

Character Region Awareness for Text Detection

Youngmin Baek, Bado Lee, Dongyo Han, Sangdoo Yun, and Hwalsuk Lee*
 Clova AI Research, NAVER Corp.

{youngmin.baek, bado.lee, dongyo.han, sangdoo.yun, hwalsuk.lee}@navercorp.com

Abstract

Scene text detection methods based on neural networks have emerged recently and have shown promising results. Previous methods trained with rigid word-level bounding boxes exhibit limitations in representing the text region in an arbitrary shape. In this paper, we propose a new scene text detection method to effectively detect text area by exploring each character and affinity between characters. To overcome the lack of individual character level annotations, our proposed framework exploits both the given character-level annotations for synthetic images and the estimated character-level ground-truths for real images acquired by the learned interim model. In order to estimate affinity between characters, the network is trained with the newly proposed representation for affinity. Extensive experiments on six benchmarks, including the TotalText and CTW-1500 datasets which contain highly curved texts in natural images, demonstrate that our character-level text detection significantly outperforms the state-of-the-art detectors. According to the results, our proposed method guarantees high flexibility in detecting complicated scene text images, such as arbitrarily-oriented, curved, or deformed texts.

1. Introduction

Scene text detection has attracted much attention in the computer vision field because of its numerous applications, such as instant translation, image retrieval, scene parsing, geo-location, and blind-navigation. Recently, scene text detectors based on deep learning have shown promising performance [8, 40, 21, 4, 11, 10, 12, 13, 17, 24, 25, 32, 26]. These methods mainly train their networks to localize word-level bounding boxes. However, they may suffer in difficult cases, such as texts that are curved, deformed, or extremely long, which are hard to detect with a single bounding box. Alternatively, character-level awareness has many advantages when handling challenging texts by linking the successive characters in a bottom-up manner. Unfortunately,

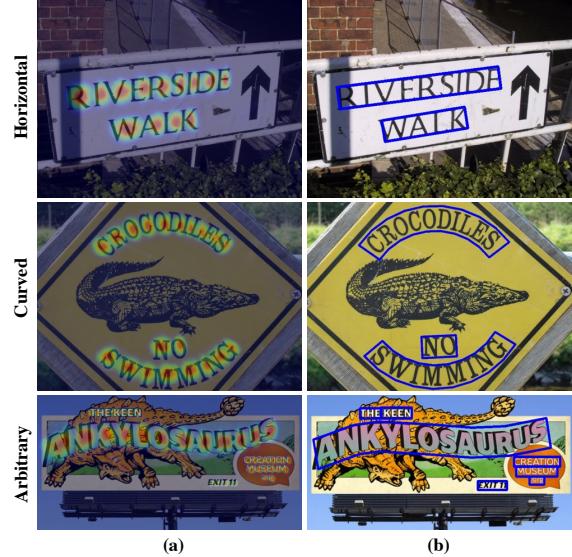


Figure 1. Visualization of character-level detection using CRAFT.
 (a) Heatmaps predicted by our proposed framework. (b) Detection results for texts of various shape.

most of the existing text datasets do not provide character-level annotations, and the work needed to obtain character-level ground truths is too costly.

In this paper, we propose a novel text detector that localizes the individual character regions and links the detected characters to a text instance. Our framework, referred to as CRAFT for *Character Region Awareness For Text detection*, is designed with a convolutional neural network producing the character *region score* and *affinity score*. The *region score* is used to localize individual characters in the image, and the *affinity score* is used to group each character into a single instance. To compensate for the lack of character-level annotations, we propose a weakly-supervised learning framework that estimates character-level ground truths in existing real word-level datasets.

Figure. 1 is a visualization of CRAFT’s results on various shaped texts. By exploiting character-level region awareness, texts in various shapes are easily represented. We demonstrate extensive experiments on ICDAR datasets [15, 14, 28] to validate our method, and the experi-

*Corresponding author.

ments show that the proposed method outperforms state-of-the-art text detectors. Furthermore, experiments on MSRA-TD500, CTW-1500, and TotalText datasets [36, 38, 3] show the high flexibility of the proposed method on complicated cases, such as long, curved, and/or arbitrarily shaped texts.

2. Related Work

The major trend in scene text detection before the emergence of deep learning was bottom-up, where hand-crafted features were mostly used – such as MSER [27] or SWT [5]– as a basic component. Recently, deep learning-based text detectors have been proposed by adopting popular object detection/segmentation methods like SSD [20], Faster R-CNN [30], and FCN [23].

Regression-based text detectors Various text detectors using box regression adapted from popular object detectors have been proposed. Unlike objects in general, texts are often presented in irregular shapes with various aspect ratios. To handle this problem, TextBoxes [18] modified convolutional kernels and anchor boxes to effectively capture various text shapes. DMPNet [22] tried to further reduce the problem by incorporating quadrilateral sliding windows. In recent, Rotation-Sensitive Regression Detector (RSDD) [19] which makes full use of rotation-invariant features by actively rotating the convolutional filters was proposed. However, there is a structural limitation to capturing all possible shapes that exist in the wild when using this approach.

Segmentation-based text detectors Another common approach is based on works dealing with segmentation, which aims to seek text regions at the pixel level. These approaches that detect texts by estimating word bounding areas, such as Multi-scale FCN [7], Holistic-prediction [37], and PixelLink [4] have also been proposed using segmentation as their basis. SSTD [8] tried to benefit from both the regression and segmentation approaches by using an attention mechanism to enhance text related area via reducing background interference on the feature level. Recently, TextSnake [24] was proposed to detect text instances by predicting the text region and the center line together with geometry attributes.

End-to-end text detectors An end-to-end approach trains the detection and recognition modules simultaneously so as to enhance detection accuracy by leveraging the recognition result. FOTS [21] and EAA [10] concatenate popular detection and recognition methods, and train them in an end-to-end manner. Mask TextSpotter [25] took advantage of their unified model to treat the recognition task as a semantic segmentation problem. It is obvious that training with the recognition module helps the text detector be more robust to text-like background clutters.

