

Robust Text Detection in Natural Scene Images by Generalized Color-enhanced Contrasting Extremal Region and Neural Networks

Lei Sun^{1,2*}, Qiang Huo², Wei Jia^{1,2*}, Kai Chen^{1,2*}

¹Department of Electronics Science and Technology
University of Science and Technology of China
Hefei, China

²Microsoft Research
Beijing, China
{v-lesun, qianghuo, v-kachen, v-weiji}@microsoft.com

Abstract—This paper presents a robust text detection approach based on generalized color-enhanced contrasting extremal region (CER) and neural networks. Given a color natural scene image, six component-trees are built from its grayscale image, hue and saturation channel images in a perception-based illumination invariant color space, and their inverted images, respectively. From each component-tree, generalized color-enhanced CERs are extracted as character candidates. By using a “divide-and-conquer” strategy, each candidate image patch is labeled reliably by rules as one of five types, namely, Long, Thin, Fill, Square-large and Square-small, and classified as text or non-text by a corresponding neural network, which is trained by an ambiguity-free learning strategy. After pruning non-text components, repeating components in each component-tree are pruned by using color and area information to obtain a component graph, from which candidate text-lines are formed and verified by another set of neural networks. Finally, results from six component-trees are combined, and a post-processing step is used to recover lost characters and split text lines into words as appropriate. Our proposed method achieves 85.72% recall, 87.03% precision, and 86.37% F-score on ICDAR-2013 “Reading Text in Scene Images” test set.

Keywords— text detection; natural scene image; generalized color-enhanced contrasting extremal region; neural networks

I. INTRODUCTION

Text detection in natural scene images has become a crucial task and received significant attention recently due to the great success of smart phones and large demands in content-based image search or understanding. Although many approaches (e.g., [1], [2], [3]) have been proposed, this problem remains largely unsolved, e.g., the winning team in ICDAR-2013 “Reading Text in Scene Images” competition achieved only a localization recall of about 66% [4]. The difficulties mainly come from diversities of texts (e.g., languages, font, size, color, orientation, noise, illumination, low contrast, occlusion and so on) as well as the complexity of the backgrounds [5].

Existing text detection methods can be categorized into three groups: sliding window based methods (e.g., [6], [7], [8]), connected component (CC) based methods (e.g., [1], [2], [5], [9]) and hybrid methods (e.g., [3], [10]). Among them, the extremal-region (ER) based methods, which belong to the connected component based methods, won the first places in both ICDAR-2011 and ICDAR-2013 competitions ([11], [4]). Despite their superior performance, several open problems need be addressed for ER methods. First, some text objects in images are not ERs whose pixels have either higher or lower intensity than its outer boundary pixels [12], and cannot be extracted by ER methods directly. Second, ER methods can

extract not only text components but also tremendous non-text components, including many ambiguous components. There is room for improvement in efficient and effective non-text pruning, and new methods are needed to address the ambiguity problem. Third, as discussed in [2], the repeating components in the hierarchical structure as shown in Fig. 1, affect severely candidate text-line formation. Fourth, building candidate text line is itself a challenging problem, especially when layout or background is complex. In this paper, generalized color-enhanced contrasting extremal region (CER) is proposed to address the first problem and an efficient and effective neural network based approach is proposed to address the other three problems. The main contribution of our paper is introduced briefly below and described in detail in subsequent sections.

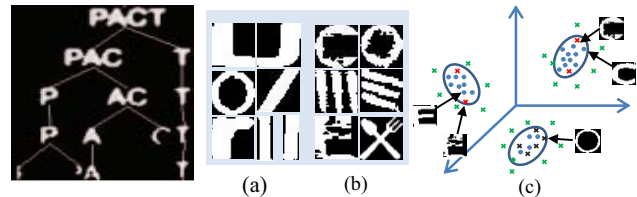


Fig. 1 The repeating MSER problem [2]. Fig. 2 Ambiguous samples and their distribution for illustrative purposes.

