

Detecting Texts of Arbitrary Orientations in Natural Images

Cong Yao, Xin Zhang, Xiang Bai, *Member, IEEE*, Wenyu Liu, *Member, IEEE*, Yi Ma, *Senior Member, IEEE*, and Zhuowen Tu, *Member, IEEE*

Abstract

Texts in a natural image directly carry rich high-level semantic information about a scene, which can be used to assist a wide variety of applications, such as image understanding, image indexing and search, geolocation or navigation, and human computer interaction. However, most existing text detection and recognition systems are designed for horizontal or near-horizontal texts. With the increasingly popular computing-on-the-go devices, detecting texts of arbitrary orientations from images taken by such devices under less controlled conditions has become an increasingly important and yet challenging task. In this paper, we propose a new algorithm to detect texts of arbitrary orientations in natural images. Our algorithm is based on a two-level classification scheme and utilize two sets of features specially designed for capturing both intrinsic and orientation-invariant characteristics of texts. To better evaluate the proposed method and compare it with other existing algorithms, we generate a more extensive and challenging dataset, which includes various types of texts in diverse real-world scenes. We also propose a new evaluation protocol, which is more suitable for benchmarking algorithms designed for texts of varying orientations. Experiments on conventional benchmarks and the new dataset demonstrate that our system compares favorably with the state-of-the-art algorithms when handling horizontal texts and achieves significantly enhanced performance on texts of arbitrary orientations in complex natural scenes.

Index Terms

C. Yao, X. Bai, and W. Liu are with the Department of Electronics and Information Engineering, Huazhong University of Science and Technology, Wuhan, 430074, China. Email: yaocong2010@gmail.com {xbai, liuwu}@hust.edu.cn.

X. Zhang is with the Department of Computer Science and Technology, Tsinghua University, Beijing, 100084, China. Email: xinzhang1111@gmail.com.

Y. Ma is with Microsoft Research Asia. Email: mayi@microsoft.com.

Z. Tu is with Microsoft Research Asia and also University of California at Los Angeles, Los Angeles, CA 90095, USA. Email: ztu@loni.ucla.edu.

22 Text detection, natural image, arbitrary orientations, feature design, evaluation protocol.

23 I. INTRODUCTION

24 Texts in a natural image typically carry very rich and poignant high-level semantic information about
25 a scene. Their existence is also ubiquitous in urban environments, e.g. traffic signs, billboards, business
26 name cards, and license plates. Effective text detection and recognition systems have been very useful in a
27 variety of applications such as robot navigation [1], image search [2], and human computer interaction [3].
28 The popularity of smart phones and ubiquitous computing devices have also made the acquisition and
29 transmission of text data more convenient and efficient than ever. Thus, automatically detecting and
30 recognizing texts from images casually captured by such devices has become an ever important task in
31 computer vision.

32 In this paper, we tackle the problem of text detection in natural images, which remains a challenging
33 task although it has been extensively studied in the past decades [4]–[16]. The difficulty of automatic
34 text detection mainly stems from two aspects: (1) diversity of text appearances and (2) complexity of
35 cluttered backgrounds. On one hand, texts, unlike conventional objects (e.g. cars and horses), typically
36 consist of hundreds and even thousands of different instances and they exhibit significant variations in
37 shapes and appearances: different texts may have different sizes, colors, fonts, languages, and orientations,
38 even within the same scene. On the other hand, many other man-made objects (such as windows and
39 railings) in the scene often exhibit similar regular structures as those of texts. Sometimes even natural
40 objects (such as grasses and leaves) may happen to distribute in a similar way as a cluster of texts. Such
41 ambiguities have made reliable text detection in natural images an extremely challenging task.

42 In the literature, most of the existing methods [6], [9], [17] have focused on detecting horizontal or
43 near-horizontal texts, as we will see in a survey of related work in Sec. (I-A). Obviously, the requirement
44 of horizontal texts severely limits the applicability of such methods in scenarios where images are taken
45 casually with a mobile device. Detecting texts with arbitrary orientations in complex natural scenes
46 remains a challenge for most practical text detection and recognition systems [18], [19]. In this work,
47 we aim to build an effective and efficient system for detecting texts of arbitrary orientations in complex
48 natural scenes (see Fig. 1).

49 Most conventional text detection methods rely on features that are primarily designed for horizontal
50 texts (such as those used in [7]). Thus, when such methods are applied to images that have texts of
51 arbitrary orientations, their performance usually drops drastically. To remedy this situation, in this paper,
52 we introduce two additional sets of rotation-invariant features for text detection. To further reduce false



Fig. 1. Detected texts in natural images.

positives produced by only using such low-level features, we have also designed a two-level classification scheme that can effectively discriminate texts from non-texts. Hence, by combining the strength of rotation-invariant local features and well trained text classifiers, our system is able to effectively detect texts of arbitrary orientations with very few false positives.

The proposed method is mostly bottom-up (data-driven) but with additional prior knowledge about texts imposed in a top-down fashion. Pixels are first grouped into connected components, corresponding to strokes or characters; connected components are then linked to form chains, corresponding to words or sentences. The connected components and chains are verified by the orientation-invariant features and discriminative classifiers. With this strategy, our method is able to combine the strength of both prior knowledge about texts (such uniform stroke width) and automatically learned classifiers from labeled training data. In this way, we can strike a good balance between systematic design and machine learning, which is shown to be advantageous over either heavy black-box learning [7] or purely heuristic design [9].

To evaluate the effectiveness of our system, we have conducted extensive experiments on both conventional benchmarks and some new (more extensive and challenging) dataset. Compared with the state-of-the-art text detection algorithms, our system performs competitively in the conventional setting of horizontal texts. We have also tested our system on a challenging dataset of 500 natural images containing texts of various orientations in complex backgrounds (see Fig. 6 and Fig. 11). On this dataset, our system works significantly better than the existing systems, with an F-measure about 0.6, more than twice that of the closest competitor.

In summary, the work presented in this paper offers the following contributions to the problem of natural image text detection:

- We have proposed a text detection algorithm which is able to detect text in complex natural scenes, for a wide range of variations in color, size, font, orientation, and language.
- We have adopted a two-level classification scheme, which incorporates both bottom-up and top-down

77 procedures, leading to high precision and recall concurrently.

- 78 • We have designed two sets of features, component level features and chain level features, which
79 can capture the intrinsic properties of texts and are robust to variations, such as rotation and scale
80 change.
- 81 • We have collected a natural image database containing 500 images with texts of different colors,
82 sizes, fonts, orientations and languages under various real-world scenarios (office, mall, street, etc.),
83 which can serve as a benchmark dataset for evaluating text detection algorithms.

84 In addition, we have devised a new evaluation protocol for benchmarking text detection algorithms
85 designed for texts of arbitrary orientations.

86 We have presented a preliminary version of our work in [20]. This paper extends that article [20] with:
87 (1) more technical details, (2) more extensive and thorough experimental evaluations, and (3) further
88 integration with an OCR system.

89 *A. Related Work*

90 There have been a large number of systems dealing with text detection in natural images and videos [4]–
91 [16], [21]–[23]. Comprehensive surveys can be found in [24], [25]. Existing approaches to text detection
92 can be roughly divided into three categories: texture-based, component-based, and hybrid methods.

93 *1) Three categories of existing approaches:*

94 **Texture-based methods** (e.g. [6], [7], [21]) treat text as a special type of texture and make use of its
95 textural properties, such as local intensities, spatial variance, filter responses, and wavelet coefficients.
96 Generally, these methods are computation demanding as all locations and scales are exhaustively scanned.
97 Moreover, these algorithms mostly only detect horizontal texts.

98 In an early work, Zhong *et al.* [26] proposed a method for text localization in color images. Horizontal
99 spatial variance was used to roughly localize texts and color segmentation was performed within the
100 localized areas to extract text components. The system of Wu *et al.* [27] adopted a set of Gaussian
101 derivatives to segment texts. Rectangular boxes surrounding the corresponding text strings were formed,
102 based on certain heuristic rules on text strings, such as height similarity, spacing and alignment. The
103 above steps were applied to an image pyramid and the results were fused to make final detections. Li *et al.*
104 [28] presented a system for detecting and tracking texts in digital video. In this system, the mean
105 and the second- and third-order central moments of wavelet decomposition responses are used as local
106 features. Zhong *et al.* [29] proposed to localize candidate caption text regions directly in the discrete
107 cosine transform (DCT) compressed domain using the intensity variation information encoded in the DCT

108 domain. The method proposed by Gllavata *et al.* [21] utilized the distribution of high-frequency wavelet
109 coefficients to statistically characterize text and non-text areas. The K -means algorithm was then adopted
110 to classify text areas in the image.

111 Different from the methods surveyed above, in which filter responses or transform domain coefficients
112 are used, the algorithm of Kim *et al.* [6] relies merely on the intensities of the raw pixels. A support
113 vector machine (SVM) classifier is trained to generate probability maps, in which the positions and
114 extents of texts are searched using adaptive mean shift. Lienhart and Wernicke [30] used complex-valued
115 edge orientation maps computed from the original RGB image as features and trained neural network to
116 distinguish between text and non-text patterns. This method is able to achieve text detection and tracking
117 with sub-pixel accuracy.

