CURVE TEXT DETECTION WITH LOCAL SEGMENTATION NETWORK AND CURVE CONNECTION

Zhao Zhou, Shufan Wu, Shuchen Kong, Yingbin Zheng, Hao Ye, Luhui Chen, Jian Pu

Videt Tech.

ABSTRACT

Curve text or arbitrary shape text is very common in real-world scenarios. In this paper, we propose a novel framework with the *local segmentation network* (LSN) followed by the *curve connection* to detect text in horizontal, oriented and curved forms. The LSN is composed of two elements, *i.e.*, proposal generation to get the horizontal rectangle proposals with high overlap with text and text segmentation to find the arbitrary shape text region within proposals. The curve connection is then designed to connect the local mask to the detection results. We conduct experiments using the proposed framework on two real-world curve text detection datasets and demonstrate the effectiveness over previous approaches.

Index Terms— Scene text detection, curved text, convolutional neural networks.

1. INTRODUCTION

Text is probably the most important way for humans to communicate and express information. With the ubiquitous image capture devices, a huge amount of scene text images are produced, which brings a great demand for automatic text content analysis. The extracted text information also helps to understand the context of the whole images. As one of the most fundamental task for text analysis, scene text detection is to identify text regions of the given scene text images, which is also an important prerequisite for many multimedia tasks, such as image understanding and video analysis.

With the development of convolutional neural networks, there have been many attempts on text detection in natural scenes and great progress are achieved in recent years. The early attempts to detect text are with annotations of horizontal texts (*e.g.*, in [1]) and the approaches for arbitrary-oriented scene text detection are also proposed (*e.g.*, in [2, 3, 4]). However, in the real-world scenarios, there are still many text regions with irregular shapes, such as the curve words or logos, and a few of such sample images are shown in Fig. 1. It is still very challenging to detect these regions with different shapes.

In order to detect text in horizontal, oriented and curved forms, we propose in this paper a novel framework incorporating the local segmentation network (LSN) and the curve connection and the pipeline of text detection is illustrated in



Fig. 1. Example text images in natural scenes (Row 1) and text detection results by different methods: CTPN [1] (Row 2), EAST [2] (Row 3), and the proposed method (Row 4).

Fig. 2. The LSN is designed with two functionality, i.e., proposal generation and text segmentation. We chose the ResNet-50 [5] as the backbone of LSN. The horizontal rectangle proposals are generated based on multiple feature maps for the adaption to the different scales of the text. These detected proposals are small local regions but with high overlap with ground-truth text. Text segmentation is used to find the arbitrary shape text region within proposals by an ROI-Align [6] to produce the same size of the features from different anchors and a segmentation subnet to fine-tune the text area. The curve connection is then designed with center line generation and text polygon generation to connect the local mask to the detection results. We demonstrate a considerable improvement for the curve connection over those regressed directly from LSN. To evaluate our proposed framework, we report the evaluations on the recent proposed curve text detection datasets, i.e., CTW1500 [7] and Total-Text [8], and we compare with several recent approaches. Notably, we achieve an improvement over the state-of-the-art TextSnake [9] while only the training images in each benchmark are employed.

The contributions can be summarized as follows.

- We propose a novel neural network architecture combing horizontal rectangle proposals and arbitrary shape text segmentation to perform curve text detection. Our framework can generate multi-scale proposals where only a small number of curve annotations are needed.
- We also propose novel strategies for the curve connection of local text segments to improve the performance of long text words and arbitrary shape text detection.
- We apply our framework to two real-world text detection datasets and find that it achieves the state-of-the-art curve text detection performance.

The rest of this paper is organized as follows. Section 2 briefly reviews the literature of scene text detection. Our proposed framework is introduced in Section 3 and the details are described in Section 4 and 5. In Section 6, we demonstrate the quantitative study on the benchmarks. Finally, We conclude our work in Section 7.

2. RELATED WORK

Scene text detection has drawn growing attention from computer vision communities in recent years. With the astonishing development of object detection, the state-of-the-art frameworks such as Faster-RCNN [10] and SSD [11] have been widely applied to text detection field. However, compared with object detection that only needs general localization for objects, scene text detection requires precise positioning for characters. Therefore, many remarkable methods based on object detection for scene text detection have been proposed. These methods focus on more precise positioning for characters and can be roughly classified into two categories, *i.e.*, the anchor-based methods and the link-based methods.