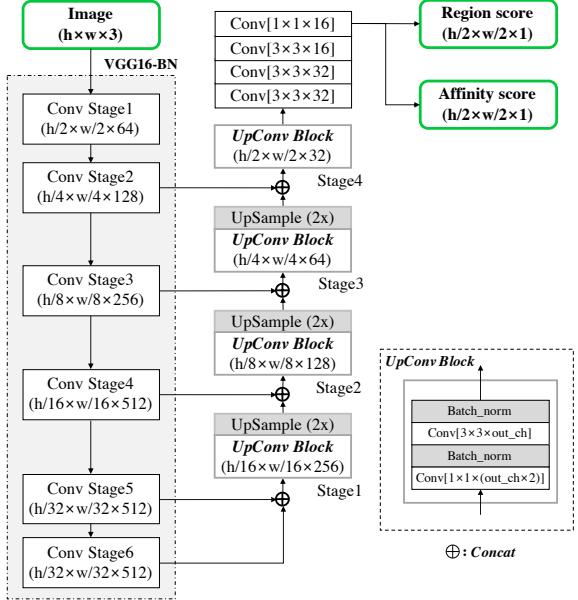


Figure 2. Schematic illustration of our network architecture.

Most methods detect text with words as its unit, but defining the extents to a word for detection is non-trivial since words can be separated by various criteria, such as meaning, spaces or color. In addition, the boundary of the word segmentation cannot be strictly defined, so the word segment itself has no distinct semantic meaning. This ambiguity in the word annotation dilutes the meaning of the ground truth for both regression and segmentation approaches.

Character-level text detectors Zhang et al. [39] proposed a character level detector using text block candidates distilled by MSER [27]. The fact that it uses MSER to identify individual characters limits its detection robustness under certain situations, such as scenes with low contrast, curvature, and light reflection. Yao et al. [37] used a prediction map of the characters along with a map of text word regions and linking orientations that require character level annotations. Instead of an explicit character level prediction, Seglink [32] hunts for text grids (partial text segments) and associates these segments with an additional link prediction. Even though Mask TextSpotter [25] predicts a character-level probability map, it was used for text recognition instead of spotting individual characters.

This work is inspired by the idea of WordSup [12], which uses a weakly supervised framework to train the character-level detector. However, a disadvantage of Wordsup is that the character representation is formed in rectangular anchors, making it vulnerable to perspective deformation of characters induced by varying camera viewpoints. Moreover, it is bound by the performance of the backbone structure (i.e. using SSD and being limited by the number of anchor boxes and their sizes).

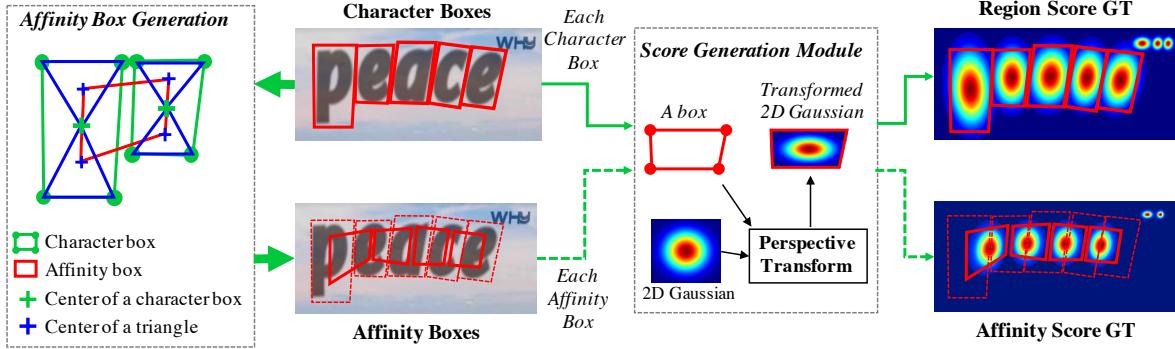


Figure 3. Illustration of ground truth generation procedure in our framework. We generate ground truth labels from a synthetic image that has character level annotations.

3. Methodology

Our main objective is to precisely localize each individual character in natural images. To this end, we train a deep neural network to predict character regions and the affinity between characters. Since there is no public character-level dataset available, the model is trained in a weakly-supervised manner.

3.1. Architecture

A fully convolutional network architecture based on VGG-16 [34] with batch normalization is adopted as our backbone. Our model has skip connections in the decoding part, which is similar to U-net [31] in that it aggregates low-level features. The final output has two channels as score maps: the *region score* and the *affinity score*. The network architecture is schematically illustrated in Fig. 2.

3.2. Training

3.2.1 Ground Truth Label Generation

For each training image, we generate the ground truth label for the *region score* and the *affinity score* with character-level bounding boxes. The *region score* represents the probability that the given pixel is the center of the character, and the *affinity score* represents the center probability of the space between adjacent characters.

Unlike a binary segmentation map, which labels each pixel discretely, we encode the probability of the character center with a Gaussian heatmap. This heatmap representation has been used in other applications, such as in pose estimation works [1, 29] due to its high flexibility when dealing with ground truth regions that are not rigidly bounded. We use the heatmap representation to learn both the *region score* and the *affinity score*.

Fig. 3 summarizes the label generation pipeline for a synthetic image. Computing the Gaussian distribution value directly for each pixel within the bounding box is very time-

consuming. Since character bounding boxes on an image are generally distorted via perspective projections, we use the following steps to approximate and generate the ground truth for both the *region score* and the *affinity score*: 1) prepare a 2-dimensional isotropic Gaussian map; 2) compute perspective transform between the Gaussian map region and each character box; 3) warp Gaussian map to the box area.

For the ground truths of the *affinity score*, the affinity boxes are defined using adjacent character boxes, as shown in Fig. 3. By drawing diagonal lines to connect opposite corners of each character box, we can generate two triangles – which we will refer to as the upper and lower character triangles. Then, for each adjacent character box pair, an affinity box is generated by setting the centers of the upper and lower triangles as corners of the box.

The proposed ground truth definition enables the model to detect large or long-length text instances sufficiently, despite using small receptive fields. On the other hand, previous approaches like box regression require a large receptive field in such cases. Our character-level detection makes it possible for convolutional filters to focus only on intra-character and inter-character, instead of the entire text instance.

