

Scene-Text-Detection Method Robust against Orientation and Discontiguous Components of Characters

Rei Endo, Yoshihiko Kawai, Hideki Sumiyoshi, Masanori Sano
Science and Technology Research Laboratories, NHK, Tokyo, Japan
{endo.r-mm, kawai.y-1k, sumiyoshi.h-di, sano.m-fo}@nhk.or.jp

Abstract

Scene-text detection in natural-scene images is an important technique because scene texts contain location information such as names of places and buildings, but many difficulties still remain regarding practical use. In this paper, we tackle two problems of scene-text detection. The first is the discontiguous component problem in specific languages that contain characters consisting of discontiguous components. The second is the multi-orientation problem in all languages. To solve these two problems, we propose a connected-component-based scene-text-detection method. Our proposed method involves our novel neighbor-character search method using a synthesizable descriptor for the discontiguous-component problems and our novel region descriptor called the rotated bounding box descriptors (RBBs) for rotated characters. We also evaluated our proposed scene-text-detection method by using the well-known MSRA-TD500 dataset that includes rotated characters with discontiguous components.

1. Introduction

Text in images contains valuable semantic information for various content-based applications [1, 6, 24]. A representative one is location information, and texts of images taken in town often contain clues to identify location. Traffic signs are located in many places in the city, and they include names of streets and places. Buildings and shops have put up signboards with their names at the entrance. By recognizing such texts in images, we can use information of the texts to identify locations. These location information can be used for city navigation and autonomous driving. In addition, since information can be obtained only from image data, it can be used in indoor or underground where GPS can not be used. Therefore, a text recognition method is useful for location identification from images.

Text-recognition methods have been extensively investigated and researchers have proposed effective ones. In particular, an optical-character-recognition method for recognizing printed text in scanned documents is one of the most

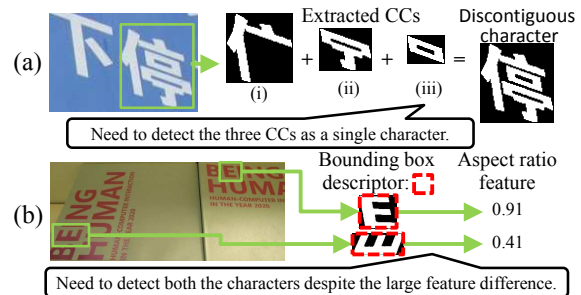


Figure 1. (a) Discontiguous-component and (b) multi-orientation problems. These problems confuse scene-text-detection methods using geometrical features.

successful [11]. However, it is still difficult for computers to detect and recognize *scene text*, which is text in natural-scene images because scene text has various interference factors such as appearance variation (e.g. changes in character size and font), language, orientation, distortion, noise, occlusion, and complex background [11]. In this paper, we tackle two problems caused by interference factors in scene-text detection. One is the discontiguous-component problem and the other is the multi-orientation problem.

The discontiguous-component problem is caused by characters with discontiguous components (hereafter called discontiguous characters). A discontiguous character consists of multiple connected components (CCs), as shown in Fig. 1. Geometrical features, which are important clues for text detection (e.g. size and aspect), calculated from a single CC of a discontiguous character, greatly differ from that of a contiguous character, and the difference easily leads to misdetection. While “i” and “j” are the only discontiguous characters in English, Chinese and Japanese have a large number of discontiguous characters. The multi-orientation problem is caused by various text orientations of scene text. Most general geometrical features of characters are not invariant to the orientation of the characters, as shown in Fig. 1. Therefore, it is necessary to solve the two problems for detecting rotated scene text written in languages such as Chinese and Japanese. Note that *text* means a word or sentence consisting of more than one character.

Many scene-text-detection methods have been proposed,

but most target horizontal English texts. Several methods detect texts, including discontinuous characters, based on the uniformity of text alignment and size [16]. However, if they are used to detect rotated text, including discontinuous characters, a large number of false detections will occur. With the existing methods based on geometrical features affected by rotation, it is difficult to correctly eliminate false detections. However, other methods estimate the text orientation [25, 28]. These methods also use the uniformity of text alignment so that when there are texts containing a discontinuous character, orientation estimation is difficult due to the uncertainty of whether a small region, *i.e.*, discontinuous component, forms a character, a part of a character, or just noise. Furthermore, we found no method for solving the two problems at the same time in our survey.

In this paper, we propose a CC-based scene-text-detection method robust against rotated discontinuous characters. Our method uses our novel neighbor-character search method with synthesizable descriptors to solve the discontinuous-component problem. The synthesizable descriptor is a particular class of descriptors that can be combined using the simple operation which lower calculation cost. Our method calculates a descriptor of a discontinuous character from only those of the individual CCs constituting the character by using synthesizable descriptors. To solve the multi-orientation problem, our method uses our novel region descriptors called the rotated bounding box descriptors (RBBs), which are synthesizable descriptors. RBBs represent coordinates of bounding boxes around regions rotated through multiples of a fixed angles, and are used to accurately calculate geometrical features of the regions. By using our neighbor-character search method with RBBs and other synthesizable descriptors, our method can detect rotated scene text that includes discontinuous characters.

The rest of paper is organized as follows. We give a brief review on related work on CC-based methods in Sec. 2 and describe our proposed scene-text-detection method in Sec. 3. We explain the experimental results and provide discussion in Sec. 4 followed by conclusions in Sec. 5.

2. Related Work

Existing scene-text-detection methods are roughly classified into two types: sliding-window based [3, 9, 10, 19] and CC based [2, 5, 11, 13, 15, 16, 21, 22, 23, 25, 26, 27]. Sliding-window-based methods use a sliding window to search for character or text regions in an image then use machine-learning techniques to identify text. Although these methods can generally detect discontinuous characters, it is difficult for them to handle large variations in character size and aspect. Because sliding-windows detect only characters the same size and aspect as those included in the training data. CC-based methods extract CCs as character candidates and form groups of character candidates with similar features. These methods can handle the variation, but existing CC-based methods do not take into

account of the discontinuous component problem. We focus on CC-based methods to handle various characters and describe related work. Maximally stable extremal region (MSER)-based methods have been particularly successful among CC-based methods.

A large number of MSER-based methods perform well in scene-text detection [5, 15, 16, 25, 27]. Because they can extract extremal regions, which are a kind of CC, that are robust against size variation, low resolution, and complex background[12]. Neumann and Matas used a multi-stage character classifier to efficiently reduce the number of false character candidates [15]. Yin *et al.* reduced the number of CCs until the overlapping region was eliminated by using the aspect ratio and stability of CCs [27]. However, their methods do not take into account the presence of discontinuous characters. Unlike English, languages such as Chinese and Japanese contain many discontinuous characters. In some cases, continuous characters can unexpectedly divided into multiple CCs due to the limitation of CC extraction algorithms. Neumann and Matas’s method can detect such characters using an estimation of a text’s base line [16]. However, it is difficult to estimate a text’s base line when the text is rotated.

