

Pyramid Mask Text Detector

Jingchao Liu^{1,2*}, Xuebo Liu^{1*}, Jie Sheng¹, Ding Liang¹, Xin Li³, Qingjie Liu²
¹SenseTime

²Beihang University

³The Chinese University of Hong Kong

liu.siqi@buaa.edu.cn, liuxuebo@sensetime.com, jie.sheng0112@gmail.com,
 liangding@sensetime.com, lixin@se.cuhk.edu.hk, qingjie.liu@buaa.edu.cn

Abstract

Scene text detection, an essential step of scene text recognition system, is to locate text instances in natural scene images automatically. Some recent attempts benefiting from Mask R-CNN formulate scene text detection task as an instance segmentation problem and achieve remarkable performance. In this paper, we present a new Mask R-CNN based framework named Pyramid Mask Text Detector (PMTD) to handle the scene text detection. Instead of binary text mask generated by the existing Mask R-CNN based methods, our PMTD performs pixel-level regression under the guidance of location-aware supervision, yielding a more informative soft text mask for each text instance. As for the generation of text boxes, PMTD reinterprets the obtained 2D soft mask into 3D space and introduces a novel plane clustering algorithm to derive the optimal text box on the basis of 3D shape. Experiments on standard datasets demonstrate that the proposed PMTD brings consistent and noticeable gain and clearly outperforms state-of-the-art methods. Specifically, it achieves an *F-measure* of 80.13% on ICDAR 2017 MLT dataset.

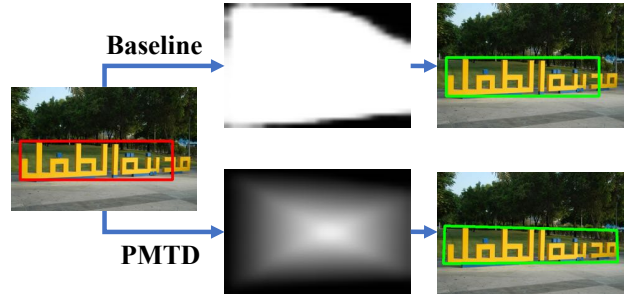
1. Introduction

Scene text detection has attracted growing research interests in the computer vision community, due to its numerous practical applications in scene understanding, license plate recognition, autonomous driving, and document analysis, etc. Recently, many works [41, 30, 14] view scene text detection as an instance segmentation problem, and several Mask R-CNN [6] based methods were proposed and achieved remarkable performance. However, there are several drawbacks in these works:

Over-simplified supervision: The common observation that most text areas in natural scenes are quadrilateral is supposed to be useful for text detection. However, the Mask R-CNN based methods, aiming for differentiating the text region from the background region rather than generating a text mask of a specific shape, ignore the consideration of



(a) Examples with imprecise segmentation labels. The area within green box denotes the manually annotated text instance. Many background pixels not belonging to the text instance are mislabeled as the foreground pixels, especially at the border of the text box, which may hurt the performance of the Mask R-CNN based methods.



(b) The red box is the predicted bounding box and the green box refers to the predicted text box. The existing Mask R-CNN based methods suffer from the errors of bounding box detection while PMTD can regress more accurate text box with the help of the informative soft text mask.

Figure 1: Imprecise segmentation labels and imprecise bounding box are detrimental to previous methods.

such kind of information and therefore can not take advantage of the given label.

Imprecise segmentation labels: Converting the quadrilateral text area into a pixel-level binary supervision signals for semantic segmentation enables directly applying Mask R-CNN to scene text detection. However, the quality of the generated pixel-level labels is unsatisfactory. As shown in Fig.1(a), many background pixels not belonging to the text region are incorrectly regarded as the foreground pixels. Trained on noisy data, the semantic segmentation based text detector is prone to generate mistakes.

Error propagation: The Mask R-CNN based methods firstly predict text bounding boxes and then perform semantic segmentation within the bounding box. Such strategy is usually reasonable for simple scenes but rather fragile when the predicted bounding box fails to cover the whole text re-

*indicates equal contribution.

gion. The reason is because determining the text box with only the text region inside the bounding box tends to exclude the outside part (See Fig. 1(b)). In other words, the errors from the object detection may be propagated to the process of finding text box, leading to performance degradation of scene text detection. We also observe that the effect of the error propagation will be amplified with the increasing of the IoU threshold for true positive text instances (quantitative results and qualitative analysis are detailed in Sec. 4.3).

In this paper, we propose the Pyramid Mask Text Detector (PMTD) to address the above problems. As depicted in Fig. 2, instead of pixel-level binary classification as done in the existing Mask R-CNN based methods, we propose to perform “soft” semantic segmentation between the text region and the background region. Explicitly, we assign a soft pyramid label (i.e., a real value between 0 and 1) for each pixel within text instance. The value of the soft pyramid label is determined by the distance to the boundary of the text box, which implicitly encoding the shape and location information into the training data. By fitting such soft text mask, the quadrilateral property of the text instance is naturally considered during training. Besides, introducing the location-aware segmentation labels reduces the impact of mislabeled pixels near the boundary of the text box.

