




# TextSnake: A Flexible Representation for Detecting Text of Arbitrary Shapes

Shangbang Long<sup>1,2</sup> , Jiaqiang Ruan<sup>1,2</sup>, Wenjie Zhang<sup>1,2</sup>, Xin He<sup>2</sup>,  
Wenhao Wu<sup>2</sup>, and Cong Yao<sup>2</sup>  

<sup>1</sup> Peking University, Beijing, China

{longlongsb, jiaqiang.ruan, zhang.wen.jie}@pku.edu.cn

<sup>2</sup> Megvii (Face++) Technology Inc., Beijing, China

{hexin, wwh}@megvii.com, yaocong2010@gmail.com

**Abstract.** Driven by deep neural networks and large scale datasets, scene text detection methods have progressed substantially over the past years, continuously refreshing the performance records on various standard benchmarks. However, limited by the representations (axis-aligned rectangles, rotated rectangles or quadrangles) adopted to describe text, existing methods may fall short when dealing with much more free-form text instances, such as curved text, which are actually very common in real-world scenarios. To tackle this problem, we propose a more flexible representation for scene text, termed as *TextSnake*, which is able to effectively represent text instances in horizontal, oriented and curved forms. In TextSnake, a text instance is described as a sequence of ordered, overlapping disks centered at symmetric axes, each of which is associated with potentially variable radius and orientation. Such geometry attributes are estimated via a Fully Convolutional Network (FCN) model. In experiments, the text detector based on TextSnake achieves state-of-the-art or comparable performance on Total-Text and SCUT-CTW1500, the two newly published benchmarks with special emphasis on curved text in natural images, as well as the widely-used datasets ICDAR 2015 and MSRA-TD500. Specifically, TextSnake outperforms the baseline on Total-Text by more than 40% in F-measure.

**Keywords:** Scene text detection · Deep neural network · Curved text

## 1 Introduction

In recent years, the community has witnessed a surge of research interest and effort regarding the extraction of textual information from natural scenes, a.k.a. scene text detection and recognition [48]. The driving factors stem from both application prospect and research value. On the one hand, scene text detection and recognition have been playing ever-increasingly important roles in a wide range of practical systems, such as scene understanding, product search, and autonomous driving. On the other hand, the unique traits of scene text, for instance, significant variations in color, scale, orientation, aspect ratio and

pattern, make it obviously different from general objects. Therefore, particular challenges are posed and special investigations are required.



**Fig. 1.** Comparison of different representations for text instances. (a) Axis-aligned rectangle. (b) Rotated rectangle. (c) Quadrangle. (d) TextSnake. Obviously, the proposed TextSnake representation is able to effectively and precisely describe the geometric properties, such as location, scale, and bending of curved text with perspective distortion, while the other representations (axis-aligned rectangle, rotated rectangle or quadrangle) struggle with giving accurate predictions in such cases.

Text detection, as a prerequisite step in the pipeline of textual information extraction, has recently advanced substantially with the development of deep neural networks and large image datasets. Numerous innovative works [6, 9, 10, 17, 22, 28–31, 34, 36, 39, 40, 46, 47] are proposed, achieving excellent performances on standard benchmarks.

However, most existing methods for text detection shared a strong assumption that text instances are roughly in a linear shape and therefore adopted relatively simple representations (axis-aligned rectangles, rotated rectangles or quadrangles) to describe them. Despite their progress on standard benchmarks, these methods may fall short when handling text instances of irregular shapes, for example, curved text. As depicted in Fig. 1, for curved text with perspective distortion, conventional representations struggle with giving precise estimations of the geometric properties.

In fact, instances of curved text are quite common in real life [15, 43]. In this paper, we propose a more flexible representation that can fit well text of arbitrary shapes, i.e., those in horizontal, multi-oriented and curved forms. This representation describes text with a series of ordered, overlapping disks, each of which is located at the center axis of text region and associated with potentially variable radius and orientation. Due to its excellent capability in adapting for the complex multiplicity of text structures, just like a snake changing its shape to adapt for the external environment, the proposed representation is named as TextSnake. The geometry attributes of text instances, i.e., central axis points, radii and orientations, are estimated with a single Fully Convolutional Network (FCN) model. Besides ICDAR 2015 and MSRA-TD500, the effectiveness of TextSnake is validated on Total-Text and SCUT-CTW1500, which are two newly-released benchmarks mainly focused on curved text. The proposed algorithm achieves state-of-the-art performance on the two curved text datasets, while at the same time outperforming previous methods on horizontal and multi-oriented text,

even in the single-scale testing mode. Specifically, TextSnake achieves significant improvement over the baseline on Total-Text by 40.0% in F-measure.

In summary, the major contributions of this paper are three-fold: (1) We propose a flexible and general representation for scene text of arbitrary shapes; (2) Based on this representation, an effective method for scene text detection is proposed; (3) The proposed text detection algorithm achieves state-of-the-art performance on several benchmarks, including text instances of different forms (horizontal, oriented and curved).

## 2 Related Work

In the past few years, the most prominent trend in the area of scene text detection is the transfer from conventional methods [3, 24] to deep learning based methods [12, 13, 17, 29, 47]. In this section, we look back on relevant previous works. For comprehensive surveys, please refer to [41, 48]. Before the era of deep learning, SWT [3] and MSER [24] are two representative algorithms that have influenced a variety of subsequent methods [11, 42]. Modern methods are mostly based on deep neural networks, which can be coarsely classified into two categories: regression based and segmentation based.

