# Faster R-CNN Based Table Detection Combining Corner Locating

Ningning Sun, Yuanping Zhu and Xiaoming Hu

School of Computer and Information Engineering

Tianjin Normal University

Tianjin, China

Email: zhuyuanping@tjnu.edu.cn

*Abstract*—**Table detection in document images has achieved remarkable improvement. However, there is still a problem of inaccurate table boundary locating. This paper proposes Faster R-CNN based table detection combining corner locating method. Firstly, coarse table detection and corner locating are implemented through Faster R-CNN network. Secondly, those corners belonging to the same tables are grouped by coordinate matching. At the same time, unreliable corners are filtered. Finally, table boundaries are adjusted and refined by corresponding corner group. The proposed method improves the precision of table boundary locating at pixel-level. Experiment results show that our method effectively improves the precision of table detection. It achieves an F-measure of 94.9% on ICDAR2017 POD dataset. Compared with traditional Faster R-CNN method, our method increases by 2.8% in F-measure and 2.1% at pixel-level localization.**

*Keywords*-**table detection; corner locating; Faster R-CNN;**

## I. INTRODUCTION

Table detection has always been an important area of our research, which promotes the development of layout analysis. Due to external and internal factors, table detection has always faced many problems. External factors mainly include tilt, blur, noise, illumination and occlusion et al, which are inevitable. Therefore, we should take some measures to weaken the influence of external factors. The measures include tilt refining [1], deblurring [2], denoising, illumination preprocessing [3] et al. However, these technologies are still immature and need to be continuously optimized. Internal factors mainly include: (1) complex layouts, such as single-column page, double-column pages and three-column pages; (2) The diversity of tables, whether contain horizontal or vertical lines, font style, font size and content format; (3) The interference of table location, including side-by-side tables, list tables and nested tables. The complication of document images leads to table detection difficulty.

In the past few years, people have implemented table detection based on traditional methods. Traditional methods are divided into two categories, shape-based methods and texture-based methods. Shape-based methods include top-down method [4] and bottom-up method [5]. For example, RLSA (Run-Length Smoothing Algorithm) [6] is a typical shape-based method commonly used in table detection. Texture-based methods include model-based method, feature-based method [7] and multiscale-based method [8]. Traditional

methods are limited to many factors, which promotes the rapid development of deep learning based methods. Therefore, deep learning based methods are the current mainstream methods in table detection. A lot of experiments have proved that deep learning shows better in learning features. Deep learning based methods are divided into two categories, bounding-box based method and semantic segmentation method. Bounding-box based method can predict object proposal. However, bounding-box method is limit to the variable length bounding-box list. Hence, this method leads to inaccurate table locating. Semantic segmentation method classifies each pixel in the image, which typically uses a per-pixel softmax, like Mask R-CNN method [9]. Semantic segmentation method is commonly applied to natural scene images.

Different from the previous methods, we propose a new method that combines corner locating method. Corners mentioned in this paper refer to table corners, as shown in Fig. 1. Corners mainly have the following characteristics: (1) The region with a certain range centered on the vertices of table is called table corner. (2) Corners belonging to the same table are called a corner group. (3) Corners have the same size and shape except of position. They locates in top-left, top-right, bottom-left and bottom-right of the table. For convenience, they are named top-left corner, top-right corner, bottom-left
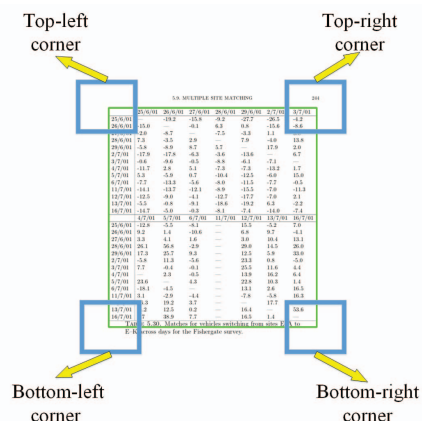