The first contribution is the generalization of our previously proposed candidate character generation method, i.e. color-enhanced CER [13]. Three new ideas are introduced. First, color-enhanced CER is generalized from perception-based HVC color space [14] to perception-based illumination invariant (PII) color space [15]. Color-enhanced CER is robust to noise but sensitive to illumination, which causes two problems: 1) If different parts in an original CER have different illumination, some parts will be lost after color enhancement; 2) If adjacent components have same color but different illumination, the grouping rule based on color will not be reliable. The new PII color space can solve both problems. Second, enhancement can be done iteratively so that the limitations of ER methods analyzed in Sec. II.C can be partially addressed. Third, CER can be enhanced not only by color, but also by other properties (e.g., stroke width [16]), to improve ER methods further.

The second contribution is to propose a “divide-and-conquer” strategy to solve the text and non-text classification problem. Classifying a connected component image patch as text or non-text is very challenging due to the ambiguity problem illustrated in Fig. 2. The ambiguity may include 1) the component itself is ambiguous (Fig. 2(a)); 2) the shapes of some non-text samples are very similar to characters (Fig. 2(b)).

*This work was done when Lei Sun, Wei Jia and Kai Chen were interns in Speech Group, Microsoft Research, Beijing, China.

A distribution of these two kinds of ambiguous samples is illustrated in Fig. 2(c). Existing solutions (e.g., [1], [2], [9]) mainly include or combine the following three steps: 1) Pre-prune non-text components with handcrafted features and classifiers (e.g., random forests, SVM); 2) Group remaining components into candidate text-lines; 3) Verify each candidate line. Although text-line information can reduce ambiguity and improve classification accuracy indeed, candidate text-line grouping itself is an open problem, especially when image layout or background is complex. So the performance of the first pre-pruning step is very crucial, as it determines the difficulty of the following candidate text line grouping problem as well as the performance of the whole system. We propose a “divide-and-conquer” strategy to address this problem. Given a binary image patch, it can be labeled reliably based on its aspect ratio and filled rate as one of five types, namely, Long, Thin, Fill, Square-large and Square-small, as illustrated in Fig. 8. For each type, a classifier can be built to answer the text and non-text question. Several benefits can be gained from this strategy: 1) Each classification problem is simpler; 2) Ambiguities are mainly restricted to Square-large sub-problem while the ambiguity issue in other sub-problems are reduced greatly; 3) The classifiers built for each sub-problem are complementary, so the characters can have a higher probability to be kept when they are included in several components with different types as shown in Fig. 1.

The third contribution is to propose an ambiguity-free learning strategy to address the ambiguity problem. The main idea is to use “bootstrap” strategy [17] to sample a small unambiguous representative active set [18] to approximate the real distribution of the samples for the corresponding image type. The difference between our method and the previous ones (e.g., [17], [19]) is that we build the character-model and ambiguity-model to filter ambiguous samples from wrongly classified samples in each “bootstrap” iteration. Three benefits can be gained from this strategy: 1) Without ambiguity, even simple classifiers can achieve very high accuracy on training set; 2) Generalization ability for positive samples is increased significantly so that a high recall rate can be achieved even if most of positive samples are only synthetic data; 3) Reduced labeling cost because it’s much easier to label the data manually due to very little confusion.

The fourth contribution is to train neural network based classifiers for four types of image patch, namely, Long, Thin, Fill, Square-large, by using the above ambiguity-free learning strategy to solve the text and non-text classification problem. Unlike previous methods (e.g., [2], [3], [5], [9]) which are mainly based on handcrafted features, our approach learns “meaningful” features directly from raw binary images, which is inspired by the recent success in deep learning [20]. There are two reasons not to use handcrafted features: 1) Due to the high degree of intra-class variation of text caused by font styles, stroke thickness, blur, distortion, noise or even part loss [9], it is very difficult to design a group of efficient and effective handcrafted features which can really distinguish text from non-text in the corresponding feature space; 2) Due to the information loss caused by feature extraction, the ambiguity problem will become more severe. Actually, we only use a two-hidden-layer neural network, instead of deep architectures (e.g., convolutional neural networks [21]) which have achieved

great success in many different fields [20]. There are two reasons: 1) Shallow architecture has lower computational cost; 2) With the specific learning strategies we proposed, the shallow architecture can achieve competitive accuracy for our task already.