Anchor-based Methods. The horizontal box in object detection was widely used in document text detection but the orientation of box is various when facing scene text detection. Ma et al. [4] proposed a multi-oriented scene text detection approach by generating six-orientations anchors at each point of the feature map. Quadrilateral anchor was introduced in [3] to detect text with tighter quadrangle. Zhou et al. [2] combined these representations together and proposed an efficient detector using a single fully convolution network with two branches. Liu et al. [7] applied 14 landmark points to represent curved text flexibly. All of these methods aforementioned focused on pursuing a tighter boundary of text but they were limited due to their templates that lack variability to cover text with extremely aspect ratio, such as long text or non-quadrilateral text.

Link-based Methods. Link-based methods are more robust when facing scene with long text or non-quadrilateral

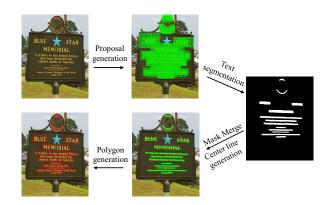


Fig. 2. Curved text detection pipeline of the proposed framework.

text. Tian *et al.* [1] introduced CTPN using LSTM [12] to link several text proposals. Shi *et al.* [13] proposed the SegLink framework that decomposes the text into segments and links and links text segments together. And Tian *et al.* [14] introduced a graph method called Min-Cost Flow to link single characters. However, these methods are based on a strong prior by restricting all the proposals to link should lie in a line. This hypothesis, in a manner, makes the problem easier and tractable but can not handle curved text.

Recently, several approaches have been proposed to detect text with form of arbitrary shapes and the curve text datasets were provided for research. [7] proposed a polygon-based curve text detector (CTD) which can directly detect curve text without empirical combination. The framework of TextSnake [9] considered a text instance as a sequence of ordered disks. To deal with the problem of separation of the close text instances, Li *et al.* [15] designed the PSENet by a progressive scale algorithm to gradually expands the predefined kernels. Different from previous methods, our local segmentation network combines the advantages of link-based methods and anchor-based methods using proposals to rough locate text and get accurate boundary by text segmentation. Our method also achieves state-of-the-art performance on the recent curve text detection datasets.

3. FRAMEWORK

The pipeline of the proposed framework is illustrated in Fig. 2. Each text can be represented by a bunch of text segments. For each segment, it has a coarse boundary represented by a square and a fine boundary represented by a mask. The former can fast classify into foreground/background proposals by a detection task, while the latter is used to get a tighter boundary of text. When faced with long text or curved text, segments will be linked together through their masks. After using the curve connection operations including the mask merging, center line generation, and polygon generation, a more smooth and robust boundary of text can be obtained.

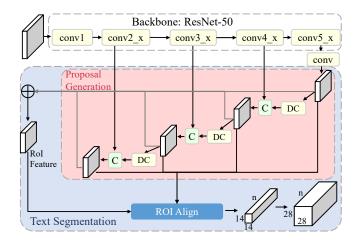


Fig. 3. The architecture of Local Segmentation Network. C and DC indicates the concatenation operation and the deconvolution layer.



Fig. 4. Sample image region with text segmentation after proposal generation.

4. LOCAL SEGMENTATION NETWORK

4.1. Network Structure

Our network is shown in Fig. 3. Here we choose ResNet50 [5] as the backbone to capture more detailed information. We remove the last fully-connected layer in ResNet50 and concatenate extra convolution layer to get deeper features with larger receptive fields. We feed the feature maps after each stage to the feature merging network. For each feature merging network, we apply a convolutional layer with 1×1 kernels to make two merge features have the same dimensions. In addition, we use 4 different stride features for the class square. These features can provide rich information and different receptive fields for robust detection. In order to get both fine and coarse feature, these 4 features are concatenated as the ROI feature. After the ROI feature, we use ROIAlign [6] to get same size feature with different size of the square. Finally, we upsample the features to predict text mask.



Fig. 5. Effect of segment merging strategy. Left: original images with ground truth; Middle: results by neighboring segment connection; Right: results by curve connection.