Recently, deep neural network (DNN)-based methods, which are hybrid methods of sliding-window-based and CC-based methods, also have been successfully applied in various computer vision tasks [1, 7, 18, 28, 29]. Although DNN-based methods require a huge amount of training data and high computational performance, they can detect characters from various images with high accuracy. Bissacco *et al.* developed a system using a deep-learning method with histogram of gradient features, and the system can detect characters in uncontrolled situations [1]. He *et al.* proposed a text-attentional convolution neural network that can learn a text region mask, character label, and binary text/non-text information in uniform schema [7]. These methods work well in English but are still affected by the discontinuous component problem. One reason is that some languages that include discontinuous characters have many types of characters. Large amounts of training data and computational resources are required to correctly classify these characters. GPU can solve the problem about computational resources, but we are also assuming use with devices that cannot use GPU such as smart-phones. Another reason is that the shapes of some CCs constituting discontinuous characters are often the same as those of contiguous characters (*e.g.* the CC (iii) in Fig. 1 is similar to "O" or zero in English.). There are difficulties in distinguishing them even when a DNN is used.

The multi-orientation problem still needs to be solved. Most existing methods use geometrical features to detect text candidates from detected characters based on the assumption that text orientation is horizontal or near-horizontal. Yin *et al.* proposed an adaptive clustering method that uses the geometrical and orientation features of

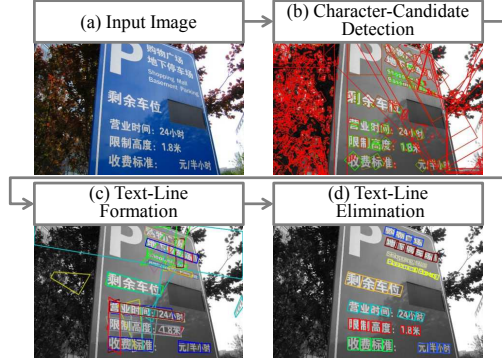


Figure 2. Overview of proposed method and sample results of each processing stage. (a) Input image. (b) Green rectangles represent CCs with high character probability, and red ones represent those with low character probability. (c) Randomly colored rectangles represent text-line candidates. (d) Rectangles represent each single character of detected text lines. Characters with same color rectangles belong to same detected text line.

character candidate pairs to detect rotated texts [25]. Zhang *et al.*'s method detects rotated texts by using geometrical features and information of a text-block, which is detected by text-block fully constitutional networks [28]. However, the geometrical features of existing methods are not accurate because they use a non-rotated bounding box or a circumscribed circle. When an elongated shaped character is oriented, the non-rotated bounding box changes greatly, as shown in Fig. 1. The radius of the circumscribed circle is used as the geometrical distance of a character pair, but the circumscribed circle is affected by character shape. For example, in detecting the text of "Alice," existing methods evaluate the distance values of "Al," "li," "ic", and "ce" under the assumption that the characters belong to the same text when their distance values are similar. However, a circumscribed circle's size of "Al" is much larger than that of "li"; therefore, the detection of "Alice" fails.

3. Proposed Scene-Text-Detection Method

Figure 2 gives an overview of our method. At the beginning of all processing, CCs are extracted from an input image and classified as characters or non-characters. Then, text-line candidates are detected using our neighbor-character search method with RBB and other synthesizable descriptors. Finally, overlapping text-line candidates are throttled into one candidate based on text probabilities. Since the key ideas regarding our method are synthesizable descriptors and RBB, which contribute to our neighbor-character search method, we describe them first then the details of each processing stage. Note that descriptors and features are clearly distinguished in this paper. A descriptor is a storing format for characteristics of a region. A feature is a value or a set of values actually used by the classifiers. For example, a bounding box is a descriptor, and the aspect ratio calculated by the bounding box is a feature.

3.1. Region Descriptors and Features

3.1.1 Synthesizable Descriptors

We propose synthesizable descriptors for efficient computation of compound region descriptors, inspired by the incrementally computable descriptor [15]. The compound region is a region consisting of multiple CCs, where a region of a discontinuous character is always a compound region. By using a synthesizable descriptor, descriptors of a compound region consisting of CCs can be simply calculated from only descriptors of the CCs. Let r_i and r_j be regions consisting of one or more CCs, r_c be the compound region of r_i and r_j , and $\psi(r)$ be a descriptor calculation function of r . When the descriptor is a synthesizable descriptor, there exists \oplus satisfying $\psi(r_c) = \psi(r_i) \oplus \psi(r_j)$, where \oplus denotes an operation that combines the descriptors of r_i and r_j . The proposed scene-text-detection method calculates descriptors of compound regions using a synthesizable descriptor to reduce calculation cost.

The following synthesizable descriptors are used with our proposed method for compound regions: RBB and horizontal continuity descriptors and other typical descriptors transformed into synthesizable descriptors. The \oplus of a descriptor that is not specifically mentioned is an addition (+).

- **RBB** ($y_{\theta_m}^{\min}(r), y_{\theta_m}^{\max}(r), x_{\theta_m}^{\min}(r), x_{\theta_m}^{\max}(r)$). RBB represents coordinates of a bounding box rotated to an arbitrary angle θ_m as shown by Fig. 3. The RBB is used to accurately calculate the geographical features of a region rotated to θ_m . Each value is defined as

$$y_{\theta_m}^{\min}(r) = \min_{r \in r} \{ \tan(\theta_m)x_r + y_r \}, \quad (1)$$

$$y_{\theta_m}^{\max}(r) = \max_{r \in r} \{ \tan(\theta_m)x_r + y_r \}, \quad (2)$$

$$x_{\theta_m}^{\min}(r) = \min_{r \in r} \{ x_r + \tan(\theta_m)(h - y_r) \}, \quad (3)$$

$$x_{\theta_m}^{\max}(r) = \max_{r \in r} \{ x_r + \tan(\theta_m)(h - y_r) \}, \quad (4)$$

where x_r and y_r are coordinates of a pixel r in a region r , h is the height of the input image, and θ_m is obtained from a number of calculation angles, N^{ori} , which is a predefined fixed value. The θ_m is defined as

$$\theta_m = \frac{\pi}{2} \times \frac{m}{N^{\text{ori}}} \quad (m = \{0, 1, \dots, N^{\text{ori}}\}). \quad (5)$$

We empirically set $N^{\text{ori}} = 6$. The \oplus of each RBB values are min, max, min and max respectively.

