

Influence of quality of pixel level annotations on text detection performance in natural images

I. Dorkić*, M. Brisinello*, R.Grbic** and M. Herceg**

* TTTech Auto Central and Eastern Europe, Osijek, Croatia

** Faculty of Electrical Engineering, Computer Science and Information Technology Osijek, Osijek, Croatia
ivan.dorkic@tttech-auto.com, matteo.brisinello@tttech-auto.com, ratko.grbic@ferit.hr, marijan.herceg@ferit.hr

Abstract - Text detection in natural images is a task that arises in many computer vision applications. State-of-the-art text detection methods are mainly based on deep neural networks designed for instance segmentation task. However, most of the available datasets for text detection do not have fine annotations at the pixel level which are required during supervised learning of such networks. Usually, a whole or reduced text bounding box is used as a segmentation mask. In this paper, a method that generates a synthetic dataset with precise annotations at the pixel level is proposed. The method is based on the available Synthtext script for generating synthetic datasets with text instances. By creating synthetic datasets with precise and coarse annotations at the pixel level we explore the efficiency of the state-of-the-art text detector TextFuseNet.

Keywords - object detection; text detection; instance segmentation; Synthtext; TextFuseNet

I. INTRODUCTION

The text represents a set of words that conveys a certain meaning to the reader. In our daily life, we often encounter text while reading newspapers, emails, information and traffic signs, street names, and so on. Text detection and recognition are for an adult person relatively simple task. Reading a text with a computer is successfully solved within domains where the structured text appears like in digitalized books, official documents, license plates, and so on. Such text is usually on a legible background, with a standard font and line spacing and horizontally aligned with no occlusions. However, text can appear in so-called unstructured form - we find it in random places in natural images. Typical examples are scenes of the streets or shops where text appears in posters, advertisements, cafe names, products, and so on. Such text instances can have non-standardized font and spacing, can be in mixed colors, multi-oriented, occluded, curved, and with a colorful background. Detection of such text in natural images is still quite challenging for computer algorithms.

Nowadays, text detection in natural images is based on deep Artificial Neural Networks (ANNs). The most effective methods [1] are based on deep Convolutional Neural Networks (CNNs) that perform instance segmentation task. To successfully learn such a network in a supervised manner, an appropriately annotated dataset is required which means that each text instance in the image is annotated at a pixel level. However, most of the

available datasets provide only a bounding region of a text instance in form of a rectangle or polygon, and usually, every pixel within such bounding region is considered a part of the text by the applied method for text detection [2]. In the rest of this paper, we call such pixel-level annotations coarse. We believe that more precise text instance annotations on a pixel level of the training dataset can potentially boost text detector efficiency.

To explore whether more precise annotations at the pixel level of the training dataset can boost text detector efficiency, we propose an approach for obtaining precisely annotated synthetic datasets for text detection. The available code from [3,4] is used to generate synthetic images where the text appears in natural images. The code is further modified to produce precise annotations on a pixel level in form of a segmentation mask or polygon for each character and whole text instance. Such datasets are then used to train state-of-the-art text detector TextFuseNet [1] to quantify the importance of precise annotating of text instances in the training dataset.

The paper is structured as follows. In Section II, a problem of text detection and text annotation is presented, followed by an overview of existing datasets for text detection. After that, an overview of the most recent solutions for text detection is given. A process of creating datasets with precise and coarse annotations on a pixel level is given in Section III along with the process of text detector training. The obtained results on the created datasets and the accompanying discussion are given in Section IV. At the very end of the paper, a conclusion is given.

II. RELATED WORK

Most text detection methods, at least the most effective ones, are based on deep neural networks for instance segmentation task like the well-known Mask RCNN [5]. These networks are further modified to take into account text appearance in natural images and are trained in a supervised manner based on available text detection datasets. Most of the datasets do not provide the fine annotations of each character or the whole text instance at the pixel level (in form of a mask or a polygon). In theory, the most precise way of specifying a text mask would result in a mask consisting of pixels within the bounding box that exactly belong to the text in the image. An example of an image with several text instances is shown in Fig. 1 a) while the difference between the commonly

used coarse approach (Fig.1 b)) and the precise labeling approach (Fig.1 c)) is clearly visible. Most of the instance segmentation methods based on Mask R-CNN solve this problem by treating all pixels within the bounding box as a segmentation mask during training.