Regression based text detection methods [17] mainly draw inspirations from general object detection frameworks. TextBoxes [17] adopted SSD [19] and added “long” default boxes and filters to handle the significant variation of aspect ratios of text instances. Based on Faster-RCNN [26], Ma *et al.* [23] devised Rotation Region Proposal Networks (RRPN) to detect arbitrary-Oriented text in natural images. EAST [47] and Deep Regression [8] both directly produce the rotated boxes or quadrangles of text, in a per-pixel manner.

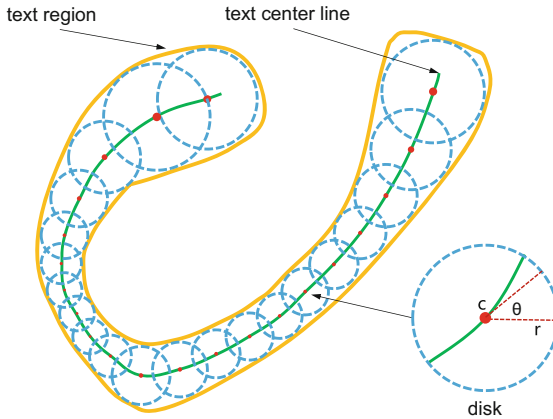
Segmentation based text detection methods cast text detection as a semantic segmentation problem and FCN [21] is often taken as the reference framework. Yao *et al.* [39] modified FCN to produce multiple heatmaps corresponding various properties of text, such as text region and orientation. Zhang *et al.* [46] first use FCN to extract text blocks and then hunt character candidates from these blocks with MSER [24]. To better separate adjacent text instances, the method of [36] distinguishes each pixel into three categories: non-text, text border and text. These methods mainly vary in the way they separate text pixels into different instances.

The methods reviewed above have achieved excellent performances on various benchmarks in this field. However, most works, except for [6, 15, 39], have not paid special attention to curved text. In contrast, the representation proposed in this paper is suitable for text of arbitrary shapes (horizontal, multi-oriented and curved). It is primarily inspired by [6, 39] and the geometric attributes of text are also estimated via the multiple-channel outputs of an FCN-based model. Unlike [39], our algorithm does not need character level annotations. In addition, it also shares a similar idea with SegLink [29], by successively decomposing text into local components and then composing them back into text instances. Analogous to [45], we also detect linear symmetry axes of text instances for text localization.

Another advantage of the proposed method lies in its ability to reconstruct the precise shape and regional strike of text instances, which can largely facilitate the subsequent text recognition process, because all detected text instances could be conveniently transformed into a canonical form with minimal distortion and background (see the example in Fig. 9).

### 3 Methodology

In this section, we first introduce the new representation for text of arbitrary shapes. Then we describe our method and training details.



**Fig. 2.** Illustration of the proposed TextSnake representation. Text region (in yellow) is represented as a series of ordered disks (in blue), each of which is located at the center line (in green, a.k.a symmetric axis or skeleton) and associated with a radius  $r$  and an orientation  $\theta$ . In contrast to conventional representations (e.g., axis-aligned rectangles, rotated rectangles and quadrangles), TextSnake is more flexible and general, since it can precisely describe text of different forms, regardless of shapes and lengths. (Color figure online)

#### 3.1 Representation

As shown in Fig. 1, conventional representations for scene text (e.g., axis-aligned rectangles, rotated rectangles and quadrangles) fail to precisely describe the geometric properties of text instances of irregular shapes, since they generally assume that text instances are roughly in linear forms, which does not hold true for curved text. To address this problem, we propose a flexible and general representation: TextSnake. As demonstrated in Fig. 2, TextSnake expresses a text instance as a sequence of overlapping disks, each of which is located at the center line and associated with a radius and an orientation. Intuitively, TextSnake is

able to change its shape to adapt for the variations of text instances, such as rotation, scaling and bending.

Mathematically, a text instance  $t$ , consisting of several characters, can be viewed as an ordered list  $S(t)$ .  $S(t) = \{D_0, D_1, \dots, D_i, \dots, D_n\}$ , where  $D_i$  stands for the  $i$ th disk and  $n$  is the number of the disks. Each disk  $D$  is associated with a group of geometry attributes, i.e.  $D = (c, r, \theta)$ , in which  $c$ ,  $r$  and  $\theta$  are the center, radius and orientation of disk  $D$ , respectively. The radius  $r$  is defined as half of the local width of  $t$ , while the orientation  $\theta$  is the tangential direction of the center line around the center  $c$ . In this sense, text region  $t$  can be easily reconstructed by computing the union of the disks in  $S(t)$ .

Note that the disks do not correspond to the characters belonging to  $t$ . However, the geometric attributes in  $S(t)$  can be used to rectify text instances of irregular shapes and transform them into rectangular, straight image regions, which are more friendly to text recognizers.