- **Horizontal Continuity** ($v_{r, \theta_m}^{\min}, v_{r, \theta_m}^{\max}$). Two vectors represent the sections in which the region of r is continuous in θ_m as shown by Fig. 4. The initial vectors of v_{r, θ_m}^{\min} and v_{r, θ_m}^{\max} are $\{x_{\theta_m}^{\min}(r)\}$ and $\{x_{\theta_m}^{\max}(r)\}$, respectively. The \oplus is shown in Algorithm 1.
- **Area** $n^{\text{ara}}(r)$ [15]. The number of pixels.
- **Total Color** ($s^{c0}(r), s^{c1}(r), s^{c2}(r)$). The totals of pixel values in HSI color space.
- **Number of Stroke-Support Pixels (SSPs)** $n^{\text{ssp}}(r)$. An SSP [16] is used to calculate the stroke width.

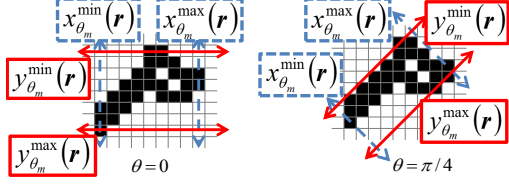


Figure 3. Rotated-bounding boxes and RBBs when $\theta = 0$ and $\pi/4$

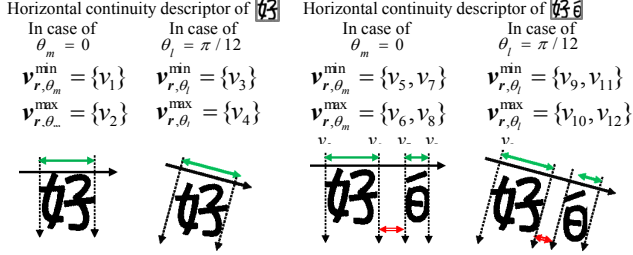


Figure 4. Examples of horizontal continuity descriptor. Green lines represent horizontal continuity descriptors, and red ones represent discontinuity features calculated by the descriptors.

- **Total Stroke Width** $s^{\text{swd}}(r)$. The total of stroke-width value of SSP.
- **Perimeter** $n^{\text{peri}}(r)$ [15]. The length of the boundary.
- **Total Smoothness Value** $s^{\text{smt}}(r)$. The total of difference values of adjacent boundary pixels' gradient direction. The details of smoothness features are given in previous studies [27].

For RBB, Eqs. 1 and 2 represent the min and max values of the vertical position, and Eqs. 3 and 4 represent the min and max values of the horizontal position in the θ_m -rotated coordinate system, as shown in Fig. 3. Therefore, RBB of θ_m allows for accurate calculation of r 's size, aspect, and position in the θ_m -rotated coordinate system. As N^{ori} increases, geometrical features with more various angles can be calculated, but the total RBB calculation cost of all CCs increases. However, RBB is incrementally computable, which is a particular class of descriptors that can be quickly calculated for all CCs extracted by MSER algorithm[15]. The descriptor's computing complexity of all CCs in an image of N pixels is $O(N)$. In our experiments, total RBB computation time of all CCs with $N^{\text{ori}} = 6$ was only 0.07 seconds longer than the time with $N^{\text{ori}} = 1$.

3.1.2 Features from Synthesizable Descriptors

Various geometrical features in the θ_m -rotated coordinate system can be obtained from RBB of θ_m . The basic features for a single region and region pair are as follows.

Some size features of r can be simply calculated from RBB of a region r . The rotated height $h_{\theta_m}(r)$ and rotated width $w_{\theta_m}(r)$ are calculated as

Algorithm 1 \oplus of horizontal continuity of r_i and r_j constituting r_c for certain θ_m

Input: $v_{r_i,\theta_m}^{\min}, v_{r_i,\theta_m}^{\max}$: Continuity descriptor of r_i

Input: $v_{r_j,\theta_m}^{\min}, v_{r_j,\theta_m}^{\max}$: Continuity descriptor of r_j

Output: $v_{r_c,\theta_m}^{\min}, v_{r_c,\theta_m}^{\max}$: Continuity descriptor of r_c

$v_{r_c,\theta_m}^{\min} \leftarrow \emptyset, v_{r_c,\theta_m}^{\max} \leftarrow \emptyset$

$q_i \leftarrow 0, q_j \leftarrow 0, v^{\min} \leftarrow 0, v^{\max} \leftarrow 0$

while $q_i < |v_{r_i,\theta_m}^{\min}|$ or $q_j < |v_{r_j,\theta_m}^{\min}|$ **do**

if $q_i \geq |v_{r_i,\theta_m}^{\min}|$ or $v_{r_i,\theta_m}^{\min}[q_i] > v_{r_j,\theta_m}^{\min}[q_j]$ **then**

$v^{\min} \leftarrow v_{r_j,\theta_m}^{\min}[q_j], v^{\max} \leftarrow v_{r_j,\theta_m}^{\max}[q_j], q_j \leftarrow q_j + 1$

else

$v^{\min} \leftarrow v_{r_i,\theta_m}^{\min}[q_i], v^{\max} \leftarrow v_{r_i,\theta_m}^{\max}[q_i], q_i \leftarrow q_i + 1$

end if

while do

if $q_i < |v_{r_i,\theta_m}^{\min}|$ and $v_{r_i,\theta_m}^{\min}[q_i] \leq v^{\max} + c_{\theta_m}$ **then**

$v^{\max} \leftarrow v_{r_i,\theta_m}^{\max}[q_i], q_i \leftarrow q_i + 1$

end if

if $q_j < |v_{r_j,\theta_m}^{\min}|$ and $v_{r_j,\theta_m}^{\min}[q_j] \leq v^{\max} + c_{\theta_m}$ **then**

$v^{\max} \leftarrow v_{r_j,\theta_m}^{\max}[q_j], q_j \leftarrow q_j + 1$

end if

if v^{\max} is not updated **then break**

end while

v_{r_c,θ_m}^{\min} push back $v^{\min}, v_{r_c,\theta_m}^{\max}$ push back v^{\max}

end while

$$h_{\theta_m}(r) = y_{\theta_m}^{\max}(r) - y_{\theta_m}^{\min}(r) + c_{\theta_m}, \quad (6)$$

$$w_{\theta_m}(r) = x_{\theta_m}^{\max}(r) - x_{\theta_m}^{\min}(r) + c_{\theta_m}, \quad (7)$$

$$c_{\theta_m} = \sqrt{\tan(\theta_m) \times \tan(\theta_m) + 1}, \quad (8)$$

where c_{θ_m} is a coefficient to prevent division by zero and is 1 when the coordinate system is not rotated ($\theta_m = 0$). The length of each coordinate system has a different scale, but a rotated aspect ratio can be calculated without normalization because the lengths of the same coordinate systems can be directly compared.