4.2. Proposal Generation

As is shown in the Fig. 4 top, a text region can be covered by several overlap square. We define each square through (x,y,l). Here x,y represents the center coordinate of the square. Actually, we consider the size of feature map extracted from network is w and h. The relation between the coordinate of square and the feature map can be formulated as the following equations:

$$x_i = i \times \text{stride}, i \in \{0, 1...w - 1\}$$

$$y_j = j \times \text{stride}, j \in \{0, 1...h - 1\}$$
(1)

 $l \in \{s \times k | s = 8, 16, 32, 64, k = 2, 2.5, 3, 3.5\}$ is the width of the square that can cover nearly all the scale of text. The network only produce 2 channels for text classification. Each predict feature point has 4 different anchor sizes.

4.3. Text Segmentation

In order to get a tight boundary of the text region, we predict a text mask for each positive square after the text classification stage. ROIAlign is used to get the same-size feature from various size of ROI features. Then the extracted features are concatenated and upsampled to get the final mask result. The result of the processing of text segmentation above is shown in Fig. 4 bottom.

4.4. Loss Function

To train the LSN, the loss function is formulated as:

$$L = \lambda_1 L_{\text{cls}} + \lambda_2 L_{\text{segment}}, \tag{2}$$

where $L_{\rm cls}$ and $L_{\rm segment}$ represent the classification loss and segmentation loss for text instances respectively. λ_1 and λ_2 balances the importance between $L_{\rm cls}$ and $L_{\rm segment}$.

For the classification term L_{cls} , the focal loss [16] is em-

Algorithm 1 Mask Merging for a given image.

```
Input: Predict mask set S; threshold s_1, s_2
Output: Merged mask queue Q
 1: Q = \{\}
 2: for each predict mask p \in S do
        if p_{i,j} > s_1 then
 3:
            p_{i,j} = 0
 4:
        else
 5:
            p_{i,j} = 1
 6:
        end if
 7:
        overlap list L = \{\}
 8:
        for each item mask m_i \in Q do
 9:
            merged mask m = m_i \cup p
10:
11:
            overlap ratio r = area_m / \min(area_m, area_p)
            if r > s_2 then
12:
                insert i into L
13:
14:
            end if
        end for
15:
16:
        m = p
        for each i \in L do
17:
            m = m \cup Q_{L[i]}
18:
            delete Q_{L[i]}
19:
        end for
20:
21:
        insert m to Q
22: end for
23: return Q
```

ployed,

$$L_{\text{cls}} = \sum_{i \in 1, 2, 3, 4} \text{FocalLoss}(p_{cls}(i), g_{cls}(i))$$

$$\text{FocalLoss}(p_t) = -\alpha_t (1 - p_t)^{\gamma} log(p_t)$$
(3)

where $p_{cls}(i)$ and $g_{cls}(i)$ is the i_{th} feature predicts and ground truth, and α_t and γ is hyper-parameters in focal loss.

For the segmentation term L_{segment} , we select $smooth_{L_1}$ loss to get the score of each pixel in square region and the definition is as follows:

$$L_{\text{segment}} = smooth_{L_1}(p_{segments}, g_{segments})$$
 (4)

$$smooth_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1\\ |x| - 0.5 & \text{otherwise} \end{cases}$$
 (5)

In our experiments, the hyper-parameters λ_1 and λ_2 are set to 1, the weights α_t and γ in focal loss are set 0.25 and 2 respectively.

5. CURVE CONNECTION

As the text segments are generated, a merging strategy for the segments is needed to have the overall detection results. A straightforward approach is to connect the neighboring segments. However, the boundary cannot be properly detected,

Algorithm 2 Text Polygon Generation.

```
Input: Merged mask set M for a given image
Output: Text polygon set P for the image
 1: P = \{\}
 2: for each proposal mask m \in M do
        Polygon p = []
 3:
        choose n positive pixels
 4:
        use principal curve to find curve center line
 5:
        choose 7 points p_i (i = 0...6) to represent center line
 6:
        for i \in [0, 6] do
 7:
           generate circumscribed rectangle of the area be-
 8:
    tween p_i and p_{i+1}
 9:
           regard rectangle points as polygon points
10:
           insert the left two points into p
11:
        end for
        insert the p into P
12:
13: end for
14: return P
```

resulting in a large precision loss regarding the performance (see Fig. 5). In this section, we design the curve connection with the components of mask merging and text polygon generation.