During the test phase, with the extended z -axis characterizing the value of the pixel-level segmentation output, we reinterpret the 2D predicted text mask into a set of 3D points. A plane clustering algorithm is proposed to regress the optimal pyramid from these 3D points. Specifically, launched with four initialized supporting planes of a pyramid, the plane clustering algorithm iteratively groups the nearest points for each supporting planes, and then updates the supporting planes by the clustered points. After the iterations, an accurate bounding pyramid is obtained and its bottom face is regarded as the output text box. Since it is not the boundary pixels but the supporting plane that gets involved in finding the text box, the error propagation issue can be alleviated, and more accurate text box can be obtained.

Our pipeline is shown in Fig.3. As there exist differences between text detection datasets and object detection datasets, such as the different distribution of aspect ratios and scales, we tailor-make a Mask R-CNN based baseline for text detection, which outperforms all previous methods on ICDAR 2017 MLT dataset. Furthermore, the proposed PMTD raises the F-measure to 80.13%. The main contributions of this paper are three-fold:

- We propose the Pyramid Mask Text Detector for scene text detection, and extensive experiments demonstrate its state-of-the-art performance on several benchmark datasets.
- We propose to perform “soft” segmentation between

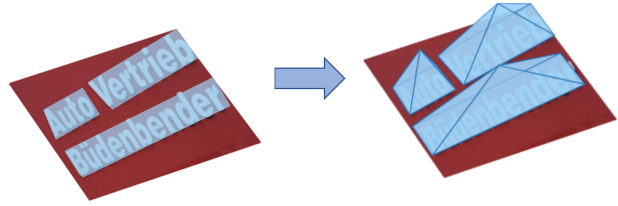


Figure 2: Previous methods aim to find $\{0, 1\}$ label for each pixel while PMTD assigns a soft pyramid label of the value $\in [0, 1]$.

text region and non-text region, incorporating the shape and location information into the model training and alleviating the inaccuracy labeling for the instance boundary.

- We introduce a novel plane clustering algorithm to find better text box with the 3D coordinate, which predicts more accurate text box and improves the robustness to imprecise bounding box predictions.

2. Related Work

Scene text detection has received significant attention over the past few years, and numerous deep learning based methods [5, 45, 22, 32, 46, 10, 20, 29, 41] have been reported in the literature. Comprehensive reviews and detailed analyses can be found in survey papers [47, 40, 43].

Earlier text detection works including [13, 16, 15] are among the first deep neural network based methods. They usually consist of multiple stages, such as candidate aggregation, word partition and false positive removal by post-processing filtering. Huang *et al.* [13] first apply the MSERs operator on the input image to generate some text candidates, then use a CNN classifier to generate a confidence map which was later used for constructing text-lines. Jaderberg *et al.* [16] train a strongly supervised character classifier to generate text saliency map, then combines bounding boxes at multiple scales and undergoes filtering and non-maximal suppression. In a later work [15], they leverage a CNN for bounding box regression and a random forest classifier for reducing the number of false-positive detections.

Recent works [5, 39, 23, 45] regard text words or lines as objects and adapt the pipeline of general object detection, e.g., Faster R-CNN [36], SSD [25] and YOLO [35] into text detection. They regress the offsets from a proposal region or a single pixel in the feature map to a horizontal rectangle and obtain good performance with well-designed modifications on horizontal text detection. Gupta *et al.* [5] improves over the YOLO network and Fully Convolutional Networks (FCN) [28] for text prediction densely, while further adopts the filter and regression steps for removing the false positives. TextBoxes [23] modifies SSD by using irregular convolutional kernels and long default anchors according to the characteristic of scene text. Built on top of Faster R-CNN, CTPN [39] develops a vertical anchor mechanism that pre-

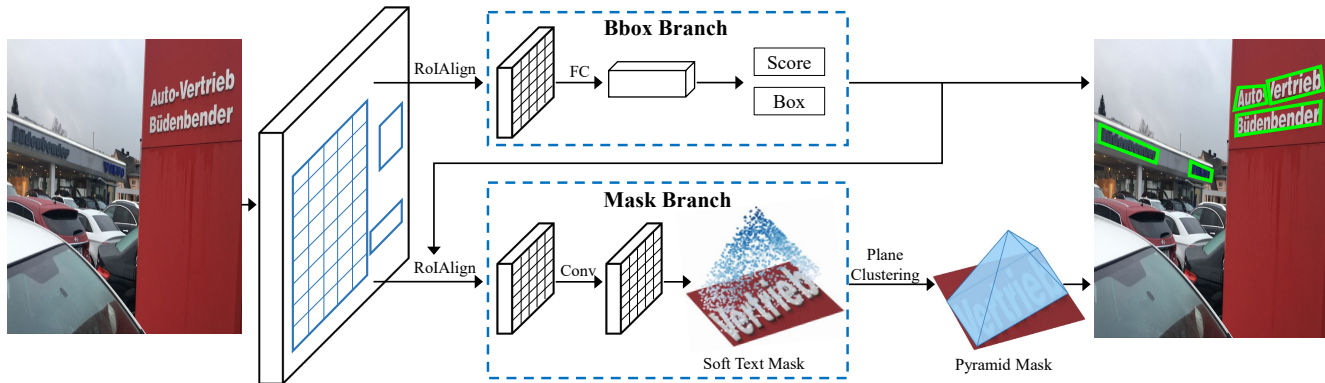


Figure 3: Overall architecture of PMTD.

dicts location and text/non-text score of each fixed-width proposal simultaneously, then connects the sequential proposals by a recurrent neural network. FEN [45] improves text recall with a feature enhancement RPN and hyper feature generation for text detection refinement.