118 The method of Weinman *et al.* [31] used a rich representation that captures important relationships
119 between responses to different scale- and orientation-selective filters. To improve the performance, con-
120 ditional random field (CRF) was used to exploit the dependencies between neighboring image region
121 labels. Based on the observation that areas with high edge density indicate text regions, text detection
122 in [32] was carried out in a sequential multi-resolution paradigm. At each level, local thresholding was
123 applied in the edge map, to highlight text areas and suppress other areas with different contrast properties;
124 hysteresis edge recovery was then performed to bring back lower-contrast edge pixels belonging to texts.

125 To speed up the text detection procedure, Chen *et al.* [7] proposed a efficient text detector, which is a
126 cascade Adaboost classifier. The weak classifiers are trained on a set of informative features, including
127 mean and variance of intensity, horizontal and vertical derivatives, and histograms of intensity gradient.
128 The detector is applied to sub-regions of the image (at multiple scales) and outputs candidate text regions.

129 Recently, Wang *et al.* [10] present a method for spotting words in natural images. They first perform
130 character detection for every letter in an alphabet and then evaluate the configuration scores for the words
131 in a specified list to pick out the most probable one. In this method, character detection is approached
132 using a nearest neighbor classifier in a sliding window fashion.

133 **Component-based methods** (e.g. [4], [9], [14], [33]) first extract candidate text components through
134 various ways (e.g. color reduction [4], [14] and Maximally Stable Extremal Region detection [11],
135 [22]) and then eliminate non-text components using heuristic rules or trained classifier, based on their
136 geometrical properties. Component-based methods are usually more efficient than texture-based methods
137 because the number of candidate components is relatively small. These methods are more robust to the
138 variations of texts, such as changes of font, scale and orientation. Moreover, the detected text components
139 can be directly used for character recognition. Due to these advantages, recent progresses in text detection

140 and recognition in natural images have been largely advanced by this category of methods [9], [11], [14],
141 [22], [33]–[35].

142 In [4], color reduction and multivalued image decomposition are performed to partition the input image
143 into multiple foreground regions. Connected component analysis is applied to these foreground regions,
144 followed by a text identification module, to filter out non-text regions.

145 The great success of sparse representation in face recognition [36] and image denoising [37] has inspired
146 numerous researchers in the community. The authors of [38] and [12] apply classification procedure to
147 candidate text components, using learned discriminative dictionaries.

148 The MSER-based methods [11], [22], [35] have attracted much attention from the community, because
149 of the excellent characteristics of MSERs (Maximally Stable Extremal Regions) [39]. MSERs can be
150 computed efficiently (near linear complexity) and are robust to noise and affine transformation. In [11],
151 MSERs are detected and taken as candidate text components. Neumann *et al.* [35] modified the original
152 MSER algorithm to take region topology into consideration, leading to superior detection performance.
153 Chen *et al.* [22] also proposed an extension to MSER, in which the boundaries of MSERs are enhanced
154 via Canny edge detection [40], to cope with image blur.

155 Epshtein *et al.* [9] proposed a novel image operator, called Stroke Width Transform (SWT), which
156 transforms the image data from containing color values per pixel to containing the most likely stroke
157 width. Based on SWT and a set of heuristic rules, this algorithm can reliably detect horizontal texts from
158 natural images.

159 While most existing algorithms are designed for horizontal or near-horizontal texts, Yi *et al.* [14] and
160 Shivakumara *et al.* [16] consider the problem of detecting multi-oriented texts in images or video frames.
161 After extracting candidate components using gradient and color based partition, the line grouping strategy
162 in [14] aggregates the components into text strings. The text strings can be in any direction. However,
163 the method of [14] relies on a large set of manually defined rules and thresholds. In [16], candidate
164 text regions are identified by K -means clustering in the Fourier-Laplacian domain. The region clusters
165 are divided into separate components using skeletonization. Even though this method can detect texts of
166 arbitrary directions, it only detects text blocks, rather than characters or words, because of the grouping
167 mechanism adopted in this method.

168 Finally, **hybrid methods** (e.g. [13], [41]) are a mixture of texture-based and component-based methods.
169 In [41], edge pixels of all possible text regions are extracted, using an elaborate edge detection method;
170 the gradient and geometrical properties of region contours are verified to generate candidate text regions,
171 followed by a texture analysis procedure to separate true text regions from non-text regions.

172 Unlike [41], the hybrid method proposed by Pan *et al.* [13] extracts candidate text components from
173 probability maps at multiple scales. The probability maps are estimated by a classifier, which is trained
174 using a set of texture features (HOG features [42]) computed in predefined patterns. A conditional random
175 field (CRF) model, combining unary component properties and binary contextual relationships, is utilized
176 to discriminate text components from non-text components.

177 *2) Our strategy:*

178 We have drawn two observations about the current text detection algorithms: (1) methods that are
179 purely based on learning (nearly black-box) [7] by training classifiers on a large amount of data can reach
180 certain but limited level of success (system [7] obtained from the authors produces reasonable results on
181 horizontal English texts but has poor performances in general cases); (2) systems that are based on smart
182 features, such as Stroke Width Transform (SWT) [9], are robust to variations in texts but they involve a
183 lot tuning and are still far from producing all satisfactory results, especially for non-horizontal texts.

184 In this paper, we adopt SWT and also design various new features that are intrinsic to texts and robust
185 to variations; a two-level classification scheme is devised to moderately utilize training to remove sensitive
186 parameter tuning by hand. We observe significant improvement over the existing approaches in dealing
187 with real-world scenes.

188 Though widely used in the community, the ICDAR datasets [43]–[45] only contain mostly horizontal
189 English texts. In [14], a dataset with texts of different directions is released, but it includes only 89
190 images without enough diversity in the texts and backgrounds. Here we collect a new dataset with 500
191 images of indoor and outdoor scenes. In addition, the evaluation methods used in [46] and the ICDAR
192 competitions [43]–[45] are mainly designed for horizontal texts. Hence, we propose a new protocol that
193 is more suitable to handle texts of arbitrary orientations (see Sec. III).

194 *B. Proposed Approach*

195 The proposed algorithm consists of four stages: (1) component extraction, (2) component analysis, (3)
196 candidate linking, and (4) chain analysis, which can be further categorized into two procedures, bottom-
197 up grouping and top-down pruning, as shown in Fig. 2. In the bottom-up grouping procedure, pixels are
198 first grouped into connected components and later these connected components are aggregated to form
199 chains; in the top-down pruning procedure non-text components and chains are successively identified
200 and eliminated. These two procedures are applied alternately when detecting text in images.

201 **Component extraction:** At this stage, edge detection is performed on the original image and the edge
202 map is fed to the SWT [9] module to produce an SWT image. Neighboring pixels in the SWT image

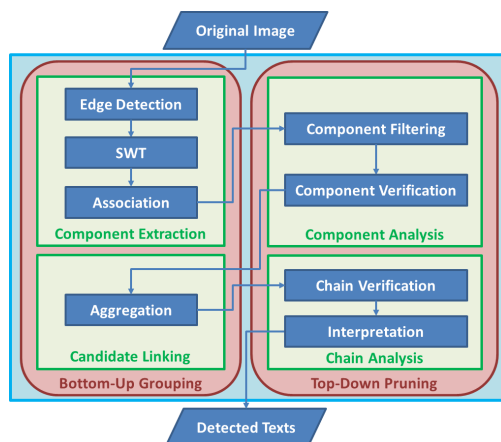


Fig. 2. Pipeline of the proposed approach.

203 are grouped together recursively to form connected components using a simple association rule.

204 **Component analysis:** Many components extracted at the component extraction stage are not parts of
 205 texts. The component analysis stage is aimed to identify and filter out those non-text components. First,
 206 the components are filtered using a set of heuristic rules that can distinguish between obvious spurious
 207 text regions and true text regions. Next, a component level classifier is applied to prune the non-text
 208 components that are hard for the simple filter.

209 **Candidate linking:** The remaining components are taken as character candidates¹. The first step of the
 210 candidate linking stage is to link the character candidates into pairs. Two adjacent candidates are grouped
 211 into a pair if they have similar geometric properties and colors. At the next step, the candidate pairs are
 212 aggregated into chains in a recursive manner.

213 **Chain analysis:** At the chain analysis stage, the chains determined at the former stage are verified by a
 214 chain level classifier. The chains with low classification scores (probabilities) are discarded. The chains
 215 may be in any direction, so a candidate might belong to multiple chains; the interpretation step is aimed
 216 to dispel this ambiguity. The chains that pass this stage are the final detected texts.

217 The remainder of this paper is organized as follows. Section II presents the details of the proposed
 218 method, including the algorithm pipeline and the two sets of features. Section III introduces the proposed
 219 dataset and evaluation protocol. The experimental results and discussions are given in Section IV.
 220 Section V concludes the paper and points out potential directions for future research.

¹In fact, components do not necessarily correspond to characters, because a single character in some languages may consist of several strokes; however, we still call them characters (or character candidates) hereafter for simplicity.