The remainder of this paper is organized as follows. In Section II, the proposed approach is described in detail. In Section III, experimental evaluation is presented. The paper is concluded in Section IV.

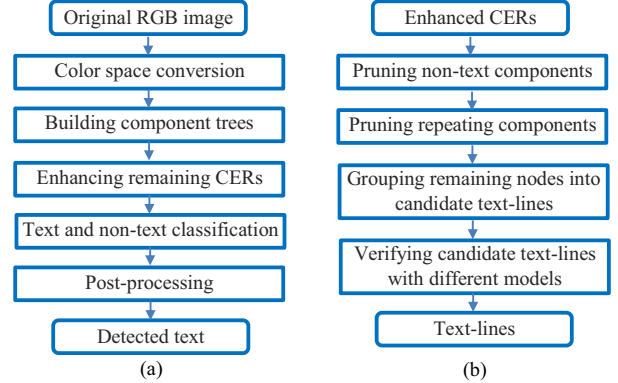


Fig. 3 Flowchart of our proposed approach: (a) whole system; (b) details of text and non-text classification module.

II. METHODOLOGY

A. Overview

As illustrated in Fig. 3, the proposed approach mainly includes the following five modules:

- 1) **Color space conversion.** The original color image is converted from RGB color space to PII color space. The transformation equation is presented in Section II.B.
- 2) **Building component trees.** Six component-trees are built from grayscale image, hue and saturation channel images in the PII color space, and their inverted images, respectively. Then the CER criterion and other simple rules are used to pre-prune obvious non-text components on each tree as in [22].
- 3) **Enhancing remaining CERs.** The dominant color of each remaining CER is estimated firstly. Pixels with dissimilar color to the dominant color are removed from each CER.
- 4) **Text and non-text classification.** A set of two-hidden-layer neural networks are trained to solve this problem. The components in the types of Long, Square-large, Thin and Fill are pruned by the corresponding neural networks first, grouped into candidate text-lines next, and then verified by the candidate text-line verification models. The classification method for the components of Square-small type is different due to two reasons: 1) Small components are easier to be grouped into text-lines; 2) the number of components of Square-small type is much larger than the number of other types, and majority of them are non-text components. So we group them into candidate text-lines first and verify them by the corresponding component model. In this way, the computation for classification is reduced a lot while the accuracy is not affected. Before candidate text-line formation, color and area information are used to prune repeating

components in the component-tree to obtain a component graph.

5) **Post-processing.** Text line information is used to recover wrongly classified text components and prune outliers. The extraction results from six component-trees are combined to get the final result.

B. Color Space Conversion

The perception-based illumination invariant color space is a new color space in which the associated metric approximates perceived distances and the color displacements capture relationships that are robust to spectral changes in illumination [15]. Let \vec{x} and $F(\vec{x})$ denote the tristimulus values of a sensor represented in XYZ coordinates and the tristimulus values in the new color space, respectively. Based on the assumptions and proofs in [15], the relationship between \vec{x} and $F(\vec{x})$ can be represented by Equation (1) as follows:

$$F(\vec{x}) = A(\hat{\ln}(B\vec{x})) \quad (1)$$

where A and B are invertible 3×3 matrices and $\hat{\ln}$ denotes the component-wise natural logarithm. A and B are learned from different datasets and their values are as follows:

$$B = \begin{bmatrix} 9.465229 \times 10^{-1} & 2.946927 \times 10^{-1} & -1.313419 \times 10^{-1} \\ -1.179179 \times 10^{-1} & 9.929960 \times 10^{-1} & 7.371554 \times 10^{-3} \\ 9.230461 \times 10^{-2} & -4.645794 \times 10^{-2} & 9.946464 \times 10^{-1} \end{bmatrix} \quad (2)$$

$$A = \begin{bmatrix} 2.707439 \times 10 & -2.280783 \times 10 & -1.806681 \\ -5.646736 & -7.722125 & 1.286503 \times 10 \\ -4.163133 & -4.579428 & -4.576049 \end{bmatrix} \quad (3)$$

Let F_1, F_2 denote two colors in PII color space, then the color distance can be simply computed as

$$d(F_1, F_2) = \|F_1 - F_2\| \quad (4)$$

where $\|\cdot\|$ denotes the usual l_2 -norm.

