

Text recognition on images from social media

Presenter: Belyaeva Oksana Vladimirovna

Authors: Akopyan M. S., Belyaeva O. V., Plechov T. P., Turdakov D. Y.

September, 2019

Optical character recognition (OCR) — translation of handwritten, printed text from images into an editable representation of characters on a computer (for example, in a text editor).

OCR on social networks is used for:

- analysis of emotional color;
- analysis of user requirements and wishes, etc.

Input data - images from social networks, in many cases characterized by:

- the presence of a complex background;
- low quality;
- distortion as a result of poor shooting conditions (taken from a smartphone);

OCR systems **do not work well** on such data.

- Development and implementation of a text information extraction tool;
 - ▶ Analysis of existing OCR systems;
 - ▶ Analysis and selection of image pre-processing methods to increase the accuracy of OCR systems;

OCR systems	Availability	Support	Advanced processing features
Finreader OCR	-	ABBYY	Camera OCR
OCROPUS	+	-	-
Tesseract OCR	+	Google	-

Demotivators



Certificates



Scanned

Государственный Думой Федерального Собрания Российской Федерации, избранной Государственной думой Российской Федерации и избранной Государственной думой Российской Федерации,

деловыми и государственными учреждениями (заведениями), учреждениями и организациями и от которых осуществляется Президентом Российской Федерации,

(в ред. Указа Президента РФ от 08.07.2013 № 613)

деловыми и иных организациях, отвечающих на основании федерального закона, указов Президента Российской Федерации,

(в ред. Указа Президента РФ от 08.07.2013 № 613)

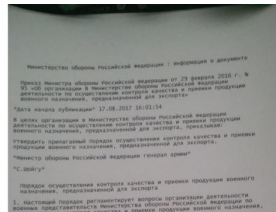
официальными деловыми и на основании приказа директора и организации, отвечающих на основании федерального закона, указов Президента Российской Федерации,

(в ред. Указа Президента РФ от 08.07.2013 № 613)

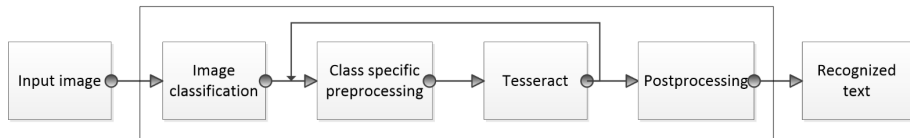
б) курьер (курьер) и несамостоятельно действующий, замещающий должность, указанные в инструкции к соответствующему органу.

2. Установить, что заместителями Председателя Правительства Российской Федерации - Членами Совета Правительства Российской Федерации на

Smartphone



Text Extraction Pipeline



- 1 Step: Image classification
- 2 Step: Preprocessing
- 3 Step: Tesseract OCR
- 4 Step: Delete characters with low confidence

1 Step - Image Classification

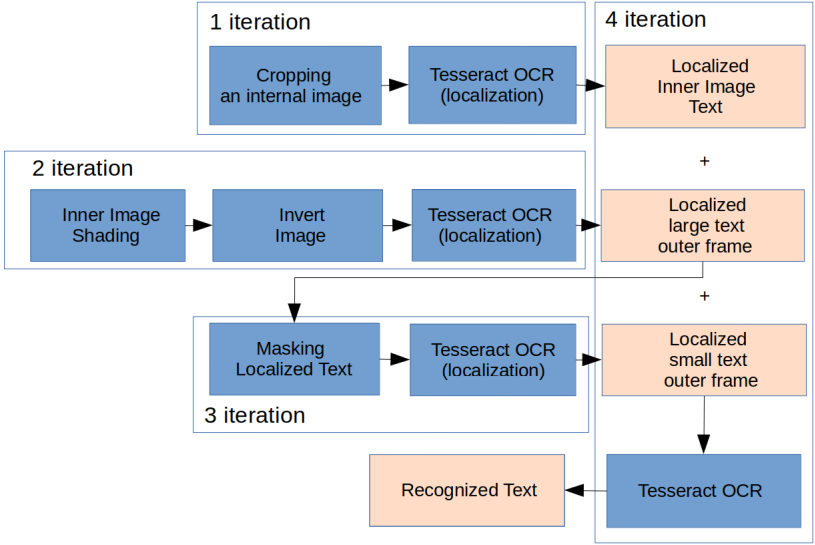
Power training and testing dataset: 11600 images.

Classifier	Accuracy
ResNet50	70
MobileNet	75
Gradient Boosting	95

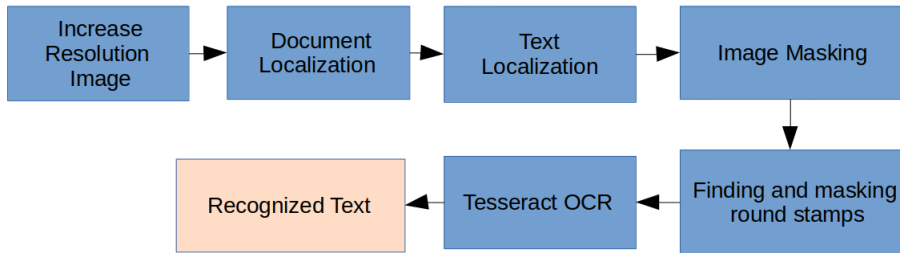
2 Step. Preprocessing methods

- Image Resolution Enhancement (IRE) - restoration of small-sized details (characters).
- Text Localization - the calculation of text blocks.
- Document Localization - calculating the position of a document in an image.

Demotivators class preprocessing



Certificates class preprocessing



Preprocessing classes Scanned, Smartphone

- Scanned images are of high quality and do not require pre-processing;
- The method of resolution enhancement is applied to images of the Smartphone class, and perspective distortion is searched and corrected.

Accuracy metric

$$\text{Character_accuracy} = \frac{\text{Number_of_characters} - \text{errors}}{\text{Number_of_characters}}$$

- Number of characters - number of characters in the document;
- Errors - minimum number of character edits.

#	Classification	Document localization	IRE	Text localization	Accuracy	GPU(CPU) seconds
1	-	-	-	-	29.5%	--(2)
2	manual	manual	NN2	NN3	87%	3.5(54)
3	manual	NN1	NN2	NN3	86%	5(58)
4	manual	NN1	-	NN3	78%	2.9(5.3)
5	manual	NN1	-	-	72%	2.7(4.6)
6	GradBoost	NN1	NN2	NN3	84.5%	5(58)

Thanks for attention