3.2.2 Weakly-Supervised Learning

Unlike synthetic datasets, real images in a dataset usually have word-level annotations. Here, we generate character boxes from each word-level annotation in a weakly-supervised manner, as summarized in Fig. 4. When a real image with word-level annotations is provided, the learned interim model predicts the character region score of the cropped word images to generate character-level bounding boxes. In order to reflect the reliability of the interim model’s prediction, the value of the confidence map over each word box is computed proportional to the number of the detected characters divided by the number of the ground truth characters, which is used for the learning weight dur-

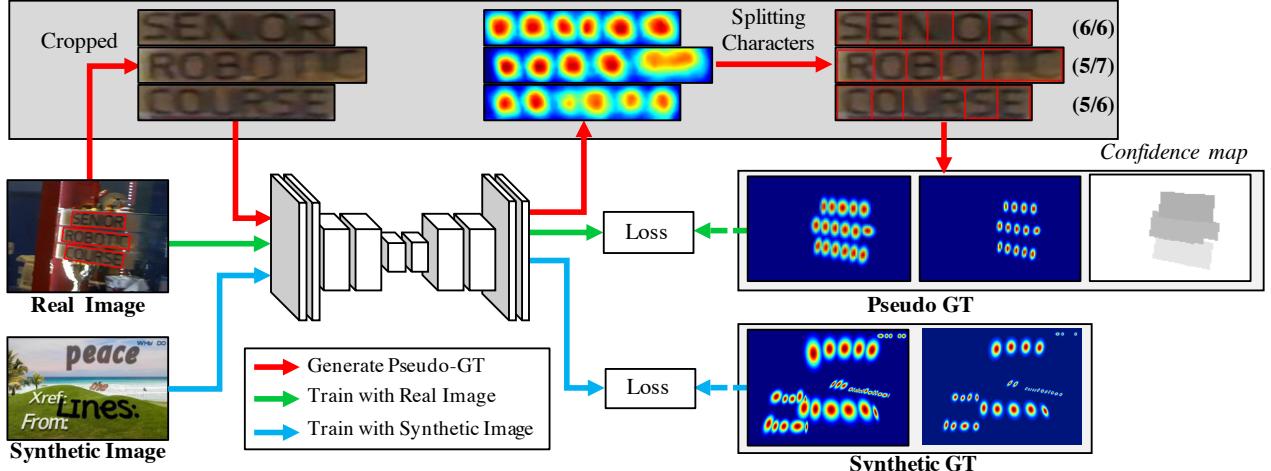


Figure 4. Illustration of the overall training stream for the proposed method. Training is carried out using both real and synthetic images in a weakly-supervised fashion.

ing training.

Fig. 6 shows the entire procedure for splitting the characters. First, the word-level images are cropped from the original image. Second, the model trained up to date predicts the *region score*. Third, the watershed algorithm [35] is used to split the character regions, which is used to make the character bounding boxes covering regions. Finally, the coordinates of the character boxes are transformed back into the original image coordinates using the inverse transform from the cropping step. The pseudo-ground truths (pseudo-GTs) for the *region score* and the *affinity score* can be generated by the steps described in Fig. 3 using the obtained quadrilateral character-level bounding boxes.

When the model is trained using weak-supervision, we are compelled to train with incomplete pseudo-GTs. If the model is trained with inaccurate region scores, the output might be blurred within character regions. To prevent this, we measure the quality of each pseudo-GTs generated by the model. Fortunately, there is a very strong cue in the text annotation, which is the *word length*. In most datasets, the transcription of words is provided and the length of the words can be used to evaluate the confidence of the pseudo-GTs.

For a word-level annotated sample w of the training data, let $R(w)$ and $l(w)$ be the bounding box region and the word length of the sample w , respectively. Through the character splitting process, we can obtain the estimated character bounding boxes and their corresponding length of characters $l^c(w)$. Then the confidence score $s_{conf}(w)$ for the sample w is computed as,

$$s_{conf}(w) = \frac{l(w) - \min(l(w), |l(w) - l^c(w)|)}{l(w)}, \quad (1)$$

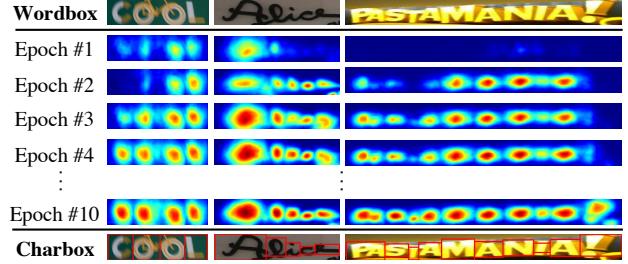


Figure 5. Character region score maps during training.

and the pixel-wise confidence map S_c for an image is computed as,

$$S_c(p) = \begin{cases} s_{conf}(w) & p \in R(w), \\ 1 & \text{otherwise,} \end{cases} \quad (2)$$

where p denotes the pixel in the region $R(w)$. The objective L is defined as,

$$L = \sum_p S_c(p) \cdot (||S_r(p) - S_r^*(p)||_2^2 + ||S_a(p) - S_a^*(p)||_2^2), \quad (3)$$

where $S_r^*(p)$ and $S_a^*(p)$ denote the pseudo-ground truth *region score* and *affinity map*, respectively, and $S_r(p)$ and $S_a(p)$ denote the predicted *region score* and *affinity score*, respectively. When training with synthetic data, we can obtain the real ground truth, so $S_c(p)$ is set to 1.