C. Generalized Color-enhanced CER

Although extremal-region has many good properties [12] and more than 94% of text components are extremal-regions in ICDAR-2003 testing set ([22], [23]), it still has limitations for candidate character generation. As illustrated in Fig. 4, the gray region in Fig. 4(a) and the characters in red bounding boxes are not extremal-regions, because they do not satisfy the strict definition of extremal-region, i.e. the minimum (maximum) intensity value of all the pixels in this region is larger (smaller) than the maximum (minimum) intensity value of all the pixels on the boundary [12].

Extracting extremal-regions on more planes in different color space can address partially this issue. To solve this problem better, generalized color-enhanced CER is proposed here. As illustrated in Fig. 5, the algorithm consists of 3 steps: 1) Estimate the dominant color of remaining pixels in the CER; 2) Extract the pixels with similar color to the dominant color to compose an enhanced CER; 3) If the number of remaining pixels is less than a threshold, the algorithm stops; otherwise, repeat steps 1) and 2). To estimate the dominant color, the remaining pixels are sorted in descending order according to their pixel value in the corresponding color channel firstly. Then the color average of the top 50% of the sorted pixels in PII color space is calculated as the dominant

color. Next, if the color distance between a pixel and the dominant color is less than a threshold T_1 (T_1 is 300 in our system), it is considered to have the similar color to the dominant color. In ICDAR-2013 testing set, there are very few text components which are not extremal-regions on the six component-trees we built previously, and all of them are listed in Fig. 6. So we can simplify the algorithm and iterate only once. Furthermore, CER can also be enhanced by other information (e.g., stroke width [16]) so that the above limitations may be addressed further, which could be our future work.

Generalized color-enhanced CER not only keeps the good properties of the color-enhanced CER [24], i.e. robust to superfluous pixels, but also more robust to illumination. The comparison results are illustrated in Fig. 7, and it can be seen that the new method has better performance. More results can be seen in Fig. 12.

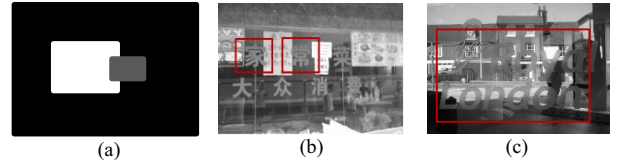


Fig. 4 Limitations of extremal-region based method.

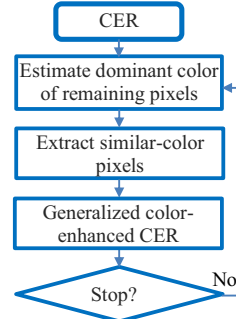


Fig. 5 Generalized color-enhanced CER.

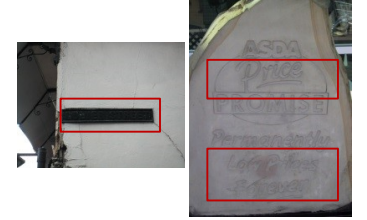


Fig. 6 Non-extremal-region text components in ICDAR-2013 dataset.



Fig. 7 Illustration of enhanced results: (a) Raw image; (b) Color-enhanced CER; (c) Generalized color-enhanced CER.

D. Pruning Non-text Components

Let w_i , h_i , Ar_i and Fr_i denote the width, height, aspect ratio and filled rate of each remaining component C_i respectively, and C_i is labeled as one of five types according to the rules in Table 1. Five types of image patches are illustrated in Fig. 8. Next, C_i is resized to the size of its corresponding type. If the aspect ratio of the component in Long type exceeds 3.5, it will be split into several regions. Generally speaking, components in Thin type contain only one thin character like “l”, “t” or “f”, and components in Long type contain two or more characters. Most components in Fill,

Square-large and Square-small types contain one or two characters. After the division, ambiguities in the Long, Fill and Square-small types are reduced greatly due to more characters information and lower variance respectively, and ambiguities in Thin type are mainly caused by the characters “I” and “l”. Ambiguities are mainly restricted to the Square-large type, so next we will take it for example to introduce the proposed ambiguity-free learning strategy which is illustrated in Fig. 9.

