# TextField: Learning A Deep Direction Field for Irregular Scene Text Detection

Yongchao Xu, Yukang Wang, Wei Zhou, Yongpan Wang, Zhibo Yang, Xiang Bai, *Senior Member, IEEE*

*Abstract*—Scene text detection is an important step of scene text reading system. The main challenges lie on significantly varied sizes and aspect ratios, arbitrary orientations and shapes. Driven by recent progress in deep learning, impressive performances have been achieved for multi-oriented text detection. Yet, the performance drops dramatically in detecting curved texts due to the limited text representation (*e.g.*, horizontal bounding boxes, rotated rectangles, or quadrilaterals). It is of great interest to detect curved texts, which are actually very common in natural scenes. In this paper, we present a novel text detector named TextField for detecting irregular scene texts. Specifically, we learn a direction field pointing away from the nearest text boundary to each text point. This direction field is represented by an image of two-dimensional vectors and learned via a fully convolutional neural network. It encodes both binary text mask and direction information used to separate adjacent text instances, which is challenging for classical segmentation-based approaches. Based on the learned direction field, we apply a simple yet effective morphological-based post-processing to achieve the final detection. Experimental results show that the proposed TextField outperforms the state-of-the-art methods by a large margin (28% and 8%) on two curved text datasets: Total-Text and SCUT-CTW1500, respectively, and also achieves very competitive performance on multi-oriented datasets: IC-DAR 2015 and MSRA-TD500. Furthermore, TextField is robust in generalizing to unseen datasets. The code is available at https://github.com/YukangWang/TextField.

*Index Terms*—Scene text detection, multi-oriented text, curved text, deep neural networks

(a) Horizontal box    (b) Rotated rectangle    (c) Quadrilateral

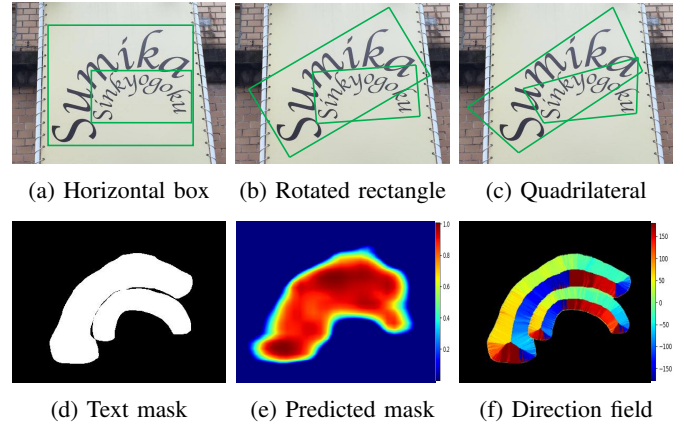(d) Text mask    (e) Predicted mask    (f) Direction field

Fig. 1: Different text representations. Classical relatively simple text representations in (a-c) fail to accurately delimit irregular texts. The text instances in (e) stick together using binary text mask representation in (d), requiring heavy post-processing to extract text instances. The proposed direction field in (f) is able to precisely describe irregular text instances.

## I. INTRODUCTION

Scene text frequently appears on many scenes and carries important information for many applications, such as product search [1], scene understanding [2], [3], and autonomous driving [4], [5]. Scene text reading is thus of great importance. As compared to general object detection, scene text detection, the prerequisite step of scene text recognition, faces particular challenges [6] due to significantly varied aspect ratios and sizes (usually small), uncontrollable lighting conditions, arbitrary orientations and shapes. To cope with these challenges, traditional methods [7]–[15] tend to involve complete pipelines and resort to specifically engineered features. The traditional pipeline usually consists of candidate character/word generation [16], [17], candidate filtering [18] and grouping [10], [13]. Each module requires careful parameter tuning and specifical heuristic rules designing to make it work properly. It is thus

Yongchao Xu, Yukang Wang, Wei Zhou, and Xiang Bai are with the School of Electronic Information and Communications, Huazhong University of Science and Technology (HUST), Wuhan, 430074, China. Yongpan Wang and Zhibo Yang are with Alibaba Group. Email: {yongchaoxu, wangyk, weizhou, xbai}@hust.edu.cn, yongpan@taobao.com, zhibo.yzb@alibaba-inc.com. (*Corresponding author: Xiang Bai.*)

difficult to optimize the whole pipeline, and also results in low detection speed.

Thanks to recent development of object detection [19]–[21] and segmentation [22] with deep learning, scene text detection has witnessed a great progress [23]–[45]. They can be roughly divided into three categories: 1) Regression-based methods. Scene text is a specific type of object. Many recent methods [23], [24], [26], [28]–[32] adapt the general object detection framework to detect texts by directly regressing horizontal/oriented rectangles or quadrilaterals, which enclose texts. Some other methods attempt to regress text parts [33], [34], [36], [42] or corners [35] followed by a linking or combination process. 2) Segmentation-based methods. Scene text detection can also be regarded as text instance segmentation. Several methods [37]–[41] rely on fully convolutional network to segment text areas. A heavy post-processing is usually involved to extract text instances from the segmented text areas. 3) Hybrid methods. Some other methods [25], [27] predict text score maps via segmentation and then obtain bounding boxes via regression.

The popular regression-based methods and existing hybrid methods [23]–[36] achieve excellent performances on standard benchmarks. Yet, they have a strong bottleneck which assumes a text instance has a linear shape, thus adopting relatively simple text representation in terms of horizontal/oriented rectangles or quadrilaterals. Their performances drop significantly

Fig. 2: Some irregular scene text detection results (enclosed by green contours) on some challenging images.

for detecting text of irregular shapes, *e.g.*, curved text. Besides, as depicted in Fig. 1(a-c), the traditional simple text representations do not achieve precise text delimitation providing texts' geometrical properties, which are useful for the subsequent recognition [46], [47]. Segmentation-based methods [37]–[41] may not suffer from this problem. Yet, as depicted in Fig. 1(e), though the predicted text region is a good estimation of text areas, it is rather difficult to separate close text instances. Indeed, many efforts of segmentation-based methods focus on how to separate segmented text regions into text instances.

In real-world scenarios, curved texts appear frequently [48] and can be easily found in real life scenes such as bottles, spherical objects, clothes, logos, signboards. In two recently released datasets (Total-Text [41] and SCUT-CTW1500 [49]) for scene text detection, around $40\%$ text instances are curved texts.

In this paper, we propose a novel text detector deemed TextField for detecting texts of arbitrary shapes and orientations. Inspired by component tree representation [16], [50]–[52] that links neighboring pixels following their intensity order to form candidate characters, we propose to learn a deep direction field, which is similar to the notion of flux image [53], to link neighboring pixels and form candidate text parts. The learned direction information is further used to group text parts into text instances. For that, the text areas are translated into text direction field first, pointing away from the nearest text boundary to each text point. Specifically, this direction field is encoded by an image of two-dimensional vectors for network training. For text areas, the field is defined as a unit vector encoding the direction, and for non-text areas, the direction field is set to $(0, 0)$. Thus, the magnitude information provides the text mask, while the direction information facilitates the post-processing of separating predicted text areas into text instances. An example of such direction field is given in Fig. 1(f). We adopt a fully convolutional network to directly regress the direction field. The candidate text pixels are then obtained by thresholding the magnitude. The direction information is used to extract text instances from candidate text pixels via some morphological operators. This results in detections with precise delimitation of irregular scene texts. Several examples are given in Fig. 2. The proposed TextField significantly outperforms other methods by 28% and 8% in F-measure on Total-Text [41] and SCUT-CTW1500 [49], respectively, while achieving very competitive performances on two widely adopted multi-oriented text datasets.

The main contributions of this paper are three folds: 1) We propose a novel direction field which can represent scene texts of arbitrary shapes. This direction field encodes both binary text mask and direction information facilitating the subsequent text grouping process. 2) Based on the direction field, we present a text detector named TextField, which efficiently detects irregular scene texts. 3) The proposed TextField significantly outperforms state-of-the-art methods on two curved text datasets and achieves competitive performances on two widely adopted multi-oriented text datasets.

The rest of this paper is organized as follows. We shortly review some related works on scene text detection in Section II. The proposed method is then detailed in Section III, followed by extensive experimental results in Section IV. Finally, we conclude and give some perspectives in Section V.

## II. RELATED WORKS

Scene text detection has been extensively exploited recently. We first review some representative methods in Section II-A. A comprehensive review of recent scene text detectors can be found in [6], [54]. The comparison of the proposed TextField with some related works is depicted in Section II-B.

### A. Scene text detection

Scene text detection methods can be roughly classified into specifically engineered and deep learning-based methods. Before the era of deep learning, scene text detector pipelines usually consist of text component extraction and filtering, component grouping, and candidate filtering. The key step is extracting text components based on some engineered features. Maximally Stable Extremal Regions (MSER) [16] and Stroke Width Transform (SWT) [17] are two representative works for text component extraction. Many traditional methods [7], [9], [11]–[13] are based on these two algorithms. Other examples of this type are [55], [56]. Most recent methods shift to deep neural networks to extract scene texts. In general, they can be roughly summarized into regression-based, segmentation-based, and hybrid methods. For the regression-based ones, they can be further divided into two categories based on the target to regress: proposal-based and part-based methods.

**Proposal-based methods:** Proposal-based methods are mainly inspired by recent object detection pipelines. TextBoxes [24] directly adapts SSD [20] for scene text detection by using long default boxes and convlutioinal filters to cope with the significantly varied aspect ratios.

