# EFFICIENT SCENE TEXT DETECTION WITH TEXTUAL ATTENTION TOWER

*L Zhang[⋆] YF Liu[⋆] H Xiao[‡] L Yang[⋆] GM Zhu[⋆] SA Shah[†] M Bennamoun[†] and PY Shen[⋆]*

[⋆]Xidian University, School of computer science and technology, China
[‡]OrionStar Ltd., China
[†]University of Western Australia, Australia

## ABSTRACT

Scene text detection has received attention for years and achieved an impressive performance across various benchmarks. In this work, we propose an efficient and accurate approach to detect multi-oriented text in scene images. The proposed feature fusion mechanism allows us to use a shallower network to reduce the computational complexity. A self-attention mechanism is adopted to suppress false positive detections. Experiments on public benchmarks including ICDAR 2013, ICDAR 2015 and MSRA-TD500 show that our proposed approach can achieve better or comparable performances with fewer parameters and less computational cost.

***Index Terms***— Scene text detection, multi-oriented text, textual attention tower

## 1. INTRODUCTION

Oriented scene text detection is one of the most challenging computer vision tasks. The primary task is to spot text objects in different types of scenes. Text object may differ in many aspects, such as font type, texture, orientation. Moreover, bounding a single text object with an up-right detection box may lead to a low IoU and detection quality. Several approaches [1, 2, 3] have already presented impressive successes on various public benchmarks and competitions.

The key of text detection is designing features to distinguish text from backgrounds. Recently, Convolutional Neural Networks (CNN) based methods such as EAST [1] and IncepText [4] have achieved the state-of-the-art performance for text detection. Like other computer vision tasks, deeper networks provide better performances. EAST initially adopts PVANET [5] and VGG-16 [6], the subsequent approaches used ResNet [7] and then ResNeXt [8].

Although several text detection frameworks have been designed, many of the recently proposed models mainly focus on detection precision. These approaches achieve high precision by complex models and high computational cost, but
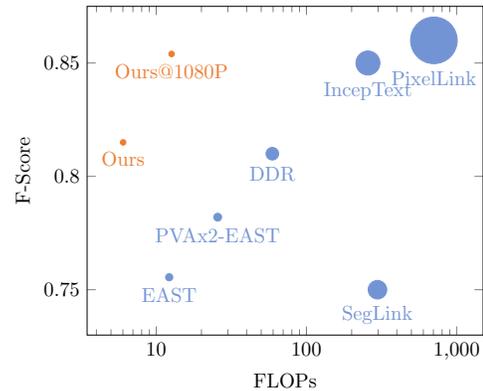
**Fig. 1**: Performance versus floating point operations (FLOPs) on ICDAR 2015 text localization challenge. **Ours** denotes our model evaluated at 720P resolution, **Ours@1080P** denotes our model evaluated at 1080P resolution. The area of each node denotes the total parameters in its network.

performance increase is relatively limited. To address these limitations, we design a novel, computationally efficient and extendable network structure to perform competitive detection compared with the former approaches.

Our main contribution can be summarized as follows:

- We propose a novel, efficient *Textual Attention Tower* (TAT) structure.

- We evaluate our proposed method on ICDAR 2013, ICDAR 2015 and MSRA TD-500 datasets. As shown in Fig. 1, the proposed module achieves a significant decrease in the computational cost and a higher accuracy compared with the state-of-the-art models.

## 2. PROPOSED METHOD

### 2.1. Architecture Overview

An overview of our framework is illustrated in Fig. 2. The MobileNetV2 architecture is used as the base network. In order to further reduce the computational cost, only the first seven residual blocks of the MobileNetV2 are used.
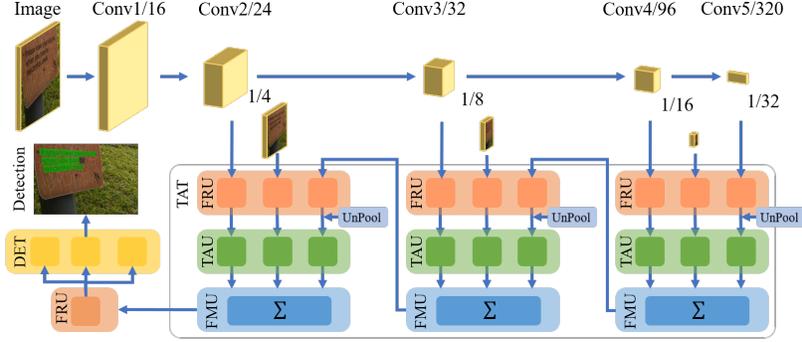
**Fig. 2**: Overview of our proposed model. TAT stands for Textual Attention Tower, FRU stands for Feature Refine Unit, TAU stands for Textual Attention Unit, FMU stands for Feature Mixup Unit, DET stands for the Detection Branch. Conv i/c stands for the feature maps with $c$ output channels extracted from the $i$-th stage of the MobileNetV2. 1/n stands for the input image down-sampled to 1/n of the original scale.