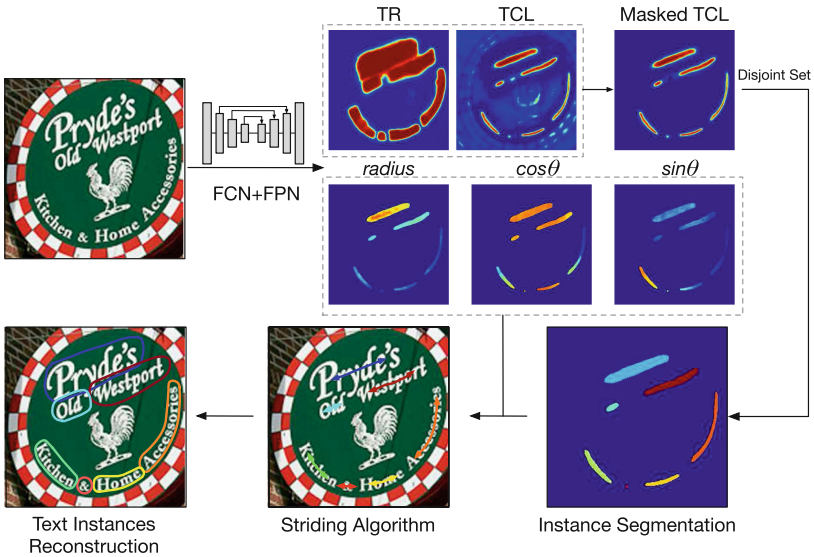
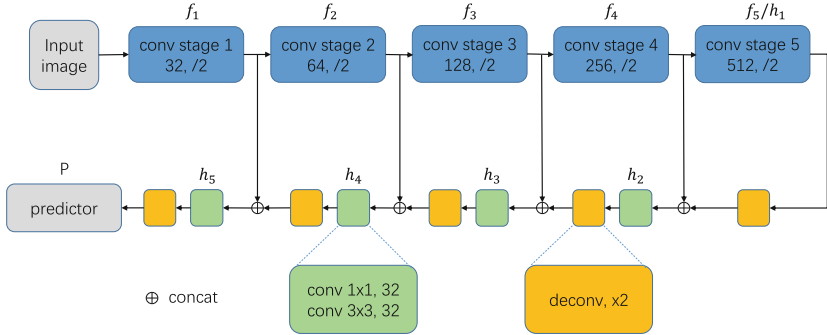


Fig. 3. Method framework: network output and post-processing

### 3.2 Pipeline

In order to detect text with arbitrary shapes, we employ an FCN model to predict the geometry attributes of text instances. The pipeline of the proposed method is illustrated in Fig. 3. The FCN based network predicts score maps of text center line (TCL) and text regions (TR), together with geometry attributes, including  $r$ ,  $\cos\theta$  and  $\sin\theta$ . The TCL map is further masked by the TR map since TCL is naturally part of TR. To perform instance segmentation, disjoint set is utilized, given the fact that TCL does not overlap with each other. A striding

algorithm is used to extract the central axis point lists and finally reconstruct the text instances.



**Fig. 4.** Network Architecture. Blue blocks are convolution stages of VGG-16. (Color figure online)

### 3.3 Network Architecture

The whole network is shown in Fig. 4. Inspired by FPN [18] and U-net [27], we adopt a scheme that gradually merges features from different levels of the stem network. The stem network can be convolutional networks proposed for image classification, e.g. VGG-16/19 [33] and ResNet [7]. These networks can be divided into 5 stages of convolutions and a few additional fully-connected (FC) layers. We remove the FC layers, and feed the feature maps after each stage to the feature merging network. We choose VGG-16 as our stem network for the sake of direct and fair comparison with other methods.

As for the feature merging network, several stages are stacked sequentially, each consisting of a merging unit that takes feature maps from the last stage and corresponding stem network layer. Merging unit is defined by the following equations:

$$h_1 = f_5 \quad (1)$$

$$h_i = conv_{3 \times 3}(conv_{1 \times 1}[f_{6-i}; UpSampling_{\times 2}(h_{i-1})]), \text{ for } i = 2, 3, 4, 5 \quad (2)$$

where  $f_i$  denotes the feature maps of the  $i$ -th stage in the stem network and  $h_i$  is the feature maps of the corresponding merging units. In our experiments, upsampling is implemented as deconvolutional layer as proposed in [44].

After the merging, we obtain a feature map whose size is  $\frac{1}{2}$  of the input images. We apply an additional upsampling layer and 2 convolutional layers to produce dense predictions:

$$h_{final} = UpSampling_{\times 2}(h_5) \quad (3)$$

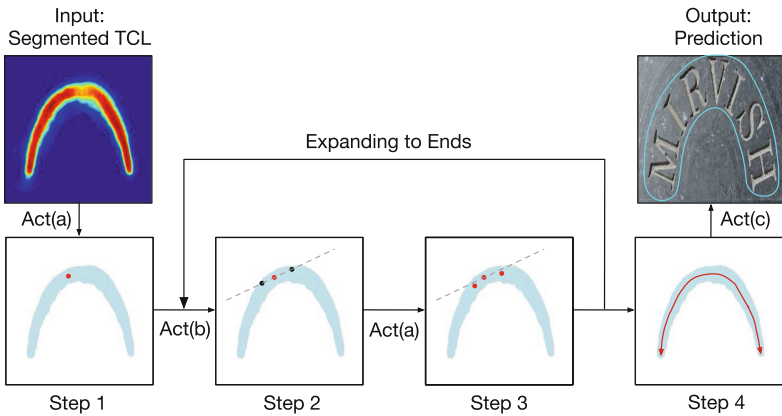
$$P = conv_{1 \times 1}(conv_{3 \times 3}(h_{final})) \quad (4)$$

where  $P \in \mathcal{R}^{h \times w \times 7}$ , with 4 channels for logits of TR/TCL, and the last 3 respectively for  $r$ ,  $\cos\theta$  and  $\sin\theta$  of the text instance. As a result of the additional upsampling layer,  $P$  has the same size as the input image. The final predictions are obtained by taking softmax for TR/TCL and regularizing  $\cos\theta$  and  $\sin\theta$  so that the squared sum equals 1.