Fig. 1: The green bounding-box is a table and the blue bounding-boxes are corners.
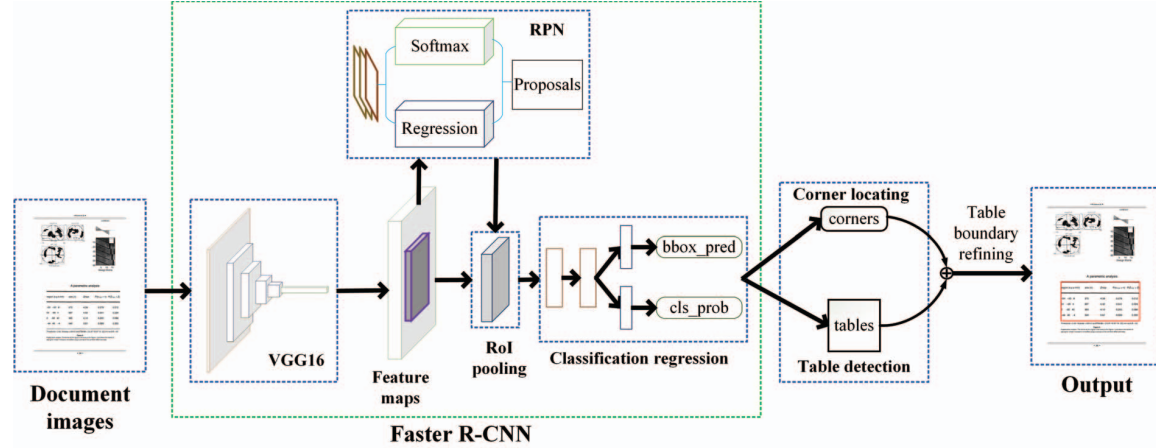
Fig. 2: The framework of Faster R-CNN based table detection combining corner locating.

corner and bottom-right corner, respectively.

Our method can be implemented in two steps. Firstly, table detection is implemented through Faster R-CNN [10]. Secondly, table boundaries are refined by reliable corners. Compared with the previous bounding-box based method, our method has two advantages: (1) For most tables, there exist corner groups for refining table boundary. (2) By refining table boundaries, the precision of table detection has significantly improved. Details will be introduced in Section 3. We have verified our method on a standard layout analysis dataset, ICDAR2017 Competition of Page Object Detection dataset [11]. Results show that the proposed method is superior to the widely used Faster R-CNN and SSD [12] method.

## II. RELATED WORK

Table detection has always been the important research. Therefore, people have done a plenty of work in it. T. Kieninger et al. [13] [14] have completed one of the ground-breaking work on table detection. However, the method they proposed could only be processed for specific documents. F. Cesarini et al. [15] have proposed a top-down table detection method . The region surrounded by the lines is identified as a table candidate area by detecting horizontal and vertical lines of the document images. Basilios et al. [16] further studies on the basis of Cesari by detecting horizontal lines, vertical lines and intersect points. However, their methods were poorly versatile and not suitable for tables without horizontal and vertical lines. Inspired by [15] [16], we think of refining table boundaries in table detection. Traditional methods have made great progress in table detection, which promotes the development of table detection. At the same time, it is accompanied with many problems, such as partial detection, over-segmented tables, under-segmented tables, and missed tables.

At present, the mainstream table detection methods are deep learning based methods. Deep learning has become a part of leading systems in many fields, especially in the field of layout analysis. R-CNN (region CNN) [17] is a pioneering work using deep learning for object detection. R-CNN is accompanied with the emergence of double counting, which increases the amount of calculation and reduces the performance of the network. Therefore, R. Girshick refers to SPP-net (Spatial Pyramid Pooling) [18] and proposes Fast R-CNN [19] to improve the speed. On this basis, Faster R-CNN method is proposed. Faster R-CNN exponentially reduces computation complexity. Given the fact that Faster R-CNN is superior in both speed and accuracy, we choose Faster R-CNN for table detection.