TextBoxes++ [30] extends TextBoxes by regressing quadrilaterals instead of horizontal bounding boxes. Ma *et al.* [29] attempt to solve the multi-oriented text detection by adopting Rotated Regional Proposal Network (RRPN) in the pipeline of faster r-cnn. Quadrilateral sliding windows are adopted in [23] to detect multi-oriented texts. Wordsup [28] explores word annotations for character-based text detection. SSTD [26] introduces the attention mechanism by FCN to suppress background interference, improving accurate detection of small texts. In [31], Liao *et al.* propose to apply rotation-invariant and sensitive features for text/non-text box classification and regression, respectively, boosting long multi-oriented text detection. Wang et al. [32] propose instance transformation network by considering geometry-aware information for scene text detection.

**Part-based methods:** Some other regression-based methods tend to regress text parts while predicting the linking relationship between them. In [33], the authors propose a Connectionist Text Proposal Network (CTPN) by first predicting vertical text parts, then adopting a recurrent neural network to link text parts. Shi *et al.* present a network named SegLink [34] to first detect text parts named text segments while predicting the linking relationship between neighboring text segments. A novel framework named Markov Clustering Network (MCN) is proposed in [36]. In this work, the authors propose to regard an image as a stochastic flow graph, where the flows are strong between text nodes (*i.e.*, text pixels) but weak for the others. Then a Markov clustering process is applied to form text instances from the predicted flow graph. In [35], *Lyu et al.* propose to first regress four corners of text boxes, followed by a combination of corners and Non-Maximum Suppression (NMS) process to achieve accurate multi-oriented text localization.

**Segmentation-based methods:** Segmentation-based approaches regard text detection as a text area segmentation problem, which is usually achieved via Fully Convolutional Neural Network (FCN). They mainly differ in how to post-process the predicted text regions into words or text lines. In [37], Zhang *et al.* adopted an FCN to estimate text blocks, on which candidate characters are extracted using MSER. Then they use traditional grouping and filtering strategies to achieve multi-oriented text detection. In addition to text block (word or line) prediction, Yao *et al.* [38] also propose to predict both individual characters and the orientation of text boxes via an FCN in a holistic fashion. Then a grouping process based on the three estimated properties of text yields the text detection. Ch'ng *et al.* [41] fine-tune DeconvNet [57] to achieve curved text detection. In [40], the authors consider the text detection as an instance segmentation problem using multi-scale image inputs. They adopt an FCN to predict text blocks, followed by two CNN branches predicting text lines and instance-aware segmentations from the estimated text blocks. Wu *et al.* [39] introduce text border in addition to text/non-text segmentation, which results in a three-class semantic segmentation, facilitating the separation of adjacent text instances. Xue *et al.* further improve [39] by exploiting bootstrapping techniques and designing semantics-aware text border detection technique for accurate text localization.

**Hybrid methods:** It is also worth to mention that some other methods leverage segmentation to classify text/non-text pixels and then localize texts via bounding box regression. For example, East [25] and Deep regression [27] both perform per-pixel rotated rectangle or quadrilateral estimation.

### B. Comparison with related works

**TextField Versus Traditional component-based methods:** Traditional methods rely on engineered features to extract text components, and heuristic grouping rules to form text instances. Each module requires careful parameter tuning, resulting in sub-optimal performance and slow runtime of the whole pipeline. The proposed TextField leverages deeply learned direction field which encodes both text mask and direction information facilitating subsequent text grouping process. The whole pipeline is more effective in both performance and runtime.

**TextField Versus Proposal-based and hybrid** methods: The proposal-based and hybrid scene text detection methods are mainly inspired by recent object detection pipelines, which have relatively less flexible text representations. They usually regress text instances in form of horizontal/oriented rectangles or quadrilaterals, having limited ability in detecting irregular texts (*e.g.*, curved texts). TextField does not suffer from this limitation. Benefiting from the proposed direction field, TextField is able to accurately detect texts of irregular shapes.

**TextField Versus Part-based methods:** Part-based methods decompose the text instances into text parts, then attempt to link the neighboring text parts. They enjoy a more flexible representation, and can somehow alleviate the problem of relatively simple text representation inherited in proposal-based methods. Yet, driven by the employed linking or combination strategy, these methods usually produce multi-oriented text detections. The proposed direction field is versatile in representing multi-oriented and curved texts, making TextField perform equally well in detecting any irregular texts.

**TextField Versus Segmentation-based methods:** Due to the significantly varied sizes and aspect ratios, most segmentation-based methods are built upon semantic segmentation, followed by a heavy post-processing step to separate the predicted text areas into text instances. In addition to text mask, some information such as text border, text line, text box orientation, or linking relationship between neighboring pixels is also predicted to ease the separation of adjacent texts. Yet, such additional information either limits the method to multi-oriented text detection or also faces similar problem with text semantic segmentation in separating adjacent texts. TextField directly regresses the direction field which encodes both text mask and direction information that points away from text boundary, thus allowing simple separation of adjacent texts. In this sense, TextField is more elegant and efficient in detecting irregular texts.

It is worth to mention that direction information has also been diversely exploited in some other applications [58]–[60], which involve different definitions or usages.
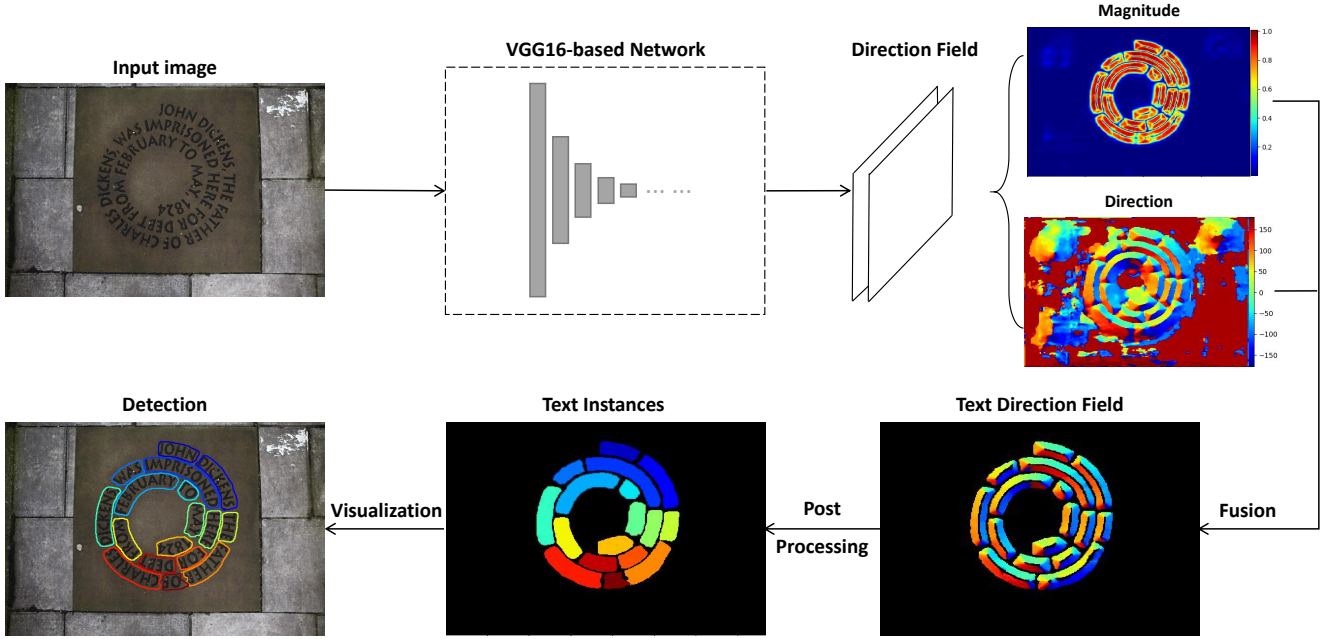
Fig. 3: Pipeline of the proposed method. Given a test image, the network predicts a novel direction field in terms of a two-channel map, which can be regarded as an image of two-dimensional vectors. To better show the predicted direction field, we calculate and visualize its magnitude and direction information. Text instances are then obtained based on these information via the proposed post-processing using some morphological tools.



Fig. 4: Illustration of the proposed direction field. Given a training image and its text annotation, a binary text mask can be easily generated. For each text pixel $p$, we find its nearest non-text pixel $N_p$. Then, a two-dimensional unit vector that points away from $N_p$ to $p$ is defined as the direction field on $p$. For non-text pixels, the direction field is set to $(0,0)$. On the right, we visualize the direction information of the text direction field.