To avoid wrong labels, positive (text) samples in initial active set are all synthetic data including about 1000 different kinds of font libraries with various transformations (e.g., affine transformation, translation, and salt-and-pepper noise). Negative (non-text) samples are extracted from ICDAR-2011 training set and enriched by Oxford Buildings dataset [25] and ImageNet dataset [26]. To save labeling efforts, images containing no text are selected, and ambiguous negative samples are separated from other negative samples. Then an initial model (neural network) is trained on the synthetic single-character positive samples and unambiguous negative samples. This model will be used as character-model later. The ambiguity-model is trained on the separated ambiguous samples and other unambiguous negative samples.

Next, the initial model is enhanced by real data with a bootstrap strategy [17]. It is straightforward to enrich positive samples, i.e. directly add wrongly classified positive samples into active set. Because initial model cannot handle the components with two or more characters, an English document dataset is used for enhancement first. Besides, several public datasets like ICDAR-2011 training set, Street-View-Text dataset [27], NEOCR dataset [28], KAIST scene text database [29], MSRA-TD500 database [5], as well as about several thousand text-intensive web images (e.g., posters, book covers) are used to enrich the positive active set. To enrich negative samples, ImageNet dataset is used. Generalized color-enhanced CERs are extracted from new images in the dataset, and only the samples classified as text are retained and filtered by the character-model and ambiguity-model in turn. The survived data are then checked by human labelers and added to the active set to retrain the model. As this process repeats, the model performance keeps improving. So far, we have used about 1.5 million images in ImageNet to enrich our negative active set.

The same learning strategy is used to train the classifier for the Thin type. Thanks to the “divide-and-conquer” strategy, there are much fewer ambiguities for Long and Fill types such that the original bootstrap strategy alone is enough for the Long and Fill types.

Besides, we let the models complement each other so that the variance for each model is reduced further. For example, we let the model for the Long type focus on learning the texture of multi-characters, and the model for Square-large type focus on learning the shape of one or two characters. Furthermore, we allow certain overlap between Long and Square-large types, i.e., if a component with an aspect ratio less than 2 is pruned by the Long classifier, it will be verified again by the Square-large classifier. In this way, the training

data for each model is reduced greatly while the accuracy is maintained.

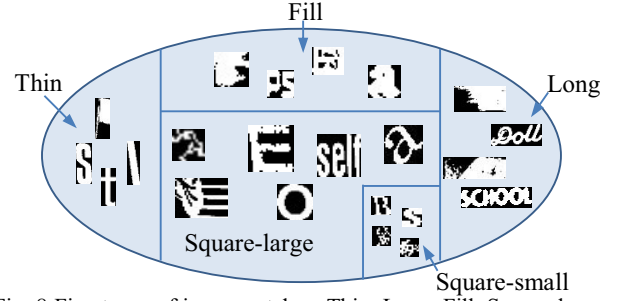


Fig. 8 Five types of image patches: Thin, Long, Fill, Square-large and Square-small.

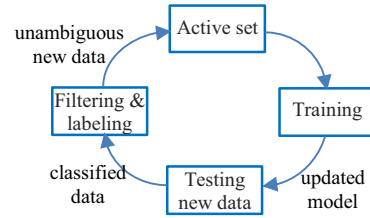


Fig. 9 Ambiguity-free learning.

Table 1 Rules for labeling image patch type.