The key of reducing computational cost and parameter size is our Textual Attention Tower (TAT) architecture, which is designed to fuse the extracted feature maps. To avoid the degeneration of low level features in deep CNN, we use the down-sampled input images as extra channels of the intermediate feature maps.

The detection branches in our proposed method are denoted as *DET*. Inspired by [1, 9], we use rotated box (RBOX) to describe text regions. Thus the *DET* branch is simply $1 \times 1$ convolutions to map final feature to detections.

### 2.2. Textual Attention Tower

The Textual Attention Tower (TAT) is designed to fuse the feature maps from different stages. As we adopt segmantation based methodology to regress the geometric information of text regions, detecting text regions can be seen as two simple subtasks: text/non-text prediction and distance regression. Both of these tasks needs large receptive field, and can perform well on tensors with less channels. As shown in Fig. 2, the TAT has three main parts, Feature Refine Unit (FRU), Textual Attention Unit (TAU) and Feature Mixup Unit (FMU).

**Feature Refine Unit** The FRU is a "bottleneck" residual block [7] to refine the feature maps and reduce its number of channels. Regardless of the number of input channels, we set the number of output channels of all the FRU modules to 32. As shown in Fig. 2, we adopt a dedicated FRU module for each input feature map, and for each down-sampled image, we adopt two cascaded FRU modules to extract the low-level feature maps.

**Textual Attention Unit** The TAU is a spatial self-attention module designed to encode the global context information. The key idea of TAU is to gather global context information to support the inference at the current position. We adopt dilated convolution as the basic operation to enlarge the receptive field. (**a**) The first part of a TAU module is a standard convolution block $c$ to reduce the number of channels of the input feature map $f$, which ensures that the sub-

sequent operations can be performed at a low computational cost. (**b**) The second part is a context encoder, consisting of four dilated convolutional blocks $e_1$, $e_2$, $e_3$ and $e_4$. Each encoder $e_i$ has an isolated depth-wise convolution layer with a different dilation rate $r = 2i - 1$. Appropriate padding is configured to ensure that the output of $e_i$ has the same spatial scale as the input feature maps. The detailed configuration of dilated convolutions used in TAU modules is illustrated on Fig. 3. (**c**) The last part of a TAU module is a convolutional decoder block $dec$, which accepts the concatenated feature maps from all encoders and decodes it to the spatial attention map. Therefore, the TAU can be formulated as

$$\textbf{TAU}(f) = x \otimes \sigma(dec(\langle e_1(c(x)), \cdots, e_4(c(x))\rangle)) \quad (1)$$

where $\sigma$ is the sigmoid nonlinearity and $\otimes$ is the element-wise multiplication with broadcast semantic.

**Feature Mixup Unit** The FMU is a simple, element-wise operator to finally fuse all these features. In all of our experiments, we adopt the element-wise addition as our FMU function.

### 2.3. Loss Function

Our loss function can be formulated as

$$L = \lambda_c L_c + \lambda_d L_d + \lambda_r L_r, \quad (2)$$

where $L_c$ is the classification loss, $L_d$ is the distance regression loss and $L_r$ stands for the rotation regression loss. $\lambda_c$, $\lambda_d$ and $\lambda_r$ are coefficients to balance the different loss terms. We set $\lambda_c$ to 1, $\lambda_d$ to 2 and $\lambda_r$ to 20 in all of our experiments.