A. A brief overview of available datasets for text detection

The ICDAR 2013 dataset [2] consists of 229 training images and 233 testing images. All images are of different resolutions. It is a set in which all the images are photographed with a camera and the text on them is natural, i.e. it is not in focus. It is a standard reference dataset for horizontal text detection.

Synthtext in the wild [3] is a synthetically generated dataset in which words are placed in images of natural scenes. The text instances are placed in an image by considering the outline of the scene and the depth map of the scene. In this way, the text was added to sufficiently large areas. It consists of 800,000 images with approximately 8 million words artificially added to the images. The images are of different resolutions. Each text instance is annotated with word-level and character-level bounding boxes.

The COCO_Text dataset [6] contains 63,686 images with 145,859 text instances, where words are annotated by using polygons. It is the largest real-world dataset of images with text appearance in a natural form. It consists mostly of horizontal and multi-oriented text instances which are rotated by a certain angle. A small part of the images also contains curved text.

The Multi-Lingual scene Text (MLT-2019) dataset [7] contains 20,000 images. The images are of different resolutions. The text in the images is displayed in its natural form. Annotation of the dataset contains word level and character level text bounding boxes along with the corresponding transcription and language class. It contains text in 10 different languages.

The Incidental Scene Text (ICDAR 2015) dataset [8] contains 1,670 images and 17,548 annotated regions. This dataset contains images in which the text is not in the focus. Therefore, most of the text instances in the images are out of focus and blurry. Text instances are annotated with word-level bounding boxes in form of quadrilaterals.

B. A brief overview of text detection methods

In [9], a method called Mask TextSpotter is proposed, which is an end-to-end neural network for spotting text with arbitrary shapes. Mask TextSpotter uses a simple end-to-end learning procedure that can achieve the detection and recognition of text directly from two-dimensional space using semantic segmentation. This method is efficient in detecting instances of irregular shape text, such as curved text. It is evaluated on four English datasets and one multi-language dataset, achieved results are much better than other similar methods.

Most existing solutions for text detection in images are based on neural networks primarily intended for instance segmentation. In [1], a method called TextFuseNet is proposed, which is mainly based on the architecture of Mask R-CNN and TextSpotter, which uses richer fused features to detect text. TextFuseNet collects and fuses text features from different levels using a multi-path fusion architecture that can efficiently match and fuse different representations. TextFuseNet detects text at the word and character level. During training, TextFuseNet uses a more precise polygon segmentation mask. The text instance annotation is a "narrowed" bounding box compared to the approach of most other methods, where everything inside the bounding box is considered a segmentation mask. The detection success measured with F1 score on the ICDAR2013 dataset is 94.3%, and on the ICDAR2015 dataset it is 92.1%

It appears that one problem in building efficient object detectors is the lack of large data sets that have precisely annotated text at the pixel level. In [10], a method for annotating text on images at the level of pixels is proposed. The previously mentioned COCO_Text data set was used. From the images in that set, 1,000,000 text images per border of the rectangular bounding box are cropped. These images were then used to train a deep neural network that is used for text segmentation, that is, the background is separated from the text itself using this network. After the training phase, two thresholds were determined. If the probability is less than 0.3, that part of the image is considered the background, if the probability is greater than 0.7, the image contains text, and if the probability is between these two thresholds, the network is unsure of its decision. Using this method, 14,690 images were generated where text is precisely annotated at the pixel level.