The method combining corners is a novel method in object detection. A corner detection layer and an analytical sampling layer are used to replace RPN network in DeNet [20]. X. Wang have presented PLN (Point Linking Network) [21], which uses a full convolution network to regress the bounding-box and the corresponding corner/center point. P. Lyu et al [22] have developed a multi-oriented scene text detection model, which can detect multi-oriented scene text in natural scene via corner localization. CornerNet [23] model is proposed for natural scene. Compared with all existing one-step detectors, Cornernet achieves the best result in MS COCO. Inspired by the above methods, we propose a method via corner locating to refine table boundary.

## III. METHODOLOGY

### A. Architecture

Based on Faster R-CNN method, we combine corner locating method in table detection. The proposed method implements table detection and improves the precision. The structure of the network is shown in Fig. 2. Firstly, the backbone network VGG16 extract feature maps from document images. Secondly, RPN and Fast R-CNN implement table detection and corner locating. Thirdly, corners are grouped by coordinate matching. Unreliable corners are filtered through grouping and several prior rules. Finally, tables are adjusted and refined via the corresponding reliable corners. Our method improves the
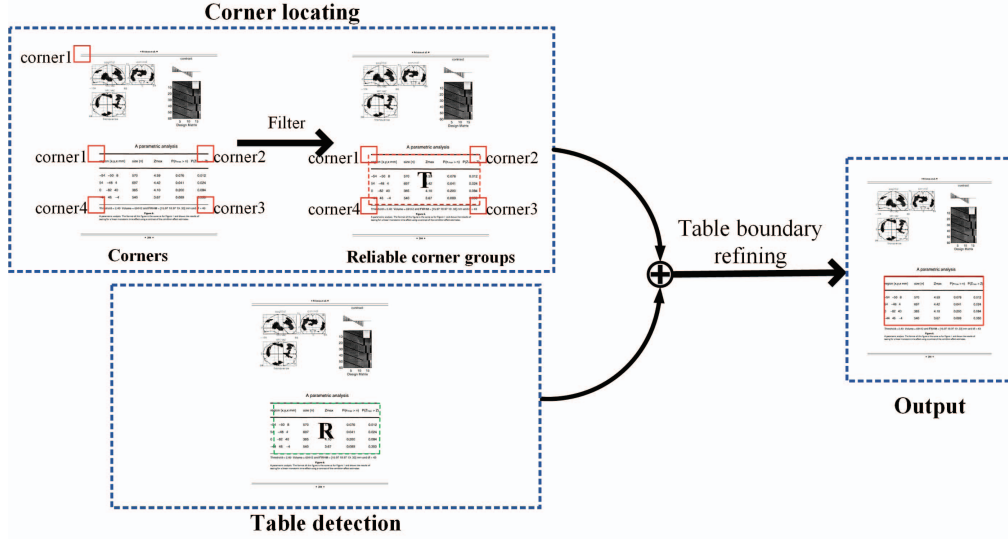
Authorized licensed use limited to: Xian Jiaotong University. Downloaded on May 19,2021 at 10:25:39 UTC from IEEE Xplore. Restrictions apply.

Fig. 3: Table boundary refining.

precision of table detection, especially at pixel-level localization.

### B. Table Detection

Table detection is implemented by two steps. Firstly, the backbone network VGG16 extracts feature maps from document images. Secondly, feature maps are sent into the RPN to implement coarse detection. Through coarse detection, object foreground and background that meet the requirements are obtained. Thirdly, feature maps and ROI regions generated by RPN are fed together to the ROI pooling layer. ROI pooling layer generates a fixed-size feature map. Finally, table detection is implemented through classification and regression. For convenience, the table detected by Faster R-CNN is called $T$.

