

₁ Einollah: Layout detection, a road into OCR's Utopia

₂ Vahid Rezanezhad¹, Clemens Neudecker¹, Mike Gerber¹, Kai Labusch¹, Robin
₃ Schaefer¹, Konstantin Barier¹ ¹

₄ ¹ Staatsbibliothek Berlin (SBB)

₅ December 15, 2020

6 Abstract

7 Many works have accomplished to segment layout of documents in order to do a better
8 Optical character recognition (OCR). Many features in documents can disturb OCR like
9 Drop capitals, curved lines, scan noises out of print space, skewed lines and so on. So, here
10 we have tried to do a perfect layout segmentation to improve OCR results with end-to-end
11 pipeline named EYNOLLAH. Beside layout detection which is implemented by pixel-wise
12 segmentation, we have also used some heuristic methods to classify marginals and to find
13 order of reading for text regions. This tool have used multiple models for variety of tasks
14 where those models are trained with Ground Truth we have in SBB. In this paper we will
15 explain how the models are trained and how the tool in general works.

16 **Keywords**— Pixel-wise segmentation, Artificial intelligence, Layout detection, Deep learning, Neu-
17 ral networks

18 1 Introduction

19 In the last years, many institutes have the idea of digitalisation in order to have data online and
20 to get the data fast. Digitalization also will helps us to provide a faster searching engine. To do
21 so it is required to transform scanned documents into a digitalized form, namely to read scanned
22 documents by machine. This can be done in two stages, first layout detection and then optical
23 character recognition.

24 In historic documents we need to recognize for example illustrations, text region and etc. Extracting
25 text regions in documents would enable us to do OCR with machine as well. Recently great im-
26 provements have been done in semantic segmentation for images. However this improvements were
27 not enough in the case of historic documents in order to do a clean and great layout detection. We
28 here tried to apply different kind of image classification and semantic segmentation to get a state
29 of art layout and textline detector. By the way for this purpose only artificial intelligence methods
30 were not enough and we had to use also some heuristic methods to boost layout detection.

31 As already mentioned in EYNOLLAH we are benefited from machine learning for semantic segmen-
32 tation. Recently many people have tried to use Convolutional neural networks (CNN) models to do
33 layout detection in order to feed optical character recognition (OCR) models. This is now possible
34 with image segmentation which uses Fully Convolutional models (FCN) , Unet , Resnet-Unet or
35 other models. All those models actually build from two parts, encoder and decoder part.

36 In SBB as Qurator team we first prepared a tool by using pixel-wise segmentation for layout and
37 textline detection and it was working quit good for our data resources in SBB. Our data in SBB
38 have good quality and scale and therefore we did not face big problems with layout detection. But
39 as soon as we received other documents from other source where the data quality were not good
40 or scale was so big we realized serious problems in the tool (?). Beside this our previous tool was
41 trained to detect only main regions like background, text regions, images and separators. And the
42 reading order of text regions were not considered. This pushed us to change our strategy and it was
43 this that only a pixel-wise segmentation could not be enough. We realized that for some elements we

44 need heuristic methods. One of main elements which is detected by heuristic method was marginals.
45 The other thing we add to our pipeline in EYNOLLAH was scale classifier. It is crucial to know
46 the scale of documents in order to do a clean layout detection. So, with all experience gained from
47 previous work and learning from other works like dhsegment tool (?), at the end of day we were
48 able to produce a tool which output state of art result in the case of layout detection and textline
49 detection.
50 Analyzing and testing thousands of documents give us this knowledge that elements of layout are
51 playing a big role in reading order of text regions. We discovered that separators and headers are
52 two main factor in reading order. So using those elements from layout detection we now had this
53 ability to find reading order too. So, EYNOLLAH alongside a layout, textline detection, enhancing,
54 deskewing and scaling can also say sth about reading order. However, in the case of wrong reading
55 order, other methods can be used for finding it. One of the methods we have thought about it and
56 we have "how to" knowledge in SBB to do it is to detect reading order semantically. It means after
57 OCR, using text mining we can arrange text regions based on their semantics.