## III. PROPOSED METHODOLOGY

### A. Overview

The proposed method relies on a fully convolutional neural network to produce a dense per-pixel direction field for detecting irregular texts. The pipeline is depicted in Fig. 3. In general, we regard the text detection problem as text instance segmentation. For that, we propose a novel direction field, aiming at segmenting texts and also separating adjacent text instances. More specifically, for a text pixel $p$, its direction field is represented by a two-dimensional unit vector that points away from its nearest text boundary pixel. This direction field is detailed in Section III-B. Benefiting from such novel representation, the proposed TextField can easily separate text instances that lie close to each other. Furthermore, such direction field is appropriate for describing text of arbitrary shapes. We adopt a VGG16-based network to learn the direction field. To preserve spatial resolution and take full advantage of multi-level information, we exploit a widely used multi-level feature fusion strategy. The network architecture is presented
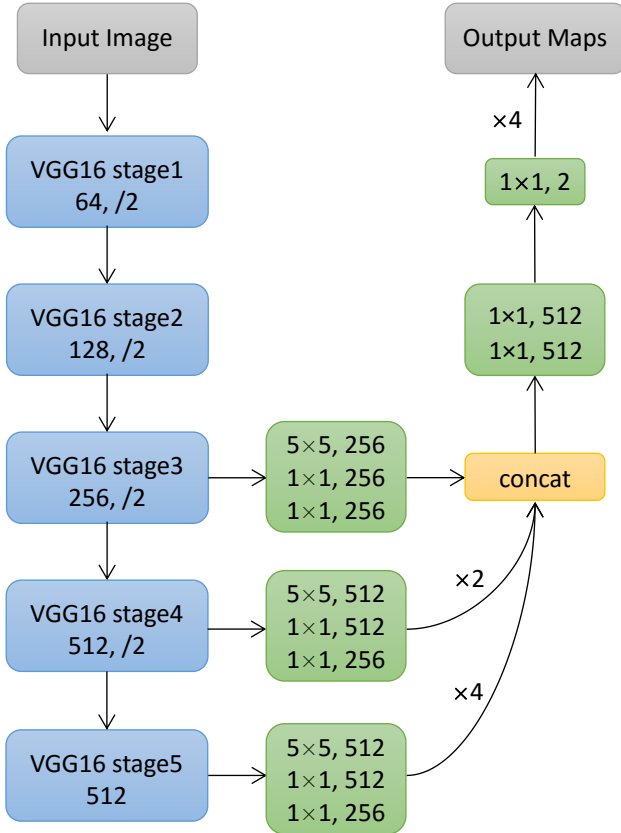
Fig. 5: Network architecture. We adopt the pre-trained VGG16 [61] as the backbone network and multi-level feature fusion to capture multi-scale text instances. The network is trained to predict dense per-pixel direction field.

for each pixel $p$ inside a text instance $T$, let $N_p$ be the nearest pixel to $p$ lying outside the text instance $T$, we then define a two-dimensional unit vector $V_{gt}(p)$ that points away from $N_p$ to the underlying text pixel $p$. This unit vector $V_{gt}(p)$ directly encodes approximately relative location of $p$ inside $T$ and highlights the boundary between adjacent text instances. For the non-text area, we represent those pixels with $(0,0)$. Formally, the proposed direction field is given by:

$$V_{gt}(p) = \begin{cases} \overrightarrow{N_p p}/\left|\overrightarrow{N_p p}\right|, & p \in \mathbb{T} \\ (0,0), & p \notin \mathbb{T} \end{cases} \quad (1)$$

where $\left|\overrightarrow{N_p p}\right|$ denotes length of the vector starting from pixel $N_p$ to $p$, and $\mathbb{T}$ stands for all the text instances in an image. In practice, for each text pixel $p$, it is simple to compute its nearest pixel $N_p$ outside the text instance containing $p$ by distance transform algorithm. Consequently, it is rather straightforward to transform a traditional text annotation to the proposed direction field.

The proposed direction field given by Eq. (1) is appropriate for detecting irregular texts. In fact, the magnitude of direction field $V$ is equivalent to binary text mask. Thus, we rely on magnitude of $V$ to differentiate text and non-text pixels. The direction information encoded in $V$ facilitates the separation of adjacent text instances (see Sec. III-E).

*C. Network architecture*

The proposed network architecture to learn the direction field for detecting irregular texts is depicted in Fig. 5. We adopt a fully convolutional neural network which mainly consists of two parts: feature extraction and multi-level feature fusion. The backbone network to extract features is the VGG16 network [61] pre-trained on ImageNet [62]. We discard the last pooling layer and its following fully connected layers. Since text sizes may vary significantly, it is difficult to detect small text instances with only coarse features. Therefore, we merge features from different stages to capture multi-scale text instances. More specifically, we exploit the feature maps from $stage3$, $stage4$, and $stage5$ of the VGG16 backbone network. These multi-level features are upsampled to the same size as the feature map from $stage3$, and are then merged together by concatenation. This is followed by three convolution layers, resulting in a two-channel map that predicts the direction field given by Eq. (1). Finally, we upsample the predicted direction field to the original size. We adopt bilinear interpolation for all the upsampling operations.

It is worth to note that the proposed method is not severely bottlenecked by the limited receptive field. In fact, the proposed direction field in Eq. (1) only relies on local clues (*i.e.*, location of the nearest text boundary). Thus, we only require a receptive field that covers the short side of text instance. Whereas, for the classical proposal-based methods, a receptive field larger than the long side of underlying text instance is usually needed. Consequently, the proposed method is more flexible in detecting irregular long texts.

in Section III-C. Some specific adaptions for the network training are given in Section III-D, including online hard negative mining and a weighted loss function for our per-pixel regression task. Both adaptions are dedicated to force our network to focus more on hard pixels and eliminate the effects caused by quantitative imbalance between foreground and background pixels. Finally, a novel post-processing based on mathematical tools (see Section III-E) is proposed to group pixels, forming detected text instances thanks to the predicted text direction field.

*B. Direction field*

As pointed out in Sec. II-B, though proposal-based and part-based text detectors have achieved impressive performances on multi-oriented text detection, they do not perform well for curved texts. Segmentation-based approaches can some-how tackle this limitation via binary text mask (of arbitrary shapes) segmentation. Yet they can hardly separate adjacent text instances. To address these issues, we propose a novel direction field for detecting irregular texts.

Instead of binary text mask involved in the segmentation-based approaches, we propose the direction field that encodes both binary text mask and direction information that can be used to separate adjacent text instances. As illustrated in Fig. 4,
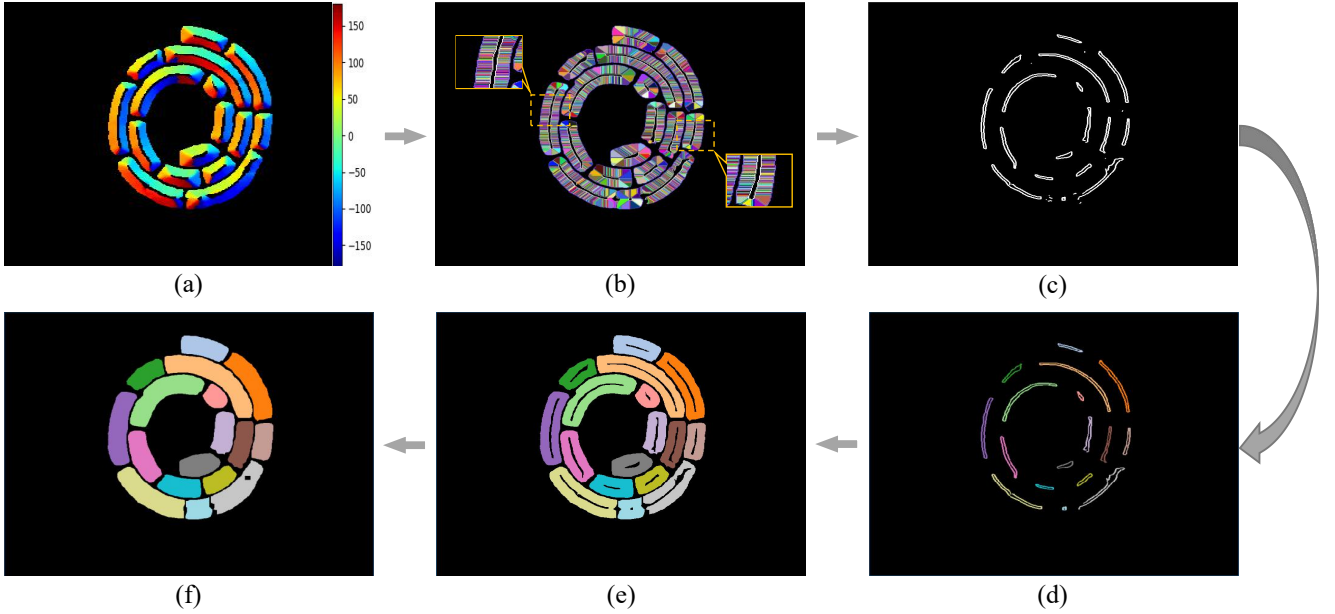
Fig. 6: Illustration of the proposed post-processing on a test image. (a): Directions on candidate text pixels; (b): Text superpixels (in different color) and their representatives (in white); (c): Dilated and grouped representatives of text superpixels; (d): Labels of filtered representatives; (e): Candidate text instances; (f) Final segmented text instances.

## D. Optimization

*1) Training objective:* We leverage the network depicted in Section III-C to regress the proposed direction field. The network parameters are optimized with an instance-balanced Euclidean loss. More specifically, the loss function to minimize is a weighted sum of the mean squared error on each pixel of the image domain $\Omega$. This is given by:

$$L = \sum_{p \in \Omega} w(p) * \|V_{gt}(p) - V_{pred}(p)\|_2, \qquad (2)$$

where $V_{pred}$ is the predicted direction field, and $w(p)$ denotes the weight coefficient of pixel $p$. Since text sizes may vary significantly in scene images, if all text pixels contribute equally to the loss function, large text instances will be dominative in the loss computation while small ones will be ignored. To tackle this problem, we adopt an instance-balanced strategy. More precisely, for an image containing $N$ text instances, the weight $w$ for a given pixel $p$ is defined as follows:

$$w(p) = \begin{cases} \dfrac{\sum_{T \in \mathbb{T}} |T|}{N * |T_p|}, & p \in \mathbb{T} \\ \\ 1, & p \notin \mathbb{T} \end{cases} \qquad (3)$$

where $|T|$ denotes the total number of pixels in a text instance $T$, and $T_p$ stands for the text instance containing pixel $p$. In this way, each text instance of any size is endowed with the same weight, contributing equally to the loss function in Eq. (2). This is consistent with current text detection system such that each text instance is equally important.

