

AON: Towards Arbitrarily-Oriented Text Recognition

Zhazhan Cheng¹ Yangliu Xu² Fan Bai³ Yi Niu¹
 Shiliang Pu¹ Shuigeng Zhou^{3*}

¹Hikvision Research Institute, China; ²Tongji University, Shanghai, China;

³Shanghai Key Lab of Intelligent Information Processing, and School of
 Computer Science, Fudan University, Shanghai, China;

{chengzhazhan; xuyangliu; niuyi; pushiliang}@hikvision.com;
 1993sallyxyl@tongji.edu.cn; {fbai17; sgzhou}@fudan.edu.cn

Abstract

Recognizing text from natural images is a hot research topic in computer vision due to its various applications. Despite the enduring research of several decades on optical character recognition (OCR), recognizing texts from natural images is still a challenging task. This is because scene texts are often in irregular (e.g. curved, arbitrarily-oriented or seriously distorted) arrangements, which have not yet been well addressed in the literature. Existing methods on text recognition mainly work with regular (horizontal and frontal) texts and cannot be trivially generalized to handle irregular texts. In this paper, we develop the arbitrary orientation network (AON) to directly capture the deep features of irregular texts, which are combined into an attention-based decoder to generate character sequence. The whole network can be trained end-to-end by using only images and word-level annotations. Extensive experiments on various benchmarks, including the CUTE80, SVT-Perspective, IIIT5k, SVT and ICDAR datasets, show that the proposed AON-based method achieves the-state-of-the-art performance in irregular datasets, and is comparable to major existing methods in regular datasets.

In the past years, there have been many works to solve scene text recognition [6, 22, 31, 39]. Although these approaches have shown promising results, most of them can effectively handle only regular texts that are often tightly-bounded, horizontal and frontal. However, in real-world



Figure 1. Examples of irregular (slant/perspective, curved and oriented etc.) texts in natural images. Subfigures (a) - (b), (c) - (d) and (e) - (f) are slant/perspective, curved and oriented images respectively.

applications, many scene texts are in irregular arrangements (e.g. arbitrarily-oriented, curved, slant and perspective etc.) as shown in Fig. 1, so most existing methods cannot be widely applied in practice.

Recently, there are two related works aiming at irregular texts: the spatial transformer network (STN) [18] - based method by [32] and the attention-based method with fully convolutional network (FCN) [23] by [39]. Shi *et al.* [32] attempted to first rectify irregular (e.g. curved or perspective distorted) texts to approximately regular texts, then recognized the rectified images with an attention-based sequence recognition network. However, in complicated (e.g. arbitrarily-oriented or serious curved) natural scenes, it is hard to optimize the STN-based method without human-labeled geometric ground truth. Besides, training STN needs sophisticated skills. For example, the

1. Introduction

Scene text recognition has attracted much research interest of the computer vision community [6, 15, 22, 27, 31, 39] because of its various applications such as road sign recognition and navigation reading for advanced driver assistant system (ADAS). Though Optical Character Recognition (OCR) has been extensively studied for several decades, recognizing texts from natural images is still a challenging task due to complicated environments (e.g. uneven lighting, blurring, perspective distortion and orientation).

*Corresponding author.

thin-plate-spline (TPS) [5]-based STN [32] should be given some initialization pattern for the fiducial points, and is not quite effective for arbitrarily-oriented scene texts. Yang *et al.* [39] introduced an auxiliary dense character detection task for encouraging the learning of visual representations with a fully convolutional network. Though the method showed better performance on irregular texts, it was carried out with an exhausting multi-task learning (MTL) strategy and relied on character-level bounding box annotations. Note that, though the attention-based model has the potential to perform 2D feature selection [38], we found in experiments that directly training attention-based model on irregular texts is difficult due to irregular character placements. This situation motivates us to explore new and more effective methods to recognize irregular scene texts.

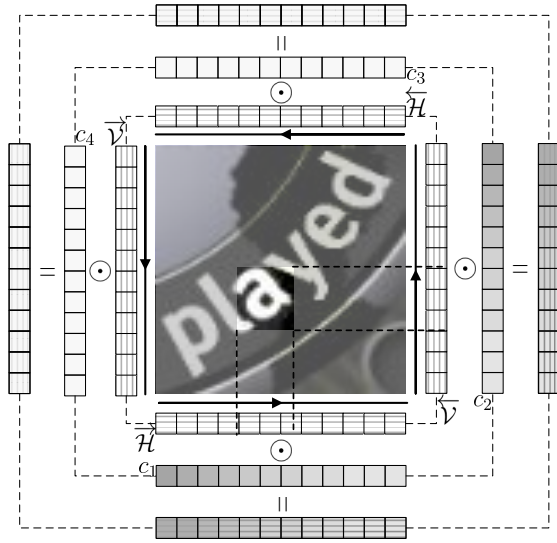


Figure 2. Illustration of character visual representation in four directions: \vec{H} : $left \rightarrow right$, \overleftarrow{H} : $right \rightarrow left$, \vec{V} : $top \rightarrow bottom$ and \overleftarrow{V} : $bottom \rightarrow top$ and four character placement clues c_1, c_2, c_3 and c_4 . Here, there are three squares connected with dashed lines. The innermost square represents the four 1D sequences of features, each comes along with an arrowed line. The middle square refers to the placement clues used for weighting the corresponding sequences of features. The outermost square stands for the weighted feature sequence by conducting Hadamard product \odot with character placement clues and horizontal/vertical features. For the character ‘a’ in the image, we can represent it by the four weighted sequences of features.