$$aspect_{\theta_m}(r) = \frac{\min(w_{\theta_m}(r), h_{\theta_m}(r))}{\max(w_{\theta_m}(r), h_{\theta_m}(r))}. \quad (9)$$

RBB of region pair r_i and r_j allows accurate calculation of rotated-region pair features: height difference h^{diff} , width difference w^{diff} , horizontal minimum distance x^{dis} , and vertical overlap y^{ovl} . In general, each character size and each interval of character pairs are similar; therefore, these features are useful for rotated text detection. These features of r_i and r_j are defined as

$$h_{\theta_m}^{\text{diff}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{|h_{\theta_m}(\mathbf{r}_i) - h_{\theta_m}(\mathbf{r}_j)|}{\max(h_{\theta_m}(\mathbf{r}_i), h_{\theta_m}(\mathbf{r}_j))}, \quad (10)$$

$$w_{\theta_m}^{\text{diff}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{|(w_{\theta_m}(\mathbf{r}_i) - w_{\theta_m}(\mathbf{r}_j))|}{\max(w_{\theta_m}(\mathbf{r}_i), w_{\theta_m}(\mathbf{r}_j))}, \quad (11)$$

$$x_{\theta_m}^{\text{dis}}(\mathbf{r}_i, \mathbf{r}_j) = \max(x_{\theta_m}^{\text{min}}(\mathbf{r}_i) - x_{\theta_m}^{\text{max}}(\mathbf{r}_j), x_{\theta_m}^{\text{min}}(\mathbf{r}_j) - x_{\theta_m}^{\text{max}}(\mathbf{r}_i)) + c_{\theta_m}, \quad (12)$$

$$y_{\theta_m}^{\text{ovl}}(\mathbf{r}_i, \mathbf{r}_j) = \min(y_{\theta_m}^{\text{max}}(\mathbf{r}_i), y_{\theta_m}^{\text{max}}(\mathbf{r}_j)) - \max(y_{\theta_m}^{\text{min}}(\mathbf{r}_i), y_{\theta_m}^{\text{min}}(\mathbf{r}_j)) + c_{\theta_m}. \quad (13)$$

We now discuss the advanced features obtained from the above synthesizable descriptors as follows. These features are used in a calculation of a character pair distance at the text-line-formation stage mentioned in Sec 3.2.2. When the values of these features are small, a probability that the regions related to the features are characters is high. Discontinuity is a feature we introduce in this paper. Other features are similar to typical features for scene-text detection, but some have been extended with RBB to be more robust against orientation.

- **Discontinuity of Single Character**

$$d^{\text{dcn}}(\mathbf{r}_i) = \begin{cases} 0 & (n^{\text{cnt}} = 1) \\ \sum_{m=0}^{N^{\text{ori}}} \sum_{q=1}^{n^{\text{cnt}}} \frac{v_{\mathbf{r}_i, \theta_m}^{\text{min}}[q] - v_{\mathbf{r}_i, \theta_m}^{\text{max}}[q-1]}{w_{\theta_m}(\mathbf{r}_i)} & (\text{else}) \end{cases}, \quad (14)$$

where $v[q]$ is the $(q+1)$ th element of vector v , n^{cnt} is $|v_{\mathbf{r}_i, \theta_m}^{\text{min}}|$. Discontinuity represents the geometrical distance between CCs constituting the same compound region as shown by Fig. 4 (b). Generally, the discontinuity of correct characters is small or zero. Our method prevents an unintended compounding of non-close CCs by using the discontinuity.

- **Color Difference of Pair**

$$d^{\text{clr}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{\sqrt{(d^{\text{c0}})^2 + (d^{\text{c1}})^2 + (d^{\text{c2}})^2}}{255}, \quad (15)$$

where $d^{\text{c0}} = (s^{\text{c0}}(\mathbf{r}_i)/n^{\text{ara}}(\mathbf{r}_i)) - (s^{\text{c0}}(\mathbf{r}_j)/n^{\text{ara}}(\mathbf{r}_j))$, $d^{\text{c1}} = (s^{\text{c1}}(\mathbf{r}_i)/n^{\text{ara}}(\mathbf{r}_i)) - (s^{\text{c1}}(\mathbf{r}_j)/n^{\text{ara}}(\mathbf{r}_j))$, $d^{\text{c2}} = (s^{\text{c2}}(\mathbf{r}_i)/n^{\text{ara}}(\mathbf{r}_i)) - (s^{\text{c2}}(\mathbf{r}_j)/n^{\text{ara}}(\mathbf{r}_j))$.

- **Stroke Width Difference of Pair**

$$d_{\theta_m}^{\text{swd}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{|swdm(\mathbf{r}_i) - swdm(\mathbf{r}_j)|}{\max(w_{\theta_m}(\mathbf{r}_i), h_{\theta_m}(\mathbf{r}_i))}, \quad (16)$$

where $swdm(\mathbf{r}) = s^{\text{swd}}(\mathbf{r})/n^{\text{ssp}}(\mathbf{r})$.

- **Vertical Overlap of Pair**

$$d_{\theta_m}^{\text{ovl}}(\mathbf{r}_i, \mathbf{r}_j) = 1 - \frac{y_{\theta_m}^{\text{ovl}}(\mathbf{r}_i, \mathbf{r}_j)}{\min(h_{\theta_m}(\mathbf{r}_i), h_{\theta_m}(\mathbf{r}_j))}. \quad (17)$$

- **Height and Width Differences of Pair**

$$d_{\theta_m}^{\text{h}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{|h_{\theta_m}(\mathbf{r}_i) - h_{\theta_m}(\mathbf{r}_j)|}{\max(h_{\theta_m}(\mathbf{r}_i), h_{\theta_m}(\mathbf{r}_j))}, \quad (18)$$

$$d_{\theta_m}^{\text{w}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{|w_{\theta_m}(\mathbf{r}_i) - w_{\theta_m}(\mathbf{r}_j)|}{\max(w_{\theta_m}(\mathbf{r}_i), w_{\theta_m}(\mathbf{r}_j))}. \quad (19)$$