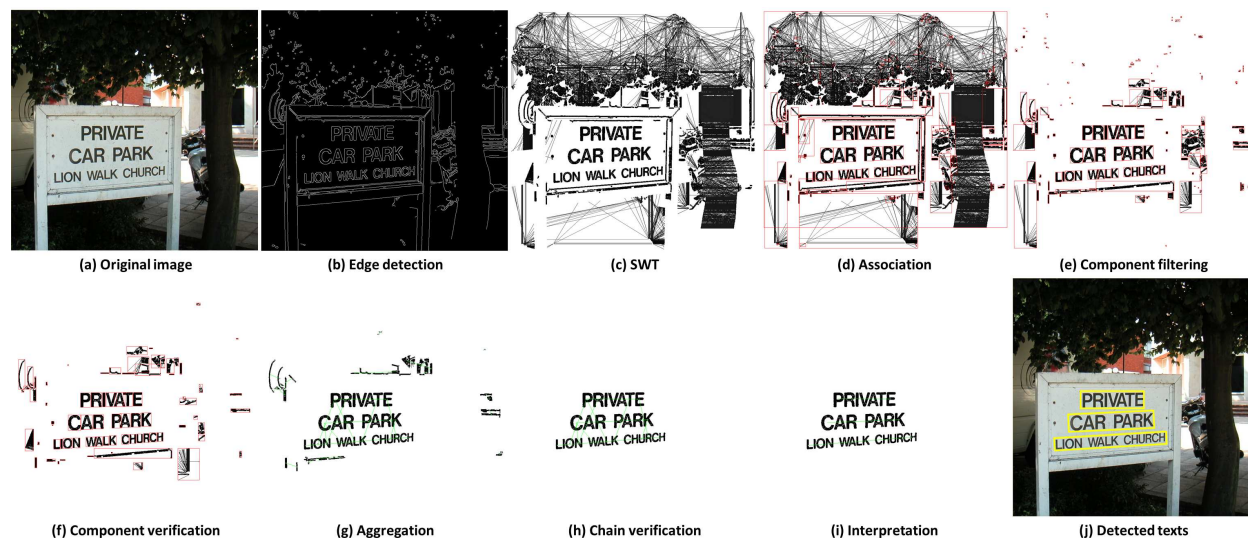


Fig. 3. Text detection process. See text for details.

II. METHODOLOGY

In this section, we present the details of the proposed algorithm. Specifically, the pipeline of the algorithm will be presented in Section II-A and the details of the features will be described in Section II-B.

A. Algorithm Pipeline

1) Component Extraction:

To extract connected components from the image, SWT [9] is adopted for its effectiveness and efficiency. In addition, it provides a way to discover connected components from edge map directly, which makes it unnecessary to consider the factors of scale and direction.

SWT runs on edge map, so we use Canny edge detector [40] to produce an edge map (Fig. 3 (b)) from the original image (Fig. 3 (a)). SWT is a local image operator which computes per pixel width of the most likely stroke containing the pixel. See [9] for details. The resulting SWT image is shown in Fig. 3 (c).

The next step of this stage is to group the pixels in the SWT image into connected components. The pixels are associated using a simple rule that the ratio of SWT values of neighboring pixels is less than 3.0. The connected components are shown in Fig. 3 (d). Note the red rectangles in the image; each rectangle contains a connected component.

TABLE I
BASIC COMPONENT PROPERTIES AND THEIR VALID RANGES.

Property	Definition	Range
width variation	$WV(c) = \frac{\sigma(c)}{\mu(c)}$	[0, 1]
aspect ratio	$AR(c) = \min\left\{\frac{w(c)}{h(c)}, \frac{h(c)}{w(c)}\right\}$	[0.1, 1]
occupation ratio	$OR(c) = \frac{q}{w(c) * h(c)}$	[0.1, 1]

237 2) *Component Analysis:*

238 The purpose of component analysis is to identify and eliminate the connected components that are
239 unlikely parts of texts. Towards this end, we devise a two-layer filtering mechanism.

240 The first layer is a filter consists of a set of heuristic rules. This filter runs on a collection of statistical
241 and geometric properties of components, which are very fast to compute. True text components usually
242 have nearly constant width and compact structure (not too thin and long), so width variation, aspect ratio
243 and occupation ratio are chosen as the basic properties to filter out obvious non-text components.

244 For a connected component c with q foreground pixels (black pixels in the SWT image), we first
245 compute its bounding box $bb(c)$ (its width and height are denoted by $w(c)$ and $h(c)$, respectively) and
246 the mean as well as standard deviation of the stroke widths, $\mu(c)$ and $\sigma(c)$. The definitions of these basic
247 properties and the corresponding valid ranges are summarized in Tab. I.

248 The components with one or more invalid properties will be taken as non-text regions and discarded.
249 This preliminary filter proves to be both effective and efficient. A large portion of obvious non-text
250 regions are eliminated after this step. Notice the difference between Fig. 3 (d) and Fig. 3 (e).

251 The second layer is a classifier trained to identify and reject the non-text components that are hard to
252 remove with the preliminary filter. A collection of component level features, which capture the differences
253 of geometric and textural properties between text components and non-text components, are used to
254 train this classifier. The criteria for feature design are: scale invariance, rotation invariance and low
255 computational cost. To meet these criteria, we propose to estimate the center, characteristic scale and
256 major orientation of each component (Fig. 4) before computing the component level features. Based on
257 these characteristics, features that are both effective and computational efficient can be obtained. The
258 details of these component level features are discussed in Sec. II-B1.

259 For a component c , the barycenter $o(c)$, major axis $L(c)$, minor axis $l(c)$, and orientation $\theta(c)$ are
260 estimated using Camshift [47] by taking the SWT image of component c as distribution map. The center,



Fig. 4. Component characteristics. The green points are the centers of the components. The radii of the pink circles represent their characteristic scales while the blue lines indicate the major orientations. The two images, which contain the same text line, are taken from different viewpoints and distances.

261 characteristic scale and major orientation of component c are defined as:

$$O(c) = o(c), \quad (1)$$

262

$$S(c) = L(c) + l(c), \quad (2)$$

263

$$\Theta(c) = \theta(c). \quad (3)$$

264 These characteristics are invariant to translation, scale and rotation to some degree (Fig. 4). As we will
 265 explain in Sec. II-B1, this is the key to the scale and rotation invariance of the component level features.

266 We train a component level classifier using the component level features. Random Forest [48] is chosen
 267 as the strong classifier. The component level classifier is the first level of the two-level classification
 268 scheme. The probability of component c , $p_1(c)$, is the fraction of votes for the positive class (text) from
 269 the trees. The components whose probabilities are lower than a threshold T_1 are eliminated and the
 270 remaining components are considered as character candidates (Fig. 3 (f)).

271 To ensure high recall in this stage, the threshold T_1 is set very low, because high threshold may filter
 272 out true text components.

273 3) Candidate Linking:

274 The character candidates are aggregated into chains at this stage. This stage also serves as a filtering
 275 step because the candidate characters cannot be linked into chains are taken as components accidentally
 276 formed by noises or background clutters, and thus are discarded.

277 Firstly, character candidates are linked into pairs. In [9], whether two candidates can be linked into a
 278 pair is determined based on the heights and widths of their bounding boxes. However, bounding boxes
 279 are not rotation invariant, so we use their characteristic scales instead. If two candidates have similar

280 stroke widths (ratio between the mean stroke widths is less than 2.0), similar sizes (ratio between their
 281 characteristic scales does not exceed 2.5), similar colors and are close enough (distance between them is
 282 less than two times the sum of their characteristic scales), they are labeled as a pair. The above parameters
 283 are optimized using the training data of the ICDAR datasets [43]–[45], however, this parameter setting
 284 turns out to be effective for all the datasets evaluated in this paper.

285 Unlike [9], which only considers horizontal linkings, the proposed algorithm allows linkings of arbitrary
 286 directions. This endows the system with the ability of detecting texts of arbitrary orientations, not limited
 287 to horizontal texts (see Fig. 1). Note that a character candidate may belong to several pairs.

288 Next, a greedy hierarchical agglomerative clustering [49] method is applied to aggregate the pairs
 289 into candidate chains. Initially, each pair constitutes a chain. Then the similarity between each couple of
 290 chains that share at least one common candidate and have similar orientations is computed; chains with
 291 the highest similarity are merged together to form a new chain. The orientation consistency $s_o(C_1, C_2)$
 292 and population consistency $s_p(C_1, C_2)$ between two chains C_1 and C_2 , which share at least one common
 293 candidate, are defined as:

$$s_o(C_1, C_2) = \begin{cases} 1 - \frac{\gamma(C_1, C_2)}{\pi/2} & \text{if } \gamma(C_1, C_2) \leq \Gamma \\ 0 & \text{otherwise} \end{cases}, \quad (4)$$

294 and

$$s_p(C_1, C_2) = \begin{cases} 1 - \frac{|n_{C_1} - n_{C_2}|}{|n_{C_1} + n_{C_2}|} & \text{if } \gamma(C_1, C_2) \leq \Gamma \\ 0 & \text{otherwise} \end{cases}, \quad (5)$$

295 where $\gamma(C_1, C_2)$ is the included angle of C_1 and C_2 while n_{C_1} and n_{C_2} are the candidate numbers of
 296 C_1 and C_2 . Γ is used to judge whether two chains have similar orientations and is empirically set to $\frac{\pi}{8}$.
 297 The similarity between two chains C_1 and C_2 is defined as the harmonic mean [50] of their orientation
 298 consistency and population consistency:

$$s(C_1, C_2) = \frac{2s_o(C_1, C_2)s_p(C_1, C_2)}{s_o(C_1, C_2) + s_p(C_1, C_2)}. \quad (6)$$

299 According to this similarity definition, the chains with proximal sizes and orientations are merged with
 300 priority. This merging process proceeds until no chains can be merged.

301 At last, the character candidates not belonging to any chain are discarded. The candidate chains after
 302 aggregation are shown in Fig. 3 (g). Each green line represents a candidate chain.