Considering that scene texts are with arbitrary orientations, works in [22, 32, 46, 10, 26, 9] make the above methods possible for multi-oriented text detection. RRPN [32] introduces inclined anchors with angle information for arbitrary-oriented text prediction and rotated RoI pooling layer to project arbitrary-oriented proposals to the feature map for a text region classifier. TextBoxes++ [22] improves TextBoxes by regressing horizontal anchors to more general quadrilaterals enclosing oriented texts. It also proposes an efficient cascaded non-maximum suppression for quadrilaterals or rotated rectangles. With dense predictions and one step post processing, EAST [46] and DDR [10] both directly produce the rotated boxes or quadrangles of text at each point in the text region. Recent text spotting methods like FOTS [26] and He *et al.* [9] show that training text detection and recognition simultaneously could greatly boost detection performance.

Except for the above regression-based methods, [3, 20, 29, 41] cast text detection as a segmentation problem. PixelLink [3] first segments out text regions by linking pixels within the same instance, then extracts text bounding boxes directly from the segmentation without location regression. TextSnake [29] employs an FCN [28] model to estimate the geometry attributes of text instances and uses a striding algorithm to extract the central axis point lists and finally reconstruct the text instances. Segmentation based methods tend to link adjacent text regions together incorrectly. To address this problem, PSENet [20] finds text kernels with various scales and proposes a progressive scale expansion algorithm to separate text instances standing close to each other accurately. SPCNET [41] views text detection as an instance segmentation problem, based on Mask R-CNN, it

proposes a text context module and a re-score mechanism to suppress false positives.

Although Curved text detection [1, 44] has attracted growing research interests recently, quadrilateral text detection is still a fundamental and challenging problem to be solved. PMTD is designed specially for quadrilateral text detection and significantly improves the state-of-the-art result from 74.3% [14] to 80.13% on ICDAR 2017 MLT.

3. Methodology

In this section, we firstly introduce a strong baseline. Then the soft pyramid label is proposed, which encodes the shape and location information into the training data. Finally, a new boundary regression algorithm, namely, plane clustering, is introduced to find the most fitting pyramid of the predicted soft text mask.

3.1. Our Baseline

Our baseline is based on Mask R-CNN with ResNet50 backbone [7]. In the training stage, we treat the axis-aligned bounding rectangle of the text region as the ground-truth bounding box and assign pixels inside text boundary to positive segmentation label. In the test stage, we firstly find all the connected areas in the predicted mask, then select the one with the maximum area, and finally obtain the output text box by finding the minimum bounding rectangle of this connected area.

We design a strong baseline by making the following three modifications:

Data augmentation: To enhance the generalization ability to various scales and aspect ratios, we apply data augmentations to enlarge scene text datasets:

1. Random horizon flip with a probability of 0.5.
2. Random resize the height and width of images to 640-2560 individually, without keeping the original aspect ratio.

3. Random select one 640×640 crop region from the resized image.

RPN Anchor: When adopting the FPN module, we can quantify the anchor by three parameters: the base scale of anchors, the feature maps where anchors searched, and the aspect ratios of anchors.

First of all, based on statistics of the data-augmented ground truth bounding box's height and width, we set the base scale of the anchor to 4×4 among all the four feature maps $\{1/4, 1/8, 1/16, 1/32\}$ uniformly.

For the anchor's aspect ratio, we calculate out five dedicated aspect ratios: $\{0.17, 0.44, 1.13, 2.90, 7.46\}$. The detail of generating aspect ratios is as follows: first, analyze the data-augmented ground truth bounding box's aspect ratio, then get the 5% quantile 0.17 and 95% quantile 7.46, finally insert three values in equal proportion between the 5% and 95% quantiles to form the final aspect ratio list.

OHEM: In the bounding box branch, we adopt the OHEM [38] to learn the hard samples. In our settings, we first sort the samples provided by RPN in the descending order of the sum of classification loss and location loss, then select the top 512 difficult samples to update the network.

3.2. Motivation

Although our baseline achieves remarkable performance, it still has the same drawbacks as other Mask R-CNN based methods, as mentioned in Sec. 1:

- These methods are not considering the common observation that most text areas in natural scenes are quadrilateral. They break down the quadrilateral structure into a pixel-wise classification problem which loses the shape information of the mask.
- Converting the quadrilateral text areas into pixel-level supervision is imprecise. Many background pixels not belonging to the text region are incorrectly regarded as the foreground pixels, as shown in Fig. 1(a). The mislabeled boundary pixels may cause an unexpectedly misjudged loss.
- Mask R-CNN based methods firstly predict bounding boxes and then predict text mask for every bounding box. The imprecisely predicted bounding box limits the mask branch to generate accurate text mask. In other words, the errors from the object detection will be propagated to the following steps, as shown in Fig. 1(b).

These problems motivate us to build pyramid mask text detector (PMTD), a new pipeline for scene text detection. The PMTD predicts a soft text mask for each text region and apply plane clustering algorithm to convert the predicted soft mask to the pyramid mask.