To directly maximize the IoU between the candidate and the ground truth, we adopt dice loss [10] as our classification loss, which can be formulated as

$$L_{cls} = 1 - 2\frac{\textbf{S}^*\hat{\textbf{S}}}{\textbf{S}^* + \hat{\textbf{S}}}, \quad (3)$$

where $\hat{\textbf{S}}$ and $\textbf{S}^*$ are the generated and the predicted score maps.
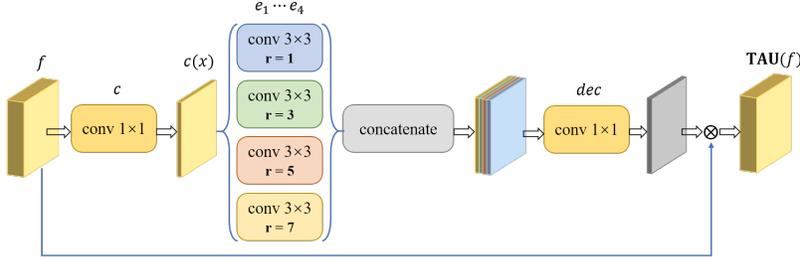
**Fig. 3**: Detailed structure of our proposed Textual Attention Unit (TAU).

For regression tasks, we adopt IoU loss, which can be formulated as

$$L_{dis} = -\log \frac{\mathbf{intersection}(\mathbf{D}^*, \hat{\mathbf{D}})}{\mathbf{union}(\mathbf{D}^*, \hat{\mathbf{D}})}, \qquad (4)$$

where $\hat{\mathbf{D}}$ and $\mathbf{D}^*$ are the generated and the regressed distance maps, **intersection** and **union** are functions to calculate the area of the intersection and union parts between $\hat{\mathbf{D}}$ and $\mathbf{D}^*$.

The rotation regression loss is constructed based on the cosine function, which can be seen as a loose variant of Smoothed-L1 loss.

$$L_r(\mathbf{R}^*, \hat{\mathbf{R}}) = 1 - \cos(\mathbf{R}^* - \hat{\mathbf{R}}) \qquad (5)$$

## 3. EXPERIMENTS

### 3.1. Benchmark Datasets and Data Augmentation

We evaluate our approach on three public benchmark datasets. All these datasets consist of scene text objects with arbitrary orientations.

**ICDAR 2013** [11] consists of 229 training images and 233 testing images of different resolutions, and most of the text instances are horizontal or near-horizontal.

**ICDAR 2015** [12] is used in the ICDAR2015 Robust Reading Competition Challenge 4. It contains 1000 images for training and 500 images for testing. The text bounding boxes have multi-orientations, and they are specified by the coordinates of their four corners in a clockwise manner.

**MSRA-TD500** [13] contains 300 training images and 200 testing images. Different from ICDAR2015, this dataset consists of text lines and separate words.

**Data Augmentation**. First, we randomly rotate each image by -15 to 15 degrees. For each image, we randomly choose a text box as a *kernel*, and then randomly expand the kernel to fit the crop size, which is set to 640 in all of all experiments. To further enrich the variety of object scales, we perform *subsampling* in the object-centrical cropping processes. During the expansion of the kernel , we first expand the kernel to $k \times 640$ where $k \sim U(0.5, 2)$, and then resize

the patch to $640 \times 640$ pixels with a bilinear interpolation. Color-space jittering and Gaussian blurring are then applied to the cropped image patches.

### 3.2. Experimental Setup

All of our models are trained on a local machine with 4 NVIDIA TITAN X Pascal GPUs. Our proposed models are trained with the ADADELTA [14] optimizer. We set the initial learning rate to 1 and the weight decay coefficient to $1 \times 10^{-5}$. The base CNN is initialized with parameters pretrained on ImageNet, the rest of the parameters in our model are initialized according to [15]. We adopt Cross-GPU Batch Normalization proposed in [16] to avoid issue that the data distribution have a significant difference with ImageNet

### 3.3. Experimental Results

Our proposed model is trained and evaluated on ICDAR2015, ICDAR2013 and MSRA TD-500.

**Oriented Text Detection on ICDAR 2015**. We conduct experiments on ICDAR2015 Challenge 4. The initialized model is optimized with an ADADELTA optimizer for 600 epochs. As shown in Table 1, when the test images are fed at original scale ($1280 \times 720$), our model achieves an F-score of 81.5; When the test images are upsampled to $1920 \times 1080$ via bilinear interpolation, our model reaches 85.4 in F-score without any fine-tuning or model ensembling, which outperforms most of the previous methods.