From the above analysis, we can see that most existing methods directly encode a text image as a 1D sequence of features and then decode them to the predicted text, which implies that any text in an image is treated in the same direction such as *from left to right* by default. However, this is not true in the wild. After carefully analyzing the typical character placement styles of natural text images,

we suggest that the visual representation of an arbitrarily-oriented character in a 2D image can be described in four directions: $left \rightarrow right$, $right \rightarrow left$, $top \rightarrow bottom$ and $bottom \rightarrow top$. Concretely, we can encode the input image to four feature sequences of four directions: *horizontal features* (\vec{H}), *reversed horizontal features* (\overleftarrow{H}), *vertical features* (\vec{V}) and *reversed vertical features* (\overleftarrow{V}), as shown in Fig. 1, and the length of each sequence is equal. The horizontal/vertical features can be extracted by down-sampling the height/width of feature maps to 1. In order to represent an arbitrarily-oriented character, a weighting mechanism can be used to combine the four feature sequences of different directions. We call the weights *character placement clues*, which are denoted as c_1, c_2, c_3 and c_4 in Fig. 1. The character placement clues can be learned from the input images with a convolutional-based network, which guides to effectively integrate the four sequences of features, and then a filter gate (FG) generates the integrated feature sequence as the character’s visual representation. Therefore, an arbitrarily-oriented character in a 2D image can be represented as the combination of horizontal and vertical features by conducting the Hadamard product with the sequences of features and the corresponding placement clues. In Fig. 1, c_1 and c_2 play the dominant role in determining the visual representation of character ‘a’. In this paper, we call the four-direction feature extraction network and the clues extraction network *arbitrary orientation network* (AON), which means that it can effectively handle arbitrarily-oriented texts.

In this paper, we develop a novel method for robustly recognizing both regular and irregular natural texts by employing the proposed *arbitrary orientation network* (AON).

Major contributions of this paper are as follows:

1. We propose the arbitrary orientation network (AON) to extract scene text features in four directions and the character placement clues.
2. We design a filter gate (FG) for fusing four-direction features with the learned placement clues. That is, FG is responsible for generating the integrated feature sequence.
3. We integrate AON, FG and an attention-based decoder into the character recognition framework. The whole network can be directly trained end-to-end without any character-level bounding box annotations.
4. We conduct extensive experiments on several public irregular and regular text benchmarks, which show that our method obtains state-of-the-art performance in irregular benchmarks, and is comparable to major existing methods in regular benchmarks.

2. Related works

In recent years, several methods have been proposed for scene text recognition. For the general information of text recognition, readers can refer to Ye and Doermann’s recent survey [41]. Basically, there are two types of scene text recognition approaches: bottom-up and top-down.

Traditional methods mostly follow the *bottom-up* pipeline: first extracting low-level features for individual character detection and recognition one by one, then integrating these characters into words based on a set of heuristic rules or a language model. For example, [27] defined a set of handcrafted features such as aspect ratio, hole area ratio etc. to train a Support Vector Machine (SVM) classifier. [35, 36] first fetched each character in the cropped word image by sliding window, then recognized it with a character classifier trained by the extracted HOG descriptors [40]. However, the performance of these methods is limited due to the low representation capability of handcrafted features. With the advancement of neural-network-based methods, many researchers developed deep neural architectures and achieved better results. [4] adopted a fully connected network of 5 hidden layers for character feature representation, then used an n-gram language model to recognize characters. [37] developed a CNN-based feature extraction framework for character recognition, and applied a non-maximum suppression method for final word predictions. [16] also proposed a CNN-based method with structured output layer for unconstrained recognition. These above methods require the segmentation of each character, which can be very challenging because of the complicated background clutter and the inadequate distance between consecutive characters. Besides, segmentation annotations require additional resource consuming.

The other approaches work in a top-down style: directly predicting the entire text from the original image without detecting the characters. [17] conducted a 90k-class classification task with a CNN, in which each class represents an English word. Consequently, the model can not recognize out-of-vocabulary words. Recent works solve this problem as a sequence recognition problem, where images and texts are separately encoded as patch sequences and character sequences, respectively. [34] extracted sequences of HOG features to represent images, and generated the character sequence with the recurrent neural network (RNN). [13] and [31] proposed the end-to-end neural networks that combines CNN and RNN for visual feature representation, then the CTC [10] Loss was combined with the RNN outputs for calculating the conditional probability between the predicted and the target sequences. [22] used a recursive CNN to learn broader contextual information, and applied the attention-based decoder for sequence generation. [6] proposed a focus mechanism to eliminate the attention drift to improve the regular text recognition performance. However, since