*2) Online hard negative mining:* In scene images, text instances usually occupy a small area of the image. Thus, the number of text pixels and non-text pixels is rather imbalanced. To alleviate this problem and to make the network training focus more on pixels which are hard to distinguish, we adopt hard mining following the online hard negative mining strategy proposed in [63]. More specifically, non-text pixels are sorted in a decreasing order of their per-pixel loss. Then only the front $\gamma * (\sum_{T \in \mathbb{T}} |T|)$ non-text pixels are reserved for backpropagation, where $\gamma$ is a given hype-parameter that denotes the ratio of non-text pixels with respect to the total number of text pixels when computing the total loss.

## E. Inference and post-processing

For a given image, the trained network predicts the direction field in terms of 2D vectors. We propose a novel post-processing pipeline using some morphological tools to obtain the final text detection results from this prediction. Precisely, as described in Section III-B, the magnitude of the predicted direction field $V_{pred}$ highlights text/non-text areas. Thus, we first threshold the magnitude image with a thresholding value $\lambda_m$ to obtain candidate text pixels $C$. It is worth to note that pixels lying around text symmetrical axis usually have low magnitude due to the cancellation of opposite direction in learning and upsampling. The text detection problem then amounts to group candidate text pixels into text instances. For that, we first segment the candidate text areas into text superpixels (*i.e.*, text parts depicted in different color in Fig. 6(b)), which are then grouped together to form candidate text instances. A last text instance filtering step is adopted to yield the final detected texts. This process is depicted in Fig. 6 and Algorithm 1, and summarized in the following:

**Text superpixel segmentation:** The magnitude information of the predicted direction field $V_{pred}$ is used to generate candidate text pixels. Then we rely on the direction information carried by $V_{pred}$ to segment the candidate text areas into text superpixels. Precisely, for each candidate text pixel $p$, the

direction information $\angle V_{pred}(p)$ is binned into one of the 8 directions, pointing to its nearest neighboring candidate text pixel denoted as $\mathcal{P}(p)$, standing for the parent of pixel $p$. Each candidate text pixel points to a unique neighboring pixel. Consequently, the parent image $\mathcal{P}$ forms a forest structure $\mathcal{F}$, partitioning the candidate text areas into text superpixels, each of which is represented by a tree $\mathcal{T} \in \mathcal{F}$. This text superpixel segmentation can be efficiently achieved by blob labeling algorithm (see line 7-15 in Algorithm 1).

**Text superpixel grouping:** Based on the segmented text superpixels represented by trees, we propose a simple grouping method to form candidate text instances. Since the proposed direction field encodes the direction away from the nearest boundary, the root pixels of all trees locate near the symmetry axis of each text instance. We consider all these root pixels as the representatives of all the text superpixels. The representatives of a text instance usually are close to each other (See Fig. 6). We apply a simple dilation $\delta$ (with $3 \times 3$ structuring element) to group the representatives of the same text instance. This is followed by a connected component labeling that forms candidate text instances. The text superpixel grouping is depicted in line 17-21 of Algorithm 1.

**Text instance filtering:** After the extraction of candidate text instances, we apply some filtering strategies to get rid of some non-text instances following their shapes and sizes. As illustrated in Fig. 6, the representative pixels of a text instance should have a symmetrical distribution of directions. Therefore, all the representative pixels of a text instance should be approximately paired in the sense of having opposite directions. Based on this observation, we count the ratio of non-paired representatives, and filter out the candidate text instances having a ratio lower than a given value $\lambda_r$ (set to 0.6). For the remaining candidate text instances, we apply a morphological closing $\phi$ (with $11 \times 11$ structuring element) to fill the inside holes. Then we also discard some noisy candidate instances whose areas are smaller than $\lambda_a$ (set to 200). The remaining candidate text instances are the final detected texts. The text instance filtering is given in line 23-27 in Algorithm 1.

Specifically, the proposed post-processing is detailed in Algorithm 1. The core body of the algorithm is the blob labeling to construct text superpixels via the forest structure encoded by the parenthood image $\mathcal{P}$. This blob labeling process can be efficiently implemented using a stack data structure $S$ and an auxiliary image $visited$. The text superpixels are labeled by the image $\mathcal{L}$. Then we identify the representative pixels $R$ by root pixels of those trees. These representative pixels are also stored by an image $\mathcal{M}$. We then apply a dilation $\delta$ with kernel $k_1 \times k_1$ ($k_1 = 3$) to group representative pixels, followed by a connected labeling $CC\_Labeling$ to form candidate text instances. We then filter out some candidate text instances by the ratio of non-paired representatives $Filter\_Unbalanced\_Text$. The label of each remaining candidate text instance is then propagated to all the pixels inside the same text superpixels. Finally, we apply a closing $\phi$ with kernel $k_2 \times k_2$ ($k_2 = 11$) to fill the holes inside each candidate text instance, followed by a removal of small candidate text instances. This post-processing gives the final detected texts encoded by $\mathcal{M}$.

---

**Algorithm 1:** Text inference with a morphological post-processing on predicted direction field $V_{pred}$. $\mathcal{M}$ is the final text instance segmentation map. See the corresponding texts in Section III-E for details.

---

**1** $Text\_Inference(V_{pred}, \lambda_m, \lambda_r, \lambda_a)$
**2** $\mathcal{M}, \mathcal{L} \leftarrow 0$, $l \leftarrow 0$, $C, R, S \leftarrow \emptyset$, $visited \leftarrow$ **False**,
　　$\mathcal{P} \leftarrow p_0$ //initialization ;
**3** *//get candidate text pixels*
**4** **foreach** $p \in \Omega$ **do**
**5** 　 **if** $|V_{pred}(p)| \geq \lambda_m$ **then** $C \leftarrow C \cup p$ ;
**6** *//blob lableing to construct trees encoded by $\mathcal{P}$*
**7** **foreach** $p \in C$ **and not** $visited(p)$ **do**
**8** 　 $S.push(p)$, $l \leftarrow l + 1$ ;
**9** 　 **while** $S \neq \emptyset$ **do**
**10** 　　 $p' \leftarrow S.pop()$, $visited(p') \leftarrow$ **True**, $\mathcal{L}(p') \leftarrow l$ ;
**11** 　　 $\mathcal{P}(p') \leftarrow \mathcal{N}_{\angle V_{pred}(p')}(p')$ ;
**12** 　　 **foreach** $q \in \mathcal{N}(p')$ **do**
**13** 　　　 **if** $q \in C$ **and not** $visited(q)$ **then**
**14** 　　　　 **if** $q = \mathcal{N}_{\angle V_{pred}(p')}(p')$ **or**
　　　　　　 $p' = \mathcal{N}_{\angle V_{pred}(q)}(q)$ **then**
**15** 　　　　　 $S.push(q)$ ;
**16** *//grouping text superpixels via their representatives*
**17** **foreach** $p \in C$ **do**
**18** 　 **if** $\mathcal{P}(p) = p_0$ **then**
**19** 　　 $R \leftarrow R \cup p$, $\mathcal{M}(p) \leftarrow 1$ ;
**20** $\mathcal{M} \leftarrow \delta_{k_1}(\mathcal{M})$ ;
**21** $\mathcal{M} \leftarrow CC\_Labeling(\mathcal{M})$ ;
**22** *//text instance filtering by the shape and size*
**23** $\mathcal{M} \leftarrow Filter\_unblanced\_Text(\mathcal{M}, R, \lambda_r)$ ;
**24** **foreach** $r \in R$ **do**
**25** 　 $\mathcal{M} \leftarrow Propagate\_Label(\mathcal{M}, \mathcal{L}, r)$ ;
**26** $\mathcal{M} \leftarrow \phi_{k_2}(\mathcal{M})$ ;
**27** $\mathcal{M} \leftarrow Filter\_Small\_Regions(\mathcal{M}, \lambda_a)$ ;
**28** **return** $\mathcal{M}$ ;