### 3.4 Inference

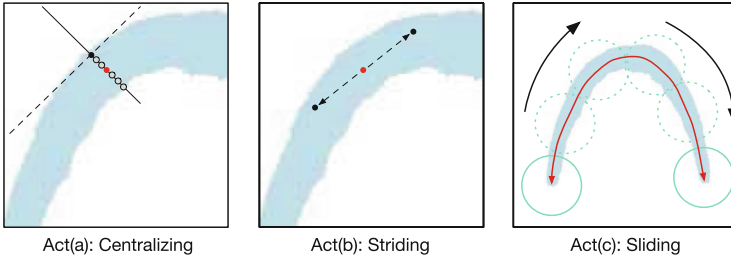
After feed-forwarding, the network produces the TCL, TR and geometry maps. For TCL and TR, we apply thresholding with values  $T_{tcl}$  and  $T_{tr}$  respectively. Then, the intersection of TR and TCL gives the final prediction of TCL. Using disjoint-set, we can efficiently separate TCL pixels into different text instances.

Finally, a striding algorithm is designed to extract an ordered point list that indicates the shape and course of the text instance, and also reconstruct the text instance areas. Two simple heuristics are applied to filter out false positive text instances: (1) The number of TCL pixels should be at least 0.2 times their average radius; (2) At least half of pixels in the reconstructed text area should be classified as TR.



**Fig. 5.** Framework of post-processing algorithm. Act(a) centralizing: relocate a given point to the central axis; Act(b) striding: a directional search towards the ends of text instances; Act(c) sliding: a reconstruction by sliding a circle along the central axis.

The procedure for the striding algorithm is shown in Fig. 5. It features 3 main actions, denoted as Act(a), Act(b), and Act(c), as illustrated in Fig. 6. Firstly, we randomly select a pixel as the starting point, and centralize it. Then, the search process forks into two opposite directions, striding and centralizing until it reaches the ends. This process would generate 2 ordered point lists in two opposite directions, which can be combined to produce the final central axis list that follows the course of the text and describes the shape precisely. Details of the 3 actions are shown below.



**Fig. 6.** Mechanisms of centralizing, striding and sliding

**Act(a) Centralizing.** As shown in Fig. 6, given a point on the TCL, we can draw the tangent line and the normal line, respectively denoted as dotted line and solid line. This step can be done with ease using the geometry maps. The midpoint of the intersection of the normal line and the TCL area gives the centralized point.

**Act(b) Striding.** The algorithm takes a stride to the next point to search. With the geometry maps, the displacement for each stride is computed and represented as  $(\frac{1}{2}r \times \cos\theta, \frac{1}{2}r \times \sin\theta)$  and  $(-\frac{1}{2}r \times \cos\theta, -\frac{1}{2}r \times \sin\theta)$ , respectively for the two directions. If the next step is outside the TCL area, we decrement the stride gradually until it's inside, or it hits the ends.

**Act(c) Sliding.** The algorithm iterates through the central axis and draw circles along it. Radii of the circles are obtained from the  $r$  map. The area covered by the circles indicates the predicted text instance.

In conclusion, taking advantage of the geometry maps and the TCL that precisely describes the course of the text instance, we can go beyond detection of text and also predict their shape and course. Besides, the striding algorithm saves our method from traversing all pixels that are related.

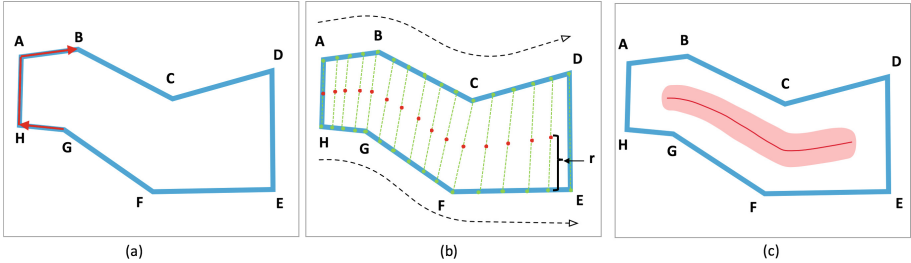
### 3.5 Label Generation

**Extracting Text Center Line.** For triangles and quadrangles, it's easy to directly calculate the TCL with algebraic methods, since in this case, TCL is a straight line. For polygons of more than 4 sides, it's not easy to derive a general algebraic method.

Instead, we propose a method that is based on the assumption that, text instances are snake-shaped, i.e. that it does not fork into multiple branches. For a snake-shaped text instance, it has two edges that are respectively the *head* and the *tail*. The two edges near the head or tail are running parallel but in opposite direction.

For a text instance  $t$  represented by a group of vertexes  $\{v_0, v_1, v_2, \dots, v_n\}$  in clockwise or counterclockwise order, we define a measurement for each edge  $e_{i,i+1}$  as  $M(e_{i,i+1}) = \cos\langle e_{i+1,i+2}, e_{i-1,i} \rangle$ . Intuitively, the two edges with  $M$  nearest to  $-1$ , e.g.  $AH$  and  $DE$  in Fig. 7, are the head and tail. After that, equal number





**Fig. 7.** Label Generation. (a) Determining text head and tail; (b) Extracting text center line and calculating geometries; (c) Expanded text center line.

of anchor points are sampled on the two sidelines, e.g.  $ABCD$  and  $HGFE$  in Fig. 7. TCL points are computed as midpoints of corresponding anchor points. We shrink the two ends of TCL by  $\frac{1}{2}r_{end}$  pixels, so that TCL are inside the TR and makes it easy for the network to learn to separate adjacent text instances.  $r_{end}$  denotes the radius of the TCL points at the two ends. Finally, we expand the TCL area by  $\frac{1}{5}r$ , since a single-point line is prone to noise.