- **Top, Bottom, and Both Alignments of Triplet**

$$d_{\theta_m}^{\text{top}}(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_k) = \quad (20)$$

$$\begin{cases} y_{\theta_m}^{\text{diff.t}}(\mathbf{r}_j, \mathbf{r}_k) & (\mathbf{r}_i = \emptyset) \\ y_{\theta_m}^{\text{diff.t}}(\mathbf{r}_j, \mathbf{r}_k) - x_{\theta_m}^{\text{ctr}}(\mathbf{r}_j, \mathbf{r}_k) \frac{y_{\theta_m}^{\text{diff.t}}(\mathbf{r}_i, \mathbf{r}_j)}{x_{\theta_m}^{\text{ctr}}(\mathbf{r}_i, \mathbf{r}_j)} & (\text{else}) \end{cases},$$

$$d_{\theta_m}^{\text{btm}}(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_j) = \quad (21)$$

$$\begin{cases} y_{\theta_m}^{\text{diff.b}}(\mathbf{r}_j, \mathbf{r}_k) & (\mathbf{r}_i = \emptyset) \\ y_{\theta_m}^{\text{diff.b}}(\mathbf{r}_j, \mathbf{r}_k) - x_{\theta_m}^{\text{ctr}}(\mathbf{r}_j, \mathbf{r}_k) \frac{y_{\theta_m}^{\text{diff.b}}(\mathbf{r}_i, \mathbf{r}_j)}{x_{\theta_m}^{\text{ctr}}(\mathbf{r}_i, \mathbf{r}_j)} & (\text{else}) \end{cases},$$

$$d_{\theta_m}^{\text{both}}(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_k) = d_{\theta_m}^{\text{top}}(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_k) + d_{\theta_m}^{\text{btm}}(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_k), \quad (22)$$

where

$$x_{\theta_m}^{\text{ctr}}(\mathbf{r}_i, \mathbf{r}_j) = \frac{x_{\theta_m}^{\text{max}}(\mathbf{r}_j) + x_{\theta_m}^{\text{min}}(\mathbf{r}_i) - x_{\theta_m}^{\text{max}}(\mathbf{r}_i) + x_{\theta_m}^{\text{min}}(\mathbf{r}_j)}{2},$$

$$y_{\theta_m}^{\text{diff.t}}(\mathbf{r}_i, \mathbf{r}_j) = y_{\theta_m}^{\text{max}}(\mathbf{r}_j) - y_{\theta_m}^{\text{max}}(\mathbf{r}_i),$$

$$y_{\theta_m}^{\text{diff.b}}(\mathbf{r}_i, \mathbf{r}_j) = y_{\theta_m}^{\text{min}}(\mathbf{r}_i) - y_{\theta_m}^{\text{min}}(\mathbf{r}_j).$$

- **Interval Similarity of Triplet**

$$d_{\theta_m}^{\text{inter}}(\mathbf{r}_i, \mathbf{r}_j, \mathbf{r}_k) = \quad (23)$$

$$\begin{cases} 0 & (\mathbf{r}_k = \emptyset) \\ \frac{\max(x_{\theta_m}^{\text{dis}}(\mathbf{r}_i, \mathbf{r}_j) - x_{\theta_m}^{\text{dis}}(\mathbf{r}_j, \mathbf{r}_k), 0)}{\max(w_{\theta_m}(\mathbf{r}_i), w_{\theta_m}(\mathbf{r}_j))} & (\text{else}) \end{cases}.$$

Interval similarity represents the geometrical distance to the nearest CC that is not included in the same text. Our method promotes a compounding of close CCs by using the interval similarity.

Note that RBB models only in-place rotation but not perspective distortion. For example, values of Eq. 18 and 19 for characters belonging to the same text with perspective distortion tends to be large resulting in the characters are determined to belong to different text. The alignment features shown in Eq. 20, 21 and 22 are important for detecting such text. Because characters belonging to the same text are regularly aligned and the values of the features are low regardless of the presence or absence of perspective distortion.

3.2. Text-Detection Flow

3.2.1 Character-Candidate Detection

With our method, CCs are extracted using an MSER algorithm and classified as character candidates or non-character candidates using a sequential classifier. The classification is broken down into three stages for better computational efficiency. The first and second stages use AdaBoost[17] and the third stage uses a DNN. The CCs that are classified as characters by the preceding identifier proceed to the next stage. The features used on the first and second stages are almost the same as those in [15].

In the first stage, the classifier uses low-computational cost features and removes only clearly non-character CCs. The difference from [15] is that a minimum-rotated aspect

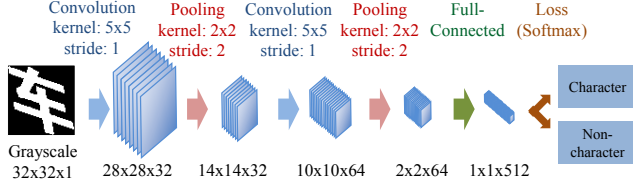


Figure 5. Structure of deep neural network

is used instead of the general aspect. The minimum-rotated aspect is defined as

$$aspect^{\min}(\mathbf{r}) = \min_{m \in \{0, 1, \dots, N_{ori}\}} (aspect_{\theta_m}(\mathbf{r})). \quad (24)$$

Here, $aspect^{\min}(\mathbf{r})$ is robust against the rotation of \mathbf{r} and improves the classification accuracy of rotated regions.

In the second stage, our method uses smoothness in addition to the features described in [15]. Smoothness is calculated by $s^{\text{smt}}(\mathbf{r})/n^{\text{peri}}(\mathbf{r})$. The second classifier outputs the class-conditional probability, and our method reduces the redundancy of detected CCs using this probability. The redundancy-reduction algorithm is similar to the linear-reduction algorithm in [27], but the difference is that the class-conditional probability is used instead of the regularized variation.

In the third stage, our method detects CCs that are characters with high probability by using the powerful discrimination performance of a DNN. Figure 5 shows the structure of the DNN used in our method. Note that all images are rotation-normalized using the angle of the minimum-aspect rotated-bounding box. After this stage, each CC has a character probability in the range of $[0, 1]$.

3.2.2 Text-Line Formation

At the text-line-formation stage, many rotated text-line candidates are detected on the basis of CCs labeled as characters with the DNN. Such a CC labeled as a character is called a base character. Our method recursively searches for neighbors of the base character using the search method to form text-line candidates. The search method is based on a character pair distance and some heuristic rules. The character pair distance represents the similarity of characters, and the distance of characters is low when the characters belong to the same text.