Image Patch Type	Rules	Resized size (pixels)	
		width	height
Thin	$Ar_i < 0.4$	12	36
Long	$Ar_i > 1.6$	28	12
Fill	$(0.4 \leq Ar_i \leq 1.6) \cap (Fr_i > 0.7)$	20	20
Square-large	$(0.4 \leq Ar_i \leq 1.6) \cap (Fr_i \leq 0.7) \cap (w_i \geq 12) \cap (h_i \geq 12)$	28	28
Square-small	others	12	12

E. Pruning Repeating Components

Due to the high accuracy of the previous pruning step (Table 3), most of the remaining nodes on the component-tree are text components, and can be grouped into many sub-trees like the structure in Fig. 1. Using only several simple rules can prune the repeating components robustly. Let S_p and S_c denote the area of node N_p and its children N_c respectively, and let S_{pb} and S_{cb} denote the area of their bounding boxes. Two simple rules are used to prune duplicate and noisy nodes: 1) If $S_c/S_p > 0.9$, N_p is removed due to duplicate; 2) If $S_c/S_p > 0.8 \cap S_{cb}/S_{pb} < 0.7$, N_p is removed because this rule implies more noise in N_p than its children. Let H_p denote the height of N_p and NUM_{sc} denote the number of N_p 's children whose height is smaller than $0.6 * H_p$. If $NUM_{sc} > 3$, N_p is removed due to the high probability to be wrongly classified. As for other nodes, they are removed if they have similar color to their parent. In this way, the hierarchical structure of the remaining sub-trees is reduced to a simple graph structure.

F. Candidate Text Line Formation

Due to the previous two steps, the candidate text line formation problem becomes very simple. The grouping

method in our previous work [24] is used. The only difference is that the color distance in PII space is used instead of the one in the HVC color space.

G. Verification

Different types of candidate text lines are verified by different verification models. If more than half of the nodes in the line are Square-small type, the candidate text-line is classified as “small-lines”, otherwise is classified as “large-lines”. The small-lines are pruned by the model built for small-line classification first, and verified by the model built for isolated components of Square-small type. If more than half of the nodes are classified as non-text, this line will be removed. To verify large-lines, the character-model and the model built for large-lines are both used. If the candidate line contains only one component, it will be verified by the character-model; otherwise, verified by the large-line model. Here, the ambiguity set is added to the negative set to train the character model.

H. Post-processing

Finally, results from all the six component-trees are combined. If the bounding boxes from different trees are overlapped, the one with more components is kept. After that, text-line information is used to recover the lost characters and split the text line into words as appropriate.

III. EXPERIMENTS

In the first experiment, the performance of the four classifiers mentioned in Section II.D is evaluated. The testing set is labeled from ICDAR-2011 testing set. All neural networks use rectified linear units as their hidden nodes and are trained by standard back propagation algorithm with a mini-batch based stochastic gradient descent for optimization. The topology of the neural networks for each type is listed in Table 2 and the classification results in Table 3. The accuracy of each neural network can reach easily to 100% on training set after 30 epochs. The performance of the component model used for Square-large type is shown in Fig. 10, and no overtraining is observed. Examples of weights from the first hidden layer of this model are printed out in Fig. 11. It can be seen that some “meaningful features” can be learned directly from raw data indeed. The top row images in Fig. 11 look like characters and can be understood as “character filter”, while the bottom row are noisy and can be treated as “non-text filter”.

In the second experiment, we take the Square-large type for example to compare the performance of the deep and shallow architectures on our task. The results can be seen in Table 4, and the performance of the deep architecture is only slightly better than the shallow model.

In the third experiment, we compare the performance of the proposed ambiguity-free learning strategy to the traditional bootstrap strategy. The result is listed in Table 5. It is observed that the proposed ambiguity-free learning strategy has much better generalization ability for positive samples while the accuracy on negative samples is only slightly lower.

In the fourth experiment, the proposed method is compared with the methods in ICDAR-2013 Robust Reading competition [30]. Examples of several difficult cases which cannot be

solved properly by other methods, but can by our method are shown in Fig. 12. A comparison of benchmark results are listed in Table 6. The recall and F-score of our proposed method are much higher than that of other methods. Actually, even higher recall and precision rates of our method could be reported if a better word segmentation method was implemented. Because our text detection system is designed to extract text lines only for follow-up OCR (e.g., [31]), we need to cut each text line into words to use the evaluation tool of ICDAR-2013 robust reading competition to obtain the benchmark results which can be compared directly with other published methods. The simple rules we implemented so far are not robust enough to segment the extracted text-lines into words. Due to the poor performance of word segmentation, both precision and recall rates of our method are under-estimated.