**Calculating  $r$  and  $\theta$ .** For each points on TCL: (1)  $r$  is computed as the distance to the corresponding point on sidelines; (2)  $\theta$  is computed by fitting a straight line on the TCL points in the neighborhood. For non-TCL pixels, their corresponding geometry attributes are set to 0 for convenience.

### 3.6 Training Objectives

The proposed model is trained end-to-end, with the following loss functions as the objectives:

$$L = L_{cls} + L_{reg} \quad (5)$$

$$L_{cls} = \lambda_1 L_{tr} + \lambda_2 L_{tcl} \quad (6)$$

$$L_{reg} = \lambda_3 L_r + \lambda_4 L_{sin} + \lambda_5 L_{cos} \quad (7)$$

$L_{cls}$  in Eq. 5 represents classification loss for TR and TCL, and  $L_{reg}$  for regression loss of  $r$ ,  $\cos\theta$  and  $\sin\theta$ . In Eq. 6,  $L_{tr}$  and  $L_{tcl}$  are cross-entropy loss for TR and TCL. Online hard negative mining [32] is adopted for TR loss, with the ratio between the negatives and positives kept to 3:1 at most. For TCL, we only take into account pixels inside TR and adopt no balancing methods.

In Eq. 7, regression loss, i.e.  $L_r$ ,  $L_{sin}$  and  $L_{cos}$ , are calculated as Smoothed-L1 loss [4]:

$$\begin{pmatrix} L_r \\ L_{cos} \\ L_{sin} \end{pmatrix} = \text{SmoothedL1} \begin{pmatrix} \frac{\widehat{r} - r}{r} \\ \widehat{\cos\theta} - \cos\theta \\ \widehat{\sin\theta} - \sin\theta \end{pmatrix} \quad (8)$$

where  $\hat{r}$ ,  $\widehat{\cos\theta}$  and  $\widehat{\sin\theta}$  are the predicted values, while  $r$ ,  $\cos\theta$  and  $\sin\theta$  are their ground truth correspondingly. Geometry loss outside TCL are set to 0, since these attributes make no sense for non-TCL points.

The weights constants  $\lambda_1$ ,  $\lambda_2$ ,  $\lambda_3$ ,  $\lambda_4$  and  $\lambda_5$  are all set to 1 in our experiments.

## 4 Experiments

In this section, we evaluate the proposed algorithm on standard benchmarks for scene text detection and compare it with previous methods. Analyses and discussions regarding our algorithm are also given.

### 4.1 Datasets

The datasets used for the experiments in this paper are briefly introduced below:

**SynthText** [5] is a large scale dataset that contains about 800K synthetic images. These images are created by blending natural images with text rendered with random fonts, sizes, colors, and orientations, thus these images are quite realistic. We use this dataset to pre-train our model.

**TotalText** [15] is a newly-released benchmark for text detection. Besides horizontal and multi-Oriented text instances, the dataset specially features *curved text*, which rarely appear in other benchmark datasets, but are actually quite common in real environments. The dataset is split into training and testing sets with 1255 and 300 images, respectively.

**CTW1500** [43] is another dataset mainly consisting of curved text. It consists of 1000 training images and 500 test images. Text instances are annotated with polygons with 14 vertexes.

**ICDAR 2015** is proposed as the Challenge 4 of the 2015 Robust Reading Competition [14] for incidental scene text detection. There are 1000 images for training and 500 images for testing. The text instances from this dataset are labeled as word level quadrangles.

**MSRA-TD500** [38] is a dataset with multi-lingual, arbitrary-oriented and long text lines. It includes 300 training images and 200 test images with text line level annotations. Following previous works [22,47], we also include the images from HUST-TR400 [37] as training data when fine-tuning on this dataset, since its training set is rather small.

For experiments on ICDAR 2015 and MSRA-TD500, we fit a minimum bounding rectangle based on the output text area of our method.

### 4.2 Data Augmentation

Images are randomly rotated, and cropped with areas ranging from 0.24 to 1.69 and aspect ratios ranging from 0.33 to 3. After that, noise, blur, and lightness are randomly adjusted. We ensure that the text on the augmented images are still legible, if they are legible before augmentation.



**Fig. 8.** Qualitative results by the proposed method. Top: detected text contours (in *yellow*) and ground truth annotations (in *green*). Bottom: combined score maps for TR (in *red*) and TCL (in *yellow*). From left to right in column: image from ICDAR 2015, TotalText, CTW1500 and MSRA-TD500. Best viewed in color. (Color figure online)

### 4.3 Implementation Details

Our method is implemented in Tensorflow 1.3.0 [1]. The network is pre-trained on SynthText for one epoch and fine-tuned on other datasets. We adopt the Adam optimizer [16] as our learning rate scheme. During the pre-training stage, the learning rate is fixed to  $10^{-3}$ . During the fine-tuning stage, the learning rate is set to  $10^{-3}$  initially and decays exponentially with a rate of 0.8 every 5000 iterations. During fine-tuning, the number of iterations is decided by the sizes of datasets. All the experiments are conducted on a regular workstation (CPU: Intel(R) Xeon(R) CPU E5-2650 v3 @ 2.30 GHz; GPU: Titan X; RAM: 384 GB). We train our model with the batch size of 32 on 2 GPUs in parallel and evaluate our model on 1 GPU with batch size set as 1. Hyper-parameters are tuned by grid search on training set.