As training is performed, the CRAFT model can predict characters more accurately, and the confidence scores $s_{conf}(w)$ are gradually increased as well. Fig. 5 shows the character region score map during training. At the early stages of training, the region scores are relatively low for unfamiliar text in natural images. The model learns the ap-

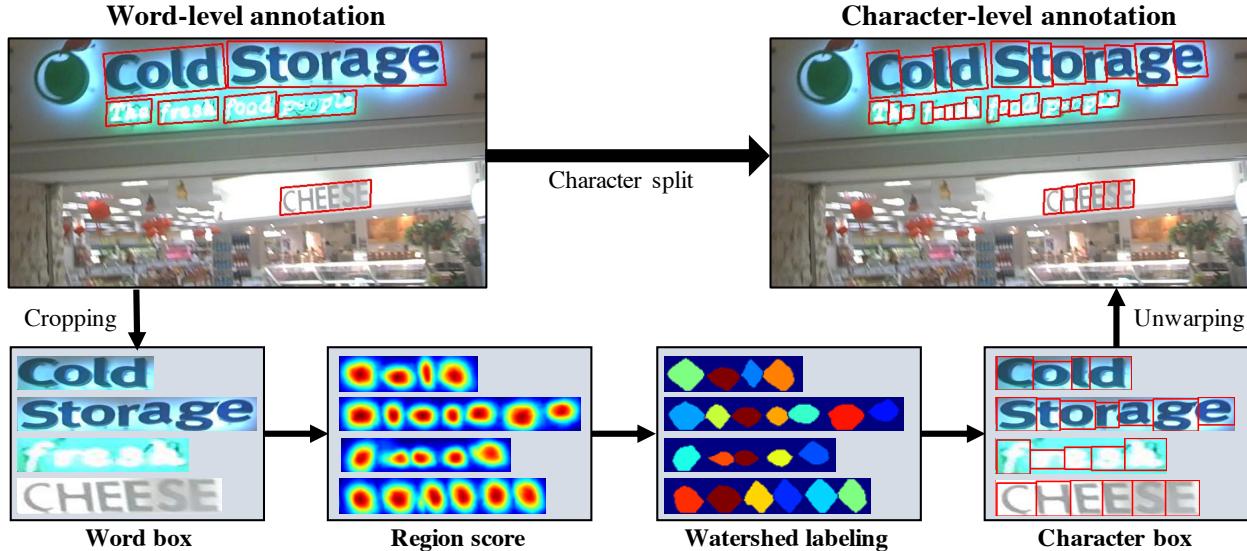


Figure 6. Character split procedure for achieving character-level annotation from word-level annotation: 1) crop the word-level image; 2) predict the region score; 3) apply the watershed algorithm; 4) get the character bounding boxes; 5) unwarp the character bounding boxes.

pearances of new texts, such as irregular fonts, and synthesized texts that have a different data distribution against that of the SynthText dataset.

If the confidence score $s_{conf}(w)$ is below 0.5, the estimated character bounding boxes should be neglected since they have adverse effects when training the model. In this case, we assume the width of the individual character is constant and compute the character-level predictions by simply dividing the word region $R(w)$ by the number of characters $l(w)$. Then, $s_{conf}(w)$ is set to 0.5 to learn unseen appearances of texts.

3.3. Inference

At the inference stage, the final output can be delivered in various shapes, such as word boxes or character boxes, and further polygons. For datasets like ICDAR, the evaluation protocol is word-level intersection-over-union (IoU), so here we describe how to make word-level bounding boxes *QuadBox* from the predicted S_r and S_a through a simple yet effective post-processing step.

The post-processing for finding bounding boxes is summarized as follows. First, the binary map M covering the image is initialized with 0. $M(p)$ is set to 1 if $S_r(p) > \tau_r$ or $S_a(p) > \tau_a$, where τ_r is the region threshold and τ_a is the affinity threshold. Second, Connected Component Labeling (CCL) on M is performed. Lastly, *QuadBox* is obtained by finding a rotated rectangle with the minimum area enclosing the connected components corresponding to each of the labels. The functions like *connectedComponents* and *minAreaRect* provided by OpenCV can be applied for this purpose.

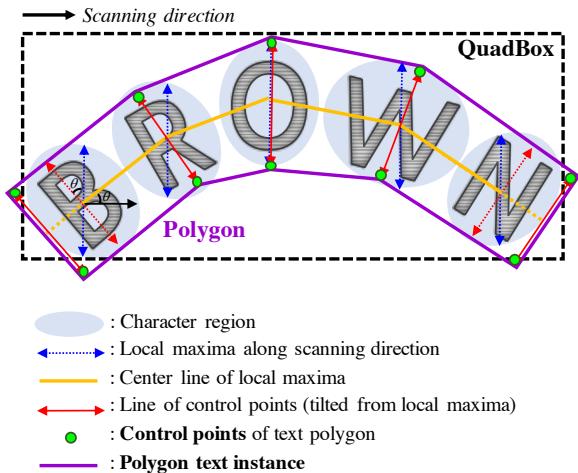


Figure 7. Polygon generation for arbitrarily-shaped texts.

Note that an advantage of CRAFT is that it does not need any further post-processing methods, like Non-Maximum Suppression (NMS). Since we have image blobs of word regions separated by CCL, the bounding box for a word is simply defined by the single enclosing rectangle. On a different note, our character linking process is conducted at a pixel-level, which differs from other linking-based methods [32, 12] relying on searching relations between text components explicitly.

Additionally, we can generate a polygon around the entire character region to deal with curved texts effectively. The procedure of polygon generation is illustrated in Fig. 7. The first step is to find the local maxima line of character regions along the scanning direction, as shown in the figure

Method	IC13(DetEval)			IC15			IC17			MSRA-TD500			FPS
	R	P	H	R	P	H	R	P	H	R	P	H	
Zhang et al. [39]	78	88	83	43	71	54	-	-	-	67	83	74	0.48
Yao et al. [37]	80.2	88.8	84.3	58.7	72.3	64.8	-	-	-	75.3	76.5	75.9	1.61
SegLink [32]	83.0	87.7	85.3	76.8	73.1	75.0	-	-	-	70	86	77	20.6
SSTD [8]	86	89	88	73	80	77	-	-	-	-	-	-	7.7
Wordsup [12]	87.5	93.3	90.3	77.0	79.3	78.2	-	-	-	-	-	-	1.9
EAST* [40]	-	-	-	78.3	83.3	80.7	-	-	-	67.4	87.3	76.1	13.2
He et al. [11]	81	92	86	80	82	81	-	-	-	70	77	74	1.1
R2CNN [13]	82.6	93.6	87.7	79.7	85.6	82.5	-	-	-	-	-	-	0.4
TextSnake [24]	-	-	-	80.4	84.9	82.6	-	-	-	73.9	83.2	78.3	1.1
TextBoxes++* [17]	86	92	89	78.5	87.8	82.9	-	-	-	-	-	-	2.3
EAA [10]	87	88	88	83	84	83	-	-	-	-	-	-	-
Mask TextSpotter [25]	88.1	94.1	91.0	81.2	85.8	83.4	-	-	-	-	-	-	4.8
PixelLink* [4]	87.5	88.6	88.1	82.0	85.5	83.7	-	-	-	73.2	83.0	77.8	3.0
RRD* [19]	86	92	89	80.0	88.0	83.8	-	-	-	73	87	79	10
Lyu et al.* [26]	84.4	92.0	88.0	79.7	89.5	84.3	70.6	74.3	72.4	76.2	87.6	81.5	5.7
FOTS [21]	-	-	87.3	82.0	88.8	85.3	57.5	79.5	66.7	-	-	-	23.9
CRAFT(ours)	93.1	97.4	95.2	84.3	89.8	86.9	68.2	80.6	73.9	78.2	88.2	82.9	8.6