58 2 Model

59 In this work for all pixel-wise segmentation models and even for enhancing we have used a resnet50-
60 unet model. In the case of scale classifier resnet50 encoder with a two dense connected layer at the
61 top is used. For encoder part (resnet50) we applied pretrained weights (weights url). The encoder
62 part, Resnet-50 does map input resolution into low resolution in return it increase features of input.
63 On other side encoder part is mapping from low resolution back into higher original resolution of
64 input image.
65 The encoder part, Resnet-50, takes advantage of pretrianed weights which are trained with image
66 classification task (imagenet).

67
68 Our used Resnet50-unet model has 38.15 M parameters where only parameters of decoder part

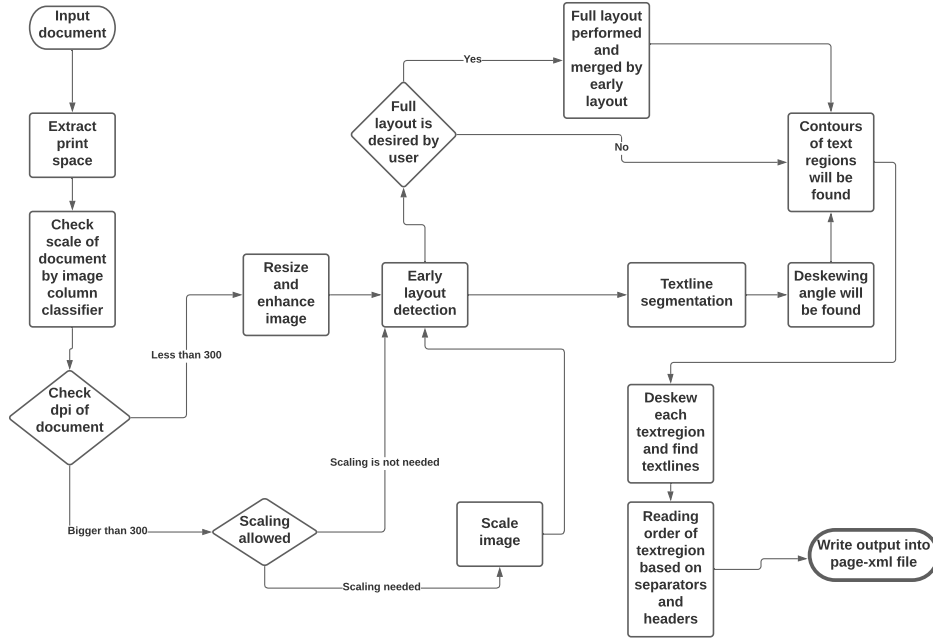


Figure 1: Algorithm flowchart

are trained fully (It means 14.71 M parameters are trained fully). The scale classifier model includes 25.6 M parameters where again only parameters of dense layers are fully trained (2.16 M parameters).

3 Algorithm stages

- Check dpi and scale of Document.
- if dpi is less than 300 and scale is lower than desired scale, document will be scaled and enhanced. During this project we found out in order to do a great layout detection we need to first scale the document to scale of our interest. When we are talking about a great layout detection it does not mean that only a higher IOU metric, however we mean that ability to isolate each textblock from others. We can get documents from different sources with different scales and quality, so scale document to the desired scale of your trained model (we

have trained our models with specific scale of documents) would lead to detect text regions separately and moreover this let the model to recognize headers much easier and better. Since, this is only the scale that differentiate headers from other types of text regions.

- extract print space in document.
- implementing first stage layout detection which its goal is to detect main regions (text regions, images, separators, background) by machine and marginals with heuristic method.
- if full layout is interest of user, this tool would go through document in order to detect headers and drop capitals too.
- textline detection will be done.
- for each text region , slope of deskewing will be calculated
- by applying heuristic method textlines in each text region will be found after deskewing.
- result will be written in a Page-xml data.

4 Scaling

As already mentioned in order to do an optimized layout detection we need to scale input into scale that we have trained our models. In our GT in SBB, our document had somehow same scale. One column documents had width between 1000 till 2000 pixels, for two columns document it was 2500 pixels in width and for three, four and five columns documents it was respectively 3000, 4000 and 5000 pixels. And finally for documents with six columns or more than it , the width was more than 6000 pixel. So, in order to get best performance (separated and isolated text regions) we need to feed model with mentioned scale. To do so we needed to first recognize number of columns in a document. For this purpose we trained an document classifier which returns the number of columns in document. We classified all documents in SBB into six classes which includes documents with

101 one, two, three, four, five columns and last class was documents with six columns and more. This
102 classification can include obviously more classes but since we have not enough newspapers with
103 more than six columns, therefore we decided to make newspapers with six columns and more that it
104 one class. This is worthy to mention again and highlight the role of scaling in a clean and sensitive
105 layout detection. So, it is recommended to boost GT for scaling purpose.