---

## IV. Experiments

The proposed method is appropriate for detecting irregular texts. In the following, we evaluate the proposed method on four public benchmark datasets: SCUT-CTW1500 [49] and Total-Text [41] which contain curved texts, ICDAR2015 Incidental Scene Text (IC15) [64] and MSRA-TD500 [65] which mainly consist of multi-oriented texts in terms of oriented rectangles or general quadrilaterals. SynthText in the Wild [66] is also adopted to pre-train the network. A short description of these datasets and adopted evaluation protocol is given in Section IV-A. Some implementation details are depicted in Section IV-B, followed by curved text detection results in Section IV-C. The experimental results on multi-oriented text detection is given in Section IV-D to demonstrate the versatility of the proposed TextField. To further demonstrate the generality of TextField, cross dataset experiments are also presented in Section IV-E. The runtime analysis and some failures cases are given in Section IV-F and Section IV-G, respectively.
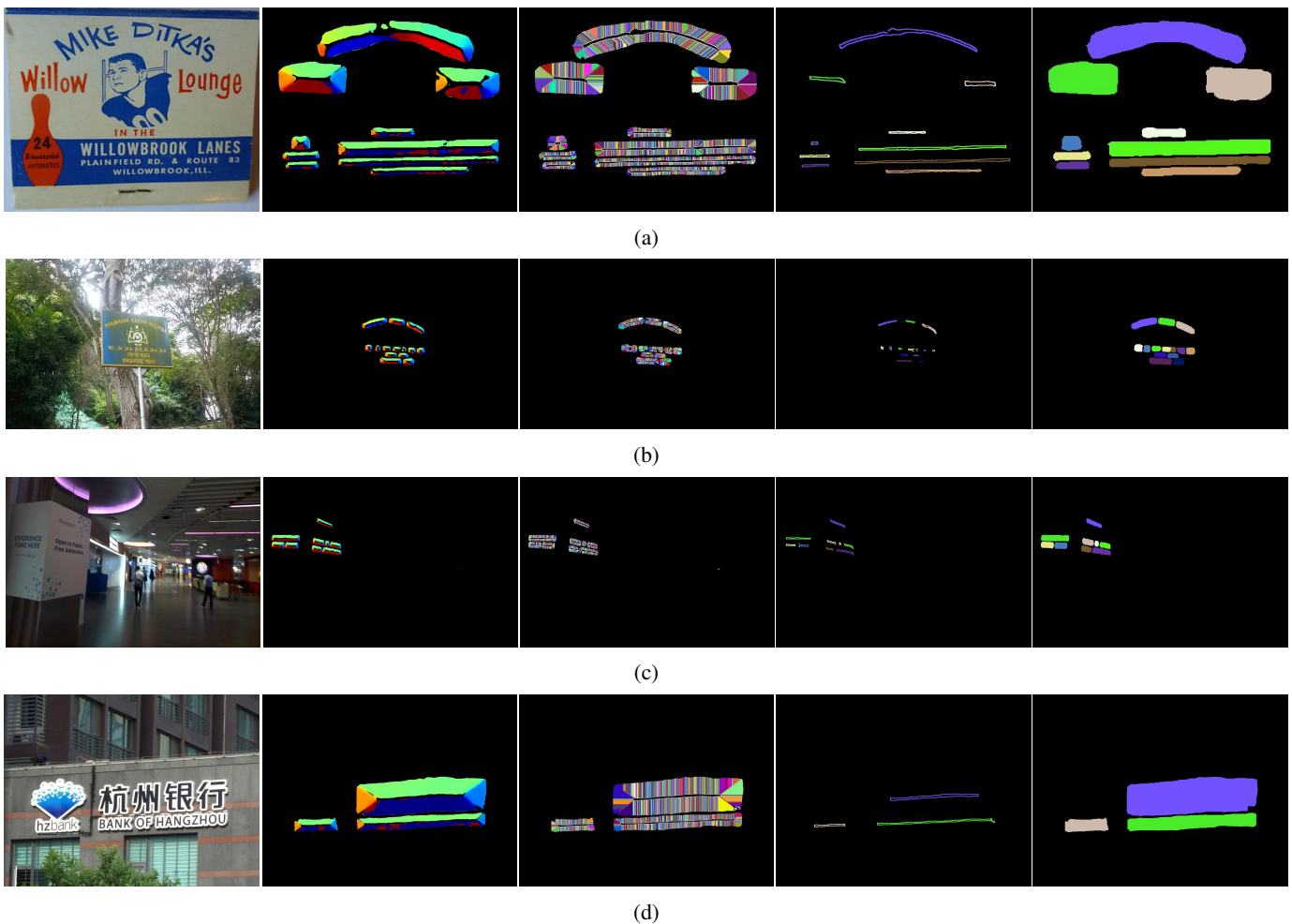
Fig. 7: Visualization of learned direction field and some involved post-processing steps on test images from SCUT-CTW500 in (a), Total-Text in (b), IC15 in (c), and MSRA-TD500 in (d), respectively. From left to right: input images, directions on candidate text pixels, text superpixels (in different color) and their representatives (in white), labels (in different color) of filtered representatives, and final segmented instances.

## A. Datasets and evaluation protocol

**SynthText in the Wild** [66]: SynthText contains 800k synthetic images generated by blending natural images with artificial text. Annotations are given in character, word, and line level. This dataset with word level annotation is used to pre-train the proposed model.

**SCUT-CTW1500** [49]: Different from classical multi-oriented text datasets, this dataset is quite challenging due to many curved texts. It consists of 1000 training images and 500 testing images. This dataset has more than 10k text annotations and at least one curved text per image. Each text instance is labeled by a polygon with 14 points. The annotation is given in line or curve level.

**Total-Text** [41]: Total-Text dataset also aims at solving the arbitrary-shaped text detection problem. It contains 1555 scene images, divided into 1255 training images and 300 testing images. This dataset contains many curved and multi-oriented texts. Annotations are given in word level with polygon-shaped bounding boxes instead of conventional rectangular bounding boxes.

**ICDAR2015 Incidental Scene Text (IC15)** [64]: This dataset is widely used to benchmark multi-oriented text detectors. It was released for the Challenge 4 of ICDAR2015 Robust Reading Competition. Different from previous datasets with text captured in relatively high resolution, scene images in this dataset are taken by Google Glasses in an incidental manner. Therefore, text in these images is of various scales, orientations, contrast, blurring, and viewpoint, making it challenging for detection. This dataset is composed of 1000 training images and 500 testing images. Annotations are provided with word-level bounding quadrilaterals.

**MSRA-TD500** [65]: This dataset is dedicated for detecting multi-lingual long texts of arbitrary orientations. It consists of 300 training images and 200 testing images, annotated at the level of text lines. Since the number of training images is rather small, similar with other methods, we also utilize the images from HUST-TR400 [67] as extra training data.

**Evaluation protocol**: We follow the standard evaluation protocol relying on $precision$, $recall$, and $f\text{-}measure$. Precisely,

they are defined as following:

$$precision = \frac{TP}{TP + FP},$$
$$Recall = \frac{TP}{TP + FN}, \quad (4)$$
$$f\text{-}measure = 2 \times \frac{precision \times recall}{precision + recall},$$

where $TP$, $FP$, and $FN$ stands for the number of correctly detected text instances, incorrect detections, and missing text instances, respectively. For a detected text instance $T$, if $T$ intersects a ground truth text instance with an IOU larger than a given thresholding value (typically set to 0.5), then the text instance $T$ is considered as a correct detection. Since there is a trade-off between $recall$ and $precision$, $f\text{-}measure$ is a common compromised measurement for performance assessment.

### B. Implementation Details

Data augmentation strategy is adopted to increase the training data and avoid over-fitting. Specifically, images are first randomly cropped with area ratios ranging from 0.1 to 1 and aspect ratios ranging from 0.3 to 3. The cropped image is then randomly rotated with 0 and $\pm 90$ degrees. Note that the randomly cropped patch is selected only when the proportion of contained texts with respect to all ground truth text areas in the original images is larger than a threshold value, randomly set to 0.1, 0.3, 0.5, and 0.7. Finally, the augmented images are resized to $384 \times 384$ or $768 \times 768$ during different training stages detailed in the following.

The proposed network is pre-trained on SynthText for one epoch, and then finetuned on SCUT-CTW1500, Total-Text, ICDAR2015 Incidental Scene Text, and MSRA-TD500, respectively. The training process is divided into three stages. In the pre-training stage, the augmented images are resized to $384 \times 384$ for the sake of training speed. The learning rate and the hyper-parameter $\gamma$ involved in online hard negative mining are set to $10^{-4}$ and 3, respectively. Then we finetune our model on each dataset for about 100 epochs with the same settings as pre-training stage. We continue to train the network for another 100 epochs by resizing the augmented images to $768 \times 768$ aiming at better handling multi-scale texts. In this last training stage, the learning rate is decayed to $10^{-5}$ and $\gamma$ is set to 6. In the whole training process, we adopt Adam [68] to optimize the network. All the experiments are conducted on Caffe [69] using a workstation with a single Titan Xp GPU.

### C. Curved text detection

The proposed TextField is appropriate to detect irregular texts. We first conduct experiments on two curved text datasets: SCUT-CTW1500 and Total-Text.

**SCUT-CTW1500**: This dataset mainly contains curved and multi-oriented texts. For each image, the annotation is given in line or curve level. The size of testing image is rather small. In testing phase, the images are resized to $576 \times 576$. The threshold parameter $\lambda_m$ is set to 0.59 for post-processing. A visualization example of the learned direction field and

TABLE I: Quantitative results of different methods evaluated on SCUT-CTW1500. * indicates the result obtained from [49].

| Methods | recall | precision | f-measure |
|---|---|---|---|
| SegLink * [34] | 0.400 | 0.423 | 0.408 |
| CTPN * [33] | 0.538 | 0.604 | 0.569 |
| EAST * [25] | 0.491 | 0.787 | 0.604 |
| DMPNet * [23] | 0.560 | 0.699 | 0.622 |
| CTD [49] | 0.652 | 0.743 | 0.695 |
| CTD+TLOC [49] | 0.698 | 0.774 | 0.734 |
| TextField (Ours) | **0.798** | **0.830** | **0.814** |

TABLE II: Quantitative results of different methods evaluated on Total-Text.

| Methods | recall | precision | f-measure |
|---|---|---|---|
| Ch'ng *et al.* [41] | 0.400 | 0.330 | 0.360 |
| Liao *et al.* [24] | 0.455 | 0.621 | 0.525 |
| TextField (Ours) | **0.799** | **0.812** | **0.806** |

some involved post-processing steps is depicted in Fig. 7(a). Some qualitative results are given in Fig. 8(a). The proposed TextField correctly detects text of arbitrary shapes with very accurate text boundaries. The quantitative results are shown in Tab. I. Compared with other state-of-the-art methods, our proposed method outperforms them by a large margin in terms of recall, precision, and f-measure. The proposed TextField achieves 81.4% F-measure, improving the state-of-the-art methods by 8.0%.