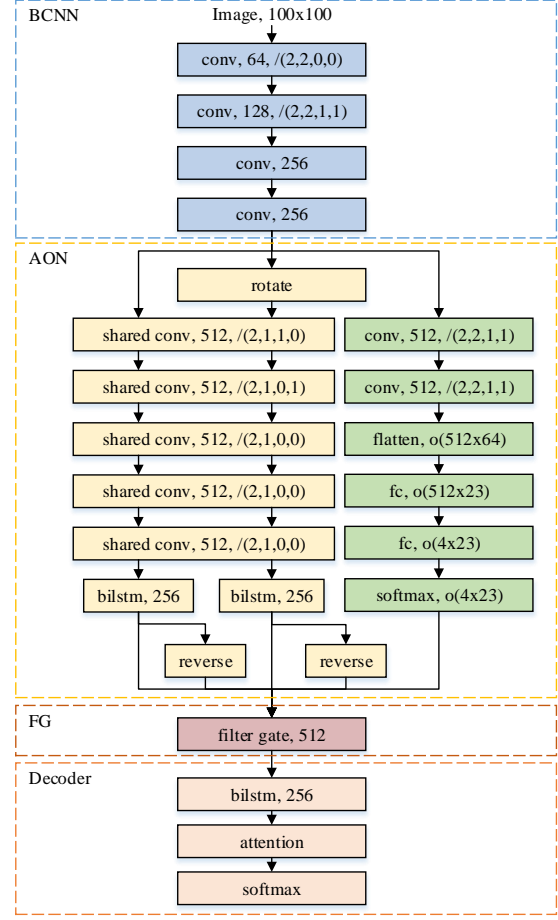


Figure 3. The network architecture of our method, which consists of four components: 1) the basal convolutional neural network (BCNN) for low-level visual representation; 2) the arbitrary orientation network (AON) for capturing the horizontal, vertical and character placement features; 3) the filter gate (FG) for combing four feature sequences with the character placement clues; 4) the attention-based decoder (Decoder) for predicting character sequence. The above four modules are shown in the blue, golden, dull-red and brown dashed boxes, respectively. Meanwhile, all convolution or shared convolution blocks have the following format: $name, c, / (s_h, s_w, p_h, p_w)$. The bilstm and filter gate blocks are represented as $name, c$. The flatten, fc (fully-connected) and softmax operations have the format: $name, o(c, l)$. Here, c, s_h, s_w, p_h, p_w, l and o represent the number of channels, stride height, stride width, pad height, pad width, length of feature maps, pooling operation and output shape, respectively. The whole network can be trained end-to-end.

a text image is encoded into a 1D-based sequence of features, these methods can not effectively handle the irregular texts such as the arbitrarily-oriented texts. In order to recognize irregular texts, [32] applied the spatial transformer network (STN) [18] for text rectification, then recognized

the rectified text images with the sequence recognition network. [39] introduced an auxiliary dense character detection task for encouraging the learning of visual representations with a fully convolutional network (FCN) [23]. In practice, training STN-based methods is extremely difficult without human-labeled geometric ground truth, especially for texts in complicated (e.g. curved, arbitrarily-oriented or perspective etc.) environments. Besides, sophisticated tricks are also required. For example, to train the thin-plate-spline (TPS) [5]-based STN [32]-based method, the initialization pattern should be given for the fiducial points. Though [39] can recognize characters in a 2D image, the method relies on the multi-task learning framework (including 3 task branches and 2 tunable super-parameters) and character-level bounding box annotations, which results in large amount of resource consuming. While obtaining better performance on irregular texts, it performs worse on regular texts.

Different from existing approaches, in this paper we first extract deep feature representations of images by using an arbitrary orientation network (AON), then use a filter gate (FG) to generate the integrated sequence of features, which are fed to an attention-based decoder for generating predicted sequences. Furthermore, we can once for all train the whole network end-to-end with only word-level annotations.

Note that in the OCR field, natural text reading systems often consist of two steps: 1) detecting each word's location in natural images and 2) recognizing text from the cropped image. In general, robust detection is helpful in recognizing texts. Therefore, several methods [20, 25, 30, 43] have been proposed for multi-oriented text detection. Though this work focuses on the recognition task, our AON-based method can directly recognize arbitrarily-oriented texts, which alleviates the pressure of text detection.

3. The Framework

The framework of whole network is shown in Fig. 3, which consists of four major components: 1) The *basal convolutional neural network* (a nomenclature for initial layer, denoted by BCNN) for extracting low-level visual features; 2) The *arbitrary orientation network* (AON) for generating four-direction sequences of features and the character placement clues; 3) The *filter gate* (FG) for combining the four sequences of features with the learned placement clues to generate the integrated feature sequence, and 4) the *attention-based decoder* for predicting character sequence.

3.1. Basal Convolutional Neural Network (BCNN)

The BCNN module is responsible for capturing the foundational visual representation of text images, and outputs a group of feature maps. BCNN can help reduce the computational cost and graphic memory. As shown in Fig. 3, we use

four convolution blocks as the foundational feature extractor. The outputs of BCNN must be square feature maps. We empirically found that higher-level feature representation as the initial state of AON can yield better performance.

3.2. Multi-Direction Feature Extraction Module

This module includes the arbitrary orientation network (AON) and the filter gate (FG), which constitute the core of the proposed method. With the extracted foundational features, we devise AON for capturing arbitrarily-oriented text features and the corresponding character placement clues. We also design FG for integrating multi-direction features by using the character placement clues. The details of AON and FG will be described in next section.

3.3. Attention-based Decoder

An attention-based decoder is a recurrent neural network (RNN) that directly generates the target sequence (y_1, \dots, y_M) from an input feature sequence $(\hat{h}_1, \dots, \hat{h}_L)$. Bahdanau et al. [3] first proposed the architecture of attention-based decoder. At the t -th step, the attention module generates an output y_t as follows:

$$y_t = \text{softmax}(W^T s_t), \quad (1)$$

where W^T is a learnable parameter, and s_t is the RNN hidden state at time t , computed by

$$s_t = \text{RNN}(y_{t-1}, g_t, s_{t-1}), \quad (2)$$

where g_t is the weighted sum of sequential feature vectors $\hat{\mathcal{H}} : (\hat{h}_1, \dots, \hat{h}_L)$, that is,

$$g_t = \sum_{j=1}^L \alpha_{t,j} \hat{h}_j, \quad (3)$$

where $\alpha_t \in \mathbb{R}^L$ is a vector of the *attention weights*, also called *alignment factors* [3]. In the computation of *attention weights*, α_t is often evaluated by scoring each element in $\hat{\mathcal{H}}$ separately and normalizing the scores as follows:

$$\alpha_t = \text{Attend}(s_{t-1}, \hat{\mathcal{H}}), \quad (4)$$

where *Attend* describes the attending process [7].