106 5 Enhancing

107 This tool is able to work with all document from any sources with verity of qualities. And one of the
108 key models which we trained in SBB was image enhancing model. If a document with low quality
109 and low scale is fed into our layout model it was so difficult to detect background (or better say the
110 space between text regions) precisely. So, after scaling we applied enhancing model to increase the
111 quality of input.

112 To train an enhancing model which is again image to image model, we needed documents with low
113 quality and GT would be corresponding image with higher quality. Since we did not have such
114 a GT in SBB and the quality of documents in SBB was good we decided with scaling down and
115 scaling up again to original size, provide low quality versions of our documents. With this method
116 we had already a great data to train our enhancing model.

117 6 Print space extraction

118 An early problem that can occur with layout detection and correspondingly affect OCR is out of
119 print space black offsets which are caused mostly by document scanning. Again DNN models help
120 us to extract print space by segmenting document by pixel-wise segmentation method. For print
121 space extraction it is necessary to feed model with full document or better to say that model should
122 not be trained in patches. Dhsegmnet project (?) has provided a big GT for print space extraction,
123 however we in SBB also annotated 800 documents more in order to boost our trained model. We



Figure 2: Print space extraction for few examples is shown

124 have shown few examples of print space extraction in ??.

125 7 Layout detection

126 7.1 Early layout detection: Only main regions

127 As already mentioned, since object detection does not meet the best performance for OCR purposes
 128 (this is checked in SBB and we have seen that object detection for text region detection does not
 129 work well for l-form text regions and in the case of documents with multi-columns), so we decided
 130 to first train models (pixel-wise segmentation) which can see the documents in patches and tried
 131 to detect text regions (or more importantly background) as separate and isolated possible. This

early and important and sensible model only detect main and mother regions in a document like background, text regions, separators and images. Thanks to sensibility of this early layout, we are able to detect marginals in a good manner with heuristic method. This is a key decision and step which help us to detect regions perfectly however this step cost longer time in the case of full layout detection.

7.2 Full layout detection

One of source of problems in a state of art OCR is of course problems occurred because of drop capitals in a document. It was a motivation for our group in SBB to care about drop capitals in layout detection too. Beside drop capitals we also have another important text type which can play key role in reading order detection. As it will be discussed in following sections, headers alongside separators are two main factor in reading order in a document. But this is not the only reason we are interested in headers. For search engines, headers are great source of information. In full layout detection we care at this time more about these two classes. After full layout is performed we merge it with early layout detection.

8 textline detection

Textline detection is a binary image segmentation which help us to detect where the textlines are located. For textline segmentation we had a great GT but the problem was all documents contain only one column. As it is imaginable, even with augmentation (scaling) we still had problem with documents with dense textlines. We could come over this problem with manually scaling the document and some parameterization in our prediction tool in order to prepare GT for bigger variety of documents including more columns in the document. Afterward we trained our model again and this time result was as it was expected. Thanks to augmentation, our model is capable of segmenting vertical lines too. However as it is