### C. Corner Locating and Filtering

For convenience, we use the same model to implement table detection and corner detection. The process of corner detection is the same as table detection. In this work, we set corner size to 80 in default. We name the corresponding corners $corner1$, $corner2$, $corner3$ and $corner4$, respectively (in a clockwise direction). At the same time, they can also be defined as $C1(x_1, y_1)$, $C2(x_2, y_2)$, $C3(x_3, y_3)$ and $C4(x_4, y_4)$.

Inspired by CornerNet [23], corners belonging to the same table are grouped by coordinate matching. For convenience, we call the rectangular region determined by corner group as:

$$R = \{C_i'(x_i', y_i') | i \epsilon \{1, 2\}\} \quad (1)$$

$$\begin{cases} x_1' = (x_1 + x_4)/2 \\ y_1' = (y_1 + y_2)/2 \\ x_2' = (x_2 + x_3)/2 \\ y_2' = (y_3 + y_4)/2 \end{cases} \quad (2)$$

Not all corners are reliable, unreliable corners can be filtered by corner grouping method. Corners satisfied the following location constraint relationships would be retained:

- $C1$ and $C2$ are on the same horizontal line approximately.
- $C3$ and $C4$ are on the same horizontal line approximately.
- $C1$ and $C4$ are on the same vertical line approximately.
- $C2$ and $C3$ are on the same vertical line approximately.

The deviation of adjacent corners greater than $z$ are rejected ($z$ is suitable threshold), defined as:

$$|(x_i - x_j)| \leq z \qquad \{i = 1, j = 4\} or \{i = 2, j = 3\} \quad (3)$$

and

$$|(y_i - y_j)| \leq z \qquad \{i = 1, j = 2\} or \{i = 3, j = 4\} \quad (4)$$

in which, $z$ is related to image size. $z = \sqrt{\theta w h}$ ($0.0002 \leq \theta \leq 0.0005$), image is of size $w \times h$. In this paper, the default value of $\theta$ is 0.0003.

### D. Table Boundary Refining

Observing table documents, we noticed that many tables contain horizontal lines but without vertical lines. Observing table detection results implemented by Faster R-CNN, We found that inaccurate $T$ were mainly caused by inaccurate table detection of left and right boundaries, while upper and lower boundaries had little effect on table detection. Therefore, only left and right boundaries of $T$ are refined through reliable corner groups. The steps are shown Fig. 3. The meanings of the two items are explained as follows:

- Search corresponding $R$ for $T$. The rule we used is IOU (Intersection-over-Union). If $IOU(R, T)$ is greater than the threshold we set (the default is 0.5), $R$ is retained.
- Refine left and right boundaries of $T$ by merging the corresponding $R$ and $T$.

Suppose $(x_1'', y_1'')$, $(x_2'', y_2'')$ are coordinates on the top-left and bottom-right vertex of the $T$, respectively. Therefore, $T$ is defined as $(x_1'', y_1'', x_2'', y_2'')$. Same as before, $R$ is defined as $(x_1', y_1', x_2', y_2')$ and refined table is defined as $(x_1^*, y_1^*, x_2^*, y_2^*)$. The relationship between them described as:

$$\begin{cases} x_1^* = \dfrac{x_1''+x_1'}{2} \\ y_1^* = y_1'' \\ x_2^* = \dfrac{x_2''+x_2'}{2} \\ y_2^* = y_2'' \end{cases} \tag{5}$$

After the above steps, table boundaries are refined by reliable corner groups. The proposed method is also applicable to other datasets.