Above, the *RNN* function in Eq. (2) represents an LSTM recurrent network. Note that the decoder is capable of generating sequences of variable lengths. Following [34], a special end-of-sequence (EOS) token is added to the target set, so that the decoder completes the generation of characters when EOS is emitted.

3.4. Network Training

We integrate the BCNN, AON, FG and attention decoder into one network, as shown in Fig. 3. Therefore, given an

input image \mathcal{I} , the loss function of the network is as follows:

$$\mathcal{L} = - \sum_t \ln P(\hat{y}_t | \mathcal{I}, \theta), \quad (5)$$

where \hat{y}_t is the ground truth of the t -th character and θ is a vector that combines all the network parameters.

3.5. Character Sequence Decoding

Decoding is the final process to generate the predicted characters. Following the decoding conventions, two processing modes are given: unconstrained (lexicon-free) mode and constrained mode. We execute unconstrained text recognition by directly selecting the most probable character. While in constrained text recognition, with respect to different types of lexicons (their sizes are denoted by “50”, “1k” and “full” respectively), we calculate the conditional probability distributions for all lexicon words, and take the one with the highest probability as the output result.

4. Technical Details of AON and FG

4.1. Arbitrary Orientation Network (AON)

We develop an *arbitrary orientation network* consisting of the *horizontal network* (HN), the *vertical network* (VN) and the *character placement clue network* (CN) for extracting horizontal, vertical and placement features respectively.

The HN encodes the foundational feature maps into a sequence of horizontal feature vectors $\mathcal{H} \in \mathbb{R}^{L \times D}$ by first performing downsampling on height directly by 5 shared convolutional blocks (described bellow) with the corresponding pooling strategy (shown in Fig. 3) to 1, and using the bidirectional LSTM to further encode the feature sequence, then generating the reversed feature sequence by conducting reverse operation (described in Eq. (6) and (7)), where L and D represent the length of $\hat{\mathcal{H}}$ and the channel number, respectively. Symmetrically, VN first rotates the square feature maps by 90 degrees, then generates the vertical feature vectors $\mathcal{V} \in \mathbb{R}^{L \times D}$ with the same procedure as HN. Here, *reversion* can accelerate training convergence, thus indirectly impacts the training of CN.

Since we describe each character sequence in four directions: *left* \rightarrow *right*, *right* \rightarrow *left*, *top* \rightarrow *bottom* and *bottom* \rightarrow *top*, \mathcal{H} and \mathcal{V} can be represented as follows:

$$\mathcal{H} = \begin{cases} \vec{\mathcal{H}} : (h_1, \dots, h_L)^T, & \text{left} \rightarrow \text{right} \\ \overleftarrow{\mathcal{H}} : (h_L, \dots, h_1)^T, & \text{right} \rightarrow \text{left} \end{cases} \quad (6)$$

$$\mathcal{V} = \begin{cases} \vec{\mathcal{V}} : (v_1, \dots, v_L)^T, & \text{top} \rightarrow \text{bottom} \\ \overleftarrow{\mathcal{V}} : (v_L, \dots, v_1)^T, & \text{bottom} \rightarrow \text{top} \end{cases} \quad (7)$$

For each text image, the CN outputs the corresponding character placement clues $\mathcal{C} \in \mathbb{R}^{4 \times L}$ as:

$$\mathcal{C} = (c_1, \dots, c_L)^T. \quad (8)$$

Here, for any $c_i \in \mathbb{R}^4$, we have $\sum_{j=1}^4 c_{ij} = 1$, where c_{ij} refers to the j -th direction’s weight. The extraction process of clues is depicted as the green blocks in Fig. 3.

In practice, we find that it is hard to train the HN and VN respectively and simultaneously. The state of each branch is easy to corrupted on orientation distribution unbalanced training datasets. Therefore, we design a shared convolution mechanism that performs the same convolutional filter operations for both horizontal and vertical process, and the shared convolution block is shown in Fig. 4. With the shared convolutional mechanism, the network is robust and easy to learn on orientation unbalanced training datasets.

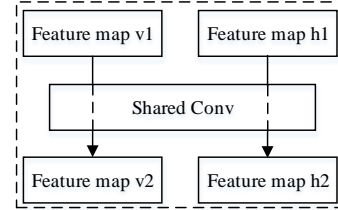


Figure 4. The shared convolution block in the dashed box provides a mechanism that multi-groups of feature maps share the same convolutional filters.

4.2. Filter Gate (FG)

With the captured four feature sequences and character placement clues, we design a filter gate to neglect the irrelevant features. Formally, given the i -th features $(\vec{\mathcal{H}}_i, \overleftarrow{\mathcal{H}}_i, \vec{\mathcal{V}}_i, \overleftarrow{\mathcal{V}}_i)$, we use the corresponding placement clue c_i to attend the appropriate features:

$$\hat{h}'_i = [\vec{\mathcal{H}}_i \ \overleftarrow{\mathcal{H}}_i \ \vec{\mathcal{V}}_i \ \overleftarrow{\mathcal{V}}_i] c_i. \quad (9)$$

Then an activation operation is performed as follows:

$$\hat{h}_i = \tanh(\hat{h}'_i). \quad (10)$$

Above, \hat{h}_i indicates the i -th element of $\hat{\mathcal{H}}$: $(\hat{h}_1, \dots, \hat{h}_L)$.

