AMD64)] on win32
Type "copyright", "credits" or "license()" for more information.
>>>
===== RESTART: C:\Users\CRT13\Desktop\ocrPIA\1.py =====
Image to be processed: 5a.jpg
(IIBIIIOIII Jam; 08 8W! JHOK th - -. 53 95mm IIIIM paAomaJ .Io Ino - pup Iaaq leonu 52 0 531215
u ., Ag! Ham 2 IO snags uImonquI aIIII momIM 1. 11 ~ , 2 IO am am mm cum aIdoad .Io suosea; san
BIIaJ mI ) UIJIP I.uo.p OIIM a Idoad uawo-M 102qu SIAaIUp pamumsp am Amp JeIIIIaIIM awII am m quup
A were 10 )IIIIUp mop oIIM aIdoad Io aAa al'I'I uowo "QM s IuIIp mau 35am, IouooIe- -ou J0 Iouoal
e-MOI .IaIIIIa 3J2 mu; s IuIIp JaIIIo uo'uIzIIIOOm am SIAaIIIIst pm;I SJaManOJO-Iw Aw quI Iaaq quoo
Ie ..%o 0,, I2 IIIIAA paAAoIIoI uaIIamaI-I 3123A IseI Jaaq aaII-IouooIe ue pauouneI .IasIaMpnI I
3801M 83808 V NO SNLLLBQ 38V SHBMVWHOUII 918 MOH ' II? 10 ill? 12 Iaaq .Iaaq quooIepu, sI IIOII
VZIII'IUAS LAN mu um mmmmm
>>>
===== RESTART: C:\Users\CRT13\Desktop\ocrPIA\1.py =====
Image to be processed: 5b.jpg
Vlllrnl HOW BIG LIQUOR'MAKERS ARE BETTING ON A SOBER FUTURE udweiser launched an aicohoI-free be
er last year; Heineken followed with a 0.0% alcohol beer this May. Microbrewers and distillers ar
e working on other drinks that are either Iow-alcoholi or no-aicohol. These new drinks will catch
the eye of people who dont drink or arent drinking at the time, whether they are designated driver
s, pregnant women, people who dont drink for religious reasons, or people who want the taste of a
weil-made beer without the intoxicating effects of a weil-made beer. But NV f iyi'ii.iiiizei?irsi
x i,j , is no-aicohoi beer beer at all, or an artificial concoction? ' y, . it starts off a
s normal beer and then the alcohol is distilled out, or removed with reverse osmosis (dontrty maki
ng it m . " with your home R0 filter though). , For mm VinePojinEater ,ugmy V " l, ' a v A
5% Value
```

References

Patents:

- [1] US6577762B1, Background Surface thresholding
<https://patents.google.com/patent/US6577762B1>
- [2] US7400768B1, Enhanced optical recognition of digitized images through selective bit-insertion.
<https://patents.google.com/patent/US7400768B1>
- [3] US9298980B1, Image preprocessing for character recognition.
<https://patents.google.com/patent/US9298980B1>
- [4] US20120063690A1, Object-Based Optical Character Recognition Pre-Processing Algorithm.
<https://patents.google.com/patent/US20120063690A1>
- [5] US7106905B2, Systems and methods for processing text-based electronic documents.
<https://patents.google.com/patent/US7106905B2>
- [6] US20130329023A1, Text recognition driven functionality.
<https://patents.google.com/patent/US20130329023A1>

Literature:

- [7] Eugene Borovikov, A survey of modern optical character recognition techniques
<https://arxiv.org/abs/1412.4183>
- [8] M Seeger, C Dance, Binarising camera images for OCR (ICDAR 2001, Proceedings of the 6th International Conference on Document Analysis and Recognition)
<http://ieeexplore.ieee.org/document/953754/>
- [9] Ranjith Unnikrishnan, Ray Smith, Combined Script and Page Orientation Estimation using the Tesseract OCR engine (ICDAR '07 Proceedings of the Ninth International Conference on Document Analysis and Recognition)
<https://dl.acm.org/citation.cfm?id=1304846>
- [10] Ray Smith, An Overview of the Tesseract OCR Engine (MOCR '09 Proceedings of the International Workshop on Multilingual OCR)
<https://dl.acm.org/citation.cfm?id=1577809>
- [11] Ray Smith, Daria Antoniva, Dar-Shyang Lee, Adapting the Tesseract Open Source OCR Engine for Multilingual OCR (MOCR '09 Proceedings of the International Workshop on Multilingual OCR)
<https://dl.acm.org/citation.cfm?id=1577804>
- [12] Zheng Zhang, CL Tan, Binarizing document image using coplanar prefilter (ICDAR 2001, Proceedings of the 6th International Conference on Document Analysis and Recognition)
<http://ieeexplore.ieee.org/document/953750/>
- [13] Zheng Zhang, CL Tan, Correcting document image warping based on regression of curved text lines (ICDAR 2003, Proceedings of the 9th International Conference on Document Analysis and Recognition)
<http://ieeexplore.ieee.org/document/1227732/>
- [14] Zheng Zhang, CL Tan, Recovery of distorted document images from bound volumes (ICDAR 2001, Proceedings of the 6th International Conference on Document Analysis and Recognition)
<http://ieeexplore.ieee.org/document/953826/>