**Total-Text**: We also evaluate the proposed TextField on Total-Text whose annotations are given in word level. This dataset mainly contains curved and multi-oriented texts. In testing, all images are resized to $768 \times 768$. The threshold parameter $\lambda_m$ is set to 0.50 for post-processing. A visualization example of the learned direction field and some involved post-processing steps is illustrated in Fig. 7(b). Some qualitative results are depicted in Fig. 8(b). From this figure, we can observe that TextField also precisely detects word level irregular texts. And TextField is able to accurately separate close text instances of arbitrary shapes. The quantitative results are given in Tab. II. The proposed TextField achieves 80.6% F-measure on this dataset, significantly outperforming other methods.

From the qualitative results depicted in Fig. 8(a-b) and quantitative results given in Tab. I and Tab. II, the proposed TextField is able to detect irregular texts in both line-level and word-level. TextField establishes new state-of-the-art results in detecting curved texts.

### D. Multi-oriented text detection

As shown in Section IV-C, the proposed TextField significantly outperforms other methods on curved text detection. To further demonstrate the ability of TextField in detecting texts of arbitrary shapes, we evaluate TextField on ICDAR2015 Incidental Scene Text and MSRA-TD500 dataset, showing that TextField also achieves very competitive results on widely adopted multi-oriented datasets. Note that for these two experiments, we fit each text instance achieved with TextField by a minimum oriented bounding rectangle.

**ICDAR2015 Incidental Scene Text**: Images in this dataset are of low resolution and contain many small text instances.

TABLE III: Comparison of methods on ICDAR2015 Incidental Scene Text. $^\dagger$ means that the base net of the model is not VGG16. $^*$ stands for multi-scale version.

| Methods | recall | precision | f-measure | FPS |
|---|---|---|---|---|
| Zhang et al. [37] | 0.430 | 0.708 | 0.536 | 0.48 |
| CTPN [33] | 0.516 | 0.742 | 0.609 | 7.1 |
| Yao et al. [38] | 0.587 | 0.723 | 0.648 | 1.61 |
| DMPNet [23] | 0.682 | 0.732 | 0.706 | - |
| SegLink [34] | 0.768 | 0.731 | 0.750 | - |
| MCN [36] | 0.800 | 0.720 | 0.760 | - |
| EAST [25] | 0.728 | 0.805 | 0.764 | 6.52 |
| SSTD [26] | 0.730 | 0.800 | 0.770 | 7.7 |
| RRPN [29] | 0.730 | 0.820 | 0.770 | - |
| ITN [32] | 0.741 | 0.857 | 0.795 | - |
| EAST $^\dagger$ [25] | 0.735 | 0.836 | 0.782 | 13.2 |
| Lyu et al. [35] | 0.707 | **0.941** | 0.807 | 3.6 |
| TextBoxes++ [30] | 0.767 | 0.872 | 0.817 | 11.6 |
| RRD [31] | 0.790 | 0.856 | 0.822 | 6.5 |
| TextField (Ours) | **0.805** | 0.843 | **0.824** | 6.0 |
| WordSup $^*$ [28] | 0.770 | 0.793 | 0.782 | 2 |
| EAST $^{\dagger *}$ [25] | 0.783 | 0.833 | 0.807 | - |
| He et al. $^{\dagger *}$ [27] | 0.800 | 0.820 | 0.810 | 1.1 |
| TextBoxes++ $^*$ [30] | 0.785 | 0.878 | 0.829 | 2.3 |
| Lyu et al. $^*$ [35] | 0.797 | **0.895** | **0.843** | 1 |
| TextField $^*$ (Ours) | **0.839** | 0.843 | 0.841 | 1.8 |

TABLE IV: Comparison of methods on MSRA-TD500. $^\dagger$ stands for the base net of the model is not VGG16.

| Methods | recall | precision | f-measure |
|---|---|---|---|
| He et al. [70] | 0.610 | 0.760 | 0.690 |
| EAST [25] | 0.616 | 0.817 | 0.702 |
| ITN [32] | 0.656 | 0.803 | 0.722 |
| Zhang et al. [37] | 0.670 | 0.830 | 0.740 |
| RRPN [29] | 0.680 | 0.820 | 0.740 |
| He et al. $^\dagger$ [27] | 0.700 | 0.770 | 0.740 |
| Yao et al. [38] | 0.753 | 0.765 | 0.759 |
| EAST $^\dagger$ [25] | 0.674 | 0.873 | 0.761 |
| Wu et al. [39] | 0.780 | 0.770 | 0.770 |
| SegLink [34] | 0.700 | 0.860 | 0.770 |
| RRD [31] | 0.730 | 0.870 | 0.790 |
| Lyu et al. [35] | 0.762 | 0.876 | 0.815 |
| MCN [36] | **0.790** | **0.880** | **0.830** |
| TextField (Ours) | 0.759 | 0.874 | 0.813 |

TABLE V: Cross-dataset evaluations of different methods on corresponding word-level and line-level datasets.

| Methods | Total-Text (train on IC15) | | |
|---|---|---|---|
| | recall | precision | f-measure |
| SegLink [34] | 0.332 | 0.356 | 0.344 |
| EAST [25] | 0.431 | 0.490 | 0.459 |
| TextField (Ours) | **0.652** | **0.615** | **0.633** |
| Methods | IC15 (train on Total-Text) | | |
| | recall | precision | f-measure |
| TextField (Ours) | **0.660** | **0.771** | **0.711** |
| Methods | SCUT-CTW1500 (train on TD500) | | |
| | recall | precision | f-measure |
| TextField (Ours) | **0.700** | **0.753** | **0.726** |
| Methods | MSRA-TD500 (train on SCUT-CTW1500) | | |
| | recall | precision | f-measure |
| TextField (Ours) | **0.758** | **0.853** | **0.803** |

Therefore, images are not resized. The original resolution of $1280 \times 720$ is used in testing. The threshold parameter $\lambda_m$ is set to 0.69 for post-processing. A visualization example of the learned direction field and some involved post-processing steps is shown in Fig. 7(c). Some detection results on this dataset are given in Fig. 8(c), where challenging texts of variant contrast and scales are correctly detected. The quantitative evaluation compared with other methods are depicted in Tab. III. The proposed TextField achieves competitive results with other state-of-the-art methods on this dataset. Following [25], [27], [30], [35], we also report the results of TextField under multi-scale evluation using $384 \times 384$, $768 \times 768$, and $1024 \times 1024$ inputs on IC15. TextField is also very competitive with other methods under multi-scale evaluation. Note that for fair comparison, we mainly compare with other methods using the same backbone network (*i.e.*, VGG16 network).

**MSRA-TD500**: This dataset contains both English and Chinese texts whose annotations are given in terms of text lines. The text scale varies significantly. In testing, we resize the images into $768 \times 768$. The threshold parameter $\lambda_m$ is set to 0.64 for post-processing. Due to the large character spacing in this dataset, we also group the detected texts with small aspect ratios before evaluating the TextField using the IC15 evaluation code. A visualization example of the learned direction field and some involved post-processing steps is depicted in Fig. 7(d). Some qualitative illustrations are shown in Fig. 8(d). The proposed TextField successfully detects long text lines of arbitrary orientations and sizes. The quantitative comparison with other methods on this dataset is given in Tab. IV. TextField also achieves competitive performance with other methods in detecting long multi-oriented texts. Specifically, TextField performs slightly worse than the methods in [35] and [36] on MSRA-TD500. Yet, the performance of TextField is much better than them on IC15 dataset.

From the qualitative results in Fig. 8(c-d) and quantitative

evaluations in Tab. III and Tab. IV, the proposed TextField is also capable to accurately detect multi-oriented texts in both line-level and word-level. This demonstrates the versatility of the proposed TextField.