303 4) Chain Analysis:

304 The candidate chains formed at the previous stage might include false positives that are random
 305 combinations of scattered background clutters (such as leaves and grasses) and repeated patterns (such

306 as bricks and windows). To eliminate these false positives, a chain level classifier is trained using the
 307 chain level features (Sec. II-B2).

308 Random Forest [48] is again used. The chain level classifier is the second level of the two-level
 309 classification scheme. The probability of chain C , $p_2(C)$, is the fraction of votes for the positive class
 310 (text) from the trees. The chains whose probabilities are lower than a threshold T_2 are eliminated.

311 To make better decisions, the total probability of each chain is also calculated. For a chain C with n
 312 candidates $c_i, i = 1, 2, \dots, n$, the total probability is defined as:

$$p(C) = \left(\frac{\sum_{i=1}^n p_1(c_i)}{n} + p_2(C) \right) / 2. \quad (7)$$

313 The chains whose total probabilities are lower than a threshold T are discarded.

314 As texts of arbitrary orientations are considered, the remaining chains may be in any direction.
 315 Therefore, a candidate might belong to multiple chains. For example, in Fig. 3 (h) the character ‘P’ in
 316 the first line is linked in three chains (note the green lines). In reality, however, a character is unlikely to
 317 belong to multiple text lines. If several chains compete for the same candidate, only the chain with the
 318 highest total probability will survive (note the difference between Fig. 3 (h) and Fig. 3 (i)).

319 The survived chains are outputted by the system as detected texts (Fig. 3 (j)). For each detected text, its
 320 orientation is calculated through linear least squares [49] using the centers of the characters; its minimum
 321 area rectangle [51] is estimated using the orientation and the bounding boxes of the characters. Word
 322 partition, which divides text lines into separate words, is also implemented in the proposed algorithm;
 323 but it is not shown in Fig. 3 since the general task of text detection does not require this step.

324 The whole algorithm described above is performed twice to handle both bright text on dark background
 325 and dark text on bright background, once along the gradient direction and once along the inverse direction.
 326 The results of two passes are fused to make final decisions. For clarity, only the results of one pass are
 327 presented in Fig. 3.

328 *B. Feature Design*

329 We design two collections of features, component level features and chain level features, for classifying
 330 text and non-text, based on the observation that it is the median degree of regularities of text rather than
 331 particular color or shape that distinguish it from non-text, which usually has either low degree (random
 332 clutters) or high degree (repeated patterns) of regularities. At character level, the regularities of text come
 333 from nearly constant width and texturelessness of strokes, and piecewise smoothness of stroke boundaries;
 334 at line level, the regularities of text are similar colors, sizes, orientations and structures of characters, and
 335 nearly constant spacing between consecutive characters.

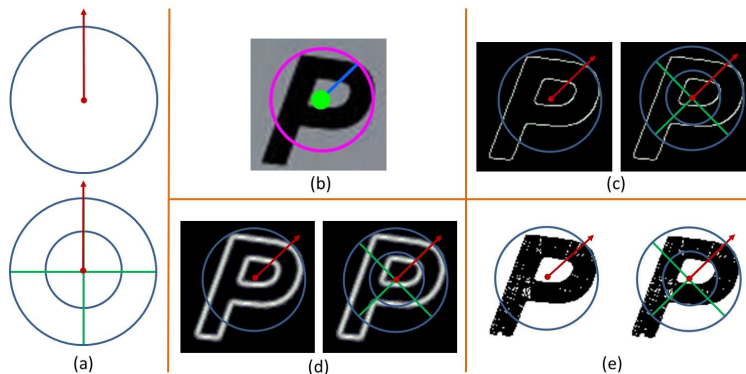


Fig. 5. Templates and calculation of scalable rotative descriptors. (a) Two templates used for computing the descriptors. The radius space and angle space are partitioned evenly in a coarse-to-fine manner. The red arrows indicate the reference orientations of the templates. (b) Component and its characteristics. (c)(d)(e) Calculation of contour shape, edge shape and occupation ratio. See text for details.

1) Component Level Features:

Inspired by Shape Context [52] and Feature Context [53], we devise two templates (Fig. 5 (a)) to capture the regularities of each component in coarse and fine granularity, respectively. The radius and orientation of the templates are not stationary, but adaptive to the component. When computing descriptors for a component, each template is placed at the center and rotated to align with the major orientation of the component; the radius is set to the characteristic scale of the component. Different cues from the sectors are encoded and concatenated into histograms. In this paper, the following cues are considered for each sector:

- **Contour shape** [54]. Contour shape is a histogram of oriented gradients. The gradients are computed on the component contour (Fig. 5 (c)).

- **Edge shape** [54]. Edge shape is also a histogram of oriented gradients; but the gradients are computed at all the pixels in the sector (Fig. 5 (d)).

- **Occupation ratio**. Occupation ratio is defined as the ratio between the number of the foreground pixels of the component within the sector and the sector area (Fig. 5 (e)).

To achieve rotation invariance, the gradient orientations are rotated by an angle $\Theta(c)$, before computing contour shape and edge shape. Then, the gradient orientations are normalized to the range $[0, \pi]$. Six orientation bins are used for computing histograms of contour shape and edge shape, to cope with different fonts and local deformations. For each cue, the signals computed in all the sectors of all the templates are concatenated to form a descriptor. We call these descriptors scalable rotative descriptors, because they are

355 computed based on templates that are scalable and rotative. Scalable rotative descriptors are similar to
 356 PHOG [55], as they both adopt spatial pyramid representation [56]. Different from the templates used for
 357 computing PHOG, our templates are circular and their scale and orientation are adaptive to the component
 358 being described. This is the key to the scale and rotation invariance of these descriptors.

359 The characteristic scale is crucial for the computation of scalable rotative descriptors because it deter-
 360 mines the scales of the templates. Too small templates may miss important information of components
 361 while too large templates will introduce noises and interferences from other components and background.
 362 The value of characteristic scale calculated using Eqn. 2 is a good trade-off in practice.

363 We found through experiments (not shown in this paper) that using finer templates can slightly improve
 364 the performance, but will largely increase the computational burden.

365 In addition, another three types of features are considered:

366 – **Axial ratio.** Axial ratio is computed by dividing the major axis of the component c with its minor
 367 axis: $XR(c) = L(c)/l(c)$.

368 – **Width variation.** This feature is the same as defined in Tab. I.

369 – **Density.** The density of component c is defined as the ratio between its pixel number q and character-
 370 istic area (here the characteristic area is $\pi \cdot S^2(c)$, not the area of the bounding box): $D(c) = q/(\pi \cdot S^2(c))$.

371 2) *Chain Level Features:*

372 Eleven types of chain level features, which are not specific to rotation and scale, are designed to
 373 discriminate text lines from false positives (mostly repeated patterns and random clutters) that cannot be
 374 distinguished by the component level features.

375 For a candidate chain C with n ($n \geq 2$) candidates $c_i, i = 1, 2, \dots, n$, the features are defined as below
 376 and summarized in Tab. II:

377 – **Candidate count.** This feature is adopted based on the observation that false positives usually have
 378 very few (for random clutters) or too many (for repeated patterns) candidates.

379 – **Average probability.** The probabilities given by the component level classifier are reliable. This
 380 feature is the average of all the probabilities ($p_1(c_i), i = 1, 2, \dots, n$) of the candidates belonging to C .

381 – **Average turning angle.** Most texts present in linear form, so for a text line the mean of the turning
 382 angles at the interior characters ($\tau(c_i), i = 2, 3, \dots, n - 1$) is very small; however, for random clutters
 383 this property will not hold. $\tau(c_i)$ is the included angle between the line $\overline{O(c_{i-1})O(c_i)}$ and $\overline{O(c_i)O(c_{i+1})}$.

384 – **Size variation.** In most cases characters in a text line have approximately equal sizes; but it's not
 385 that case for random clutters. The size of each component is measured by its characteristic scale $S(c_i)$.

386 – **Distance variation.** Another property of text is that characters in a text line are distributed uni-

387 formly, i.e. the distances between consecutive characters have small deviation. The distance between two
 388 consecutive components is the distance of their centers $O(c_{i-1})$ and $O(c_i)$.

389 – **Average direction bias.** For most text lines, the major orientations of the characters are nearly
 390 perpendicular to the major orientation of the text line. Direction bias of component c_i , $\beta(c_i)$, is the
 391 included angle between c_i and the orientation of the chain.

392 – **Average axial ratio.** Some repeated patterns (e.g. barriers) that are not texts consist of long and thin
 393 components, this feature can help differentiate them from true texts.

394 – **Average density.** On the contrary, other repeated patterns (e.g. bricks) consist of short and fat
 395 components, this feature can be used to eliminate this kind of false positives.

396 – **Average width variation.** False positives formed by foliage usually have varying widths while
 397 texts have constant widths. This feature is defined as the mean of all the width variation values of the
 398 candidates.

399 – **Average color self-similarity.** Characters in a text line usually have similar but not identical color
 400 distributions with each other; yet in false positive chains, color self-similarities [57] of the candidates are
 401 either too high (repeated patterns) or too low (random clutters). The color similarity $cs(x, y)$ is defined
 402 as the cosine similarity of the color histograms of the two candidates x and y .

403 – **Average structure self-similarity.** Likewise, characters in a text line have similar structure with each
 404 other while false positives usually have almost the same structure (repeated patterns) or diverse structures
 405 (random clutters). The structure similarity $ss(x, y)$ is defined as the cosine similarity of the edge shape
 406 descriptors of the two components x and y .