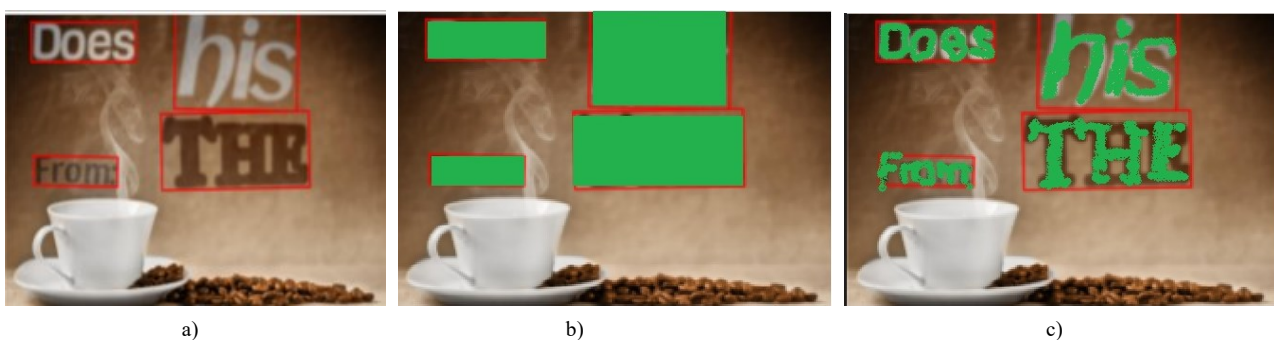


Figure 1. Examples of annotated images: a) input image with text bounding boxes, (b) coarse segmentation mask, (c) precise segmentation mask.

III. CREATING PRECISELY ANNOTATED SYNTHETIC DATASETS

This section describes the proposed approach for the precise annotating of text instances in synthetic natural images. For this purpose, we used an existing script available at [4] which was used for developing a famous Synthtext in the wild dataset. Based on created synthetic datasets we trained and tested TextFusenet text detector.

The code from [4] produces natural images where text instances are placed at suitable locations in the image. The script for each text instance provides coordinates of the text bounding box and information which image pixels correspond to the characters of the text. The example of generated text instance is shown in Fig. 2 where annotation at a pixel level provided by the script [4] is overlaid in red color. It can be noticed that generated annotation is quite imprecise. Certain pixels that belong to the background are annotated as text, so the actual shape of each character is not clearly visible from the generated mask.

To annotate the text instance more precisely at the pixel level we first extract the text instance from the image based on the provided text bounding box. Then k -means clustering algorithm [11] is applied to the extracted image to perform color segmentation. By doing so the segmentation of the image is performed based on the color similarity of the pixels. K -means clustering is performed in RGB color space. The resulting centers of k -means algorithm represent the dominant colors in the image. Since every character of each generated text instance is always in the same color, the assumption is that the pixels that belong to the text will be assigned to a single group and all other pixels (e.g. background) to other groups. In this paper, the optimal number of clusters was determined experimentally, and it is four. If less than four clusters are used, the text wouldn't always be segmented properly. Segmentation of the image shown in Fig. 2 in four clusters by k -means algorithm is shown in Fig. 3. In this case the first cluster represents pixels belonging to text (Fig. 3 a) while the remaining clusters represent the background and text border. Obviously, cluster representing text can be used as a more precise text mask than the annotation produced by the script [4] which is shown in Fig. 2. Furthermore, since script [4] provides rectangle bounding boxes on a character level also, we can easily extract precise mask for each character of a text instance from the obtained cluster representing text.

Since the result of the k -means algorithm depends on the initial initialization of the centers and the appearance of the image itself, it is not possible to know in advance in which group the pixels belonging to the text will be assigned. For this purpose, we developed CNN that performs binary classification whether the image contains text or not. In that way, we can determine which cluster actually contains pixels belonging to the text, i.e. required precise mask of the text. The CNN structure is shown in Table I. It is based on convolutional and max pooling layers. An input image is scaled to the size of 64x64 pixels. Dropout layer is used to prevent overfitting. To train such CNN, 812 different images containing text instances (like Fig. 4) were segmented with k equal to four



Figure 2. Example of generated text instance with overlaid annotation in red color produced by script [4].



Figure 3. Segmentation of the image shown in Fig. 2 using k -means algorithm in four clusters.

resulting in 3,248 images (similar to Fig. 3). The resulting images were manually labeled to class 1 or 0. It is necessary to do so because images containing text are considered positive examples, while images without text are considered negative examples. Several examples of such images are shown in Fig. 4. The images were divided into a training set, a validation set, and a test set. The training dataset contains 512 images with text and 1536 images without text. The validation set contains 50 images with text and 150 without text. The test set contains 250 images with text and 750 without text. The number of training epochs was 40 and the batch size was 32. During learning online augmentation was used such as rotation, horizontal flip, vertical flip, zoom, etc.

The evaluation of the built CNN is shown in Table II using the confusion matrix on the test dataset with 0.5 threshold applied to network output. The built CNN shows accuracy on the test of over 99%, with precision equal to 97.6% and recall equal to 98.4%.