*E. Loss Function*

The loss function consists of two parts. (1)The classification layer determines the proposal is foreground or background. (2) The regression layer is used to predict the central anchor point of the proposal, the coordinates corresponding to x, y, w and h. Faster R-CNN follows the definition of multi-task loss. The loss function consists of classification loss and regression loss, defined as:

$$L(\{p_i\}, \{t_i\}) = \frac{1}{N_{cls}} \sum_i L_{cls}(p_i, p_i^*) + \lambda \frac{1}{N_{reg}} \sum_i p_i^* L_{reg}(t_i, t_i^*) \tag{6}$$

in which, $L_{cls}(p_i, p_i^*)$ is the classification loss, expressed as a logarithmic function of two categories (foreground and background), calculated as: $L_{cls}(p_i, p_i^*) = -\log[p_i p_i^* + (1 - p_i^*)(1 - t_i^*)]$. $L_{reg}(t_i, t_i^*)$ is the regression loss, the calculation formula is: $L_{reg}(t_i, t_i^*) = R(t_i - t_i^*)$. $R$ is a smooth L1 function, $p_i^* L_{reg}$ means that the regression loss is only activated for positive anchors ($p_i^* = 1$) and is disabled for otherwise($p_i^* = 0$).

## IV. EXPERIMENT

In order to verify the effectiveness of our method, we comprehensively evaluate our method on ICDAR2017 POD dataset [11].

*A. Experimental Details*

The experiment is trained in tensorflow 1.3.0 environment, with one GTX1080Ti, 11G of GPU memory is sufficient. The network is randomly initialized under the default setting of tensorflow with no pretraining on any external dataset. The network is fine-tuned with an initial learning rate of 0.0001 and slowly reduce, minibatch from 2 to 10. Model is continued to train for 50k iterations.

The input resolution is set to $600 \times$ h when the network is trained (h $\leq$1000). The output resolution ie set to $1200 \times$ 1200. Corner size is set to a fixed value (related to image size). To reduce overfitting, we use random horizontally-flipped images for data augmentation. Not only original images but also horizontally-flipped images are used in the testing
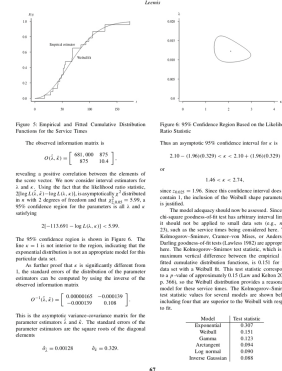


Fig. 4: ICDAR2017 POD dataset.

phase. Output contains five classes together through Faster R-CNN, which are tables and four corners. Table boundary is refined through all reliable corners. Because the initial generated anchor is restricted by the radio and size, so small objects (for example: corners) below the minimum anchor size are frequently missed in the latter process.

*B. Dataset*

ICDAR2017 POD dataset is mainly for studying layout analysis. Peking University rigorously labels such data. The official dataset includes 792 document images (contain tables), of which contains 549 images for training and 243 images for testing. The dataset mainly consists of one column and two columns. The types are divided into four categories: text, table, figure and formula. The format of the dataset we used is shown in Fig. 4.

The dataset has complex page layout, more than partial pixel loss. The situation deeply increases the difficulty of table detection. According to our observation, there are three reasons affecting table detection. (1) The interference of background, other three categories (text, figure, formula) should be treated as background. (2) The diversity of tables, tables with lines and tables without lines. (3) The complexity of the page layout, a single column, two columns, list tables and side-by-side tables.

*C. Experimental Results*

In this section, we carry out ablation studies. Refer to the evaluation metrics of the ICDAR2017 POD competition, IOU is set to 0.6 and the evaluation metrics including precision, recall and F-measure.