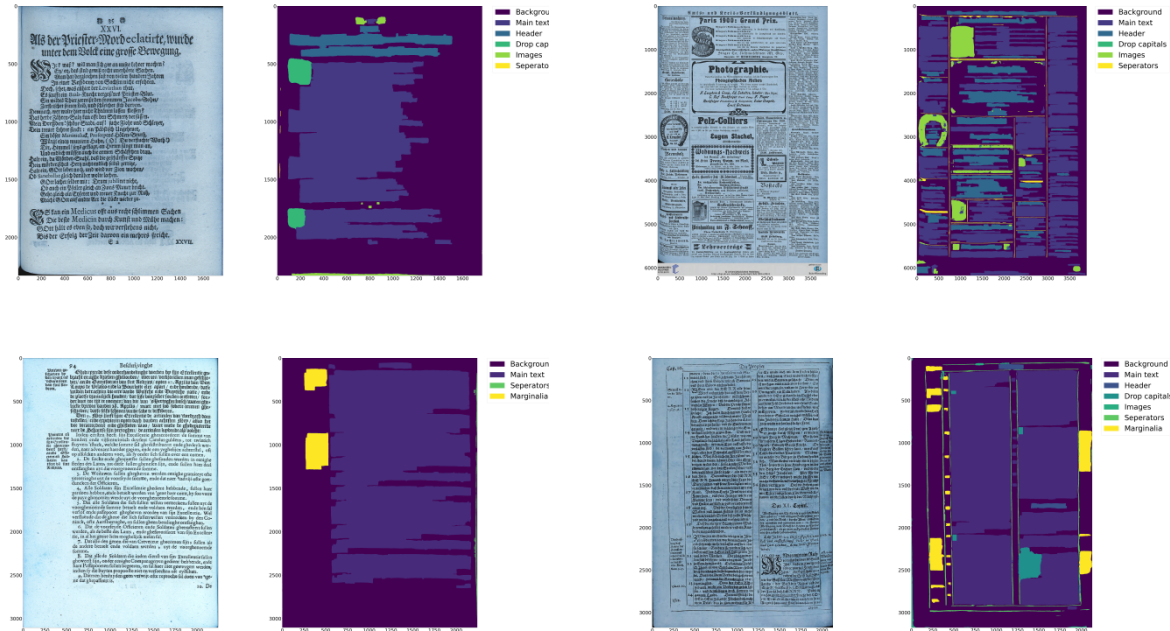


Figure 3: Full layout detection

clear in our results , textline segmentation includes artifacts for example for separators or images. But this artifacts does not affect our result so much since we care only of textlines in text regions. Textline segmentation result has a great role in deskewing and consequently detecting textlines bounding boxes. This amazing role will be explained in the next section.

9 deskewing

Skewed documents can strongly affect textline detection, in order to resolve this of course a robust deskewing is needed. We in SBB applied many different ways and algorithms to do deskewing. At the end we realized that textline segmentation result can help us in this way. The idea was that if we rotate a skewed document for a specific angle, deskewing angle, mean of variance of sum of textline segmentation in direction X should be maximum. This metric helped us to deskew documents properly. You can see this metric for some given documents in ??.



Figure 4: Pixel-wise textline segmentation

166 For historic documents we face those that are skewed locally and in a different manner from each
 167 other. So, "eynollah" does not only deskew the document but also apply it on each text region.
 168 This helped us to detect textline rectangle for documents with heterogeneous skewed regions.

169 10 Curved lines textline detection

170 11 reading order

171 see figures ????