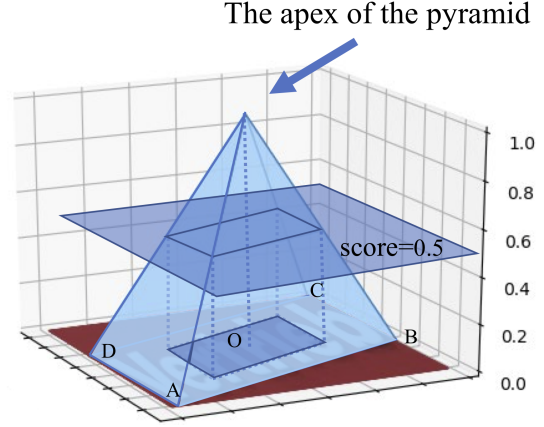


Figure 4: Generation of soft pyramid label. For a pixel in the text area, its label is the height of the pyramid.

3.3. Pyramid Label

We refine the mask's hard label of the class $\in \{0, 1\}$ to the soft pyramid label of the score $\in [0, 1]$ so that the PMTD can capture the shape and location information from the data. Specifically, we assign the center of text region as the apex of the pyramid with an ideal value $score = 1$ and the boundary of text region as the bottom edge of the side of the pyramid. We use the linear interpolation to fill each triangle side of the pyramid, as illustrated in Fig. 4.

Formally, given the four corner points $A(x_a, y_a)$, $B(x_b, y_b)$, $C(x_c, y_c)$, $D(x_d, y_d)$ of a quadrilateral, the value $score_p$ for the point $P(x_p, y_p)$ can be calculated as follows. First, the center of text region $O(x_o, y_o)$ can be obtained by:

$$x_o = (x_a + x_b + x_c + x_d)/4 \quad (1)$$

$$y_o = (y_a + y_b + y_c + y_d)/4 \quad (2)$$

For every region R_{OMN} (region between two rays OM and ON) from R_{OAB} , R_{OBC} , R_{OCD} , R_{ODA} , the \vec{OP} can be decomposed uniquely:

$$\vec{OP} = \alpha \vec{OM} + \beta \vec{ON} \quad (3)$$

$$\begin{bmatrix} x_p - x_o \\ y_p - y_o \end{bmatrix} = \begin{bmatrix} x_m - x_o & x_n - x_o \\ y_m - y_o & y_n - y_o \end{bmatrix} \begin{bmatrix} \alpha \\ \beta \end{bmatrix} \quad (4)$$

Then, α and β can be obtained by

$$\begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} x_m - x_o & x_n - x_o \\ y_m - y_o & y_n - y_o \end{bmatrix}^{-1} \begin{bmatrix} x_p - x_o \\ y_p - y_o \end{bmatrix} \quad (5)$$

The region R which P belongs to needs to satisfy the following condition:

$$\alpha \geq 0 \text{ and } \beta \geq 0 \quad (6)$$

Then the $score_p$ can be calculated by:

$$score_p = \max(1 - (\alpha + \beta), 0) \quad (7)$$

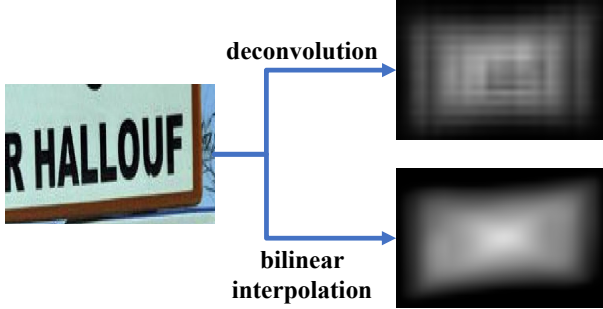


Figure 5: Deconvolution causes checkerboard pattern in our experiment, so we use bilinear interpolation for upsample to get a more accurate mask.

During the training stage, such supervision is reasonable. If one pixel locates near the center of the instance, its receptive field will be filled with positive pixels and deserves a higher score consequently. While the receptive field of the pixels near the boundary will contain much background context, and the scores of these pixels should be close to 0. In this respect, a larger receptive field is vital for PMTD to attain more precise results. So in the mask head, we replace the first four convolution layers to dilated convolution with stride 2 to enlarge receptive field.

Moreover, as mentioned in [34], deconvolution may cause the checkerboard pattern, which is harmful to the pixelwise regression, as illustrated in Fig. 5. To avoid this, we replace the deconvolution layer in the mask head to bilinear interpolation and a followed convolution layer.

We employ pixelwise L_1 loss to optimize the predicted text mask. Following the design in Mask R-CNN, the loss function of the whole network is as follows:

$$L = L_{\text{rpn}} + \lambda_1 L_{\text{cls}} + \lambda_2 L_{\text{box}} + \lambda_3 L_{\text{pyramid_mask}} \quad (8)$$

λ_1 , λ_2 and λ_3 are set to 1, 1, 5 respectively in our experiments.

Training with this new style label alleviates the pixel mislabeling problem. Taking a background pixel near the boundary as an example. Although it is mistakenly regarded as the foreground, its ground truth in our methods is still close to 0, while in previous Mask R-CNN based methods, this pixel is labeled as 1.

3.4. Plane Clustering

In this section, we will illustrate the plane clustering algorithm in details, which is an iteratively updated clustering algorithm for regressing the most fitting text box from the predicted soft text mask.