5. Performance Evaluation

We conduct extensive experiments to validate the proposed method on both irregular and regular recognition benchmarks. To be fair, we train the *HN* used in AON as the baseline model (denoted by Naive_base), which is similar to the previous works focusing on regular text recognition. We also combine *HN* with the TPS-based STN used in [32] as the STN-based control model (denoted by STN_base). All control experiments are conducted with similar training data and in similar running environment. We compare our model with not only the major existing methods (including the-state-of-the-art ones), but also the above two baseline models: Naive_base and STN_base. Furthermore, we explore the roles of HN, VN, CN and FG in AON.

5.1. Datasets

The regular and irregular benchmarks are as follows:

SVT-Perspective [28] contains 639 cropped images for testing. Images are picked from side-view angle snapshots in Google Street View, therefore one may observe severe perspective distortions. Each image is associated with a 50-word lexicon and a full lexicon.

CUTE80 (CT80 in short) [29] is collected for evaluating curved text recognition. It contains 288 cropped natural images for testing. No lexicon is associated.

ICDAR 2015 (IC15 in short) [21] contains 2077 cropped images where more than 200 irregular (arbitrarily-oriented, perspective or curved). No lexicon is associated.

IIIT5K-Words (IIIT5K in short) [26] is collected from the Internet, containing 3000 cropped word images in its test set. Each image specifies a 50-word lexicon and a 1k-word lexicon, both of which contain the ground truth words as well as other randomly picked words.

Street View Text (SVT in short) [35] is collected from the Google Street View, consists of 647 word images in its test set. Many images are severely corrupted by noise and blur, or have very low resolutions. Each image is associated with a 50-word lexicon.

ICDAR 2003 (IC03 in short) [24] contains 251 scene images, labeled with text bounding boxes. Each image is associated with a 50-word lexicon defined by Wang et al. [35]. For fair comparison, we discard images that contain non-alphanumeric characters or have less than three characters, following [35]. The resulting dataset contains 867 cropped images. The lexicons include the 50-word lexicons and the full lexicon that combines all lexicon words.

5.2. Implementation Details

Network details: The deep neural network has been detailed in Fig. 3. In our network, all images are resized to 100×100 . As for the convolutional strategy, all convolutional blocks have 3×3 size of kernels, 1×1 size of pads and 1×1 size of strides, and all pooling (max) blocks have 2×2 size of kernels. We adopt batch normalization (BN) [14] and ReLU activation right after each convolution. For the character generation task, the attention is designed with an LSTM (256 memory blocks) and 37 output units (26 letters, 10 digits, and 1 EOS symbol).

Implementation and Running Environment: We train our model on 8-million synthetic data released by Jaderberg et al. [15] and 4-million synthetic instances (excluding the images that contain non-alphanumeric characters) cropped from 80-thousand images [12] by the ADADELTA [42] optimization method. Meanwhile, we conduct data augmentation by randomly rotating each image range from 0° to 360° once. Our method is implemented under the Caffe framework [19]. The CUDA 8.0 and CUDNN v7 backend are extensively used in our implementation, so that

most modules in our method are GPU-accelerated. Our method can handle about 190/630 samples per second in the training/testing phase. The experiments are carried out on a workstation with one Intel Xeon(R) E5-2650 2.30GHz CPU, one NVIDIA Tesla P40 GPU, and 128GB RAM.

5.3. Performance on Irregular Datasets

Method	SVT-Perspective			CT80	IC15
	50	Full	None	None	None
ABBY[35]	40.5	26.1	—	—	—
Mishra <i>et al.</i> [11]	45.7	24.7	—	—	—
Wang <i>et al.</i> [37]	40.2	32.4	—	—	—
Phan <i>et al.</i> [28]	75.6	67.0	—	—	—
Shi <i>et al.</i> [31]	92.6	72.6	66.8	54.9	—
Shi <i>et al.</i> [32]	91.2	77.4	71.8	59.2	—
Yang <i>et al.</i> [39]	93.0	80.2	75.8	69.3	—
Cheng <i>et al.</i> [6]	92.6	81.6	71.5	63.9	66.2
Naive_base	92.4	83.3	70.5	75.4	67.8
STN_base	94.6	82.8	68.5	73.7	67.5
Ours	94.0	83.7	73.0	76.8	68.2

Table 1. Results on irregular benchmarks. “50” is lexicon size and “Full” indicates the combined lexicon of all images in the benchmarks. “None” means lexicon-free.

Recently, Cheng *et al.* [6] proposed FAN to improve text recognition performance, which must be trained with additional character-level bounding box annotations. Here, we also compare our method with FAN on the irregular datasets. Tab. 1 summarizes the recognition results on three irregular text datasets: SVT-Perspective, CUTE80 and ICDAR15. Comparing with the existing methods’ performance results released in the literature, we find that our method outperforms the existing methods on almost all benchmarks, except for SVT-Perspective with lexicon-free released by Yang *et al.* [39]. However, it is worthy of pointing out that Yang’s method [39] implicates its text-reading system with both word-level and character-level bounding box annotations, which is resource consuming, while our method can be easily carried out with only word-level annotations.



Figure 5. Some images rectified by TPS.

Tab. 1 also gives the results of the two baseline models. We can see that Naive_base does not recognize irregular texts well. Though theoretically TPS-based STN can handle any irregular texts, it seems not able to satisfactorily

rectify arbitrary-oriented or seriously curved texts in practice. Fig. 5 shows some rectified examples by TPS, their original images are shown in Fig. 1. We can see that except for the first and the third images, the other four images are not desirably rectified.