407 III. DATASET AND EVALUATION PROTOCOL

408 In this section, we introduce a large dataset for evaluating text detection algorithms, which contains 500
 409 natural images with real-world complexity. In addition, a new evaluation methodology which is suitable
 410 for benchmarking algorithms designed for texts of arbitrary directions is proposed.

411 A. Dataset

412 Although widely used in the community, the ICDAR datasets [43], [44] have two major drawbacks.
 413 First, most of the text lines (or single characters) in the ICDAR datasets are horizontal. In real scenarios,
 414 however, text may appear in any orientation. The second drawback is that all the text lines or characters
 415 in this dataset are in English. Therefore it is unable to use these datasets to assess systems designed for
 416 multilingual data.

TABLE II
CHAIN LEVEL FEATURES.

Feature	Definition
Candidate count	$CC(C) = n$
Average probability	$AP(C) = \frac{\sum_{i=1}^n p_1(c_i)}{n}$
Average turning angle	$ATA(C) = \begin{cases} \frac{\sum_{i=2}^{n-1} \tau(c_i)}{n-2} & \text{if } n > 2 \\ 0 & \text{otherwise} \end{cases}$
Size variation	$SV(C) = \frac{\sigma_s(c)}{\mu_s(c)}$ $\mu_s(C) = \frac{\sum_{i=1}^n S(c_i)}{n}$ $\sigma_s(C) = \sqrt{\frac{\sum_{i=1}^n (S(c_i) - \mu_s(c))^2}{n-1}}$
Distance variation	$DV(C) = \frac{\sigma_d(c)}{\mu_d(c)}$ $\mu_d(C) = \frac{\sum_{i=2}^n d(c_{i-1}, c_i)}{n-1}$ $\sigma_d(C) = \sqrt{\frac{\sum_{i=2}^n (d(c_{i-1}, c_i) - \mu_d(c))^2}{n-1}}$
Average direction bias	$ADB(C) = \frac{\sum_{i=1}^n \beta(c_i)}{n}$
Average axial ratio	$AAR(C) = \frac{\sum_{i=1}^n X R(c_i)}{n}$
Average density	$AD(C) = \frac{\sum_{i=1}^n D(c_i)}{n}$
Average width variation	$AWV(C) = \frac{\sum_{i=1}^n WV(c_i)}{n}$
Average color self-similarity	$ACS(C) = \frac{\sum_{i=1}^n CS(c_i)}{n}$ $CS(c_i) = \frac{\sum_{k \neq i} cs(c_k, c_i)}{n-1}$
Average structure self-similarity	$ASS(C) = \frac{\sum_{i=1}^n SS(c_i)}{n}$ $SS(c_i) = \frac{\sum_{k \neq i} ss(c_k, c_i)}{n-1}$

417 These two shortcomings are pointed out in [13], [14]. Two separate datasets are created: one contains
 418 non-horizontal text lines [13] and the other one is a multilingual dataset [14]. In this work, we generate
 419 a new multilingual text image dataset with horizontal as well as slant and skewed texts. We name this
 420 dataset MSRA Text Detection 500 Database (MSRA-TD500)², because it contains 500 natural images in
 421 total. These images are taken from indoor (office and mall) and outdoor (street) scenes using a packet
 422 camera. The indoor images are mainly signs, doorplates and caution plates while the outdoor images are

²http://users.loni.ucla.edu/~ztu/Download_front.htm

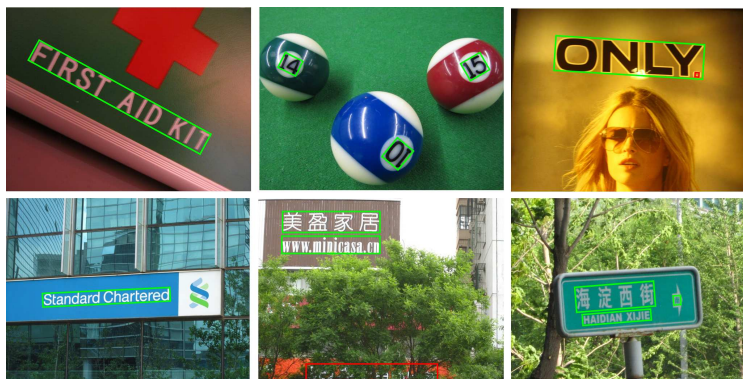


Fig. 6. Typical images from the proposed dataset along with ground truth rectangles. Notice the red rectangles. They indicate the texts within them are labeled as difficult (due to blur or occlusion).

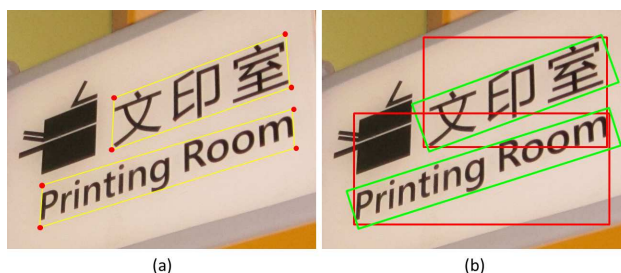


Fig. 7. Ground truth generation. (a) Human annotations. The annotators are required to locate and bound each text line using a four-vertex polygon (red dots and yellow lines). (b) Ground truth rectangles (green). The ground truth rectangle is generated automatically by fitting a minimum area rectangle using the polygon.

423 mostly guide boards and billboards in complex background. The resolutions of the images vary from
 424 1296×864 to 1920×1280 .

425 The MSRA-TD500 dataset is very challenging because of both the diversity of the texts and the
 426 complexity of the backgrounds in the images. The texts may be in different languages (Chinese, English
 427 or mixture of both), fonts, sizes, colors and orientations. The backgrounds may contain vegetation (e.g.
 428 trees and grasses) and repeated patterns (e.g. windows and bricks), which are not so distinguishable from
 429 text.

430 Some typical images from this dataset are shown in Fig. 6. It is worth mentioning that even though the
 431 purpose of this dataset is to evaluate text detection algorithms designed for texts of arbitrary orientations,
 432 horizontal and near-horizontal texts still dominate the dataset because these are the most common cases
 433 in practice.

434 The dataset is divided into two parts: training set and test set. The training set contains 300 images
 435 randomly selected from the original dataset and the rest 200 images constitute the test set. All the images
 436 in this dataset are fully annotated. The basic unit in this dataset is text line rather than word, which is
 437 used in the ICDAR dataset, because it is hard to partition Chinese text lines into individual words based
 438 on their spacings; even for English text lines, it is non-trivial to perform word partition without high
 439 level information. The procedure of ground truth generation is shown in Fig. 7.

440 B. Evaluation Protocol

441 Before presenting our novel evaluation protocol for text detection, we first introduce the evaluation
 442 method used in the ICDAR competitions as background. Under the ICDAR evaluation protocol, the
 443 performance of an algorithm is measured by F-measure, which is the harmonic mean of precision and
 444 recall. Different from the standard information retrieval measures of precision and recall, more flexible
 445 definitions are adopted in the ICDAR competitions [43], [44]. The match m between two rectangles is
 446 defined as the ratio of the area of intersection and that of the minimum bounding rectangle containing
 447 both rectangles. The set of rectangles estimated by each algorithm are called *estimates* while the set of
 448 ground truth rectangles provided in the ICDAR dataset are called *targets*. For each rectangle, the match
 449 with the largest value is found. Hence, the best match for a rectangle r in a set of rectangles R is defined
 450 as:

$$m(r; R) = \max\{m(r, r') | r' \in R\}. \quad (8)$$

451 Then, the definitions of precision and recall are:

$$precision = \frac{\sum_{r_e \in E} m(r_e; T)}{|E|}, \quad (9)$$

$$recall = \frac{\sum_{r_t \in T} m(r_t; E)}{|T|}, \quad (10)$$

452 where E and T are the sets of ground truth rectangles and estimated rectangles, respectively. The F-
 453 measure, which is a single measure of algorithm performance, is a combination of the two above measures.
 454 The relative weights of precision and recall are controlled by a parameter α , which is set to 0.5 to give
 455 equal weights to precision and recall:

$$f = \frac{1}{\frac{\alpha}{precision} + \frac{1-\alpha}{recall}}. \quad (11)$$

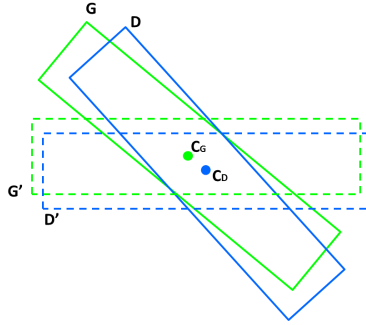


Fig. 8. Calculation of overlap ratio between detection rectangle and ground truth rectangle.

456 Minimum area rectangles [51] are used in our protocol because they (green rectangles in Fig. 7 (b))
 457 are much tighter and more accurate than axis-aligned rectangles (red rectangles in Fig. 7 (b)). However,
 458 a problem imposed by using minimum area rectangles is that it is difficult to judge whether a text line
 459 is correctly detected. As shown in Fig. 8, it is not trivial to directly compute the overlap ratio between
 460 the estimated rectangle D and the ground truth rectangle G . Instead, we calculate the overlap ratio using
 461 axis-aligned rectangles G' and D' , which are obtained by rotating G and D round their centers C_G and
 462 C_D , respectively. The overlap ratio between G and D is defined as:

$$m(G, D) = \frac{A(G' \cap D')}{A(G' \cup D')}, \quad (12)$$

463 where $A(G' \cap D')$ and $A(G' \cup D')$ denote the areas of the intersection and union of G' and D' . Obviously,
 464 the overlap ratio computed in this way is not accurate. Besides, the ground truth rectangles annotated
 465 are not accurate either, especially when the texts are skewed. Because of the imprecision of both ground
 466 truth and computed overlap ratio, the definitions of precision and recall used in the ICDAR protocol do
 467 not apply. Alternatively, we return to their original definitions.