Table 1. Results on quadrilateral-type datasets, such as ICDAR and MSRA-TD500. * denote the results based on multi-scale tests. Methods in *italic* are results solely from the detection of end-to-end models for a fair comparison. R, P, and H refer to recall, precision and H-mean, respectively. The best score is highlighted in **bold**. FPS is for reference only because the experimental environments are different. We report the best FPSs, each of which was reported in the original paper.

with arrows in blue. The lengths of the local maxima lines are equally set as the maximum length among them to prevent the final polygon result from becoming uneven. The line connecting all the center points of the local maxima is called the center line, shown in yellow. Then, the local maxima lines are rotated to be perpendicular to the center line to reflect the tilt angle of characters, as expressed by the red arrows. The endpoints of the local maxima lines are the candidates for the control points of the text polygon. To fully cover the text region, we move the two outer-most tilted local maxima lines outward along the local maxima center line, making the final control points (green dots).

4. Experiment

4.1. Datasets

ICDAR2013 (IC13) was released during the ICDAR 2013 Robust Reading Competition for focused scene text detection, consisting of high-resolution images, 229 for training and 233 for testing, containing texts in English. The annotations are at word-level using rectangular boxes.

ICDAR2015 (IC15) was introduced in the ICDAR 2015 Robust Reading Competition for incidental scene text detection, consisting of 1000 training images and 500 testing images, both with texts in English. The annotations are at the word level using quadrilateral boxes.

ICDAR2017 (IC17) contains 7,200 training images, 1,800 validation images, and 9,000 testing images with texts in 9 languages for multi-lingual scene text detection. Similar to

Method	TotalText			CTW-1500			
	R	P	H	R	P	H	
CTD+TLOC [38]	-	-	-	-	69.8	77.4	73.4
MaskSpotter [25]	55.0	69.0	61.3	-	-	-	-
TextSnake [24]	74.5	82.7	78.4	85.3	67.9	75.6	-
CRAFT(ours)	79.9	87.6	83.6	81.1	86.0	83.5	-

Table 2. Results on polygon-type datasets, such as TotalText and CTW-1500. R, P and H refer to recall, precision and H-mean, respectively. The best score is highlighted in **bold**.

IC15, the text regions in IC17 are also annotated by the 4 vertices of quadrilaterals.

MSRA-TD500 (TD500) contains 500 natural images, which are split into 300 training images and 200 testing images, collected both indoors and outdoors using a pocket camera. The images contain English and Chinese scripts. Text regions are annotated by rotated rectangles.

TotalText (TotalText), recently presented in ICDAR 2017, contains 1255 training and 300 testing images. It especially provides curved texts, which are annotated by polygons and word-level transcriptions.

CTW-1500 (CTW) consists of 1000 training and 500 testing images. Every image has curved text instances, which are annotated by polygons with 14 vertices.

4.2. Training strategy

The training procedure includes two steps: we first use the SynthText dataset [6] to train the network for 50k iter-

ations, then each benchmark dataset is adopted to fine-tune the model. Some “DO NOT CARE” text regions in ICDAR 2015 and ICDAR 2017 datasets are ignored in training by setting $s_{conf}(w)$ to 0. We use the ADAM [16] optimizer in all training processes. For multi-GPU training, the training and supervision GPUs are separated, and pseudo-GTs generated by the supervision GPUs are stored in the memory. During fine-tuning, the SynthText dataset is also used at a rate of 1:5 to make sure that the character regions are surely separated. In order to filter out texture-like texts in natural scenes, *On-line Hard Negative Mining* [33] is applied at a ratio of 1:3. Also, basic data augmentation techniques like crops, rotations, and/or color variations are applied.

Weakly-supervised training requires two types of data; quadrilateral annotations for cropping word images and transcriptions for calculating word length. The datasets meeting these conditions are IC13, IC15, and IC17. Other datasets such as MSRA-TD500, TotalText, and CTW-1500 do not meet the requirements. MSRA-TD500 does not provide transcriptions, while TotalText and CTW-1500 provide polygon annotations only. Therefore, we trained CRAFT only on the ICDAR datasets, and tested on the others without fine-tuning. Two different models are trained with the ICDAR datasets. The first model is trained on IC15 to evaluate IC15 only. The second model is trained on both IC13 and IC17 together, which is used for evaluating the other five datasets. No extra images are used for training. The number of iterations for fine-tuning is set to 25k.

4.3. Experimental Results

Quadrilateral-type datasets (ICDARs, and MSRA-TD500) All experiments are performed with a single image resolution. The longer side of the images in IC13, IC15, IC17, and MSRA-TD500 are resized to 960, 2240, 2560, and 1600, respectively. Table 1 lists the h-mean scores of various methods on ICDAR and MSRA-TD500 datasets. To have a fair comparison with end-to-end methods, we include their detection-only results by referring to the original papers. We achieve state-of-the-art performances on all the datasets. In addition, CRAFT runs at 8.6 FPS on IC13 dataset, which is comparatively fast, thanks to the simple yet effective post-processing.

For MSRA-TD500, annotations are provided at the line-level, including the spaces between words in the box. Therefore, a post-processing step for combining word boxes is applied. If the right side of one box and the left side of another box are close enough, the two boxes are combined together. Even though fine-tuning is not performed on the TD500 training set, CRAFT outperforms all other methods as shown in Table 1.