Method	IIIT5k			SVT		IC03		
	50	1k	None	50	None	50	Full	None
ABBY[35]	24.3	—	—	35.0	—	56.0	55.0	—
Wang <i>et al.</i> [35]	—	—	—	57.0	—	76.0	62.0	—
Mishra <i>et al.</i> [11]	64.1	57.5	—	73.2	—	81.8	67.8	—
Wang <i>et al.</i> [37]	—	—	—	70.0	—	90.0	84.0	—
Goel <i>et al.</i> [8]	—	—	—	77.3	—	89.7	—	—
Bissacco <i>et al.</i> [4]	—	—	—	90.4	78.0	—	—	—
Alsharif [2]	—	—	—	74.3	—	93.1	88.6	—
Almazán <i>et al.</i> [1]	91.2	82.1	—	89.2	—	—	—	—
Yao <i>et al.</i> [40]	80.2	69.3	—	75.9	—	88.5	80.3	—
Jaderberg <i>et al.</i> [16]	—	—	—	86.1	—	96.2	91.5	—
Su and Lu [33]	—	—	—	83.0	—	92.0	82.0	—
Gordo [9]	93.3	86.6	—	91.8	—	—	—	—
Jaderberg <i>et al.</i> [17]	97.1	92.7	—	95.4	80.7	98.7	98.6	93.1
Jaderberg <i>et al.</i> [16]	95.5	89.6	—	93.2	71.7	97.8	97.0	89.6
Shi <i>et al.</i> [31]	97.6	94.4	78.2	96.4	80.8	98.7	97.6	89.4
Shi <i>et al.</i> [32]	96.2	93.8	81.9	95.5	81.9	98.3	96.2	90.1
Lee <i>et al.</i> [22]	96.8	94.4	78.4	96.3	80.7	97.9	97.0	88.7
Yang <i>et al.</i> [39]	97.8	96.1	—	95.2	—	—	97.7	—
Cheng’s baseline [6]	98.9	96.8	83.7	95.7	82.2	98.5	96.7	91.5
Cheng <i>et al.</i> [6]	99.3	97.5	87.4	97.1	85.9	99.2	97.3	94.2
Naive_base	99.5	98.1	86.0	96.9	81.9	98.5	96.5	90.5
STN_base	99.5	97.8	85.9	96.3	80.7	98.5	96.2	89.2
Ours	99.6	98.1	87.0	96.0	82.8	98.5	97.1	91.5

Table 2. Results on regular benchmarks. “50” and “1k” are lexicon sizes. “Full” indicates the combined lexicon of all images in the benchmarks. “None” means lexicon-free.

5.4. Performance on Regular Datasets

AON is designed for recognizing both irregular and regular texts. Therefore, we test our method on some regular text benchmarks, the results are shown in Tab. 2. In the constrained cases, our method achieves comparable performance to the existing methods. In the unconstrained cases, our method only falls behind Cheng *et al.* [6] on the three benchmarks, and Jaderberg *et al.* [17] on IC03. For [6], two major factors lead to its high performance: a) using extra geometric annotations (location of each character) in training the attention decoder, and b) exploiting a ResNet-based feature extractor for obtaining robust feature representation. However, labelling the location of each character is extremely expensive, so it is not feasible for real applications. For fair comparison, we also gave the results of Cheng’s baseline (without the FocusNet branch) in Tab. 2, and found that our method outperforms Cheng’s baseline in most cases, which validates the superiority of our method. Though Jaderberg *et al.* [17] achieves an amazing results on IC03, their model cannot recognize out-of-vocabulary words, which limits its applicability in real world. Note that our model is trained without any character geometric information, and it performs better than the other existing

methods. As a whole, our method performs effectively in recognizing regular texts.

5.5. Deep insight into AON

Here, to further clarify the working mechanism of AON, we elaborate the roles of the major components *HN*, *VN*, *CN* and *FG* in AON, and show the placement trends of texts in some real images. These trends are generated by AON.

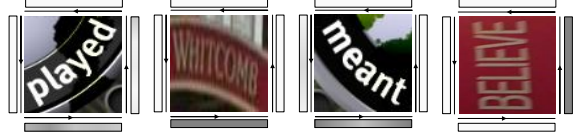


Figure 6. Illustration of learned character placement clues by AON. Each image is surrounded with four changing-gray bars with arrows of different directions. Deeper gray in the bars indicates larger weight for the corresponding directional feature.

The roles of *HN*, *VN*, *CN* and *FG* in AON. We use both horizontal sequence of features and vertical sequence of features to represent arbitrarily-oriented texts. Concretely, for horizontal/vertical texts, horizontal/vertical features are enough to represent the texts; For perspective/slant or arbitrarily-oriented text, we generate the final feature sequence by combining horizontal and vertical features. *HN* and *VN* are responsible for generating horizontal and vertical features respectively. *CN* plays an important role in learning the weights (*i.e.*, character placement clues) that are used to guide the generation of final feature sequences. *FG* is just to perform the weight-sum operation with the horizontal/vertical feature sequence and the learnt placement weights. Fig. 6 shows some examples of generated placement clues. We can see that the generated clues conform to our visual observations in the images, which validates the effectiveness of *CN* in AON.

Text placement trends generated with AON. Here we verify that the learned character placement clues can be used to generate placement trends of character sequences by positioning each character and drawing text orientations in the original images. Below is the computation process of text placement trends.