### E. Cross dataset text detection

To further demonstrate the generalization ability of the proposed TextField, we also evaluate the TextField trained on one dataset and test the trained model on a different dataset annotated in the same level (*e.g.*, word or line). Specifically, we first benchmark several classical models (trained on IC15) on Total-Text dataset. As depicted in Tab. V, The proposed TextField generalizes better on cross-dataset text detection. We then test the TextField (trained on Total-Text) on IC15 dataset, which gives acceptable results. We have also performed cross-dataset evaluations on two line-level annotated datasets: SCUT-CTW1500 and MSRA-TD500. As shown in Tab. V, for the line-based text detection, TextField also achieves very competitive results (under cross-dataset setting) with some state-of-the-art methods trained on the target dataset. Specifically, TextField trained on MSRA-TD500 containing only multi-oriented texts performs comparably with other methods properly trained on SCUT-CTW1500, a curved text dataset. Furthermore, it is worth to note that TextField trained on SCUT-CTW1500 containing mainly curved English texts also performs rather well (with a small degradation) in detecting multi-oriented Chinese texts in MSRA-TD500.
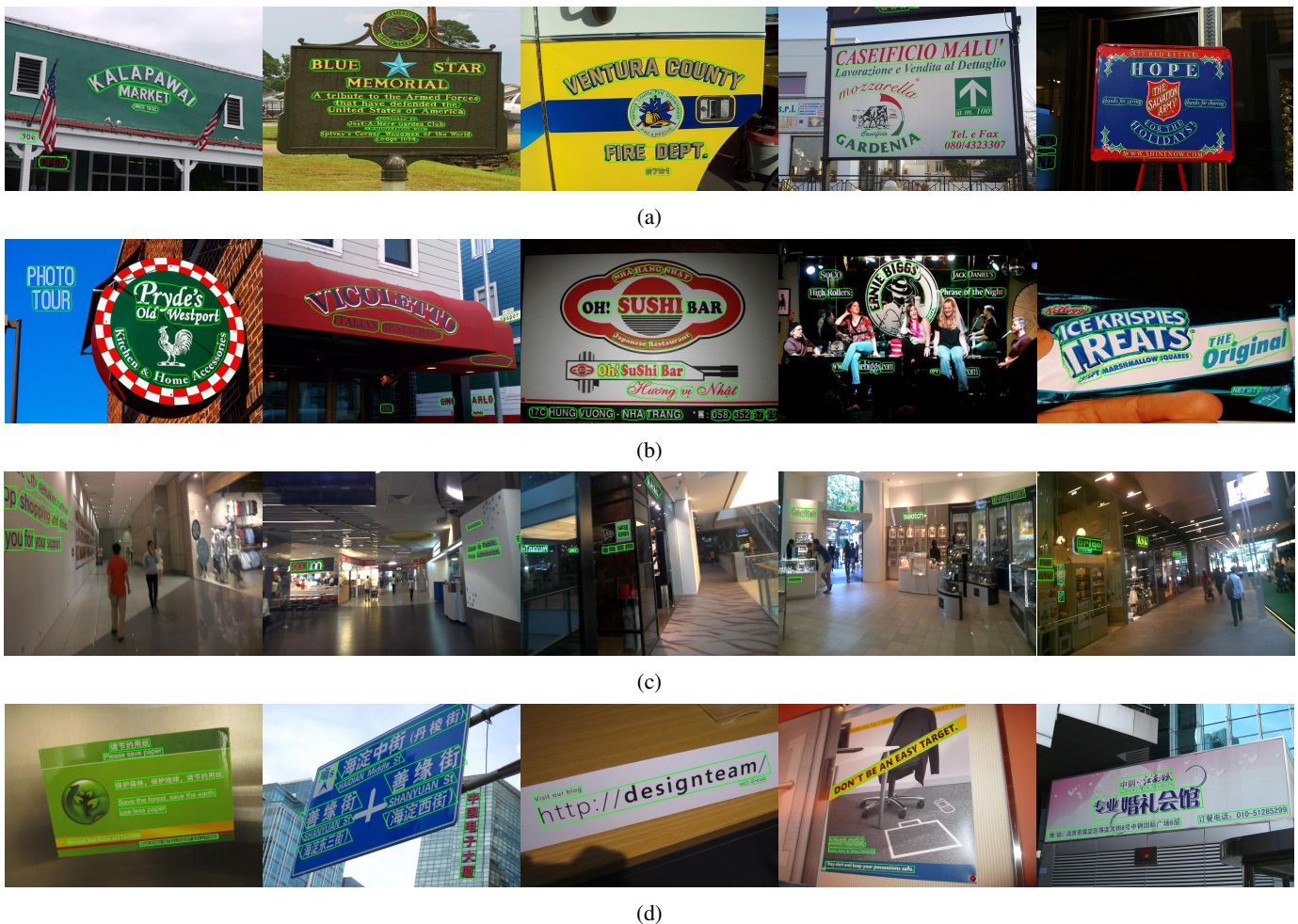
(a)



(b)



(c)



(d)

Fig. 8: Some qualitative detection results on SCUT-CTW500 in (a), Total-Text in (b), IC15 in (c), and MSRA-TD500 in (d). The arbitrary-shaped texts are correctly detected with accurate text instance boundaries.

These cross-dataset experiments demonstrate that the proposed TextField is effective in detecting irregular texts, and is also robust in generalizing to unseen datasets.

*F. Runtime*

The proposed TextField first yields the predicted direction field through the proposed network, then followed by a morphological post-processing step to achieve final text detection results. The runtime of TextField is thus decomposed into two stages: network inference and post-processing. For the network inference, using the VGG16 backbone network as depicted in Fig. 5, it takes about 130ms for a $1280 \times 720$ IC15 image and 100ms for a $768 \times 768$ MSRA-TD500 image on a Titan Xp GPU. As described in Section III-E, the post-processing is mainly composed of three steps: text superpixel segmentation, text superpixel grouping, and text instance filtering. The text superpixel segmentation could be achieved by the blob labeling algorithm which is very fast. The grouping step only involves some classical morphological operations. The text instance filtering step is also very fast thanks to the criterion incrementally computed during the grouping step. The whole post-processing stage takes about 36ms for a $1280 \times 720$ IC15

image and 24ms for a $768 \times 768$ MSRA-TD500 image on a 3.4GHz/8MB cache Intel core i7-2600, 16GB RAM. As depicted in Tab. III, the proposed TextField runs at 6.0 FPS using VGG16 backbone, which is on par with most state-of-the-art methods. Furthermore, TextField is able to accurately detect irregular texts and generalizes well to unseen datasets.

*G. Weakness*

As demonstrated in previous experiments, TextField performs well in most cases of detecting texts of arbitrary shapes. It still fails for some difficult images, such as object occlusion, large character spacing. Some failure examples are given in Fig. 9. TextField also has some false detections on some text-like areas. Note that all these difficulties are common challenges for the other state-of-the-art methods [25], [30], [34].

## V. CONCLUSION

We have presented TextField, which learns a deep direction field for irregular text detection. Specifically, we propose a novel text direction field that points away from nearest text boundary to each text point. Such two-dimensional text

Fig. 9: Some failure examples. Green contours: correct detections; Red contours: missing ground truths; Blue contours: false detections.

direction field encodes both binary text mask and direction information that facilitates the separation of adjacent text instances, which remains challenging for classical segmentation-based approaches. TextField directly regresses the direction field followed by a simple yet effective post-processing step inspired by some morphological tools. Experiments on two curved text datasets (Total-Text and SCUT-CTW1500) and two widely-used datasets (ICDAR 2015 and MSRA-TD500) demonstrate that the proposed method outperforms all state-of-the-art methods by a large margin in detecting curved texts, and achieves very competitive performances in detecting multi-oriented texts. Furthermore, based on the cross-dataset evaluations, TextField also generalizes well to unseen datasets. In the future, we would explore more robust text superpixel grouping strategy (*e.g.,* via explicitly learning the text center line) to further boost TextField, and investigate the common challenges faced by all state-of-the-art text detectors.

## REFERENCES

[1] B. Xiong and K. Grauman, "Text detection in stores using a repetition prior," in *Proc. of IEEE Winter Conference on Applications of Computer Vision*, 2016, pp. 1–9.

[2] C. Yi and Y. Tian, "Scene text recognition in mobile applications by character descriptor and structure configuration," *IEEE Trans. on Image Processing*, vol. 23, no. 7, pp. 2972–2982, 2014.

[3] C. Kang, G. Kim, and S. I. Yoo, "Detection and recognition of text embedded in online images via neural context models." in *Proc. of the AAAI Conf. on Artificial Intelligence*, 2017, pp. 4103–4110.

[4] X. Rong, C. Yi, and Y. Tian, "Recognizing text-based traffic guide panels with cascaded localization network," in *Proc. of European Conference on Computer Vision*, 2016, pp. 109–121.

[5] Y. Zhu, M. Liao, M. Yang, and W. Liu, "Cascaded segmentation-detection networks for text-based traffic sign detection," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 1, pp. 209–219, 2018.

[6] Q. Ye and D. Doermann, "Text detection and recognition in imagery: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 7, pp. 1480–1500, 2015.

[7] K. Wang and S. Belongie, "Word spotting in the wild," in *Proc. of European Conference on Computer Vision*, 2010, pp. 591–604.

[8] Y.-F. Pan, X. Hou, C.-L. Liu *et al.*, "A hybrid approach to detect and localize texts in natural scene images," *IEEE Trans. on Image Processing*, vol. 20, no. 3, pp. 800–813, 2011.

[9] L. Neumann and J. Matas, "Real-time scene text localization and recognition," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2012, pp. 3538–3545.

[10] B. Bai, F. Yin, and C. L. Liu, "Scene text localization using gradient local correlation," in *Proc. of International Conference on Document Analysis and Recognition*, 2013, pp. 1380–1384.

[11] W. Huang, Z. Lin, J. Yang, and J. Wang, "Text localization in natural images using stroke feature transform and text covariance descriptors," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2013, pp. 1241–1248.

[12] W. Huang, Y. Qiao, and X. Tang, "Robust scene text detection with convolution neural network induced mser trees," in *Proc. of European Conference on Computer Vision*, 2014, pp. 497–511.

[13] X.-C. Yin, W.-Y. Pei, J. Zhang, and H.-W. Hao, "Multi-orientation scene text detection with adaptive clustering," *IEEE Trans. Pattern Anal. Mach. Intell.*, no. 9, pp. 1930–1937, 2015.

[14] S. Lu, T. Chen, S. Tian, J.-H. Lim, and C.-L. Tan, "Scene text extraction based on edges and support vector regression," *International Journal on Document Analysis and Recognition*, vol. 18, no. 2, pp. 125–135, 2015.

[15] S. Tian, Y. Pan, C. Huang, S. Lu, K. Yu, and C. Lim Tan, "Text flow: A unified text detection system in natural scene images," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2015, pp. 4651–4659.