### 4.4 Experiment Results

**Experiments on Curved Text (Total-Text and CTW1500).** Fine-tuning on these two datasets stops at about  $5k$  iterations. Thresholds  $T_{tr}$ ,  $T_{tcl}$  are set to (0.4, 0.6) and (0.4, 0.5) respectively on Total-Text and CTW1500. In testing, all images are rescaled to  $512 \times 512$  for Total-Text, while for CTW1500, the images are not resized, since the images in CTW1500 are rather small (The largest image is merely  $400 \times 600$ ). For comparison, we also evaluated the models of EAST [47] and SegLink [29] on Total-Text and CTW1500. The quantitative results of different methods on these two datasets are shown in Table 1.

The superior performances of our method on Total-Text and CTW1500 verify that the proposed representation can handle well curved text in natural images.

**Experiments on Incidental Scene Text (ICDAR 2015).** Fine-tuning on ICDAR 2015 stops at about  $30k$  iterations. In testing, all images are resized to  $1280 \times 768$ .  $T_{tr}$ ,  $T_{tcl}$  are set to (0.4, 0.9). For the consideration that images in

**Table 1.** Quantitative results of different methods evaluated on Total-Text and CTW1500. Note that EAST and SegLink were not fine-tuned on Total-Text. Therefore their results are included only for reference. Comparative results on CTW1500 are obtained from [43].

Datasets	Total-text			CTW1500		
Method	Precision	Recall	F-measure	Precision	Recall	F-measure
SegLink [29]	30.3	23.8	26.7	42.3	40.0	40.8
EAST [47]	50.0	36.2	42.0	78.7	49.1	60.4
DeconvNet [25]	33.0	40.0	36.0	-	-	-
DMPNet [20]	-	-	-	69.9	56.0	62.2
CTD [43]	-	-	-	74.3	65.2	69.5
CTD+TLOC [43]	-	-	-	<b>77.4</b>	69.8	73.4
<b>TextSnake</b>	<b>82.7</b>	<b>74.5</b>	<b>78.4</b>	67.9	<b>85.3</b>	<b>75.6</b>

ICDAR 2015 contains many unlabeled small texts, predicted rectangles with the shorter side less than 10 pixels or the area less than 300 are filtered out.

The quantitative results of different methods on ICDAR 2015 are shown in Table 2. With only *single-scale* testing, our method outperforms most competitors (including those evaluated in multi-scale). This demonstrates that the proposed representation TextSnake is general and can be readily applied to multi-oriented text in complex scenarios.

**Experiments on Long Straight Text Lines (MSRA-TD500).** Fine-tuning on MSRA-TD500 stops at about  $10k$  iterations. Thresholds for  $T_{tr}$ ,  $T_{tcl}$  are (0.4, 0.6). In testing, all images are resized to  $1280 \times 768$ . Results are shown in Table 2. The F-measure (78.3%) of the proposed method is higher than that of the other methods.

## 4.5 Analyses and Discussions

**Precise Description of Text Instances.** What distinguishes our method from others is its ability to predict a precise description of the shape and course of text instances (see Fig. 8).

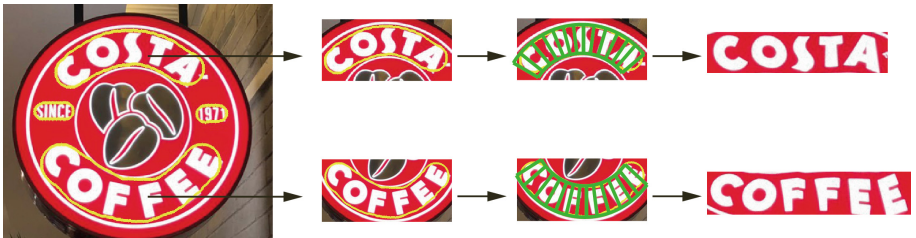
We attribute such ability to the TCL mechanism. Text center line can be seen as a kind of skeletons that prop up the text instance, and geo-attributes providing more details. Text, as a form of written language, can be seen as a stream of signals mapped onto 2D surfaces. Naturally, it should follow a course to extend.

Therefore we propose to predict TCL, which is much narrower than the whole text instance. It has two advantages: (1) A slim TCL can better describe the course and shape; (2) TCL, intuitively, does not overlap with each other, so that instance segmentation can be done in a very simple and straightforward way, thus simplifying our pipeline.

**Table 2.** Quantitative results of different methods on ICDAR 2015 and MSRA-TD500. \*stands for multi-scale, †indicates that the base net of the model is not VGG16.

Datasets	ICDAR 2015			MSRA-TD500			FPS
Method	Precision	Recall	F-measure	Precision	Recall	F-measure	
Zhang <i>et al.</i> [46]	70.8	43.0	53.6	83.0	67.0	74.0	0.48
Yao <i>et al.</i> [39]	72.3	58.7	64.8	76.5	<b>75.3</b>	75.9	1.61
SegLink [29]	73.1	76.8	75.0	86.0	70.0	77.0	-
EAST [47]	80.5	72.8	76.4	81.7	61.6	70.2	6.52
WordSup * [9]	79.3	77.0	78.2	-	-	-	2
EAST * † [47]	83.3	78.3	80.7	<b>87.3</b>	67.4	76.1	<b>13.2</b>
He <i>et al.</i> * † [8]	82.0	80.0	81.0	77.0	70.0	74.0	1.1
PixelLink [2]	<b>85.5</b>	<b>82.0</b>	<b>83.7</b>	83.0	73.2	77.8	3.0
<b>TextSnake</b>	84.9	80.4	82.6	83.2	73.9	<b>78.3</b>	1.1

Moreover, as depicted in Fig. 9, we can exploit local geometries to sketch the structure of the text instance and transform the predicted curved text instances into canonical form, which may largely facilitate the recognition stage.