We know that the alignment factors α_t produced by the attention module indicate the probability distributions over the input sequence of features for generating the glimpse vector g_t . And the four character placement clues $\mathcal{C} = [c_1, c_2, c_3, c_4]$ imply the importance of four extracted feature sequences for representing characters. With \mathcal{C} and α_t , we roughly divide the input image into $L \times L$ patches and calculate the character position distribution dis by $dis = \mathcal{C} \odot \alpha_t$, where $dis = (d_1, d_2, d_3, d_4) \in \mathbb{R}^{4 \times L}$. We further normalize each element by $norm(d_{ij}) = \frac{d_{ij}}{\sum_{i=1}^2 \sum_{j=1}^L d_{ij}}$ for $i \in (1, 2)$, and by $norm(d_{ij}) = \frac{d_{ij}}{\sum_{i=3}^4 \sum_{j=1}^L d_{ij}}$ for $i \in (3, 4)$.

Here, $norm$ indicates the normalization operation.

For a character at position (x, y) , we first compute the horizontal coordinate x with $[d_1, d_2]$ by $x = \sum_{j=1}^L \sum_{i=1}^2 j \times norm(d_{ij})$, where $i \in (1, 2)$ and $j \in (1, 2, \dots, L)$. Similarity, we compute the vertical coordinate y with $[d_3, d_4]$ by $y = \sum_{j=1}^L \sum_{i=3}^4 j \times norm(d_{ij})$.

To visualize the placement trends of texts in the input images, we mark the coordinate (x, y) on each input image as the corresponding character's position, and consecutively connect the last character's position and the current character's position with an arrow to describe the text's placement trend. Fig. 7 shows some examples of generated text placement trends of real images. We can see that the trends formed by the connected arrows basically conform to our visual observations, which again shows that our method is effective in estimating the orientations of texts in images.



Figure 7. The visualization of generated placement trends for perspective, curved and oriented images, shown in the 1st, 2nd and 3rd row, respectively. The curves formed by connected red arrows indicate text placement trends. All texts in the images are correctly recognized by our method.

6. Discussions

The necessity of CN in AON. As the attention-based decoder is able to select features for generating characters. It is natural to suspect whether CN is necessary. To answer this, we have two experiments without CN : 1) Concatenating horizontal and vertical feature sequence along the channel axis. We find the model converges slowly and cannot achieve state-of-the-art performance, because three quarters' information in the final feature sequence is superfluous. 2) Concatenating horizontal and vertical along the temporal axis. We get results of averagely about 4% lower than that of AON on all benchmarks. The above experiments show that CN is important in AON.

Impact of aspect ratio. We studied the impact of aspect ratio by experiments, but did not observed obvious negative impact for images with a large aspect ratio. Compared to the previous works [31, 32], the enlarging/shrinking operation in height does not obviously affect recognition results of horizontal texts with a large aspect ratio.

Integrating with only two directional feature sequence. It is not reasonable to integrate only two directional sequences of features. For example, as shown in Fig. 1, if we integrate the right-left and down-top directional sequences of features to generate the final feature sequence, the visual features of 'p' and 'd' will be frame-wisely mixed up due to the weighting mechanism of FG.

The computational cost of AON. Computational cost and the number of parameters are major concerns in resource-constrained scenarios such as embedded computer systems. Comparing to the *Naive_base* model, the introducing of AON increases parameters and computational cost (twice of *Naive_base*). However, the *STN_base* needs triple parameters and computational cost of *Naive_base*.

7. Conclusion

In this work, we propose a novel method to recognize arbitrarily oriented texts by 1) devising an arbitrary orientation network to extract visual features of characters in four directions and the character placement clues, 2) using a filter gate mechanism to combine the four-direction sequences of features, and 3) employing an attention-based decoder for generating character sequence. Different from most existing methods, our method can effectively recognize both irregular and regular texts from images. Experiments over both regular and irregular benchmarks validate the superiority of the proposed method. In the future, we plan to extend the proposed idea to other related tasks.

8. Acknowledgement

Fan Bai and Shuigeng Zhou were partially supported by the Program of Science and Technology Innovation Action of Science and Technology Commission of Shanghai Municipality (STCSM) under grant No. 17511105204.

References

- [1] J. Almazán, A. Gordo, A. Fornés, and E. Valveny. Word Spotting and Recognition with Embedded Attributes. *IEEE TPAMI*, 36(12):2552–2566, 2014. 7
- [2] O. Alsharif and J. Pineau. End-to-end text recognition with hybrid hmm maxout models. In *ICLR*, 2014. 7
- [3] D. Bahdanau, K. Cho, and Y. Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. In *ICLR*, 2015. 4
- [4] A. Bissacco, M. Cummins, Y. Netzer, and H. Neven. PhotoOCR: Reading Text in Uncontrolled Conditions. In *ICCV*, pages 785–792, 2013. 3, 7
- [5] F. L. Bookstein. Principal warps: Thin-plate splines and the decomposition of deformations. *IEEE TPAMI*, 11(6):567–585, 1989. 2, 4
- [6] Z. Cheng, F. Bai, Y. Xu, G. Zheng, S. Pu, and S. Zhou. Focusing Attention: Towards Accurate Text Recognition in Natural Images. In *ICCV*, pages 5076–5084, 2017. 1, 3, 6, 7