The built CNN was implemented along with k -means algorithm in the script [4] to produce precise annotation at the pixel level for each text instance during the generation process (on-the-fly). Based on the resulting mask, a polygon or RLE record is determined for each text instance and saved to the JSON file containing annotations in COCO style [12]. Although the RLE notation is somewhat more precise, in this paper we used the polygon annotations for pixel-level segmentation since the default implementation of TextFuseNet requires segmentation annotation in form of polygons. Since the script from [4] provides a bounding box for each character of the text instance, we simply annotated each character at pixel level by extracting only part of the text mask that intersects with the character bounding box.

The code from [4] script uses a database of 8000 images but during a single run over 7000 images with

TABLE I. CNN STRUCTURE FOR BINARY CLASSIFICATION OF IMAGES CONTAINING TEXT OR NOT.

Layer	Input size	Output size
InputLayer	(64,64,3)	(64,64,3)
Conv2D	(64,64,3)	(64,64,32)
Conv2D	(64,64,32)	(62,62,32)
MaxPooling2D	(62,62,32)	(31,31,32)
Dropout	(31,31,32)	(31,31,32)
Conv2D	(31,31,32)	(31,31,64)
Conv2D	(31,31,64)	(29,29,64)
MaxPooling2D	(29,29,64)	(14,14,64)
Dropout	(14,14,64)	(14,14,64)
Flatten	(14,14,64)	(12544)
Dense	(12544)	(512)
Dropout	(512)	(512)
Dense	(512)	(2)

TABLE II. CONFUSION MATRIX OBTAINED ON TEST DATASET

		Predicted	
		Text	Not text
True	Text	244	6
	Not text	4	746



Figure 4. Examples of training images obtained by k -means clustering, a) - b) images containing text, c) - d) images without text.

random text instances are generated because some images do not always provide a convenient place to add text instance of a certain size and length. A typical example of the generated image is visible in Fig. 5 and a single text instance with precise and coarse annotations at the pixel level is shown in Fig. 6 a) and b).

In the end, six different training datasets were generated with varying numbers of training images (around 1000, 3000, and 5000) and with different types of annotations at the pixel level (precise or coarse). The generated datasets are summarized in Table III.

IV. RESULTS AND DISCUSSION

This section presents the results of TextFuseNet detectors training and testing on the generated datasets. A PC with a Linux operating system (Ubuntu 20.04. LTS), NVIDIA GeForce RTX 3080Ti graphics card, and Intel Core i9-11900F processor with 32 GB of RAM was used for training and testing. The training comprised of 60 epochs and the model with the lowest validation loss was selected for further evaluation.

All built detectors were evaluated on the same test data set with respect to text instances detection. COCO evaluation metrics are used for built detectors evaluation: average precision over 10 IoU thresholds from 0.5 to 0.95



Figure 5. Examples of the generated image with precise annotations using modified *Syntex* script.

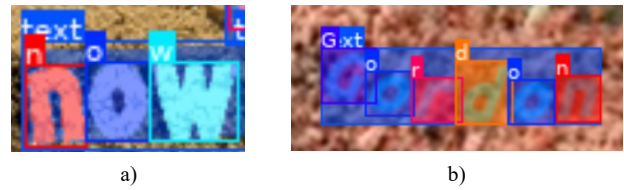


Figure 6. Examples of text instance with a) precise annotation, b) coarse annotation.

TABLE III. GENERATED SYNTHETIC DATASETS

Dataset name	Number of images in each dataset
Precise 1k	1029
Precise 3k	3265
Precise 5k	5099
Coarse 1k	1029
Coarse 3k	3265
Coarse 5k	5099
Validation	1038
Test	1032

with a step size of 0.05 - AP , average precision at IoU equal to 0.5 - AP^{50} , average precision at IoU equal to 0.75 - AP^{75} , average precision for small objects with an area less than 32^2 - AP^S , average precision for medium objects with an area between 32^2 and 96^2 - AP^M , and average precision for large objects with an area greater than 96^2 - AP^L .

The evaluation results of TextFuseNet detectors which were trained on precisely and coarsely annotated datasets with a varying number of training images are presented in Table IV and Table V. In both tables, it can be observed that metrics slightly increase with the increase of the number of training images which was expected. AP^{50} values are only slightly lower than AP^{75} indicating good localization properties of all detectors. All models have high values of average precision for medium and large objects but struggle with small ones. If we compare AP^{50} value of detectors trained on precisely and coarsely annotated datasets, then it can be noticed that the obtained values are only slightly in favor of more precise