468 Similar to the evaluation method for the PASCAL object detection task [58], in our protocol detections
 469 are considered true or false positives based on the overlap ratio between the estimated minimum area
 470 rectangles and the ground truth rectangles. If the included angle of the estimated rectangle and the ground
 471 truth rectangle is less than $\pi/8$ and their overlap ratio exceeds 0.5, the estimated rectangle is considered
 472 a correct detection. Multiple detections of the same text line are taken as false positives. The definitions
 473 of precision and recall are:

$$precision = \frac{|TP|}{|E|}, \quad (13)$$

$$recall = \frac{|TP|}{|T|}, \quad (14)$$

474 where TP is the set of true positive detections while E and T are the sets of estimated rectangles and
 475 ground truth rectangles.

476 Moreover, to accommodate difficult texts (too small, occluded, blurry, or truncated) that are hard for
 477 text detection algorithms, we introduce an elastic mechanism which can tolerate detection misses of
 478 difficult texts. The basic criterion of this elastic mechanism is: *if the difficult texts are detected by an*
 479 *algorithm, it counts; otherwise, the algorithm will not be punished.* Accordingly, the annotations of the
 480 images in the proposed dataset should be changed. Each text line considered to be difficult is given an
 481 additional “difficult” label (Fig. 6). Thus the ground truth rectangles can be categorized into two sub
 482 sets: ordinary sub set T_o and difficult sub set T_d ; ditto, the true positives TP can also be categorized into
 483 ordinary sub set TP_o , which is the set of rectangles matched with T_o , and ordinary sub set TP_d , which
 484 is the set of rectangles matched with T_d . After incorporating the elastic mechanism, the definitions of
 485 precision and recall become:

$$precision = \frac{|TP_o| + |TP_d|}{|E|} = \frac{|TP|}{|E|}, \quad (15)$$

$$recall = \frac{|TP_o| + |TP_d|}{|T_o| + |T_d|} = \frac{|TP|}{|T_o| + |T_d|}. \quad (16)$$

486 IV. EXPERIMENTS AND DISCUSSIONS

487 We implemented the proposed algorithm in C++ and evaluated it on a common server (2.53GHz CPU,
 488 48G RAM and Windows 64-bit OS). 200 trees are used for training the component level classifier and
 489 100 trees for the chain level classifier. The threshold values are: $T_1 = 0.1$, $T_2 = 0.3$ and $T = 0.4$. We
 490 found empirically that the text detectors under this parameter setting work well for all the datasets used in
 491 this paper. We believe better performances could be achieved by tuning the parameters for each dataset.

492 A. Results on Horizontal Texts

493 In order to compare the proposed algorithm with existing methods designed for horizontal texts, we
 494 evaluated the algorithm on the standard dataset used in the ICDAR 2003 Rubust Reading Competition [43]
 495 and the ICDAR 2005 Text Locating Competition [44]. This dataset contains 509 fully annotated text



Fig. 9. Detected texts in images from the ICDAR test set.

TABLE III

PERFORMANCES OF DIFFERENT TEXT DETECTION METHODS EVALUATED ON THE ICDAR TEST SET.

Algorithm	Precision	Recall	F-measure
TD-ICDAR	0.68	0.66	0.66
Epshtein <i>et al.</i> [9]	0.73	0.60	0.66
Yi <i>et al.</i> [14]	0.71	0.62	0.62
Becker <i>et al.</i> [44]	0.62	0.67	0.62
Chen <i>et al.</i> [7]	0.60	0.60	0.58
Zhu <i>et al.</i> [44]	0.33	0.40	0.33
Kim <i>et al.</i> [44]	0.22	0.28	0.22
Ezaki <i>et al.</i> [44]	0.18	0.36	0.22

496 images. 258 images from the dataset are used for training and 251 for testing. We trained a text detector
 497 (denoted by TD-ICDAR) on the training images.

498 Some detected texts of the proposed algorithm are presented in Fig. 9. Our algorithm can handle several
 499 types of challenging scenarios, e.g. variations in text font, color and size, as well as repeated patterns and
 500 background clutters. The quantitative comparison of different methods evaluated on the ICDAR test set
 501 is shown in Tab. III. Our algorithm compares favorably with the state-of-the-art algorithms when dealing
 502 with horizontal texts.

503 It is noted that existing algorithms seem to converge in performance (with F-measure around 0.66) on
 504 the ICDAR dataset. This might be due to three reasons: (1) the ICDAR evaluation method is different
 505 from the conventional methods for object detection (e.g. the PASCAL method). The ICDAR evaluation

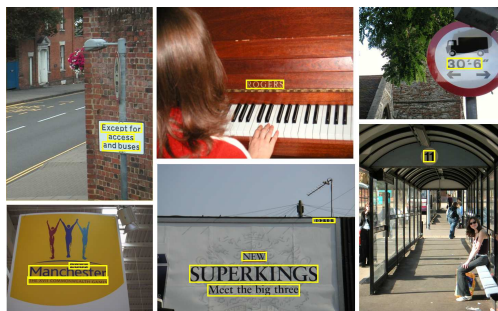


Fig. 10. Detected texts in images from the ICDAR 2011 test set.

TABLE IV

PERFORMANCES OF DIFFERENT TEXT DETECTION METHODS EVALUATED ON THE ICDAR 2011 DATASET [45].

Algorithm	Precision	Recall	F-measure
TD-ICDAR2011	0.7759	0.5549	0.6471
Kim <i>et al.</i> [45]	0.8298	0.6247	0.7128
Yi <i>et al.</i> [45]	0.6722	0.5809	0.6232
Yang <i>et al.</i> [45]	0.6697	0.5768	0.6198
Neumann <i>et al.</i> [45]	0.6893	0.5254	0.5963
Shao <i>et al.</i> [45]	0.6352	0.5352	0.5809
Guyomard <i>et al.</i> [45]	0.6297	0.5007	0.5578
Lee <i>et al.</i> [45]	0.5967	0.4457	0.5103
Sun <i>et al.</i> [45]	0.3501	0.3832	0.3659
Hanif <i>et al.</i> [45]	0.5505	0.2596	0.3419

506 method actually requires pixel-level accuracy (see Eqn. 9 and Eqn. 10), which is rigorous for detection
 507 algorithms, considering that the ground truth is given in the form of rough rectangles. (2) The ICDAR
 508 evaluation method requires word partition, that is, dividing text lines into individual words. This limits
 509 the scores of text detection algorithms either; because it is non-trivial to perform word partition without
 510 high level information. Moreover, the definitions of “word” are not consistent among different images.
 511 (3) Most algorithms assume that in the image a word or text line consists of at least two characters.
 512 However, in the ICDAR dataset some images contain single characters. In these images, most existing
 513 algorithms will fail to detect the single characters.

514 The ICDAR 2011 Robust Reading Competition Challenge 2 [45] was held to track the recent progress
 515 in the field of scene text detection and recognition. Due to the problems with the dataset used in the
 516 previous ICDAR competitions (for example, imprecise bounding boxes and inconsistent definitions of



Fig. 11. Detected texts in images from the proposed dataset. Yellow rectangles: true positives, pink rectangles: false negatives, red rectangles: false positives. Best viewed in color.

TABLE V

PERFORMANCES OF DIFFERENT TEXT DETECTION METHODS EVALUATED ON THE PROPOSED DATASET.

Algorithm	Precision	Recall	F-measure
TD-MSRA	0.63	0.63	0.60
TD-ICDAR	0.53	0.52	0.50
Epshtein <i>et al.</i> [9]	0.25	0.25	0.25
Chen <i>et al.</i> [7]	0.05	0.05	0.05

517 “word”), the dataset in the ICDAR 2011 competition is extended and relabeled [45]. Moreover, the
 518 evaluation method proposed by Wolf *et al.* [59] was adopted as the standard for performance evaluation,
 519 to replace the previous evaluation protocol which is unable to handle the cases of one-to-many and
 520 many-to-many matches and thus consistently underestimates the capabilities of text detection algorithms.

521 To enable direct and fair comparison, we trained a text detector (denoted by TD-ICDAR2011) using
 522 the training set of the ICDAR 2011 competition dataset, performed text detection on the test images
 523 and measured the performance using the method of Wolf *et al.* [59]. Fig. 10 illustrates several detection
 524 examples of our method on this dataset. The quantitative results of different text detection methods
 525 evaluated on the ICDAR 2011 dataset are shown in Tab. IV. The proposed algorithm achieves the second
 526 highest F-measure on this dataset.

527 B. Results on Texts of Arbitrary Orientations

528 Besides the ICDAR datasets, we also trained a text detector (denoted by TD-MSRA) on the proposed
 529 dataset and compared it with the systems of Chen *et al.* [7] and Epshtein *et al.* [9]. Detection examples

TABLE VI
PERFORMANCES OF DIFFERENT TEXT DETECTION METHODS EVALUATED ON THE ORIENTED SCENE TEXT DATABASE (OSTD) [14].