**Fig. 9.** Text instances transformed to canonical form using the predicted geometries.

**Generalization Ability.** To further verify the generalization ability of our method, we train and fine-tune our model on datasets *without* curved text and evaluate it on the two benchmarks *featuring* curved text. Specifically, we fine-tune our models on ICDAR 2015, and evaluate them on the target datasets. The models of EAST [47], SegLink [29], and PixelLink [2] are taken as baselines, since these two methods were also trained on ICDAR 2015.

As shown in Table 3, our method still performs well on curved text and significantly outperforms the three strong competitors SegLink, EAST and PixelLink, without fine-tuning on curved text. We attribute this excellent generalization ability to the proposed flexible representation. Instead of taking text as a whole, the representation treats text as a collection of local elements and integrates them together to make decisions. Local attributes are kept when formed into a

**Table 3.** Comparison of cross-dataset results of different methods. The following models are fine-tuned on ICDAR 2015 and evaluated on Total-Text and CTW1500. Experiments for SegLink, EAST and PixelLink are done with the open source code. The evaluation protocol is DetEval [35], the same as Total-Text.

Datasets	Total-text			CTW1500		
Methods	Precision	Recall	F-measure	Precision	Recall	F-measure
SegLink [29]	35.6	33.2	34.4	33.0	2.4	30.5
EAST [47]	49.0	43.1	45.9	46.7	37.2	41.4
PixelLink [2]	53.5	52.7	53.1	50.6	42.8	46.4
<b>TextSnake</b>	<b>61.5</b>	<b>67.9</b>	<b>64.6</b>	<b>65.4</b>	<b>63.4</b>	<b>64.4</b>

whole. Besides, they are independent of each other. Therefore, the final predictions of our method can retain most information of the shape and course of the text. We believe that this is the main reason for the capacity of the proposed text detection algorithm in hunting text instances with various shapes.

## 5 Conclusion and Future Work

In this paper, we present a novel, flexible representation for describing the properties of scene text with arbitrary shapes, including horizontal, multi-oriented and curved text instances. The proposed text detection method based upon this representation obtains state-of-the-art or comparable performance on two newly-released benchmarks for curved text (Total-Text and SCUT-CTW1500) as well as two widely-used datasets (ICDAR 2015 and MSRA-TD500) in this field, proving the effectiveness of the proposed method. As for future work, we would explore the direction of developing an end-to-end recognition system for text of arbitrary shapes.

## References

1. Abadi, M., et al.: TensorFlow: a system for large-scale machine learning. *OSDI* **16**, 265–283 (2016)
2. Deng, D., Liu, H., Li, X., Cai, D.: PixelLink: detecting scene text via instance segmentation. In: *Proceedings of AAAI* (2018)
3. Epshtein, B., Ofek, E., Wexler, Y.: Detecting text in natural scenes with stroke width transform. In: *Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2963–2970. IEEE (2010)
4. Girshick, R.: Fast R-CNN. In: *Proceedings of The IEEE International Conference on Computer Vision (ICCV)*, December 2015
5. Gupta, A., Vedaldi, A., Zisserman, A.: Synthetic data for text localisation in natural images. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2315–2324 (2016)

6. He, D., et al.: Multi-scale FCN with cascaded instance aware segmentation for arbitrary oriented word spotting in the wild. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 474–483. IEEE (2017)
7. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
8. He, W., Zhang, X.Y., Yin, F., Liu, C.L.: Deep direct regression for multi-oriented scene text detection. In: Proceedings of The IEEE International Conference on Computer Vision (ICCV), October 2017
9. Hu, H., Zhang, C., Luo, Y., Wang, Y., Han, J., Ding, E.: WordSup: exploiting word annotations for character based text detection. In: Proceedings of The IEEE International Conference on Computer Vision (ICCV), October 2017
10. Huang, L., Yang, Y., Deng, Y., Yu, Y.: DenseBox: unifying landmark localization with end to end object detection. arXiv preprint [arXiv:1509.04874](https://arxiv.org/abs/1509.04874) (2015)
11. Huang, W., Qiao, Y., Tang, X.: Robust scene text detection with convolution neural network induced MSER trees. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014. LNCS, vol. 8692, pp. 497–511. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-10593-2\\_33](https://doi.org/10.1007/978-3-319-10593-2_33)
12. Jaderberg, M., Simonyan, K., Vedaldi, A., Zisserman, A.: Reading text in the wild with convolutional neural networks. *Int. J. Comput. Vis.* **116**(1), 1–20 (2016)
13. Jaderberg, M., Vedaldi, A., Zisserman, A.: Deep features for text spotting. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014. LNCS, vol. 8692, pp. 512–528. Springer, Cham (2014). [https://doi.org/10.1007/978-3-319-10593-2\\_34](https://doi.org/10.1007/978-3-319-10593-2_34)
14. Karatzas, D., et al.: ICDAR 2015 competition on robust reading. In: 13th International Conference on Document Analysis and Recognition (ICDAR), pp. 1156–1160. IEEE (2015)
15. Kheng Chng, C., Chan, C.S.: Total-text: a comprehensive dataset for scene text detection and recognition. In: 14th IAPR International Conference on Document Analysis and Recognition (ICDAR) (2017)
16. Kingma, D., Ba, J.: Adam: a method for stochastic optimization. In: Proceedings of ICLR (2015)
17. Liao, M., Shi, B., Bai, X., Wang, X., Liu, W.: TextBoxes: a fast text detector with a single deep neural network. In: Proceedings of AAAI, pp. 4161–4167 (2017)
18. Lin, T.Y., Dollar, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017
19. Liu, W., et al.: SSD: single shot MultiBox detector. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) ECCV 2016. LNCS, vol. 9905, pp. 21–37. Springer, Cham (2016). [https://doi.org/10.1007/978-3-319-46448-0\\_2](https://doi.org/10.1007/978-3-319-46448-0_2)
20. Liu, Y., Jin, L.: Deep matching prior network: toward tighter multi-oriented text detection (2017)
21. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3431–3440 (2015)
22. Lyu, P., Yao, C., Wu, W., Yan, S., Bai, X.: Multi-oriented scene text detection via corner localization and region segmentation. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
23. Ma, J., et al.: Arbitrary-oriented scene text detection via rotation proposals. arXiv preprint [arXiv:1703.01086](https://arxiv.org/abs/1703.01086) (2017)