Polygon-type datasets (TotalText, CTW-1500) It is challenging to directly train the model on TotalText and CTW-1500 because their annotations are in polygonal in shape,

Method	IC13	IC15	IC17
Mask TextSpotter [25]	91.7	86.0	-
EAA [10]	90	87	-
FOTS [21]	92.8	89.8	70.8
CRAFT(ours)	95.2	86.9	73.9

Table 3. H-mean comparison with end-to-end methods. Our method is not trained in an end-to-end manner, yet shows comparable results, or even outperforms popular methods.

which complicates text area cropping for splitting character boxes during weakly-supervised training. Consequently, we only used the training images from IC13 and IC17, and fine-tuning was not conducted to learn the training images provided by these datasets. At the inference step, we used the polygon generation post-processing from the *region score* to cope with the provided polygon-type annotations.

The experiments for these datasets are performed with a single image resolution, too. The longer sides of the images within TotalText and CTW-1500 are resized to 1280 and 1024, respectively. The experimental results for polygon-type datasets are shown in Table 2. The individual-character localization ability of CRAFT enables us to achieve more robust and superior performance in detecting arbitrarily shaped texts compared to other methods. Particularly, the TotalText dataset has a variety of deformations, including curved texts as shown in Fig. 8, for which adequate inference by quadrilateral-based text detectors is infeasible. Therefore, a very limited number of methods can be evaluated on those datasets.

In the CTW-1500 dataset’s case, two difficult characteristics coexist, namely annotations that are provided at the line-level and are of arbitrary polygons. To aid CRAFT in such cases, a small link refinement network, which we call the *LinkRefiner*, is used in conjunction with CRAFT. The input of the *LinkRefiner* is a concatenation of the *region score*, the *affinity score*, and the intermediate feature map of CRAFT, and the output is a *refined affinity score* adjusted for long texts. To combine characters, the *refined affinity score* is used instead of the original *affinity score*, then the polygon generation is performed in the same way as it was performed for TotalText. Only *LinkRefiner* is trained on the CTW-1500 dataset while freezing CRAFT. The detailed implementation of *LinkRefiner* is addressed in the supplementary materials. As shown in Table 2, the proposed method achieves state-of-the-art performance.

4.4. Discussions

Robustness to Scale Variance We solely performed single-scale experiments on all the datasets, even though the size of texts are highly diverse. This is different from the majority of other methods, which rely on multi-scale tests to handle



Figure 8. Results on the TotalText dataset. First row: each column shows the input image (top) with its respective region score map (bottom left) and affinity map (bottom right). Second row: each column only shows the input image (left) and its region score map (right).

the scale variance problem. This advantage comes from the property of our method localizing individual characters, not the whole text. The relatively small receptive field is sufficient to cover a single character in a large image, which makes CRAFT robust in detecting scale variant texts.

Multi-language issue The IC17 dataset contains Bangla and Arabic characters, which are not included in the synthetic text dataset. Moreover, both languages are difficult to segment into characters individually because every character is written cursively. Therefore, our model could not distinguish Bangla and Arabic characters as well as it does Latin, Korean, Chinese, and Japanese. In East Asian characters’ cases, they can be easily separated with a constant width, which helps train the model to high performance via weakly-supervision.

Comparison with End-to-end methods Our method is trained with the ground truth boxes only for detection, but it is comparable with other end-to-end methods, as shown in Table. 3. From the analysis of failure cases, we expect our model to benefit from the recognition results, especially when the ground truth words are separated by semantics, rather than visual cues.

Generalization ability Our method achieved state-of-the-art performances on 3 different datasets without additional fine-tuning. This demonstrates that our model is capable of

capturing general characteristics of texts, rather than overfitting to a particular dataset.

5. Conclusion

We have proposed a novel text detector called CRAFT, which can detect individual characters even when character-level annotations are not given. The proposed method provides the *character region score* and the *character affinity score* that, together, fully cover various text shapes in a bottom-up manner. Since real datasets provided with character-level annotations are rare, we proposed a weakly-supervised learning method that generates pseudo-ground truths from an interim model. CRAFT shows state-of-the-art performances on most public datasets and demonstrates generalization ability by showing these performances without fine-tuning. As our future work, we hope to train our model with a recognition model in an end-to-end fashion to see whether the performance, robustness, and generalizability of CRAFT translates to a better scene text spotting system that can be applied in more general settings.

Acknowledgements. The authors would like to thank Beomyoung Kim, Daehyun Nam, and Donghyun Kim for helping with extensive experiments.