**Comparison of all results in ICDAR2017 POD dataset.** Our method is implemented for comparative experiment. All results are reported in Table 1. Compared with ICDAR2017 POD competition results, it can be seen that our method ranks the second in table detection with F-measure of 94.9%. Faster R-CNN+CCs+CRFs method is complicated, which adds a series of pre-processing steps to improve the results. Such as Connected Component analysis (CCs), Support Vector

1317

TABLE I    Comparison of table detection results between the proposed method and ICDAR2017 POD Competition methods.

| Method Head | Precision | Recall | F-measure |
|---|---|---|---|
| SSD [12] | 0.071 | 0.959 | 0.132 |
| CNN [17] | 0.230 | 0.221 | 0.225 |
| Faster R-CNN(VGG_CNN_M_1024 [11] | 0.670 | 0.940 | 0.782 |
| Faster R-CNN+edge_based information [11] | 0.842 | 0.890 | 0.865 |
| Faster R-CNN(VGG16) | 0.924 | 0.918 | 0.921 |
| Our method | 0.943 | 0.956 | 0.949 |
| Faster R-CNN+CCs+CRFs [11] | 0.968 | 0.953 | 0.960 |



(a)                    (b)

(c)                    (d)

Fig. 5: Some table detection results on ICDAR2017 POD dataset. (a), (b), (c) and (d) show the comparison of Faster R-CNN method and our method. Color legend: Faster R-CNN(VGG16)method - green; Our method - red.

Machine (SVM) and Conditional Random Fields (CRF). Our method is simple that only combines corner locating to optimize table boundary localization. Our method would achieves better results if it combines with those methods.

**Compared with Faster R-CNN method.** As a further comparison, Faster R-CNN (VGG16) method and our method were trained for comparative experiment. The evaluation metrics precision, recall and F-measure are the same as Table 1. IOU indicates comparison at pixel-level localization. Comparative results are reported in Table 2. Compared with Faster R-CNN (VGG16) method, the proposed method significantly increased by 2.8% in F-measure and 2.1% at pixel-level localization.

To find out what differences of table detection between

TABLE II    Improvement at pixel-level Localization.

| Method Head | Precision | Recall | F-measure | IOU |
|---|---|---|---|---|
| Faster R-CNN(VGG16) | 0.924 | 0.918 | 0.921 | 0.811 |
| Our method | 0.943 | 0.956 | 0.949 | 0.832 |

Faster R-CNN (VGG16) method and our method, we give some examples in Fig. 5. Obviously, (a) and (d) implies that that our method refines right boundary of the inaccurate tables detected by Faster R-CNN method, (b) and (c) implies that our method refines left boundary of the inaccurate tables detected by Faster R-CNN method.

## V. Conclusion

In this paper, We propose a novel table detection method. Our method combines Faster R-CNN and corner locating method. Through table boundary refining via corner locating, the precision of table detection is significantly increasing. Compared with traditional Faster R-CNN method, our method increased by 2.8% in F-measure and 2.1% at pixel-level localization on the ICDAR2017 POD dataset.

The proposed method is also applicable to other datasets. Although the proposed method is quite practical, it is still accompanied with some problems. Firstly, anchor is restricted by the radio and size, which leads to small object proposals (for example: corners) higher miss detection. Secondly, corner size adopts a fixed value, which is also disadvantageous in table boundary refining. We will continue to improve our method in the follow-up work and set adaptive corner size.

## References

[1] M. Rahman, S. Gustafson, P. Irani, and S. Subramanian, "Tilt techniques: investigating the dexterity of wrist-based input," Proc. SIGCHI Conference on Human Factors in Computing Systems, ACM Press, Apr. 2009, pp. 1943-1952.

[2] P. D. Sankhe, M. Patil, and M. Margaret, "Deblurring of grayscale images using inverse and Wiener filter," Proc. International Conference & Workshop on Emerging Trends in Technology, ACM Press, 2011, pp. 145-148.

[3] G. Antal, R. Martinez, F. Csonka, M. Sbert, and L. Szirmay-Kalos, "Combining global and local global-illumination algorithms," Proc. 19th Spring Conference on Computer Graphics, ACM Press, Apr. 2003, pp. 185-192.

[4] F. Chang, S. Y. Chu, and C. Y. Chen, "Chinese document layout analysis using an adaptive regrouping strategy," Pattern Recognition, 2005, 38(2): 261-271.