5.1. Mask Merging

For each predict mask, we set a threshold s_1 to define where is the fine region of text. When the pixel predict score is higher than s_1 , we set the pixel value 1 otherwise 0. In order to merge all boxes in the same text region, we design the mask merging approach and the pseudo-code is shown in Algorithm 1. First, we use a queue Q store the regions who are separate from each other. And a new region will compare with any of the regions on Q. If the overlapping mask of these two regions take the proportion of the smaller region mask is above threshold s_2 , these regions should be merged. After comparing with all of the regions on Q, we merge all of the regions which should be merged including the new region. The last two steps is repeated until all of the regions are operated.

5.2. Text Polygon Generation

A few text regions are obtained after the mask merging process. We chose n positive pixels in each text region and use the principal curve [17] to regress the curve center line. Seven points are chosen from the center line. For each pair of the points that are adjacent in center line, we use center point of two points as a rectangle center and generate a circumscribed rectangle of the area where are the text region between the pare center points. We will get a polygon of the text region by repeating these steps. Detail of the algorithm is shown in Algorithm 2.

6. EXPERIMENTS

6.1. Datasets

We evaluate our approach with two recent proposed curve text detection benchmarks, *i.e.*, the CTW1500 [7] and the Total-Text dataset [8] dataset.

CTW1500 [7] contains 1500 images (1000 train and 500 test). Each text instance annotation is a polygon with 14 vertexes to define the text region at the level of the text line. The text instances include both inclined texts as well as the horizontal texts.

Total-Text [8] contains not only horizontal and multioriented text instances but also the curved text. The dataset is consists of 1255 training images and 300 testing images. The images are annotated at the level of the word by a polygon with 2N vertices $(N \in \{2, ..., 15\})$.

6.2. Implementation Details

Proposal Label Generation. Assume the feature size is $W \times H$, our approach will generate $W \times H \times 4$ default anchors. For the label of each text polygon, two rules are defined to judge whether the default anchor is positive: 1. the center point of the default anchor is in ground truth and the height of default anchor is not great than 1.8 times of ground truth polygon height; 2. one of top two points in default box outside of ground truth polygon and the button is same as top. If both conditions are met, we consider it to be a positive anchor, otherwise, it is negative.

Text Segment Mask Generation. In order to separate the text instances that are very close to each other when generating text masks, we chose 50% region of ground truth polygon as a strong true region and set their scores to 1. The rest 50% regions are considered as the weak true region with a score of 0.1. The other regions are with the background score of 0.

Data Augmentation. The images are randomly rotated with 30 degrees, aspect ratios are set from 0.33 to 3 and randomly reverse image left and right. After these steps, random crop invert and blur are randomly adjusted. In order to ensure the cropped image have complete ground truth polygon, we only randomly crop the most left box to image width and apply it to the other three directions.

Model Training. Our method is implemented in Pytorch. We use Adam optimizer as our learning rate scheme. During the training stage, we chose 200 positive anchors for each feature as an input of ROI-Align, and at the testing stage we set $s_3 = 0.4$ as the positive square threshold and the maximum number of the positive square to ROI-Align is 2000.

Table 1. Quantitative results of different methods evaluated on CTW1500. LSN+CC indicates the full framework, while LSN is the results without the curve connection.

~~	BIT IS the results without the curit connection.					
	Method	Precision	Recall	F-measure		
	CTD [7]	74.3	65.2	69.5		
	CTD+TLOC [7]	77.4	69.8	73.4		
	SLPR [18]	80.1	70.1	74.8		
	TextSnake [9]	67.9	85.3	75.6		
	LSN	69.0	75.7	72.2		
	LSN+CC	83.2	78.8	80.8		

Table 2. Quantitative results of different methods evaluated on Total-Text

Method	Precision	Recall	F-measure
Total-Text [8]	40.0	33.0	36.0
Mask TextSpotter [19]	69.0	55.0	61.3
TextSnake [9]	82.7	74.5	78.4
LSN+CC	82.4	76.9	79.5

6.3. Results and Comparison

CTW1500. To evaluate our framework, we compare with different curve text detection methods: CTD and CTD+TLOC [7] by 14 landmark points to represent curved text, SLPR [18] with sliding line point regression, and TextSnake [9] which represents curve regions as a sequence of ordered, overlapping disks. Table 1 shows the results for text detection. It clearly shows that our proposed framework outperforms state-of-the-art methods on CTW1500. The improvement of LSN+CC over LSN demonstrate that adding the curve connection can not only achieve a significant improvement of precision but also helps the recall. Some text detection results are illustrated in Fig. 6.