As a reverse process of generating a pyramid label from text region, we will first construct the pyramid from the text mask, then take the bottom edge of the pyramid as the output text box. Hence, the critical point is to parameterize and rebuild the pyramid.

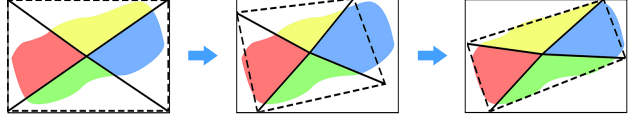


Figure 6: Illustration of the plane clustering algorithm. Every positive point in the predicted soft mask is assigned to one of four different colors, which indicates the supporting plane the point belongs to. The dashed lines are the intersection of supporting planes and the bottom plane, which form the predicted text box together. The supporting planes are refined from left to right.

Formally, the pyramid is composed of four supporting planes and one base plane. In the context of pyramid mask, we can convert the predicted soft mask into a point set of (x, y, z) , in which the (x, y) denotes the location and z stands for the predicted score of this pixel. The base plane is formulated as the plane $z = 0$, and each supporting plane can be uniquely determined by the equation $Ax + By + Cz + D = 0, C = 1$. Consequently, the task of the plane clustering algorithm is reduced to find the optimal parameter $\{A, B, D\}$ for each supporting plane, see Algorithm 1 for details.

Algorithm 1 Plane Clustering

input:

point : location = (x, y) , predicted_score = z
points = [*point_num* = $H * W$, (x, y, z)]

output:

plane : $Ax + By + Cz + D = 0, C = 1$
planes = [*plane_num* = 4, (A, B, D)]

function PLANE CLUSTERING(*points*)

$P \leftarrow \text{SELECTPOSITIVE}(\text{points})$

$\text{apex}.(x, y) \leftarrow \text{MEAN}(P).(x, y)$

$\text{apex}.z \leftarrow 1$

$\text{planes} \leftarrow \text{INITPLANES}(\text{apex})$

while *iter* < *max_iter* and REJECT(*residuals*) **do**

$G \leftarrow \emptyset \times \text{plane_num}$

for $p \in P$ **do**

$\text{plane}_i \leftarrow \text{NEARESTPLANE}(p, \text{planes})$

$G[\text{plane}_i] \leftarrow G[\text{plane}_i] \cup \{p\}$

end for

$\text{planes}, \text{residuals} \leftarrow \text{RLS}(G)$

end while

return *planes*

end function

In the initialization stage, the positive points set P is built by the condition $z > 0.1$. Then the apex of the initial pyramid is assigned as the center of P , with an ideal score $z = 1$. The four vertexes of the pyramid in the bottom face are initialized as four corner points of the predicted text bounding

Method	Precision	Recall	F-measure
FOTS [26]	80.95	57.51	67.25
FOTS* [26]	81.86	62.30	70.75
Lyu <i>et al.</i> [31]	83.80	55.60	66.80
Lyu <i>et al.</i> * [31]	74.30	70.60	72.40
PSENet [20]	77.01	68.40	72.45
Pixel-Anchor [21]	79.54	59.54	68.10
Pixel-Anchor* [21]	83.90	65.80	73.76
SPCNET [41]	66.90	73.40	70.00
SPCNET* [41]	68.60	80.60	74.10
Huang <i>et al.</i> [14]	80.00	69.80	74.30
Baseline	84.72	70.37	76.88
PMTD	85.15	72.77	78.48
PMTD*	84.42	76.25	80.13

Table 1: Comparison with other results on ICDAR 2017 MLT. * means multi scale testing.

box, shown in the left image in Fig.6.

After initializing the pyramid, an iterative updating scheme is implemented for clustering points, which is shown in Fig.6. In the assignment step, we partition each point to the nearest plane, and in the update step, we employ the robust least square algorithm (RLS) [11] to regress four supporting planes from the clustered points respectively, which is robust to the noise in the predicted text mask.

When the iteration reaches the max iteration or the regression residuals returned by RLS is small enough, the final quadrangular pyramid is obtained. Then the text box can be calculated out by the intersection of four supporting planes and the plane $z = 0$. In our experiment, the max iteration and the residual threshold is assigned to 10 and $1e-4$ respectively.

Thanks to the more informative soft text mask, the plane clustering algorithm takes advantage of the whole soft mask’s information to regress the most fitting pyramid. As the final text box is obtained from the supporting plane rather than the boundary pixels, PMTD is robust to imprecise bounding boxes, and naturally regress more accurate text boundary, detailed in Sec. 4.3.

4. Experiments

In this section, We evaluate our approach on ICDAR 2017 MLT [33], ICDAR 2015 [18] and ICDAR 2013 [19]. Experiment results demonstrate that the proposed PMTD brings consistent and noticeable gain, and clearly outperforms the state-of-the-art methods. Furthermore, the ablation study shows PMTD is robust to imprecise bounding box predictions and predicts more accurate text boxes.