**Horizontal Text Detection**. We fine-tune our model on ICDAR 2013 training set for 200 epochs using the ADADELTA optimizer. Test images in ICDAR 2013 have different resolution. We resize all the testing images to $800 \times 600$ pixels. As shown in Table.2, our model outperforms most of the existing methods in term of F-score.

**Long Text Detection**. The main challenge in MSRA TD-500 is long text detection. We initialize our model with parameters trained on ICDAR 2015, and then optimize this model with an ADADELTA optimizer in the TD-500 training set. To fit the testing images with our proposed method, we down-sample all testing images to $960 \times 720$ pixels. As

**Table 1**: Results on ICDAR2015 Challenge 4 Incidental Scene Text Localization task. MS means multi-scale testing, 1080P means test at 1920 × 1080.

| Method | Recall | Precision | F-score |
|---|---|---|---|
| SegLink [17] | 73.1 | 76.8 | 75.0 |
| EAST [1] | 71.4 | 80.6 | 75.7 |
| EAST MS [1] | 78.3 | 83.3 | 80.7 |
| PixelLink [18] | 82.0 | 85.5 | 83.7 |
| IncepText [4] | 80.6 | **90.5** | 85.3 |
| Ours | 77.8 | 85.8 | 81.5 |
| Ours 1080P | **83.2** | 87.7 | **85.4** |

**Table 2**: Evaluation results on ICDAR 2013.

| Method | recall | precision | F-score |
|---|---|---|---|
| TextBoxes++ [3] | 74.0 | 86.0 | 80.0 |
| PixelLink [18] | **83.6** | 86.4 | 84.5 |
| SegLink [17] | 83.0 | 87.7 | **85.3** |
| Ours | 74.3 | 90.3 | 82.1 |
| Ours MS | 78.6 | **92.9** | 85.2 |

shown in Table 3, on TD-500, our proposed method outperforms segmentation based methods, but cannot surpass the region proposal based frameworks such as IncepText. The main reason is that without the object-level supervision information, segmentation based methods usually fail to separate text lines and the surrounding text-like regions.

**Table 3**: Evaluation results on TD-500.

| Method | Recall | Precision | F-score |
|---|---|---|---|
| DDR [19] | 70.0 | 77.0 | 74.0 |
| EAST [1] | 67.4 | 87.3 | 76.1 |
| SegLink [17] | 70.0 | 86.0 | 77.2 |
| PixelLink [18] | 73.2 | 83.0 | 77.8 |
| IncepText [4] | **79.0** | **87.5** | **83.0** |
| Ours | 75.3 | 81.4 | 78.2 |

**Comparison of Computational Complexity**. To further compare the computational complexity for our proposed method with the existing methods. we compare the theoretical FLOPs per pixel of every method over their relative f-score achieved on ICDAR 2015. Since the multi-scale testing strategies and model-ensembling require significantly more computational resources, and the increments of FLOPs depend on the implementation, only performances achieved by a single model without multi-scale testing are included. As shown in Table 4, our proposed models outperform existing models in terms of computational complexity, and can still achieve a competitive performance.

**Table 4**: Comparison on per pixel computational complexity and corresponding relative F-scores.

| Method | FLOPs | F-score |
|---|---|---|
| PixelLink [18] | 765.65K | 83.7 |
| EAST-VGG16 [1] | 310.62K | 76.4 |
| SegLink [17] | 322.42K | 75.0 |
| IncepText [4] | 278.53K | 85.3 |
| DDR [19] | 64.34K | 81.0 |
| EAST-PVA [1] | 13.23K | 75.7 |
| Ours | 6.65K | 81.5 |
| Ours@1080P | **6.65K** | **85.4** |

**Effectiveness of TAT Module**. Table 5 summarized more detailed results of our models with different settings on IC-DAR 2015. We choose the best-performing model proposed EAST-PVAx2 in [1] as baseline in this comparison, which is listed in the first line in Table 5. When PVAx2 is replaced with full sized MobileNetV2, the effectiveness of MobileNetV2 itself performs significant performance improvement, and reduced the computational cost by 25% approximately comparing with EAST. For the other four model configurations, the last two convolutional blocks in MobileNetV2 are omitted. With FRUs significantly reduce computational cost, and TAUs improves the detection precision, our model achieves better performance than EAST with about one quarter FLOPs.