[5] J. Xi, J. Hu, and L. Wu, "Page segmentation of Chinese newspapers," Pattern Recognition, 2002, 35(12): 2695-2704.

[6] F. M. Wahl, K. Y. Wong, and R. G. Casey, "Block segmentation and text extraction in mixed text/image documents," Computer Graphics and Image Processing, 1982, 20(4):375-390.

[7] J. L. Chenxe, "A simplified approach to the HMM based texture analysis and its application to document segmentation," Pattern Recognition Letters, 1997, 18(10): 993-1007.

[8] H. Choi, and R. G. Baraniuk, "Multiscale document segmentation using wavelet-domain hidden Markov models," Proc. SPIE Document Recognition and Retrieval VII, Dec. 1999, vol. 3967, pp. 234-248.

[9] K. He, G. Gkioxari, P. Dollr, and R. Girshick, "Mask R-CNN," Proc. IEEE Symp. International Conference on Computer Vision(ICCV), 2017, pp. 2961-2969.

[10] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," Proc. International Conference on Neural Information Processing Systems, 2015, pp. 91-99.

[11] L. Gao, X. Yi, Z. Jiang, L. Hao, and Z. Tang, "ICDAR2017 Competition on Page Object Detection," Proc. 14th International Conference on Document Analysis and Recognition (ICDAR), IEEE Press, Vol. 1, Nov. 2017, pp. 1417-1422.

[12] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, and C. Y. Fu, et al. "SSD: Single Shot MultiBox Detector," Proc. European Conference on Computer Vision(ECCV), Springer Press, Cham, Oct. 2016, pp. 21-37.

[13] T. Kieninger, and A. Dengel, "Table recognition and labeling using intrinsic layout features," Proc. International Conference on Advances in Pattern Recognition, London, Springer Press, 1999, pp. 307-316.

[14] T. Kieninger, and A. Dengel, "Applying the T-RECS table recognition system to the business letter domain," Proc. Sixth International Conference on Document Analysis and Recognition, IEEE Press, 2001, pp. 518-522.

[15] F. Cesarini, S. Marinai, L. Sarti, and G. Soda, "Trainable table location in document images," Proc. Object recognition supported by user interaction for service robots, IEEE Press, 2002, Vol. 3, pp. 236-240.

[16] B. Gatos, D. Danatsas, I. Pratikakis, and S. J. Perantonis, "Automatic table detection in document images," Proc. International Conference on Pattern Recognition and Image Analysis. Springer Press, Heidelberg, Berlin, Aug. 2005, pp. 609-618.

[17] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation," Proc. IEEE Symp. Conference on Computer Vision and Pattern Recognition(CVPR), 2014, pp. 580-587.

[18] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," Proc. IEEE Symp. Transactions on Pattern Analysis and Machine Intelligence, 2015, 37(9): 1904-1916.

[19] R. Girshick, "Fast R-CNN," Proc. IEEE Symp. International Conference on Computer Vision(ICCV), 2015, pp. 1440-1448.

[20] L. Tychsen-Smith, and L. Petersson, "Denet: Scalable real-time object detection with directed sparse sampling," Proc. IEEE Symp. International Conference on Computer Vision(ICCV), 2017, pp. 428-436.

[21] X. Wang, K. Chen, Z. Huang, C. Yao, and W. Liu, "Point linking network for object detection," 2017, arXiv preprint arXiv:1706.03646.

[22] P. Lyu, C. Yao, W. Wu, S. Yan, and X. Bai, "multi-oriented scene text detection via corner localization and region segmentation," Proc. IEEE Symp. Conference on Computer Vision and Pattern Recognition(CVPR), 2018, pp. 7553-7563.

[23] H. Law, and J. Deng, "Cornernet: Detecting objects as paired keypoints," Proc. European Conference on Computer Vision(ECCV), 2018, pp. 734-750.