4.1. Datasets

ICDAR 2017 MLT is a multi-oriented, multi-scripting, and multi-lingual scene text dataset. It consists of 7200 training images, 1800 validation images, and 9000 test im-

Method	Precision	Recall	F-measure
SegLink [37]	73.10	76.80	75.00
SSTD [8]	80.00	73.00	77.00
WordSup [12]	79.33	77.03	78.16
EAST* [46]	83.27	78.33	80.72
R2CNN [17]	85.62	79.68	82.54
DDR [10]	82.00	80.00	81.00
Lyu <i>et al.</i> * [31]	89.50	79.70	84.30
RRD* [24]	88.00	80.00	83.80
TextBoxes++* [22]	87.80	78.50	82.90
PixelLink [3]	85.50	82.00	83.70
FOTS [26]	91.00	85.17	87.99
IncepText* [42]	89.40	84.30	86.80
TextSnake [29]	84.90	80.40	82.60
FTSN [2]	88.60	80.00	84.10
SPCNET [41]	88.70	85.80	87.20
PSENet [20]	89.30	85.22	87.21
Baseline	85.84	90.55	88.14
PMTD	91.30	87.43	89.33

Table 2: Comparison with other results on ICDAR 2015. * means multi scale testing. For PMTD, we only report single scale testing result.

ages. The text regions are annotated by four vertices of the quadrilateral. It is one of the largest and most challenging scene text detection datasets.

ICDAR 2015 is another multi-oriented text detection dataset only for English, which includes 1000 training images and 500 testing images. Similar to ICDAR 2017 MLT, the text region is also annotated as a quadrilateral.

ICDAR 2013 is a dataset that points at the horizontal text in the natural scene. This dataset consists of 229 training images and 233 testing images.

4.2. Comparisons with Other Methods

In this section, we compare PMTD with state-of-the-art methods on standard datasets. As shown in Tab. 1, 2, 3, our method outperforms others in all datasets.

ICDAR 2017 MLT: ImageNet [4] pre-trained ResNet50 is adapted to initialize network parameter. We train our model using ICDAR 2017 MLT training and validation images for 160 epochs. We use SGD as our optimizer with batch size 64. The initial learning rate is 0.08 and decays to one-tenth of the previous at the 80th and 128th epoch. During the training stage, images are cropped to 640×640 patches as described in Sec. 3.1. Results are shown in Tab. 1. For single scale testing, with resizing images’ long side to 1600, PMTD achieves an F-measure of 78.48%. We also resize the long side to 1600 and 2560 for multi-scale testing, and it achieves 80.13% F-measure, which outperforms the state-of-the-art method by 5.83%. Qualitative results are shown in Fig. 7.

ICDAR 2015: For ICDAR 2015, we use the pre-trained model from ICDAR 2017 MLT, and finetune another 40



Figure 7: Detection results of PMTD. Best viewed in color.

Method	ICDAR13 Eval	DetEval
CTPN [39]	85.00	86.00
SegLink [37]	-	85.30
TextBoxes* [23]	85.00	86.00
SSTD [8]	87.00	88.00
WordSup [12]	-	90.34
R2CNN [17]	87.73	-
DDR [10]	-	86.00
MCN [27]	88.00	-
Lyu <i>et al.</i> * [31]	88.00	-
RRD* [24]	89.00	-
TextBoxes++* [22]	88.00	89.00
PixelLink* [3]	-	88.10
FEN* [45]	91.60	92.30
FOTS* [26]	92.50	92.82
SPCNET [41]	92.10	-
Baseline	91.73	92.25
PMTD	93.40	93.59

Table 3: Comparison with other results on ICDAR 2013. * means multi scale testing. For PMTD, we only report single scale testing result.

epochs using ICDAR 2015 training data. Learning rate is set to 0.0008 and unchanged during training. For testing, images' long side are resized to 1920. As shown in Tab. 2, PMTD outperforms all other methods and achieves 1.19% higher F-measure than our baseline, which demonstrates the proposed PMTD brings consistent gain on different datasets.

ICDAR 2013: Similar to ICDAR 2015, we finetune 40 epochs on ICDAR 2017 MLT pre-trained model using IC-

Method	Precision	Recall	F-measure
Baseline	84.72	70.37	76.88
Baseline+DC+BU	85.17	70.75	77.29
PMTD	85.15	72.77	78.48

Table 4: Results of our models with different settings on ICDAR 2017 MLT dataset. PMTD clearly outperforms our baseline.

DAR 2013 training data, with fixed 0.0008 learning rate. We resize images' long size to 960 during testing. As shown in Tab. 3, PMTD surpasses all previous methods once more and gains 1.67% improvement to the baseline on this dataset.

4.3. Ablation Study

In this section, we conduct a series of comparative experiments. Experiment results show that our method achieves better performance and predicts more accurate text boxes.

Better performance: We first compare the performance of the baseline and PMTD on ICDAR 2017 MLT dataset. Results are shown in Tab. 4.

Baseline: Mask R-CNN baseline as described in Sec. 3.1, is a solid baseline that significantly outperforms state-of-the-art methods.

Baseline+DC+BU: As described in Sec. 3.3, we use dilated convolution (DC) and bilinear upsampling (BU) to predict more accurate soft text mask. We also use these two parts for our baseline to measure their gains. Experiment result indicates that dilated convolution and bilinear upsampling only increase baseline by 0.41%.

PMTD: Our proposed method. It achieves an im-

IoU	Matched number			F-measure		
	Baseline	PMTD	Relative improve	Baseline	PMTD	Relative improve
0.5	1784	1816	1.79%	88.14%	89.33%	1.35%
0.6	1696	1729	1.95%	83.60%	84.79%	1.42%
0.7	1443	1556	7.83%	70.44%	75.31%	6.91%
0.8	799	962	20.40%	38.36%	45.32%	18.14%
0.9	107	157	46.73%	5.14%	6.73%	30.93%

Table 5: Number of true positives and F-measure under different IoU threshold on ICDAR 2015. PMTD outperforms baseline significantly when IoU threshold is high.