Table 2 Topology of neural networks for each type of image patch.

Square-large	Fill	Thin	Long
784:400:200:2	400:256:128:2	432:256:128:2	336:256:128:2

Table 3 Result of non-text components pruning.

Image Patch Type	Training set size	Testing set			
		Text set size	Acc. (%)	Non-text set size	Acc. (%)
Square-large	2,548,734	5729	94.0	39653	97.1
Long	492,544	442	91.5	32284	99.6
Fill	205,000	867	94.1	13637	98.7
Thin	45,568	254	96.1	12459	97.2

Table 4 Comparison of deep and shallow architectures.

Architecture	Topology	Acc. on pos (%)	Acc. on neg (%)
Shallow	784:400:200:2	94.0	97.1
Deep	784:1000×5:2	94.4	97.6

Table 5 Comparison of ambiguity-free learning and bootstrap.

Methods	Acc. on pos (%)	Acc. on neg (%)
Ambiguity-free learning strategy	94.0	97.1
Traditional bootstrap strategy	85.2	98.3

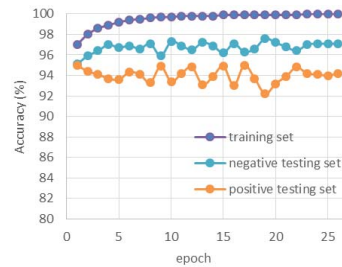


Fig. 10 Results of the component model for Square-large type.



Fig. 11 Examples of weights from the first hidden layer of the model for Square-large type.

Table 6 Comparison of benchmark results on ICDAR-2013 “Reading Text in Scene Images” test set.

Methods	Recall	Precision	F-score
Our proposed method	85.72%	87.03%	86.37%
USTB_TexStar	66.45%	88.47%	75.89%
TextSpotter	64.84%	87.51%	74.49%
UMD_IntergratedDiscrimination	62.26%	89.17%	73.33%
CASIA_NLPR	68.24%	78.89%	73.18%



Fig. 12 Examples of text detection results of our method.

IV. CONCLUSION

In this paper, a novel approach based on generalized color-enhanced contrasting extremal region (CER) and neural networks is proposed to solve the difficult problem of automatic text detection from natural scene images. Several open problems of extremal region based methods are solved properly. The proposed method has achieved 85.72% recall, 87.03% precision, and 86.37% F-score on ICDAR-2013 "Reading Text in Scene Images" test set.