Algorithm	Precision	Recall	F-measure
TD-MSRA	0.77	0.73	0.74
TD-ICDAR	0.71	0.69	0.68
Yi <i>et al.</i> [14]	0.56	0.64	0.55
Epshtein <i>et al.</i> [9]	0.37	0.32	0.32
Chen <i>et al.</i> [7]	0.07	0.06	0.06

530 of the proposed algorithm on this dataset are shown in Fig. 11. The proposed algorithm is able to detect
 531 texts of large variation in natural scenes, with the presence of vegetations and buildings. The images in
 532 the last row of Fig. 11 are some typical cases where our algorithms failed to detect the texts or gave false
 533 positives. The misses (pink rectangles) are mainly due to strong highlights, blur and low resolution; the
 534 false positives (red rectangles) are usually caused by windows, trees, or signs that are very alike text.

535 The performances are measured using the proposed evaluation protocol and shown in Tab. V. Compared
 536 with the competing algorithms, the proposed method achieves significantly enhanced performance when
 537 detecting texts of arbitrary orientations. The performances of other competing algorithms are not presented
 538 because of unavailability of their codes/executables. The average processing time of our algorithm on
 539 this dataset is 7.2s and that of Epshtein *et al.* is 6s. Our algorithm is a bit slower, but with the advantage
 540 of being able to detect texts of arbitrary orientations.

541 In [14], a dataset called Oriented Scene Text Database (OSTD), which contains texts of various
 542 orientations, is released. This dataset contains 89 images of logos, indoor scenes and street views. We
 543 perform text detection on all the images in this dataset. The quantitative results are presented in Tab. VI.
 544 Our method outperforms [14] on the Oriented Scene Text Database (OSTD), with an improvement of
 545 0.19 in F-measure.

546 From Tab. V and Tab. VI, we observe that even TD-ICDAR (only trained on horizontal texts) achieves
 547 much better performance than other methods on non-horizontal texts. It demonstrates the effectiveness
 548 of the proposed features.

549 C. Results on Texts of Different Languages

550 To further verify the ability of the proposed algorithm to detect texts of different languages, we
 551 collected a multilingual text image database from the Internet. The database contains 94 natural images



Fig. 12. Detected texts in various languages. The images are collected from the Internet.

TABLE VII

PERFORMANCES OF DIFFERENT TEXT DETECTION METHODS EVALUATED ON TEXTS OF DIFFERENT LANGUAGES.

Algorithm	Precision	Recall	F-measure
TD-MSRA	0.73	0.64	0.66
Epshtein <i>et al.</i> [9]	0.58	0.65	0.59
Chen <i>et al.</i> [7]	0.06	0.08	0.07

552 with texts of various languages, including both oriental and western languages, such as Japanese, Korean,
 553 Arabic, Greek, and Russian. We applied TD-MSRA to all the images in this database. Fig. 12 shows
 554 some detected texts in images from this database. The algorithms of Epshtein *et al.* [9] and Chen *et*
 555 *al.* [7] were adopted as baselines. The quantitative results of these algorithms are presented in Tab. VII.
 556 The proposed algorithm and the method of Epshtein *et al.* [9] both give excellent performance on this
 557 benchmark. Though TD-MSRA is only trained on Chinese and English texts, it can effortlessly generalize
 558 to texts in different languages. This indicates that the proposed algorithm is quite general and it can serve
 559 as a multilingual text detector.

560 D. Special Consideration on Single Characters

561 As pointed out in Sec. IV-A, most existing algorithms cannot handle single characters, since they assume
 562 that in the image a word or text line consists of at least two characters. To overcome this limitation, we
 563 modified the proposed algorithm to handle single characters. In the candidate linking stage, we no longer
 564 simply discard all single character candidates but retain the character candidates with high probabilities
 565 ($p_1(c) > 0.8$) instead, even if they do not belong to any chain. After this modification, the proposed
 566 algorithm is able to detect obvious single characters in natural images. Fig. 13 depicts some detected
 567 single characters by the proposed algorithm.



Fig. 13. Detected single characters in images (from the ICDAR dataset [43], [44]).

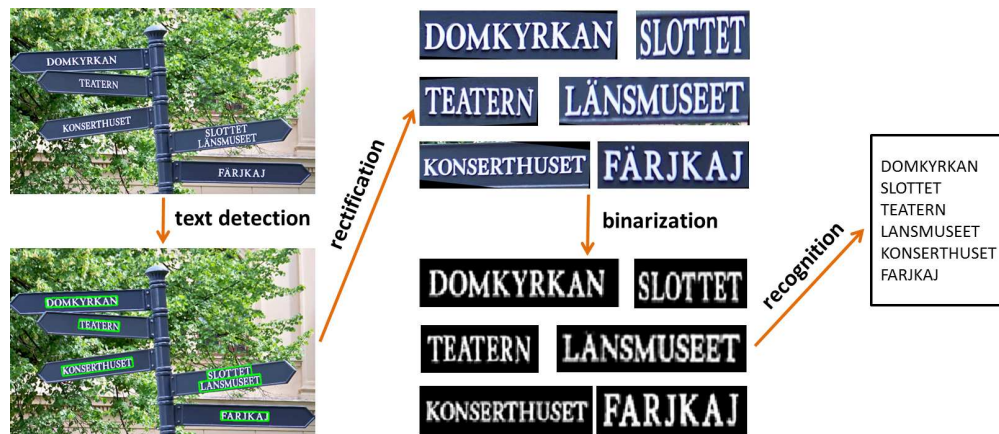


Fig. 14. Pipeline of our end-to-end scene text recognition system.

568 *E. End-to-End Scene Text Recognition*

569 Based on the proposed text detection algorithm, we developed an end-to-end scene text recognition
 570 system. The pipeline of this system is demonstrated in Fig. 14. After applying text detection to the
 571 original image, the detected text regions are rectified by a powerful low rank structure recovery technique
 572 (TILT [60]) and binarized, then these regions are fed into an off-the-shelf OCR software to produce the
 573 final recognition result. To evaluate the proposed end-to-end scene text recognition system, we generated

TABLE VIII
 END-TO-END SCENE TEXT RECOGNITION PERFORMANCES.

System	Precision	Recall	F-measure
Ours	0.58	0.51	0.53
Epshtein <i>et al.</i> [9]	0.57	0.49	0.51
Direct OCR	0.13	0.10	0.11

574 a dataset of 80 natural images with English texts and Arabic numbers. Majority of the images are from
575 the MSRA-TD500 database and the rest images are from the Internet.

576 For comparison, we tested the end-to-end text recognition system of Epshtein *et al.* [9] on this dataset.
577 To demonstrate how text detection can help effectively extract text information from natural images,
578 we also performed character recognition directly on the original images (denoted by Direct OCR). The
579 quantitative performances are computed at character level and shown in Tab. VIII. As can be seen, applying
580 OCR directly on natural images gives poor performance, because of the variation of texts and complex
581 background clutters. In contrast, both our scene text recognition system and that of Epshtein *et al.* [9]
582 achieve much higher performance on natural images. This suggests that text detection and rectification
583 are crucial steps when extracting texts from natural images.

584 V. CONCLUSIONS AND FUTURE WORK

585 We have presented a text detection system that detects texts of arbitrary directions in complex natural
586 scenes. Our system compares favorably with the state-of-the-art algorithms when handling horizontal
587 texts and achieves significantly enhanced performance on texts of arbitrary orientations. Furthermore, we
588 have proposed a multilingual database with horizontal as well as non-horizontal texts and specifically
589 designed an evaluation protocol for benchmarking algorithms designed for texts of arbitrary orientations.

590 As one might have noticed, the component level features are actually character descriptors that can
591 distinguish among different characters, thus they can be adopted to recognize characters. We plan to
592 make use of this property and develop an unified framework for text detection and character recognition
593 in the future.

594 ACKNOWLEDGMENT

595 REFERENCES

- 596 [1] G. N. DeSouza and A. C. Kak, "Vision for mobile robot navigation: A survey," *IEEE Trans. PAMI*, vol. 24, no. 2, pp.
597 237–267, 2002.
- 598 [2] S. Tsai, H. Chen, D. Chen, G. Schroth, R. Grzeszczuk, and B. Girod, "Mobile visual search on printed documents using
599 text and low bit-rate features," in *Proc. of ICIP*, 2011.
- 600 [3] B. Kisaanin, V. Pavlovi, and T. S. Huang, *Real-time vision for human-computer interaction*. Springer, 2005.
- 601 [4] A. Jain and B. Yu, "Automatic text location in images and video frames," *Pattern Recognition*, vol. 31, no. 12, pp.
602 2055–2076, 1998.
- 603 [5] Y. M. Y. Hasan and L. J. Karam, "Morphological text extraction from images," *IEEE Trans. Image Processing*, vol. 9,
604 no. 11, pp. 1978–1983, 2000.