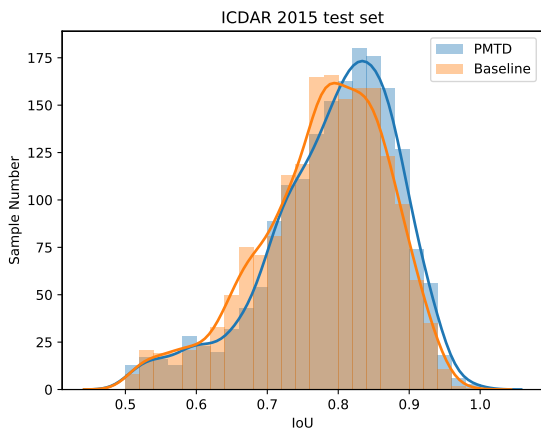


Figure 8: Distribution of predicted text boxes with different IoU. PMTD clearly predicts more accurate text boxes than baseline.

provement of 1.6% compared with the baseline. And in Sec. 4.2, experiments on ICDAR 2013 and ICDAR 2015 show PMTD brings consistent gains.

More accurate prediction: As mentioned in Sec. 1, with the help of the informative soft text mask, the plane clustering algorithm can regress more accurate text boundary and is more robust to imprecise predicted bounding box. However, the current evaluation cannot clearly reflect these two advantages due to the moderate evaluation (only require $IoU \geq 0.5$ on ICDAR 2015 and ICDAR 2017 MLT). So we evaluate PMTD and baseline under a higher IoU threshold.

Experiments are constructed on ICDAR 2015 test set for the absence of the label for ICDAR 2017 MLT test set. Results are summarized in Tab. 5. We can see PMTD outperforms baseline by a larger margin when the IoU threshold is 0.8. Especially, PMTD increases F-measure by 18.14% under 0.8 IoU threshold. Distribution of the true positive samples in different IoU are also illustrated in Fig. 8, which indicates a denser distribution in the high IoU interval for the PMTD.

Qualitative results are also shown in Fig. 7 and Fig. 9. From Fig. 7 we can see that PMTD can predict satisfactory text boxes, especially for text regions with strange shapes such as trapezoids and curves. Moreover, thanks to the gra-



Figure 9: PMTD is robust to imprecise predicted bounding box. From left to right: imprecise bounding box, predicted soft text mask, regression text box. It is worth noting that the soft text mask contains gradient information, which helps plane clustering algorithm to regress text box correctly.

dient information provided by the soft text mask, PMTD shows the robustness to imprecise predicted bounding boxes as shown in the Fig. 9.

5. Conclusion

In this work, we presented the Pyramid Mask Text Detector (PMTD), which encodes the shape and location information into the supervision and predicts a soft text mask for each text instance. A plane clustering algorithm is introduced to find the most fitting pyramid mask of the predicted soft text mask. Experiments on standard datasets demonstrate the effectiveness of our method.

References

- [1] C. K. Ch'ng and C. S. Chan. Total-text: A comprehensive dataset for scene text detection and recognition. In *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, pages 935–942. IEEE, 2017. 3
- [2] Y. Dai, Z. Huang, Y. Gao, Y. Xu, K. Chen, J. Guo, and W. Qiu. Fused text segmentation networks for multi-oriented scene text detection. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 3604–3609. IEEE, 2018. 6
- [3] D. Deng, H. Liu, X. Li, and D. Cai. Pixellink: Detecting scene text via instance segmentation. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 3, 6, 7
- [4] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. 2009. 6