References

- [1] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, pages 1302–1310. IEEE, 2017. 3
- [2] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *PAMI*, 40(4):834–848, 2018. 11
- [3] C. K. Ch’ng and C. S. Chan. Total-text: A comprehensive dataset for scene text detection and recognition. In *ICDAR*, volume 1, pages 935–942. IEEE, 2017. 2
- [4] D. Deng, H. Liu, X. Li, and D. Cai. Pixellink: Detecting scene text via instance segmentation. In *AAAI*, 2018. 1, 2, 6
- [5] B. Epshtain, E. Ofek, and Y. Wexler. Detecting text in natural scenes with stroke width transform. In *CVPR*, pages 2963–2970. IEEE, 2010. 2
- [6] A. Gupta, A. Vedaldi, and A. Zisserman. Synthetic data for text localisation in natural images. In *CVPR*, pages 2315–2324, 2016. 6
- [7] D. He, X. Yang, C. Liang, Z. Zhou, G. Alexander, I. Ororbia, D. Kifer, and C. L. Giles. Multi-scale fcn with cascaded instance aware segmentation for arbitrary oriented word spotting in the wild. In *CVPR*, pages 474–483, 2017. 2
- [8] P. He, W. Huang, T. He, Q. Zhu, Y. Qiao, and X. Li. Single shot text detector with regional attention. In *ICCV*, volume 6, 2017. 1, 2, 6
- [9] T. He, W. Huang, Y. Qiao, and J. Yao. Accurate text localization in natural image with cascaded convolutional text network. *arXiv preprint arXiv:1603.09423*, 2016. 11
- [10] T. He, Z. Tian, W. Huang, C. Shen, Y. Qiao, and C. Sun. An end-to-end textspotter with explicit alignment and attention. In *CVPR*, pages 5020–5029, 2018. 1, 2, 6, 7
- [11] W. He, X.-Y. Zhang, F. Yin, and C.-L. Liu. Deep direct regression for multi-oriented scene text detection. In *CVPR*, pages 745–753, 2017. 1, 6
- [12] H. Hu, C. Zhang, Y. Luo, Y. Wang, J. Han, and E. Ding. Wordsup: Exploiting word annotations for character based text detection. In *ICCV*, 2017. 1, 2, 5, 6
- [13] Y. Jiang, X. Zhu, X. Wang, S. Yang, W. Li, H. Wang, P. Fu, and Z. Luo. R2cnn: rotational region cnn for orientation robust scene text detection. *arXiv preprint arXiv:1706.09579*, 2017. 1, 6
- [14] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, et al. Icdar 2015 competition on robust reading. In *ICDAR*, pages 1156–1160. IEEE, 2015. 1
- [15] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazan, and L. P. De Las Heras. Icdar 2013 robust reading competition. In *ICDAR*, pages 1484–1493. IEEE, 2013. 1
- [16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 7
- [17] M. Liao, B. Shi, and X. Bai. Textboxes++: A single-shot oriented scene text detector. *Image Processing*, 27(8):3676–3690, 2018. 1, 6
- [18] M. Liao, B. Shi, X. Bai, X. Wang, and W. Liu. Textboxes: A fast text detector with a single deep neural network. In *AAAI*, pages 4161–4167, 2017. 2
- [19] M. Liao, Z. Zhu, B. Shi, G.-s. Xia, and X. Bai. Rotation-sensitive regression for oriented scene text detection. In *CVPR*, pages 5909–5918, 2018. 2, 6
- [20] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016. 2
- [21] X. Liu, D. Liang, S. Yan, D. Chen, Y. Qiao, and J. Yan. Fots: Fast oriented text spotting with a unified network. In *CVPR*, pages 5676–5685, 2018. 1, 2, 6, 7
- [22] Y. Liu and L. Jin. Deep matching prior network: Toward tighter multi-oriented text detection. In *CVPR*, pages 3454–3461, 2017. 2
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. 2
- [24] S. Long, J. Ruan, W. Zhang, X. He, W. Wu, and C. Yao. Textsnake: A flexible representation for detecting text of arbitrary shapes. *arXiv preprint arXiv:1807.01544*, 2018. 1, 2, 6
- [25] P. Lyu, M. Liao, C. Yao, W. Wu, and X. Bai. Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. *arXiv preprint arXiv:1807.02242*, 2018. 1, 2, 6, 7
- [26] P. Lyu, C. Yao, W. Wu, S. Yan, and X. Bai. Multi-oriented scene text detection via corner localization and region segmentation. In *CVPR*, pages 7553–7563, 2018. 1, 6
- [27] J. Matas, O. Chum, M. Urban, and T. Pajdla. Robust wide-baseline stereo from maximally stable extremal regions. *Image and Vision Computing*, 22(10):761–767, 2004. 2
- [28] N. Nayef, F. Yin, I. Bidid, H. Choi, Y. Feng, D. Karatzas, Z. Luo, U. Pal, C. Rigaud, J. Chazalon, et al. Icdar2017 robust reading challenge on multi-lingual scene text detection and script identification-rrc-mlt. In *ICDAR*, volume 1, pages 1454–1459. IEEE, 2017. 1
- [29] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, pages 483–499. Springer, 2016. 3
- [30] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: towards real-time object detection with region proposal networks. *PAMI*, (6):1137–1149, 2017. 2
- [31] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241. Springer, 2015. 3
- [32] B. Shi, X. Bai, and S. Belongie. Detecting oriented text in natural images by linking segments. In *CVPR*, pages 3482–3490. IEEE, 2017. 1, 2, 5, 6
- [33] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, pages 761–769, 2016. 7
- [34] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 3
- [35] L. Vincent and P. Soille. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. *PAMI*, (6):583–598, 1991. 4

- [36] C. Yao, X. Bai, W. Liu, Y. Ma, and Z. Tu. Detecting texts of arbitrary orientations in natural images. In *CVPR*, pages 1083–1090. IEEE, 2012. [2](#)
- [37] C. Yao, X. Bai, N. Sang, X. Zhou, S. Zhou, and Z. Cao. Scene text detection via holistic, multi-channel prediction. *arXiv preprint arXiv:1606.09002*, 2016. [2, 6](#)
- [38] L. Yuliang, J. Lianwen, Z. Shuaiteao, and Z. Sheng. Detecting curve text in the wild: New dataset and new solution. *arXiv preprint arXiv:1712.02170*, 2017. [2, 6, 11](#)
- [39] Z. Zhang, C. Zhang, W. Shen, C. Yao, W. Liu, and X. Bai. Multi-oriented text detection with fully convolutional networks. In *CVPR*, pages 4159–4167, 2016. [2, 6](#)
- [40] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang. East: an efficient and accurate scene text detector. In *CVPR*, pages 2642–2651, 2017. [1, 6](#)

A. LinkRefiner for CTW-1500 dataset

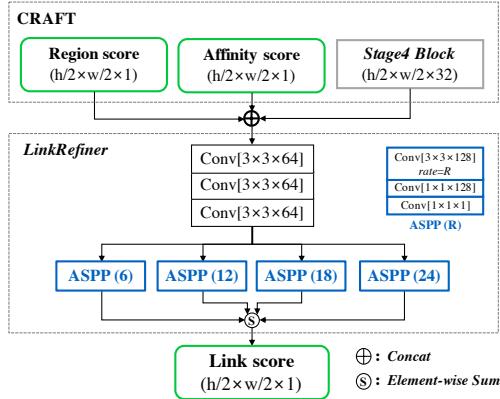


Figure 9. Schematic illustration of *LinkRefiner* architecture.

CTW-1500 dataset [38] provides polygon-only annotations without text transcriptions. Furthermore, annotations of CTW-1500 are provided at the line-level and does not consider spaces as separation cues. This is far from our assumption of affinity, which is that the score for affinity is zero for characters with a space between them.