We now define a character pair distance for calculating the similarity of two characters. In a certain θ_m -rotated coordinate system, the distance from character region \mathbf{r}_i to another character region \mathbf{r}_j , $d(\mathbf{r}_i, \mathbf{r}_j)$ is defined as $\mathbf{w}^T \mathbf{x}_{\mathbf{r}_i, \mathbf{r}_j, \theta_m}$, where \mathbf{w} is a parameter of weight vector. Then, $\mathbf{x}_{\mathbf{r}_i, \mathbf{r}_j, \theta_m}$ is a vector of ten features calculated from Eq. 14 to Eq. 23 in Sec. 3.1.2. We determine parameter \mathbf{w} by using the distance metric learning (DML) technique [27, 25], which is a semi-supervised clustering technique and learns a distance metric that satisfies given labels. Yin *et al.* used the DML technique to automatically determine

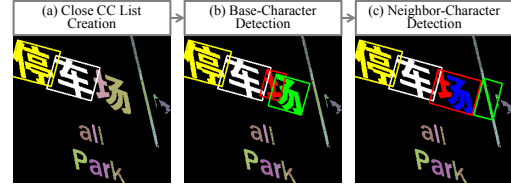


Figure 6. Flow of neighbor-character search and sample results when $\theta_m = \pi(5/12)$. (a) White region is \mathbf{r}_i , yellow one is \mathbf{r}_{i-1} , and others are in close CC list of \mathbf{r}_i . Proposed method searches neighbor character of \mathbf{r}_i . (b) Red region is initial \mathbf{r}_{i+1} and green one is \mathbf{r}_{i+2} , which is \mathbf{r}_{i+1} 's horizontal nearest CC in CC list. $d(\mathbf{r}_i, \mathbf{r}_{i+1})$ is large because \mathbf{r}_{i+1} and \mathbf{r}_{i+2} are very close and $d_{\theta_m}^{\text{inter}}(\mathbf{r}_i, \mathbf{r}_{i+1}, \mathbf{r}_{i+2})$ is large. (c) Combination of red and blue regions is \mathbf{r}_i 's neighbor character, which is detected at end.

weights and threshold to distinguish characters from non-characters. We also used this technique to determine \mathbf{w} , but only the threshold is determined manually.

Our search method uses the pair distance to detect a neighbor characters of \mathbf{r}_i . The process is performed based on four regions including \mathbf{r}_i , the rest of which are a nearest character region of the \mathbf{r}_i , referred to as \mathbf{r}_{i-1} , and two undetermined regions, referred to as \mathbf{r}_{i+1} and \mathbf{r}_{i+2} . \mathbf{r}_{i+1} represents a neighbor-character candidate of \mathbf{r}_i . \mathbf{r}_{i-1} and \mathbf{r}_{i+2} are used to determine whether \mathbf{r}_{i+1} should be included in the text-line candidate. At calculating $\mathbf{x}_{\mathbf{r}_i, \mathbf{r}_{i+1}, \theta_m}$, the input of Eq. 14 is (\mathbf{r}_{i+1}) , that of Eqs. 20, 21 and 22 are $(\mathbf{r}_{i-1}, \mathbf{r}_i, \mathbf{r}_{i+1})$, that of Eq. 23 is $(\mathbf{r}_i, \mathbf{r}_{i+1}, \mathbf{r}_{i+2})$, that of other Eqs. are $(\mathbf{r}_i, \mathbf{r}_{i+1})$. Figure 6 show examples of the four regions and the flow of neighbor-character search. Our search method searches for a neighbor character of \mathbf{r}_i by selecting and evaluating multiple-neighbor-character candidates as \mathbf{r}_{i+1} . For selecting \mathbf{r}_{i+1} and \mathbf{r}_{i+2} , our search method constructs a close CC list of \mathbf{r}_i , as shown in Fig. 6 (a). The close CC list of \mathbf{r}_i is a set of \mathbf{r}_{cc} that satisfies the following conditions:

$$d_{\theta_m}^{\text{clr}}(\mathbf{r}_{cc}, \mathbf{r}_i) < th_{cc}^{\text{clr}}, \quad d_{\theta_m}^{\text{swd}}(\mathbf{r}_{cc}, \mathbf{r}_i) < th_{cc}^{\text{swd}}, \quad (25)$$

$$x_{\theta_m}^{\text{dis}}(\mathbf{r}_{cc}, \mathbf{r}_i) < th_{cc}^{\text{hor}}, \quad d_{\theta_m}^{\text{ovl}}(\mathbf{r}_{cc}, \mathbf{r}_i) < th_{cc}^{\text{ver}}. \quad (26)$$

The close CC list is used to narrow the neighbor-character search range. Each threshold, th_{cc}^{clr} , th_{cc}^{swd} , th_{cc}^{hor} and th_{cc}^{ver} , is determined as a value in which all correct samples of training data can be detected. Next, our search method selects the initial \mathbf{r}_{i+1} from the close CC list based on heuristic rules. The initial \mathbf{r}_{i+1} is the CC with the lowest $x_{\theta_m}^{\text{dis}}(\mathbf{r}_i, \mathbf{r}_{i+1})$ that satisfies conditions $y_{\theta_m}^{\text{ovl}}(\mathbf{r}_i, \mathbf{r}_{i+1}) > 0$ and $d_{\theta_m}^{\text{h}}(\mathbf{r}_i, \mathbf{r}_{i+1}) > 0.5$. On the other hand, \mathbf{r}_{i+2} is the horizontal nearest CC that satisfies $y_{\theta_m}^{\text{ovl}}(\mathbf{r}_{i+1}, \mathbf{r}_{i+2}) > 0$ in the CC list. As a result, our search method determines the four regions for calculating $\mathbf{x}_{\mathbf{r}_i, \mathbf{r}_{i+1}, \theta_m}$.

Our search method searches for a combination with lower pair distance from the close CC list. The initial \mathbf{r}_{i+1} becomes the initial neighbor candidate $\mathbf{r}_{i+1}^{\text{cand}}$, and $d(\mathbf{r}_i, \mathbf{r}_{i+1}^{\text{cand}})$ is calculated. Then, temporary compound regions of $\mathbf{r}_{i+1}^{\text{cand}}$ and other CCs of the CC list are cre-

ated and pair distances are calculated. Next, the temporary compound region with the lowest distance to r_i is selected as r_{i+1}^{lowest} , and r_{i+1}^{cand} is replaced with r_{i+1}^{lowest} when $d(r_i, r_{i+1}^{\text{cand}}) \geq d(r_i, r_{i+1}^{\text{lowest}})$. Repeating the above operation until $d(r_i, r_{i+1}^{\text{cand}}) < d(r_i, r_{i+1}^{\text{lowest}})$, and the last r_{i+1}^{cand} is selected as the detected neighbor character of r_i , r^{nbr} . Finally, the method will add the detected character to the text-line candidate when $d(r_i, r^{\text{nbr}}) < th^{\text{dml}}$ and carry out neighbor-character search again. r_i and r_{i-1} are replaced with r^{nbr} and r_i in the next step.

At this stage, the probability of the text-line candidate is calculated for the next text-line-elimination stage. Let t_{θ_m} be a text-line candidate whose orientation angle is θ_m and $p^{\text{txt}}(t_{\theta_m})$ be the probability of t_{θ_m} . The initial value of $p^{\text{txt}}(t_{\theta_m})$ is zero, and $(th^{\text{dml}} - d(r_i, r^{\text{nbr}}))$ is added to $p^{\text{txt}}(t_{\theta_m})$ when r^{nbr} is joined to t_{θ_m} . High $p^{\text{txt}}(t_{\theta_m})$ denotes that features of character candidates constituting t_{θ_m} are very similar.