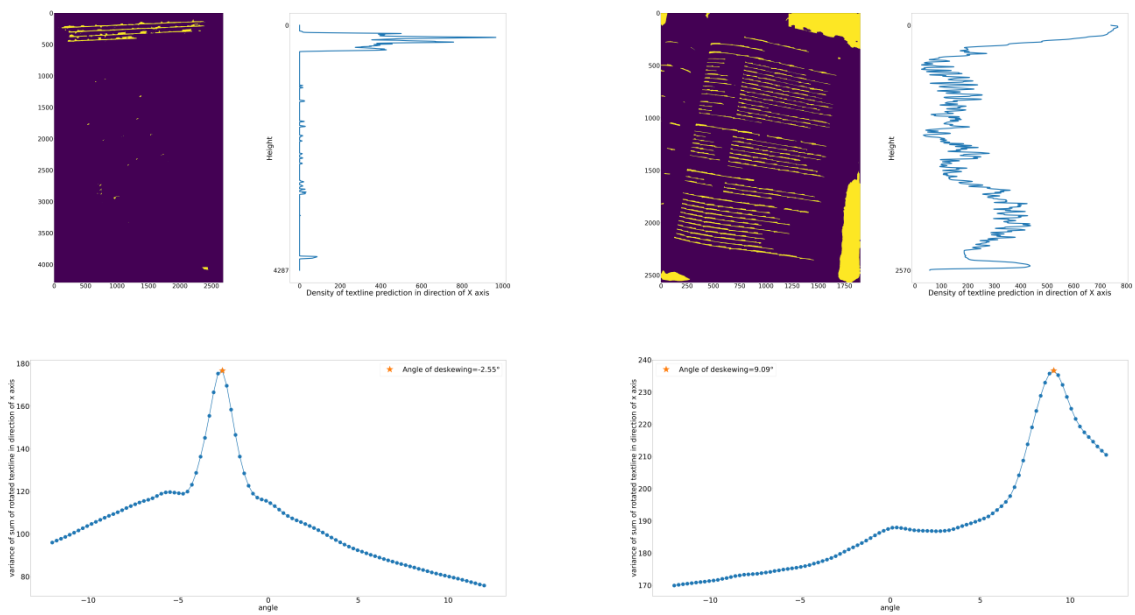


Figure 5: Deskewing procedure

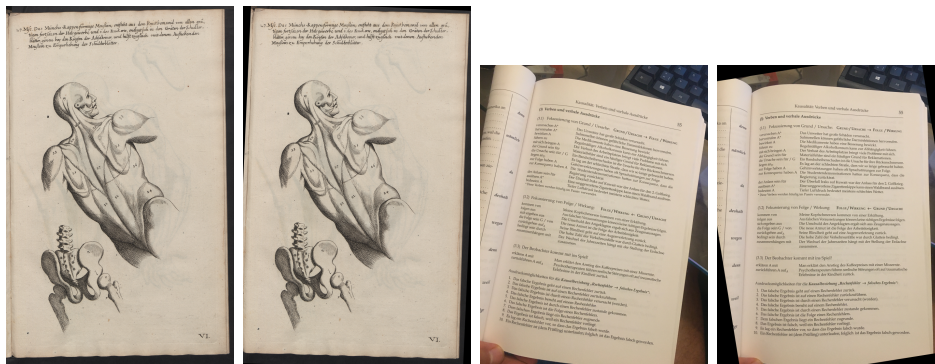


Figure 6: Deskewed documents



Figure 7: Enhancing. Left figures show org low quality images and on the right side enhanced images can be seen



Figure 8: curved lines detection



Figure 9: Drop capital integration in textlines



Figure 10: Headers role in reading order



Figure 11: Vertical textlines

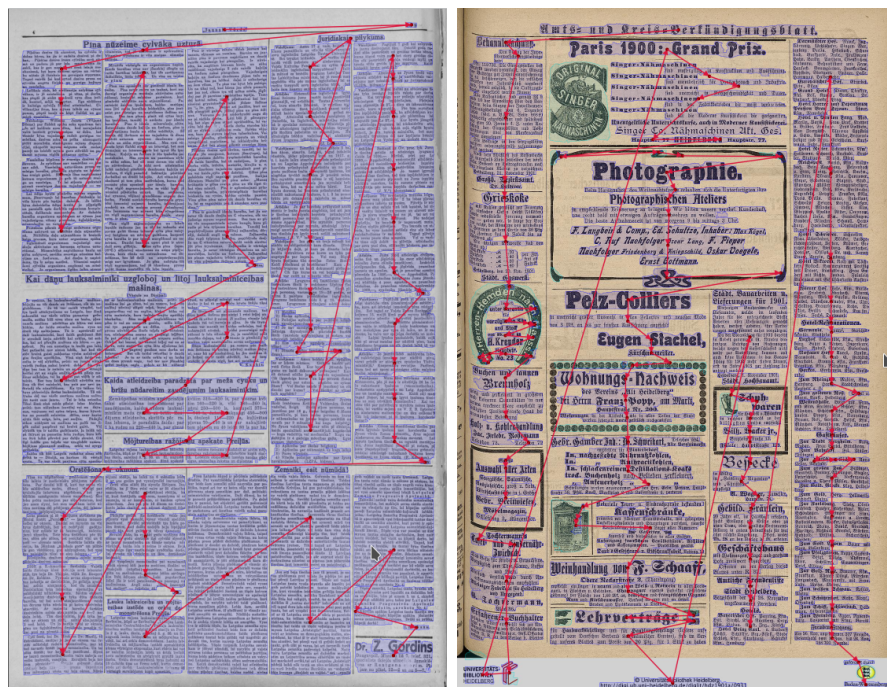


Figure 12: Reading Order