- [5] A. Gupta, A. Vedaldi, and A. Zisserman. Synthetic data for text localisation in natural images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2315–2324, 2016. [2](#)
- [6] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. [1](#)
- [7] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [3](#)
- [8] P. He, W. Huang, T. He, Q. Zhu, Y. Qiao, and X. Li. Single shot text detector with regional attention. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3047–3055, 2017. [6, 7](#)
- [9] T. He, Z. Tian, W. Huang, C. Shen, Y. Qiao, and C. Sun. An end-to-end textspotter with explicit alignment and attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5020–5029, 2018. [3](#)
- [10] W. He, X.-Y. Zhang, F. Yin, and C.-L. Liu. Deep direct regression for multi-oriented scene text detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 745–753, 2017. [2, 3, 6, 7](#)
- [11] P. W. Holland and R. E. Welsch. Robust regression using iteratively reweighted least-squares. *Communications in Statistics-theory and Methods*, 6(9):813–827, 1977. [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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 4940–4949, 2017. [6, 7](#)
- [13] W. Huang, Y. Qiao, and X. Tang. Robust scene text detection with convolution neural network induced msr trees. In *European Conference on Computer Vision*, pages 497–511. Springer, 2014. [2](#)
- [14] Z. Huang, Z. Zhong, L. Sun, and Q. Huo. Mask r-cnn with pyramid attention network for scene text detection. *arXiv preprint arXiv:1811.09058*, 2018. [1, 3, 6](#)
- [15] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision*, 116(1):1–20, 2016. [2](#)
- [16] M. Jaderberg, A. Vedaldi, and A. Zisserman. Deep features for text spotting. In *European conference on computer vision*, pages 512–528. Springer, 2014. [2](#)
- [17] 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. [6, 7](#)
- [18] 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 *Document Analysis and Recognition (ICDAR), 2015 13th International Conference on*, pages 1156–1160. IEEE, 2015. [6](#)
- [19] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. Gomez-Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazan, and L. P. de las Heras. Icdar 2013 robust reading competition. In *Document Analysis and Recognition (ICDAR), 2013 12th International Conference on*, pages 1484–1493. IEEE, 2013. [6](#)
- [20] X. Li, W. Wang, W. Hou, R.-Z. Liu, T. Lu, and J. Yang. Shape robust text detection with progressive scale expansion network. *arXiv preprint arXiv:1806.02559*, 2018. [2, 3, 6](#)
- [21] Y. Li, Y. Yu, Z. Li, Y. Lin, M. Xu, J. Li, and X. Zhou. Pixel-anchor: A fast oriented scene text detector with combined networks. *arXiv preprint arXiv:1811.07432*, 2018. [6](#)
- [22] M. Liao, B. Shi, and X. Bai. Textboxes++: A single-shot oriented scene text detector. *IEEE transactions on image processing*, 27(8):3676–3690, 2018. [2, 3, 6, 7](#)
- [23] M. Liao, B. Shi, X. Bai, X. Wang, and W. Liu. Textboxes: A fast text detector with a single deep neural network. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017. [2, 7](#)
- [24] M. Liao, Z. Zhu, B. Shi, G.-s. Xia, and X. Bai. Rotation-sensitive regression for oriented scene text detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5909–5918, 2018. [6, 7](#)
- [25] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. [2](#)
- [26] X. Liu, D. Liang, S. Yan, D. Chen, Y. Qiao, and J. Yan. Fots: Fast oriented text spotting with a unified network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5676–5685, 2018. [3, 6, 7](#)
- [27] Z. Liu, G. Lin, S. Yang, J. Feng, W. Lin, and W. L. Goh. Learning markov clustering networks for scene text detection. *arXiv preprint arXiv:1805.08365*, 2018. [7](#)
- [28] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. [2, 3](#)
- [29] S. Long, J. Ruan, W. Zhang, X. He, W. Wu, and C. Yao. Textsnake: A flexible representation for detecting text of arbitrary shapes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 20–36, 2018. [2, 3, 6](#)
- [30] 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. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 67–83, 2018. [1](#)
- [31] P. Lyu, C. Yao, W. Wu, S. Yan, and X. Bai. Multi-oriented scene text detection via corner localization and region segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7553–7563, 2018. [6, 7](#)
- [32] J. Ma, W. Shao, H. Ye, L. Wang, H. Wang, Y. Zheng, and X. Xue. Arbitrary-oriented scene text detection via rotation proposals. *IEEE Transactions on Multimedia*, 20(11):3111–3122, 2018. [2, 3](#)
- [33] N. Nayef, F. Yin, I. Bizid, 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 *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, pages 1454–1459. IEEE, 2017. [6](#)
- [34] A. Odena, V. Dumoulin, and C. Olah. Deconvolution and checkerboard artifacts. *Distill*, 2016. [5](#)
- [35] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016. [2](#)
- [36] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. [2](#)
- [37] B. Shi, X. Bai, and S. Belongie. Detecting oriented text in natural images by linking segments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2550–2558, 2017. [6](#), [7](#)
- [38] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 761–769, 2016. [4](#)
- [39] Z. Tian, W. Huang, T. He, P. He, and Y. Qiao. Detecting text in natural image with connectionist text proposal network. In *European conference on computer vision*, pages 56–72. Springer, 2016. [2](#), [7](#)
- [40] S. Uchida. Text localization and recognition in images and video. *Handbook of Document Image Processing and Recognition*, pages 843–883, 2014. [2](#)
- [41] E. Xie, Y. Zang, S. Shao, G. Yu, C. Yao, and G. Li. Scene text detection with supervised pyramid context network. *arXiv preprint arXiv:1811.08605*, 2018. [1](#), [2](#), [3](#), [6](#), [7](#)
- [42] Q. Yang, M. Cheng, W. Zhou, Y. Chen, M. Qiu, and W. Lin. Inceptext: a new inception-text module with deformable psroi pooling for multi-oriented scene text detection. *arXiv preprint arXiv:1805.01167*, 2018. [6](#)
- [43] Q. Ye and D. Doermann. Text detection and recognition in imagery: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 37(7):1480–1500, 2015. [2](#)
- [44] L. Yuliang, J. Lianwen, Z. Shuaitao, and Z. Sheng. Detecting curve text in the wild: New dataset and new solution. *arXiv preprint arXiv:1712.02170*, 2017. [3](#)
- [45] S. Zhang, Y. Liu, L. Jin, and C. Luo. Feature enhancement network: A refined scene text detector. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. [2](#), [3](#), [7](#)
- [46] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang. East: an efficient and accurate scene text detector. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 5551–5560, 2017. [2](#), [3](#), [6](#)
- [47] Y. Zhu, C. Yao, and X. Bai. Scene text detection and recognition: Recent advances and future trends. *Frontiers of Computer Science*, 10(1):19–36, 2016. [2](#)