- [7] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio. Attention-Based Models for Speech Recognition. In *NIPS*, pages 577–585, 2015. 4
- [8] V. Goel, A. Mishra, K. Alahari, and C. V. Jawahar. Whole is Greater than Sum of Parts: Recognizing Scene Text Words. In *ICDAR*, pages 398–402, 2013. 7
- [9] A. Gordo. Supervised mid-level features for word image representation. In *CVPR*, pages 2956–2964, 2015. 7
- [10] A. Graves, S. Fernández, F. Gomez, and J. Schmidhuber. Connectionist Temporal Classification : Labelling Unsegmented Sequence Data with Recurrent Neural Networks. In *ICML*, pages 369–376. ACM, 2006. 3
- [11] A. Graves, A. r. Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In *ICASSP*, pages 6645–6649, 2013. 6, 7
- [12] A. Gupta, A. Vedaldi, and A. Zisserman. Synthetic Data for Text Localisation in Natural Images. In *CVPR*, pages 2315–2324, 2016. 6
- [13] P. He, W. Huang, Y. Qiao, C. C. Loy, and X. Tang. Reading Scene Text in Deep Convolutional Sequences. In *AAAI*, pages 3501–3508, 2016. 3
- [14] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456, 2015. 6
- [15] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Synthetic Data and Artificial Neural Networks for Natural Scene Text Recognition. *arXiv preprint arXiv:1406.2227*, 2014. 1, 6
- [16] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Deep Structured Output Learning for Unconstrained Text Recognition. In *ICLR*, 2015. 3, 7
- [17] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Reading Text in the Wild with Convolutional Neural Networks. *IJCV*, 116(1):1–20, 2016. 3, 7
- [18] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial Transformer Networks. *NIPS*, pages 2017–2025, 2015. 1, 3
- [19] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional Architecture for Fast Feature Embedding. In *ACM-MM*, pages 675–678, 2014. 6
- [20] 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. 4
- [21] D. Karatzas, L. Gomez-Bigorda, A. Nicolaou, S. Ghosh, A. Bagdanov, M. Iwamura, J. Matas, L. Neumann, V. R. Chandrasekhar, S. Lu, F. Shafait, S. Uchida, and E. Valveny. ICDAR 2015 competition on Robust Reading. In *ICDAR*, pages 1156–1160, 2015. 6
- [22] C. Y. Lee and S. Osindero. Recursive Recurrent Nets with Attention Modeling for OCR in the Wild. In *CVPR*, pages 2231–2239, 2016. 1, 3, 7
- [23] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015. 1, 4
- [24] S. M. Lucas, A. Panaretos, L. Sosa, A. Tang, S. Wong, and R. Young. ICDAR 2003 robust reading competitions. In *ICDAR*, pages 682–687, 2003. 6
- [25] J. Ma, W. Shao, H. Ye, L. Wang, H. Wang, Y. Zheng, and X. Xue. Arbitrary-Oriented Scene Text Detection via Rotation Proposals. *arXiv preprint arXiv:1703.01086*, 2017. 4
- [26] A. Mishra, K. Alahari, and C. V. Jawahar. Scene Text Recognition using Higher Order Language Priors. In *BMVC*, pages 1–11, 2012. 6
- [27] L. Neumann and J. Matas. Real-time scene text localization and recognition. In *CVPR*, pages 3538–3545, 2012. 1, 3
- [28] T. Quy Phan, P. Shivakumara, S. Tian, and C. Lim Tan. Recognizing text with perspective distortion in natural scenes. In *ICCV*, pages 569–576, 2013. 6
- [29] A. Risnumawan, P. Shivakumara, C. S. Chan, and C. L. Tan. A robust arbitrary text detection system for natural scene images. *Expert Systems with Applications*, 41(18):8027–8048, 2014. 6
- [30] B. Shi, X. Bai, and S. Belongie. Detecting Oriented Text in Natural Images by Linking Segments. *arXiv preprint arXiv:1703.06520*, 2017. 4
- [31] B. Shi, X. Bai, and C. Yao. An End-to-End Trainable Neural Network for Image-based Sequence Recognition and Its Application to Scene Text Recognition. *IEEE TPAMI*, preprint, 2016. 1, 3, 6, 7, 8
- [32] B. Shi, X. Wang, P. Lyu, C. Yao, and X. Bai. Robust Scene Text Recognition with Automatic Rectification. In *CVPR*, pages 4168–4176, 2016. 1, 2, 3, 4, 5, 6, 7, 8
- [33] B. Su and S. Lu. Accurate Scene Text Recognition Based on Recurrent Neural Network. In *ACCV*, pages 35–48, 2015. 7
- [34] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to Sequence Learning with Neural Networks. In *NIPS*, pages 3104–3112, 2014. 3, 4
- [35] K. Wang, B. Babenko, and S. Belongie. End-to-end scene text recognition. In *ICCV*, pages 1457–1464, 2011. 3, 6, 7
- [36] K. Wang and S. Belongie. Word Spotting in the Wild. In *ECCV*, pages 591–604. Springer, 2010. 3
- [37] T. Wang, D. J. Wu, A. Coates, and A. Y. Ng. End-to-end text recognition with convolutional neural networks. In *ICPR*, pages 3304–3308, 2012. 3, 6, 7
- [38] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *ICML*, pages 2048–2057, 2015. 2
- [39] X. Yang, D. He, Z. Zhou, D. Kifer, and C. L. Giles. Learning to Read Irregular Text with Attention Mechanisms. In *IJCAI*, pages 3280–3286, 2017. 1, 2, 4, 6, 7
- [40] C. Yao, X. Bai, B. Shi, and W. Liu. Strokelets: A Learned Multi-scale Representation for Scene Text Recognition. In *CVPR*, pages 4042–4049, 2014. 3, 7
- [41] Q. Ye and D. Doermann. Text Detection and Recognition in Imagery: A Survey. *IEEE TPAMI*, 37(7):1480–1500, 2015. 3
- [42] M. D. Zeiler. ADADELTA: An Adaptive Learning Rate Method. *CoRR*, abs/1212.5701, 2012. 6
- [43] 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. 4