Total-Text. Table 2 demonstrate the comparison of our full framework with state-of-the-art curve text detection approaches. We compare with Total-Text [8], Mask TextSpotter [19], and TextSnake [9]. The substantial performance gains over the published works confirm the effectiveness of using the local segmentation network and then curve connection for the text detection task. Some detection results obtained on the benchmark are illustrated in Fig. 7.

7. CONCLUSIONS

We proposed a text detection framework for the arbitrary shape scene text. Initial curve proposals were generated with the text segmentation on the horizontal rectangle local-level proposals by the local segmentation network. By combining the text segments sent from the previous step, the proposals can be refined in terms of the text polygon with the curve connection algorithms. Experimental comparisons with the state-of-the-art approaches on CTW1500 and Total-Text showed



Fig. 6. Text detection results on CTW1500.

the effectiveness of the proposed LSN and curve connection for the text detection task.

8. REFERENCES

- [1] Zhi Tian, Weilin Huang, Tong He, Pan He, and Yu Qiao, "Detecting text in natural image with connectionist text proposal network," in *ECCV*, 2016, pp. 56–72.
- [2] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang, "East: an efficient and accurate scene text detector," in *CVPR*, 2017, pp. 2642–2651.
- [3] Yuliang Liu and Lianwen Jin, "Deep matching prior network: Toward tighter multi-oriented text detection," in *CVPR*, 2017, pp. 3454–3461.
- [4] Jianqi Ma, Weiyuan Shao, Hao Ye, Li Wang, Hong Wang, Yingbin Zheng, and Xiangyang Xue, "Arbitrary-oriented scene text detection via rotation proposals," *IEEE Transactions on Multimedia*, vol. 20, no. 11, pp. 3111–3122, 2018.
- [5] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in CVPR, 2016, pp. 770–778.
- [6] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick, "Mask r-cnn," in ICCV, 2017, pp. 2980–2988.
- [7] Yuliang Liu, Lianwen Jin, Shuaitao Zhang, and Sheng Zhang, "Detecting curve text in the wild: New dataset and new solution," arXiv:1712.02170, 2017.
- [8] Chee Kheng Ch'ng and Chee Seng Chan, "Total-text: A comprehensive dataset for scene text detection and recognition," in *ICDAR*, 2017, pp. 935–942.
- [9] Shangbang Long, Jiaqiang Ruan, Wenjie Zhang, Xin He, Wenhao Wu, and Cong Yao, "Textsnake: A flexible representation for detecting text of arbitrary shapes," in *ECCV*, 2018, pp. 20–36.



Fig. 7. Text detection results on Total-Text.

- [10] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in NIPS, 2015, pp. 91–99.
- [11] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg, "Ssd: Single shot multibox detector," in ECCV, 2016.
- [12] Alex Graves and Jürgen Schmidhuber, "Framewise phoneme classification with bidirectional lstm and other neural network architectures," *Neural Networks*, vol. 18, no. 5-6, pp. 602–610, 2005.
- [13] Baoguang Shi, Xiang Bai, and Serge Belongie, "Detecting oriented text in natural images by linking segments," in *CVPR*, 2017, pp. 3482–3490.
- [14] Shangxuan Tian, Yifeng Pan, Chang Huang, Shijian Lu, Kai Yu, and Chew Lim Tan, "Text flow: A unified text detection system in natural scene images," in *ICCV*, 2015.
- [15] Xiang Li, Wenhai Wang, Wenbo Hou, Ruo-Ze Liu, Tong Lu, and Jian Yang, "Shape robust text detection with progressive scale expansion network," *arXiv:1806.02559*, 2018.
- [16] Tsung-Yi Lin, Priyal Goyal, Ross Girshick, Kaiming He, and Piotr Dollár, "Focal loss for dense object detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [17] Trevor Hastie and Werner Stuetzle, "Principal curves," *Journal of the American Statistical Association*, vol. 84, no. 406, pp. 502–516, 1989.
- [18] Yixing Zhu and Jun Du, "Sliding line point regression for shape robust scene text detection," *arXiv:1801.09969*, 2018.
- [19] Pengyuan Lyu, Minghui Liao, Cong Yao, Wenhao Wu, and Xiang Bai, "Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes," in *ECCV*, 2018.