- 605 [6] K. I. Kim, K. Jung, and J. H. Kim, "Texture-based approach for text detection in images using support vector machines
606 and continuously adaptive mean shift algorithm," *IEEE Trans. PAMI*, vol. 25, no. 12, pp. 1631–1639, 2003.
- 607 [7] X. Chen and A. Yuille, "Detecting and reading text in natural scenes," in *Proc. of CVPR*, 2004.
- 608 [8] D. Chen, J. M. Odobez, and H. Bourlard, "Text detection and recognition in images and video frames," *Pattern Recognition*,
609 vol. 37, no. 3, pp. 595–608, 2004.
- 610 [9] B. Epshtein, E. Ofek, and Y. Wexler, "Detecting text in natural scenes with stroke width transform," in *Proc. of CVPR*,
611 2010.
- 612 [10] K. Wang and S. Belongie, "Word spotting in the wild," in *Proc. ECCV*, 2010.
- 613 [11] L. Neumann and J. Matas, "A method for text localization and recognition in real-world images," in *Proc. of ACCV*, 2010.
- 614 [12] M. Zhao, S. T. Li, and J. Kwok, "Text detection in images using sparse representation with discriminative dictionaries,"
615 *IVC*, vol. 28, no. 12, pp. 1590–1599, 2010.
- 616 [13] Y. Pan, X. Hou, and C. Liu, "A hybrid approach to detect and localize texts in natural scene images," *IEEE Trans. Image
617 Processing*, vol. 20, no. 3, pp. 800–813, 2011.
- 618 [14] C. Yi and Y. Tian, "Text string detection from natural scenes by structure-based partition and grouping," *IEEE Trans.
619 Image Processing*, vol. 20, no. 9, pp. 2594–2605, 2011.
- 620 [15] A. Coates, B. Carpenter, C. Case, S. Satheesh, B. Suresh, T. Wang, and A. Y. Ng, "Text detection and character recognition
621 in scene images with unsupervised feature learning," in *Proc. of ICDAR*, 2011.
- 622 [16] P. Shivakumara, T. Q. Phan, and C. L. Tan, "A laplacian approach to multi-oriented text detection in video," *IEEE Trans.
623 PAMI*, vol. 33, no. 2, pp. 412–419, 2011.
- 624 [17] K. Wang, B. Babenko, and S. Belongie, "End-to-end scene text recognition," in *Proc. ICCV*, 2011.
- 625 [18] ABBYY, *ABBYY Mobile Products*, 2012, <http://www.abbyy.com/mobile/>.
- 626 [19] 3GVision, *3GVision's Business Card Reader Application*, 2012, <http://www.i-nigma.com/TextRecognition.html>.
- 627 [20] C. Yao, X. Bai, W. Liu, Y. Ma, and Z. Tu, "Detecting texts of arbitrary orientations in natural images," in *Proc. of CVPR*,
628 2012, to appear.
- 629 [21] J. Gllavata, R. Ewerth, and B. Freisleben, "Text detection in images based on unsupervised classification of high-frequency
630 wavelet coefficients," in *Proc. of ICPR*, 2004.
- 631 [22] H. Chen, S. Tsai, G. Schroth, D. Chen, R. Grzeszczuk, and B. Girod, "Robust text detection in natural images with
632 edge-enhanced maximally stable extremal regions," in *Proc. of ICIP*, 2011.
- 633 [23] X. Zhao, K. H. Lin, Y. Fu, Y. Hu, Y. Liu, and T. S. Huang, "Text from corners: A novel approach to detect text and
634 caption in videos," *IEEE Trans. Image Processing*, vol. 20, no. 3, pp. 790–799, 2011.
- 635 [24] K. Jung, K. Kim, and A. Jain, "Text information extraction in images and video: a survey," *PR*, vol. 37, no. 5, pp. 977–997,
636 2004.
- 637 [25] J. Liang, D. Doermann, and H. Li, "Camera-based analysis of text and documents: a survey," *IJDAR*, vol. 7, no. 2, pp.
638 84–104, 2005.
- 639 [26] Y. Zhong, K. Karu, and A. K. Jain, "Locating text in complex color images," *Pattern Recognition*, vol. 28, no. 10, pp.
640 1523–1535, 1995.
- 641 [27] V. Wu, R. Manmatha, and E. M. Riseman, "Finding text in images," in *Proc. of 2nd ACM Int. Conf. Digital Libraries*,
642 1997.
- 643 [28] H. P. Li, D. Doermann, and O. Kia, "Automatic text detection and tracking in digital video," *IEEE Trans. Image Processing*,
644 vol. 9, no. 1, pp. 147–156, 2000.

- 645 [29] Y. Zhong, H. Zhang, and A. K. Jain, "Automatic caption localization in compressed video," *IEEE Trans. PAMI*, vol. 22,
646 no. 4, pp. 385–392, 2000.
- 647 [30] R. Lienhart and A. Wernicke, "Localizing and segmenting text in images and videos," *IEEE Trans. CSVT*, vol. 12, no. 4,
648 pp. 256–268, 2002.
- 649 [31] J. Weinman, A. Hanson, and A. McCallum, "Sign detection in natural images with conditional random fields," in *Proc. of*
650 *WMLSP*, 2004.
- 651 [32] M. R. Lyu, J. Song, and M. Cai, "A comprehensive method for multilingual video text detection, localization, and extraction,"
652 *IEEE Trans. CSVT*, vol. 15, no. 2, pp. 243–255, 2005.
- 653 [33] A. Ikica and P. Peer, "An improved edge profile based method for text detection in images of natural scenes," in *Proc. of*
654 *EUROCON*, 2011.
- 655 [34] R. Minetto, N. Thome, M. Cord, J. Fabrizio, and B. Marcotegui, "Snoopertext: A multiresolution system for text detection
656 in complex visual scenes," in *Proc. of ICIP*, 2010.
- 657 [35] L. Neumann and J. Matas, "Text localization in real-world images using efficiently pruned exhaustive search," in *Proc. of*
658 *ICDAR*, 2011.
- 659 [36] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Trans.*
660 *PAMI*, vol. 31, no. 2, pp. 210–227, 2009.
- 661 [37] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Trans.*
662 *Image Processing*, vol. 15, no. 12, pp. 3736–3745, 2006.
- 663 [38] W. Pan, T. D. Bui, and C. Y. Suen, "Text detection from natural scene images using topographic maps and sparse
664 representations," in *Proc. of ICIP*, 2009.
- 665 [39] J. Matas, O. Chum, M. Urban, and T. Pajdla, "Robust wide baseline stereo from maximally stable extremal regions," in
666 *Proc. of BMVC*, 2002.
- 667 [40] J. F. Canny, "A computational approach to edge detection," *IEEE Trans. PAMI*, vol. 8, no. 6, pp. 679–698, 1986.
- 668 [41] Y. Liu, S. Goto, and T. Ikenaga, "A contour-based robust algorithm for text detection in color images," *IEICE Trans. Inf.*
669 *Syst.*, vol. E89-D, no. 3, pp. 1221–1230, 2006.
- 670 [42] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. of CVPR*, 2005.
- 671 [43] S. M. Lucas, A. Panaretos, L. Sosa, A. Tang, S. Wong, and R. Young, "Icdar 2003 robust reading competitions," in *Proc.*
672 *of ICDAR*, 2003.
- 673 [44] S. M. Lucas, "Icdar 2005 text locating competition results," in *Proc. of ICDAR*, 2005.
- 674 [45] A. Shahab, F. Shafait, and A. Dengel, "Icdar 2011 robust reading competition challenge 2: Reading text in scene images,"
675 in *Proc. of ICDAR*, 2011.
- 676 [46] X. S. Hua, W. Liu, and H. J. Zhang, "An automatic performance evaluation protocol for video text detection algorithms,"
677 *IEEE Trans. CSVT*, vol. 14, no. 4, pp. 498–507, 2004.
- 678 [47] G. R. Bradski, "Real time face and object tracking as a component of a perceptual user interface," in *Proc. of IEEE*
679 *Workshop on Applications of Computer Vision*, 1998.
- 680 [48] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- 681 [49] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*,
682 *Second Edition*. New York: Springer, 2009.
- 683 [50] C. V. Rijsbergen, *Information Retrieval, Second Edition*. London: Butterworths, 1979.

- 684 [51] H. Freeman and R. Shapira, "Determining the minimum-area encasing rectangle for an arbitrary closed curve," *Comm.*
685 *ACM*, vol. 18, no. 7, pp. 409–413, 1975.
- 686 [52] S. Belongie, J. Malik, and J. Puzicha, "Shape matching and object recognition using shape contexts," *IEEE Trans. PAMI*,
687 vol. 24, no. 4, pp. 509–522, 2002.
- 688 [53] X. Wang, X. Bai, W. Liu, and L. J. Latecki, "Feature context for image classification and object detection," in *Proc. of*
689 *CVPR*, 2010.
- 690 [54] C. Gu, J. Lim, P. Arbelaez, and J. Malik, "Recognition using regions," in *Proc. CVPR*, 2009.
- 691 [55] A. Bosch, A. Zisserman, and X. Munoz, "Representing shape with a spatial pyramid kernel," in *Proc. of CIVR*, 2007.
- 692 [56] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene
693 categories," in *Proc. of CVPR*, 2006.
- 694 [57] E. Shechtman and M. Irani, "Matching local self-similarities across images and videos," in *Proc. of CVPR*, 2007.
- 695 [58] M. Everingham, L. V. Gool, C. K. I. Williams, J. Winn, and A. Zisserman, "The PASCAL Visual Object Classes (VOC)
696 challenge," *IJCV*, vol. 88, no. 2, pp. 303–338, 2010.
- 697 [59] C. Wolf and J. M. Jolion, "Object count/area graphs for the evaluation of object detection and segmentation algorithms,"
698 *IJDAR*, vol. 8, no. 4, pp. 280–296, 2006.
- 699 [60] Z. Zhang, A. Ganesh, X. Liang, and Y. Ma, "TILT: Transform invariant low-rank textures," *IJCV*, vol. 99, no. 1, pp. 1–24,
700 2012.