3.2.3 Text-Line Elimination

After text-line formation, too many text-line candidates remain and the majority overlap, as shown in Fig. 2 (c). Our method throttles overlapping text-line candidates into one by using the simple elimination algorithm. First, our method updates $p^{\text{txt}}(t_{\theta_m})$ using the probability $p^{\text{txt}}(t_{\theta_l}^{\text{ovl}})$ of candidate $t_{\theta_l}^{\text{ovl}}$, which overlaps t_{θ_m} . When θ_m and θ_l are the same, $p^{\text{txt}}(t_{\theta_l}^{\text{ovl}})$ is added to $p^{\text{txt}}(t_{\theta_m})$. On the other hand, when θ_m and θ_l are different, $-p^{\text{txt}}(t_{\theta_l}^{\text{ovl}})$ is added to $p^{\text{txt}}(t_{\theta_m})$. Next, the candidate t_{θ_m} is removed when θ_m and θ_l are different and $p^{\text{txt}}(t_{\theta_m}) < p^{\text{txt}}(t_{\theta_l}^{\text{ovl}})$. Empirically, the number of text-line candidates with correct angles is larger than that of candidates with incorrect angles. Therefore, this elimination method works well in spite of its simplicity. Finally, the remaining text-line candidates are classified using AdaBoost with typical features [14, 23].

4. Experiments

4.1. Parameter Setting

We used the training set of the MSRA-TD500 dataset [4] to determine all parameters in our method. The MSRA-TD500 includes rotated scene text of English and Chinese. It includes a relatively large number of rotated-discontiguous characters to other datasets. For the learning of the DML technique, we labeled each character region in the training set and constructed positive and negative samples. The weight vector was determined using DML technique, and we adjusted the threshold value manually. How to learn automatically the threshold is for future work. Table 1 lists the learned parameters, which were normalized by the manually adjusted threshold value.

The training of AdaBoost and the DNN was also done using the MSRA-TD 500 training set. We determined each

Params	w^{dcn}	w^{clr}	w^{swd}	w^{ver}	w^{h}
Learned-value	1.30	1.02	0.395	0.603	0.415
Params	w^{w}	w^{top}	w^{btm}	w^{both}	w^{inter}
Learned-value	0.529	0.359	0.709	0.672	0.726

Table 1. Learned weights by using MSRA-TD500 training set

threshold, assuming that detecting one or more characters constituting a text is considered detecting characters constituting the text. Because our method detects a text even when only one character constituting the text is detected at the character-candidate detection stage. The labeling of all training data was done manually.

In order to determine the value of N^{ori} , we evaluated recall and processing time of our method while changing N^{ori} . The recall increased with increasing N^{ori} when $N^{\text{ori}} \leq 6$ because text-line candidates with more various orientation were able to be detected when N^{ori} is larger. However the recall did not change or decreased when $N^{\text{ori}} > 6$ due to an excessive detection of diagonal text candidates. Therefore we set N^{ori} to 6. On the other hand, the processing time monotonically increased in all N^{ori} . In our experiments on MSRA-TD500 training set, the overall processing processing time in the case of $N^{\text{ori}} = 6$ was only 1.2 times the time in the case of $N^{\text{ori}} = 1$, although the number of search angles for text-line candidates is 6 times. The reason is that the descriptor calculation of compound regions had been made more efficient by using synthesizable descriptors. In our experiment, our method calculated descriptors of about 31 thousands compound regions per an image, and the total pixels of all compound regions is about 430 millions. In the case of RBB calculation, the numbers of min and max calculations were reduced from the total pixels to the number of compound regions (about 1/13, 800).

4.2. Performance Evaluation on MSRA-TD500

We evaluated our method on the MSRA-TD500 [4], and examples of detection results are given in Fig. 7 and 8. Our method could detect most of the rotated scene text that included discontiguous characters. In particular, it is usually difficult to correctly detect when some texts were dense such as in Fig. 7 (a) and (b). With our method, false positives, which overlapped with the correct text, were removed at the text-line-elimination stage, and the method succeeded in the detection. However, some text failed to be detected. The upper close texts were inadvertently connected, as shown in Fig. 8 (a). A part of the text was not detected due to light reflection, as shown in Fig. 8 (b). Two close text lines with perspective distortion were not detected, as shown in Fig. 8 (c). The cause of the detection failure is that two small characters at the right end of each text line mistakenly detected as one discontiguous character in the text-line-formation stage described in Sec. 3.2.2. The our method could not detect texts with non-uniform alignment, as shown in Fig. 8 (d) and (e), because we assumed



Figure 7. Examples of successful detections. Characters with same colored rectangles belong to same detected text line.



Figure 8. Examples of failure results. Each rectangle represents detected texts.

Method	Year	Precision	Recall	F-measure
Our method		0.87	0.63	0.74
Zhang <i>et al.</i> [28]	2016	0.83	0.67	0.74
He <i>et al.</i> [7]	2016	0.76	0.61	0.69
Yin <i>et al.</i> [27]	2015	0.81	0.63	0.71
Kang <i>et al.</i> [8]	2014	0.71	0.62	0.66
Yao <i>et al.</i> [20]	2014	0.62	0.64	0.61

Table 2. Experimental results on MSRA-TD500

that the alignment was uniform. Some repeated patterns were erroneously detected as text, as shown in Fig 7 (b) and 8 (e). In our method, text-line candidates composed of regularly aligned character-candidates have high p^{txt} and are detected as text, even when the character-candidates are not characters. The number of false-positives tends to increase when the image contains many repeated patterns, because our method searches for text at various directions. Although the detection fails due to other conditions as stated as above, our method was able to detect rotated-scene text that included discontinuous characters.

The quantitative comparisons of the proposed and existing methods are shown in Table 2 our method achieved the best score on precision and F-measure, and recall was the second best score. Our method could detect all characters constituting a text even when only one character constituting the text is detected at the character-candidate detection stage. Therefore, in the sequential CC classification stage of our method, we were able to set the threshold of the classifier to very high (the threshold of DNN was 0.95 in our implementation), and improved precision. However, the recall of our method was not high because the method was designed to detect text in a certain language such as Chinese, and only 36% of the MSRA-TD500 contains Chinese. Texts of English has a large variance of horizontal width and vertical width of characters than that of Chinese; therefore, the contribution to neighbor-character search by accurate geometrical features is smaller. We evaluated the recall

for only Chinese text, which was 0.67; the same as that obtained by Zhang *et al.* [28]. The precision for only Chinese text cannot be calculated because our method detect both English and Chinese without distinction, but we consider that it is similar to the precision of Table 2. Therefore, the performance of our method is excellent for text in Chinese. Our method may also be suitable for other languages with geographical features similar to Chinese (*e.g.* Japanese).