To obtain a single-long polygon from the detected characters, we employ a shallow network for link refinement, so called *LinkRefiner*. The architecture of the *LinkRefiner* is shown in Fig. 9. The input of the *LinkRefiner* is a concatenation of the *region score*, the *affinity score*, and the intermediate feature map from the network, that is the output of *Stage4* of the original CRAFT model. Atrous Spatial Pyramid Pooling (ASPP) in [2] is adopted to ensure a large receptive field for combining distant characters and words onto the same text line.

For the ground truth of the *LinkRefiner*, lines are simply drawn between the centers of the paired control points of the annotated polygons, which is similar to the text line generation used in [9]. The width of each line is proportional to the distance between paired control points. The ground truth generation for the *LinkRefiner* is illustrated in Fig. 10. The output of the model is called the *link score*. For training, only the *LinkRefiner* is trained on the CTW-1500 training dataset, while freezing CRAFT.

After training, we have the outputs produced by the model, which are the *region score*, the *affinity score*, and the *link score*. Here, the *link score* is used instead of the original *affinity score*, and the text polygon is obtained entirely through the same process as done with TotalText. The CRAFT model localizes the individual characters, and the *LinkRefiner* model combines the characters as well as the words separated by spaces, which are required by the CTW-1500 evaluation.

The results on the CTW-1500 dataset are shown in

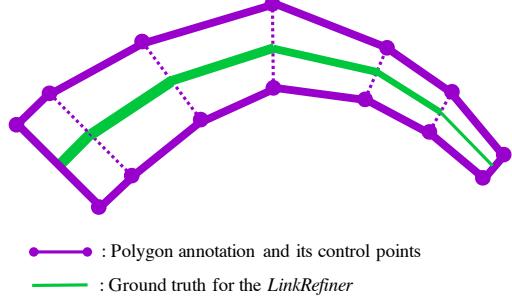


Figure 10. Ground truth generation for *LinkRefiner*.

Fig. 11. Very challenging image samples with long and curved texts are successfully detected by the proposed method. Moreover, with our polygon representation, the curved images can be rectified into straight text images, which are also shown in Fig. 11. We believe this ability for rectification can further be of use for recognition tasks.

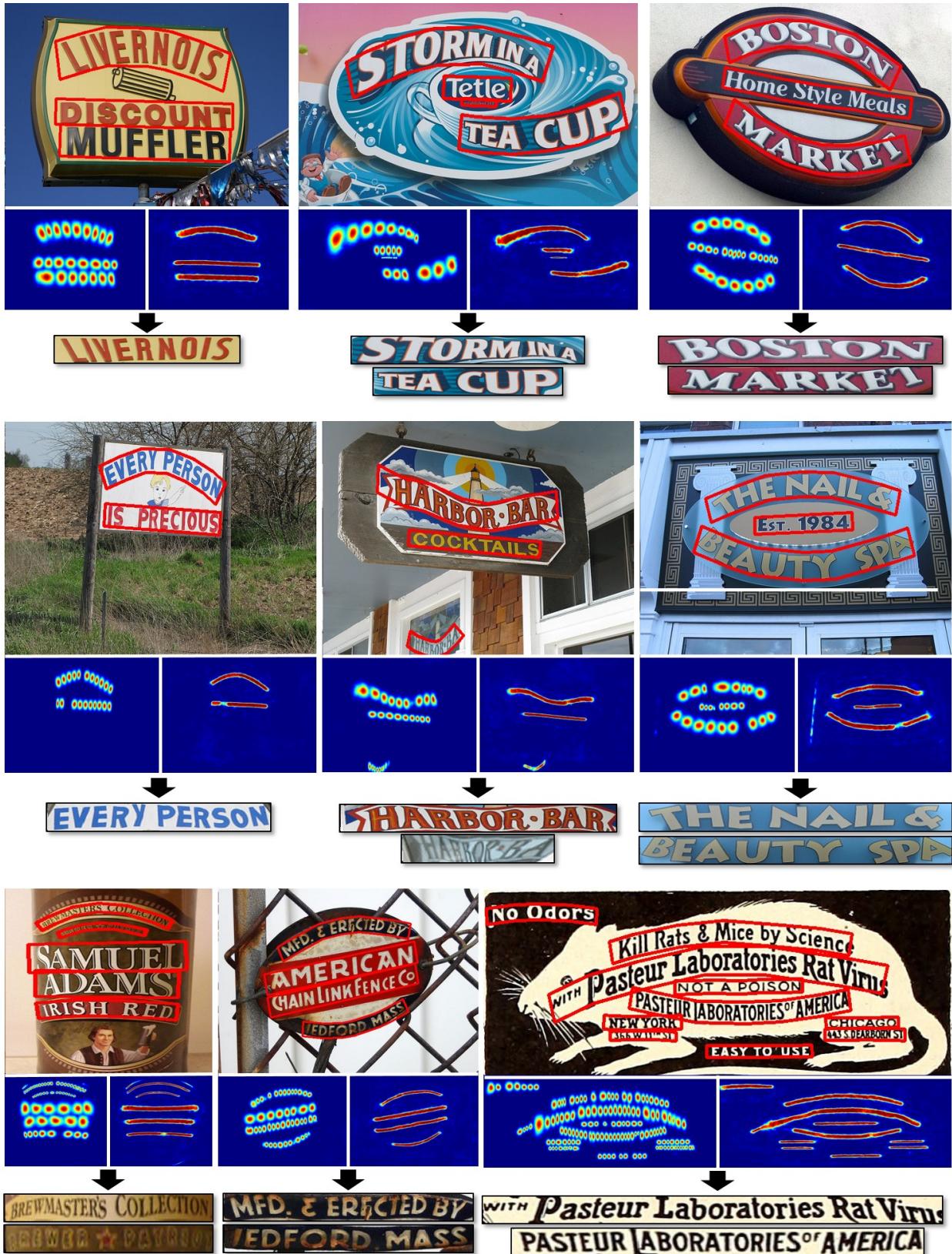


Figure 11. Results on CTW-1500 dataset. For each cluster: the input image (top), region score (middle left), link score (middle right), and the resulting rectified polygons for curved texts (bottom, below arrow) are shown. Note that the affinity scores are not rendered and are unused in the CTW-1500 dataset.