24. Neumann, L., Matas, J.: A method for text localization and recognition in real-world images. In: Kimmel, R., Klette, R., Sugimoto, A. (eds.) ACCV 2010. LNCS, vol. 6494, pp. 770–783. Springer, Heidelberg (2011). [https://doi.org/10.1007/978-3-642-19318-7\\_60](https://doi.org/10.1007/978-3-642-19318-7_60)
25. Noh, H., Hong, S., Han, B.: Learning deconvolution network for semantic segmentation, pp. 1520–1528 (2015)
26. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. In: Advances in Neural Information Processing Systems, pp. 91–99 (2015)
27. Ronneberger, O., Fischer, P., Brox, T.: U-Net: convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) MICCAI 2015. LNCS, vol. 9351, pp. 234–241. Springer, Cham (2015). [https://doi.org/10.1007/978-3-319-24574-4\\_28](https://doi.org/10.1007/978-3-319-24574-4_28)
28. Sheng, Z., Yuliang, L., Lianwen, J., Canjie, L.: Feature enhancement network: a refined scene text detector. In: Proceedings of AAAI (2018)
29. Shi, B., Bai, X., Belongie, S.: Detecting oriented text in natural images by linking segments. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017
30. Shi, B., Bai, X., Yao, C.: An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* **39**(11), 2298–2304 (2017)
31. Shi, B., Yang, M., Wang, X., Lyu, P., Yao, C., Bai, X.: ASTER: an attentional scene text recognizer with flexible rectification. *IEEE Trans. Pattern Anal. Mach. Intell.* (2018)
32. Shrivastava, A., Gupta, A., Girshick, R.: Training region-based object detectors with online hard example mining, pp. 761–769 (2016)
33. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint [arXiv:1409.1556](https://arxiv.org/abs/1409.1556) (2014)
34. Tian, S., Lu, S., Li, C.: WeText: scene text detection under weak supervision. In: Proceedings of The IEEE International Conference on Computer Vision (ICCV) (2017)
35. Wolf, C., Jolion, J.M.: Object count/area graphs for the evaluation of object detection and segmentation algorithms. *Int. J. Doc. Anal. Recognit. (IJDAR)* **8**(4), 280–296 (2006)
36. Wu, Y., Natarajan, P.: Self-organized text detection with minimal post-processing via border learning. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 5000–5009 (2017)
37. Yao, C., Bai, X., Liu, W.: A unified framework for multioriented text detection and recognition. *IEEE Trans. Image Process.* **23**(11), 4737–4749 (2014)
38. Yao, C., Bai, X., Liu, W., Ma, Y., Tu, Z.: Detecting texts of arbitrary orientations in natural images. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1083–1090. IEEE (2012)
39. Yao, C., Bai, X., Sang, N., Zhou, X., Zhou, S., Cao, Z.: Scene text detection via holistic, multi-channel prediction. arXiv preprint [arXiv:1606.09002](https://arxiv.org/abs/1606.09002) (2016)
40. Yao, C., Bai, X., Shi, B., Liu, W.: Strokelets: a learned multi-scale representation for scene text recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 4042–4049 (2014)
41. Ye, Q., Doermann, D.: Text detection and recognition in imagery: a survey. *IEEE Trans. Pattern Anal. Mach. Intell.* **37**(7), 1480–1500 (2015)
42. Yin, X.C., Yin, X., Huang, K., Hao, H.W.: Robust text detection in natural scene images. *IEEE Trans. Pattern Anal. Mach. Intell.* **36**(5), 970–983 (2014)



43. Yuliang, L., Lianwen, J., Shuaitao, Z., Sheng, Z.: Detecting curve text in the wild: new dataset and new solution. arXiv preprint [arXiv:1712.02170](https://arxiv.org/abs/1712.02170) (2017)
44. Zeiler, M.D., Krishnan, D., Taylor, G.W., Fergus, R.: Deconvolutional networks. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2528–2535 (2010)
45. Zhang, Z., Shen, W., Yao, C., Bai, X.: Symmetry-based text line detection in natural scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 2558–2567 (2015)
46. Zhang, Z., Zhang, C., Shen, W., Yao, C., Liu, W., Bai, X.: Multi-oriented text detection with fully convolutional networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 4159–4167 (2016)
47. Zhou, X., et al.: EAST: an efficient and accurate scene text detector. In: Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), July 2017
48. Zhu, Y., Yao, C., Bai, X.: Scene text detection and recognition: recent advances and future trends. *Front. Comput. Sci.* **10**(1), 19–36 (2016)