The average detection time of the our method was about 2.28 seconds (4GHz CPU without GPU), while that of Zhang *et al.*'s method is about 2.1 seconds (2GHz CPU with GPU). We were not able to directly compare them because of different environments, but our method achieved the similar speed as their method even when we did not use GPU. The detection times of our method changed depending on an input image. The minimum time was 0.442 seconds and maximum was 20.5 seconds for the MSRA-TD500. The reason for the long detection time is that our method detected many overlapping-text-line candidates when the text was dense in the input image. Improving the detection time is for future work.

5. Conclusion

We proposed an MSER-based scene-text-detection method robust against rotated discontinuous characters. RBB accurately calculates rotated geometrical features and our neighbor-character search method with synthesizable descriptors detects discontinuous characters. Our method achieved the best score with regard to precision and F-measure on the MSRA-TD500. For Chinese text, the recall was the same as the best score of the existing methods. Therefore, we conclude that our method achieves the state-of-the-art performance with respect to precision and F-measure in detecting rotated scene-texts that include discontinuous characters. For future work, we will investigate a more suitable parameter-learning method and improvement in detection time for dense text.

References

- [1] A. Bissacco, M. Cummins, Y. Netzer, and H. Neven. Photoocr: Reading text in uncontrolled conditions. In *IEEE International Conference on Computer Vision*, pages 785–792, 2013.
- [2] E. O. Boris Epshtein, Yonathan Wexler. Detecting text in natural scenes with stroke width transform. In *IEEE International Conference on Computer Vision and Pattern Recognition*, pages 2960–2970, 2010.
- [3] X. Chen and A. L. Yuille. Detecting and reading text in natural scenes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 366–373, 2004.
- [4] Cong Yao, Xiang Bai, Wenyu Liu, Yi Ma, and Zhuowen Tu. Detecting texts of arbitrary orientations in natural images. In *IEEE Conference on Computer Vision and Pattern Recognition*, volume 8, pages 1083–1090. IEEE, 2012.
- [5] L. Gomez and D. Karatzas. MSER-Based Real-Time Text Detection and Tracking. In *International Conference on Pattern Recognition*, pages 3110–3115. IEEE, 2014.
- [6] A. R. Hanson, E. Learned-Miller, and J. J. Weinman. Scene text recognition using similarity and a lexicon with sparse belief propagation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(10):1733–1746, 2009.
- [7] T. He, W. Huang, Y. Qiao, and J. Yao. Text-attentional convolutional neural network for scene text detection. *IEEE Transactions on Image Processing*, 25(6):2529–2541, 2016.
- [8] L. Kang, Y. Li, and D. Doermann. Orientation robust text line detection in natural images. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4034–4041. IEEE, 2014.
- [9] K. I. Kim, K. Jung, and J. H. Kim. Texture-based approach for text detection in images using support vector machines and continuously adaptive mean shift algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(12):1631–1639, 2003.
- [10] J.-J. Lee, P.-H. Lee, S.-W. Lee, A. Yuille, and C. Koch. Adaboost for text detection in natural scene. In *International Conference on Document Analysis and Recognition*, pages 429–434, 2011.
- [11] T. Lu, S. Palaiiahakote, C. L. Tan, and W. Liu. *Introduction to Video Text Detection*, pages 1–18. Springer London, London, 2014.
- [12] J. Matas, O. Chum, M. Urban, and T. Pajdl. Robust wide baseline stereo from maximally stable extremal regions. In *British Machine Vision Conference*, pages 36.1–36.10, 2002.
- [13] R. Minetto, N. Thome, M. Cord, N. J. Leite, and J. Stolfi. Snoopertext: A text detection system for automatic indexing of urban scenes. *Computer Vision and Image Understanding*, 122:92–104, 2014.
- [14] L. Neumann and J. Matas. Text localization in real-world images using efficiently pruned exhaustive search. In *International Conference on Document Analysis and Recognition*, pages 687–691. IEEE, 2011.
- [15] L. Neumann and J. Matas. Real-time scene text localization and recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, volume 38, pages 3538–3545. IEEE, 2012.
- [16] L. Neumann and J. Matas. Efficient scene text localization and recognition with local character refinement. In *International Conference on Document Analysis and Recognition*, pages 746–750, 2015.
- [17] R. E. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37(3):297–336, 1999.
- [18] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015.
- [19] T. Wang, D. J. Wu, A. Coates, and A. Y. Ng. End-to-end text recognition with convolutional neural networks. In *International Conference on Pattern Recognition*, pages 3304–3308, 2012.
- [20] C. Yao, X. Bai, and W. Liu. A unified framework for multi-oriented text detection and recognition. *IEEE Transactions on Image Processing*, 23(11):4737–4749, 2014.
- [21] C. Yi and Y. Tian. Localizing text in scene images by boundary clustering, stroke segmentation, and string fragment classification. *IEEE Transactions on Image Processing*, 21(9):4256–4268, 2012.
- [22] C. Yi and Y. Tian. Text extraction from scene images by character appearance and structure modeling. *Computer Vision and Image Understanding*, 117(2):182–194, 2013.
- [23] X. Yin, X.-C. Yin, H.-W. Hao, and K. Iqbal. Effective text localization in natural scene images with msr, geometry-based grouping and adaboost. In *International Conference on Pattern Recognition*, pages 725–728. IEEE Computer Society, 2012.
- [24] X.-C. Yin, H. W. Hao, J. Sun, and S. Naoi. Robust vanishing point detection for mobilecam-based documents. In *2011 International Conference on Document Analysis and Recognition*, pages 136–140, 2011.
- [25] X.-C. Yin, W. Y. Pei, J. Zhang, and H. W. Hao. Multi-orientation scene text detection with adaptive clustering. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1930–1937, 2015.
- [26] X.-C. Yin, X. Yin, K. Huang, and H.-W. Hao. Accurate and robust text detection: A step-in for text retrieval in natural scene images. In *ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1091–1092, 2013.
- [27] X.-C. Yin, X. Yin, K. Huang, and H.-W. Hao. Robust text detection in natural scene images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(5):970–983, 2014.
- [28] Z. Zhang, C. Zhang, W. Shen, C. Yao, W. Liu, and X. Bai. Multi-oriented text detection with fully convolutional networks. *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4159–4167, 2016.
- [29] S. Zhu and R. Zanibbi. A text detection system for natural scenes with convolutional feature learning and cascaded classification. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 625–632, 2016.