[16] J. Matas, O. Chum, M. Urban, and T. Pajdla, "Robust wide-baseline stereo from maximally stable extremal regions," *Image and vision computing*, vol. 22, no. 10, pp. 761–767, 2004.

[17] B. Epshtein, E. Ofek, and Y. Wexler, "Detecting text in natural scenes with stroke width transform," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2010, pp. 2963–2970.

[18] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman, "Reading text in the wild with convolutional neural networks," *International Journal of Computer Vision*, vol. 116, no. 1, pp. 1–20, 2016.

[19] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, no. 6, pp. 1137–1149, 2017.

[20] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Proc. of European Conference on Computer Vision*, 2016, pp. 21–37.

[21] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.

[22] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.

[23] Y. Liu and L. Jin, "Deep matching prior network: Toward tighter multi-oriented text detection," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 3454–3461.

[24] M. Liao, B. Shi, X. Bai, X. Wang, and W. Liu, "Textboxes: A fast text detector with a single deep neural network." in *Proc. of the AAAI Conf. on Artificial Intelligence*, 2017, pp. 4161–4167.

[25] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang, "East: an efficient and accurate scene text detector," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 2642–2651.

[26] P. He, W. Huang, T. He, Q. Zhu, Y. Qiao, and X. Li, "Single shot text detector with regional attention," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2017, pp. 3047–3055.

[27] W. He, X.-Y. Zhang, F. Yin, and C.-L. Liu, "Deep direct regression for multi-oriented scene text detection," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2017, pp. 745–753.

[28] H. Hu, C. Zhang, Y. Luo, Y. Wang, J. Han, and E. Ding, "Wordsup: Exploiting word annotations for character based text detection," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2017, pp. 4950–4959.

[29] J. Ma, W. Shao, H. Ye, L. Wang, H. Wang, Y. Zheng, and X. Xue, "Arbitrary-oriented scene text detection via rotation proposals," *IEEE Transactions on Multimedia*, 2018.

[30] M. Liao, B. Shi, and X. Bai, "Textboxes++: A single-shot oriented scene text detector," *IEEE Trans. on Image Processing*, vol. 27, no. 8, pp. 3676–3690, 2018.

[31] M. Liao, Z. Zhu, B. Shi, G. Xia, and X. Bai, "Rotation-sensitive regression for oriented scene text detection," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 5909–5918.

[32] F. Wang, L. Zhao, X. Li, X. Wang, and D. Tao, "Geometry-aware scene text detection with instance transformation network," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 1381–1389.

[33] Z. Tian, W. Huang, T. He, P. He, and Y. Qiao, "Detecting text in natural image with connectionist text proposal network," in *Proc. of European Conference on Computer Vision*, 2016, pp. 56–72.

[34] B. Shi, X. Bai, and S. Belongie, "Detecting oriented text in natural images by linking segments," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 3482–3490.

[35] P. Lyu, C. Yao, W. Wu, S. Yan, and X. Bai, "Multi-oriented scene text detection via corner localization and region segmentation," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 7553–7563.

[36] Z. Liu, G. Lin, S. Yang, J. Feng, W. Lin, and W. L. Goh, "Learning markov clustering networks for scene text detection," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 6936–6944.

[37] Z. Zhang, C. Zhang, W. Shen, C. Yao, W. Liu, and X. Bai, "Multi-oriented text detection with fully convolutional networks," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2016, pp. 4159–4167.

[38] C. Yao, X. Bai, N. Sang, X. Zhou, S. Zhou, and Z. Cao, "Scene text detection via holistic, multi-channel prediction," *arXiv preprint arXiv:1606.09002*, 2016.

[39] Y. Wu and P. Natarajan, "Self-organized text detection with minimal post-processing via border learning," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2017, pp. 5010–5019.

[40] D. He, X. Yang, C. Liang, Z. Zhou, G. Alexander, I. Ororbia, D. Kifer, and C. L. Giles, "Multi-scale fcn with cascaded instance aware segmentation for arbitrary oriented word spotting in the wild." in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 474–483.

[41] C. K. Ch'ng and C. S. Chan, "Total-text: A comprehensive dataset for scene text detection and recognition," in *Proc. of International Conference on Document Analysis and Recognition*, vol. 1, 2017, pp. 935–942.

[42] S. Tian, S. Lu, and C. Li, "Wetext: Scene text detection under weak supervision," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2017, pp. 1492–1500.

[43] C. Xue, S. Lu, and F. Zhan, "Accurate scene text detection through border semantics awareness and bootstrapping," in *Proc. of European Conference on Computer Vision*, 2018, pp. 355–372.

[44] T. He, Z. Tian, W. Huang, C. Shen, Y. Qiao, and C. Sun, "An end-to-end textspotter with explicit alignment and attention," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 5020–5029.

[45] F. Zhan, S. Lu, and C. Xue, "Verisimilar image synthesis for accurate detection and recognition of texts in scenes," in *Proc. of European Conference on Computer Vision*, 2018, pp. 249–266.

[46] B. Shi, X. Bai, and C. Yao, "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.*, vol. 39, no. 11, pp. 2298–2304, 2017.

[47] B. Shi, M. Yang, X. Wang, P. Lyu, C. Yao, and X. Bai, "Aster: An attentional scene text recognizer with flexible rectification," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2018, to appear, DOI: 10.1109/TPAMI.2018.2848939.

[48] J. Fabrizio, M. Robert-Seidowsky, S. Dubuisson, S. Calarasanu, and R. Boissel, "Textcatcher: a method to detect curved and challenging text in natural scenes," *International Journal on Document Analysis and Recognition*, vol. 19, no. 2, pp. 99–117, 2016.

[49] Y. Liu, L. Jin, S. Zhang, and S. Zhang, "Detecting curve text in the wild: New dataset and new solution," *arXiv preprint arXiv:1712.02170*, 2017.

[50] P. Salembier, A. Oliveras, and L. Garrido, "Antiextensive connected operators for image and sequence processing," *IEEE Trans. on Image Processing*, vol. 7, no. 4, pp. 555–570, 1998.

[51] L. Najman and M. Couprie, "Building the component tree in quasi-linear time," *IEEE Trans. on Image Processing*, vol. 15, no. 11, pp. 3531–3539, 2006.

[52] E. Carlinet and T. Géraud, "A comparative review of component tree computation algorithms," *IEEE Trans. on Image Processing*, vol. 23, no. 9, pp. 3885–3895, 2014.

[53] Y. Wang, Y. Xu, S. Tsogkas, X. Bai, S. Dickinson, and K. Siddiqi, "Deepflux for skeletons in the wild," *arXiv preprint arXiv:1811.12608*, 2018.

[54] Y. Zhu, C. Yao, and X. Bai, "Scene text detection and recognition: Recent advances and future trends," *Frontiers of Computer Science*, vol. 10, no. 1, pp. 19–36, 2016.

[55] L. Gomez and D. Karatzas, "Multi-script text extraction from natural scenes," in *Proc. of International Conference on Document Analysis and Recognition*, 2013, pp. 467–471.

[56] Y. Li, W. Jia, C. Shen, and A. van den Hengel, "Characterness: An indicator of text in the wild," *IEEE Trans. on Image Processing*, vol. 23, no. 4, pp. 1666–1677, 2014.

[57] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proc. of IEEE Intl. Conf. on Computer Vision*, 2015, pp. 1520–1528.

[58] X. Bai, Z. Zhang, H.-Y. Wang, and W. Shen, "Directional edge boxes: Exploiting inner normal direction cues for effective object proposal generation," *Journal of Computer Science and Technology*, vol. 32, no. 4, pp. 701–713, 2017.

[59] M. Bai and R. Urtasun, "Deep watershed transform for instance segmentation," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 2858–2866.

[60] L.-C. Chen, A. Hermans, G. Papandreou, F. Schroff, P. Wang, and H. Adam, "Masklab: Instance segmentation by refining object detection with semantic and direction features," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 4013–4022.

[61] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.

[62] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F.-F. Li, "Imagenet: A large-scale hierarchical image database," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[63] A. Shrivastava, A. Gupta, and R. Girshick, "Training region-based object detectors with online hard example mining," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2016, pp. 761–769.

[64] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu *et al.*, "Icdar 2015 competition on robust reading," in *Proc. of International Conference on Document Analysis and Recognition*, 2015, pp. 1156–1160.

[65] C. Yao, X. Bai, W. Liu, Y. Ma, and Z. Tu, "Detecting texts of arbitrary orientations in natural images," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2012, pp. 1083–1090.

[66] A. Gupta, A. Vedaldi, and A. Zisserman, "Synthetic data for text localisation in natural images," in *Proc. of IEEE Intl. Conf. on Computer Vision and Pattern Recognition*, 2016, pp. 2315–2324.

[67] C. Yao, X. Bai, and W. Liu, "A unified framework for multioriented text detection and recognition," *IEEE Trans. on Image Processing*, vol. 23, no. 11, pp. 4737–4749, 2014.

[68] D. Kinga and J. B. Adam, "A method for stochastic optimization," in *International Conference on Learning Representations (ICLR)*, vol. 5, 2015.

[69] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of ACM international conference on Multimedia*, 2014, pp. 675–678.

[70] T. He, W. Huang, Y. Qiao, and J. Yao, "Text-attentional convolutional neural network for scene text detection," *IEEE Trans. on Image Processing*, vol. 25, no. 6, pp. 2529–2541, 2016.