V. REFERENCES

- [1] L. Neumann and J. Matas, "On combining multiple segmentations in scene text recognition," in *ICDAR*, 2013, pp. 523-527.
- [2] X.C. Yin, X.W. Yin, K.Z. Huang and H.W. Hao, "Robust text detection in natural scene images," *IEEE Trans. PAMI*, 2013.
- [3] Y.F. Pan, X. Hou, and C.L. Liu, "A hybrid approach to detect and localize texts in natural scene images," *IEEE Trans. Image Processing*, vol. 20, no. 3, pp. 800-813, 2011.
- [4] D. Karatzas, et al., "ICDAR 2013 robust reading competition," in *ICDAR*, 2013, pp. 1484-1493.
- [5] C. Yao, X. Bai, W.Y. Liu, Y. Ma, and Z.W. Tu, "Detecting texts of arbitrary orientations in natural images," in *CVPR*, 2012, pp. 1083-1090.
- [6] X. Chen and A. Yuille, "Detecting and reading text in natural scenes," in *CVPR*, 2004, pp. 366-373.
- [7] A. Coates, B. Carpenter, C. Case, S. Satheesh, B. Suresh, T. Wang, D. Wu, and A. Ng, "Text detection and character recognition in scene images with unsupervised feature learning," in *ICDAR*, 2011, pp. 440-445.
- [8] R. Minetto, N. Thone, M. Cord, N.J. Leite and J. Stolfi, "T-HOG: an effective gradient-based descriptor for single line text regions," *Pattern Recognition*, vol. 46, no. 3, pp. 1078-1090, 2013.
- [9] C. Shi, C. Wang, B. Xiao, Y. Zhang, and S. Gao, "Scene text detection using graph model built upon maximally stable extremal regions," *Pattern Recognition Letters*, vol. 34, no. 2, pp. 107-116, 2013.
- [10] L. Neumann and J. Matas, "Scene text localization and recognition with oriented stroke detection," in *ICCV*, 2013, pp. 97-104.
- [11] A. Shahab, F. Shafait and A. Dengel, "ICDAR 2011 robust reading competition challenge 2: reading text in scene images," in *ICDAR*, 2011, pp. 1491-1496.
- [12] J. Matas, O. Chum, M. Urban, and T. Pajdla, "Robust wide baseline stereo from maximally stable extremal regions," in *BMVC*, 2002, pp. 384-393.
- [13] L. Sun and Q. Huo, "An improved component tree based approach to user-intention guided text extraction from natural scene images," in *ICDAR*, 2013, pp. 383-387.
- [14] A. H. Munsell, A color notation(12th ed.), Baltimore, MD: Munsell Color Company, 1971.
- [15] H.Y. Chong, S.J. Gortler and T. Zickler, "A perception-based color space for illumination-invariant image processing," in *SIGGRAPH*, 2008, pp. 1-7.
- [16] B. Epshtein, E. Ofek and Y. Wexler, "Detecting text in natural scenes with stroke width transform," in *CVPR*, 2010, pp. 2963-2970.
- [17] K.K. Sung and T. Poggio, "Example-based learning for view-based human face detection," *IEEE Trans. PAMI*, vol. 20, no. 1, pp. 39 - 51, 1998.
- [18] Z. Kalal, J. Matas, and K. Mikolajczyk, "Weighted sampling for large-scale boosting," in *BMVC*, 2008, pp. 42.1-42.10.
- [19] H.A. Rowley, S. Baluja, and T. Kanade, "Neural network-based face detection," *IEEE Trans. PAMI*, vol. 20, pp. 22-38, 1998.
- [20] Y. Bengio, A. Courville and P. Vincent, "Representation learning: a review and new perspectives," *IEEE Trans. PAMI*, vol. 35, no. 8, pp. 1798-1828, 2013.
- [21] Y. LeCun, L. Bottou, Y. Bengio and P. Haffner, "Gradient-based learning applied to document recognition," in *Proc. the IEEE* 86, 11, 1998, pp. 2278-2324.
- [22] L. Sun and Q. Huo, "A component-tree based method for user-intention guided text extraction," in *ICPR*, 2012, pp. 633-636.
- [23] L. Neumann and J. Matas, "Text localization in real-world images using efficiently pruned exhaustive search," in *ICDAR*, 2011, pp. 687-691.
- [24] L. Sun and Q. Huo, "An improved component tree based approach to user-intention guided text extraction from natural scene images," in *ICDAR*, 2013, pp. 383-387.
- [25] J. Philbin, et al, "Object retrieval with large vocabularies and fast spatial matching," in *CVPR*, 2007, pp. 1-8.
- [26] J. Deng, W. Dong, R. Socher, L. Li, K. Li, and F.F. Li, "Imagenet: a large-scale hierarchical image database," in *CVPR*, 2009, pp. 248-255.
- [27] "Street View Text Dataset," [Online]. Available: <http://vision.ucsd.edu/~kai/svt/>.
- [28] R. Nagy, A. Dicker, and M.W. Klaus, "NEOCR: a configurable dataset for natural image text recognition," in *ICDAR workshop at ICDAR*, 2011, pp. 150-163.
- [29] "KAIST scene text database," [Online]. Available: http://www.iaprtc11.org/mediawiki/index.php/KAIST_Scene_Text_Database.
- [30] "ICDAR 2013 Robust Reading Competition," [Online]. Available: <http://dag.cvc.uab.es/icdar2013competition/?ch=2&com=results>.
- [31] A. Bissacco, M. Cummins, Y. Netzer and H. Neven, "PhotoOCR: Reading Text in Uncontrolled Conditions," in *ICCV*, 2013, pp. 785-792.