**Table 5**: Effectiveness of TAT on ICDAR2015 incidental scene text location task. "M" means "MobileNetV2" and "I" means use raw input as extra feature. "P", "R", "F" represent "Precision", "Recall", "F-measure" respectively.

| M | FRU | TAU | I | R | P | F | FLOPs |
|---|---|---|---|---|---|---|---|
| | | | | 73.5 | 83.6 | 78.2 | 23.85G |
| ✓ | | | | 73.7 | **87.8** | 80.1 | 17.75G |
| ✓ | ✓ | | | 77.4 | 83.6 | 80.4 | **5.79G** |
| ✓ | ✓ | ✓ | | 77.2 | 85.8 | 81.3 | 5.85G |
| ✓ | ✓ | ✓ | ✓ | **77.8** | 85.8 | **81.5** | 6.03G |

## 4. CONCLUSION

In this paper, we propose a novel and efficient multi-oriented text detection method from natural scene images. The main idea of our design is the use of dilated convolution to keep a reasonable yet abundant information for different levels of receptive fields. Another improvement comes from using homogenous bottlenecks with the base network to refine feature maps. We achieve a better performance of our proposed technique on three public scene text benchmarks.

# 5. REFERENCES

[1] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang, "East: An efficient and accurate scene text detector," in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.

[2] Pan He, Weilin Huang, Tong He, Qile Zhu, Yu Qiao, and Xiaolin Li, "Single shot text detector with regional attention," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 3066–3074, 2017.

[3] Minghui Liao, Baoguang Shi, and Xiang Bai, "Textboxes++: A single-shot oriented scene text detector," *IEEE Transactions on Image Processing*, vol. 27, pp. 3676–3690, 2018.

[4] Qiangpeng Yang, Mengli Cheng, Wenmeng Zhou, Yan Chen, Minghui Qiu, Wei Lin, and Wei Chu, "Inceptext: A new inception-text module with deformable psroi pooling for multi-oriented scene text detection," in *IJCAI*, 2018.

[5] Kye-Hyeon Kim, Yeongjae Cheon, Sanghoon Hong, Byung-Seok Roh, and Minje Park, "Pvanet: Deep but lightweight neural networks for real-time object detection," *CoRR*, vol. abs/1608.08021, 2016.

[6] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.

[7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.

[8] Saining Xie, Ross B. Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He, "Aggregated residual transformations for deep neural networks," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5987–5995, 2017.

[9] Xuebo Liu, Ding Liang, Shi Yan, Dagui Chen, Yu Qiao, and Junjie Yan, "Fots: Fast oriented text spotting with a unified network," *CoRR*, vol. abs/1801.01671, 2018.

[10] Carole H. Sudre, Wenqi Li, Tom Vercauteren, Sébastien Ourselin, and M. Jorge Cardoso, "Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations," in *DLMIA/ML-CDS@MICCAI*, 2017.

[11] *12th International Conference on Document Analysis and Recognition, ICDAR 2013, Washington, DC, USA, August 25-28, 2013*. IEEE Computer Society, 2013.

[12] Dimosthenis Karatzas, Lluis Gomezbigorda, Anguelos Nicolaou, Suman K Ghosh, Andrew D Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Chandrasekhar, Shijian Lu, et al., "Icdar 2015 competition on robust reading," pp. 1156–1160, 2015.

[13] Cong Yao, Xiang Bai, Wenyu Liu, Yi Ma, and Zhuowen Tu, "Detecting texts of arbitrary orientations in natural images," pp. 1083–1090, 2012.

[14] Matthew D. Zeiler, "Adadelta: An adaptive learning rate method," *CoRR*, vol. abs/1212.5701, 2012.

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 1026–1034, 2015.

[16] Chao Peng, Tete Xiao, Zeming Li, Yuning Jiang, Xiangyu Zhang, Kai Jia, Gang Yu, and Jian Sun, "Megdet: A large mini-batch object detector," *CoRR*, vol. abs/1711.07240, 2017.

[17] Baoguang Shi, Xiang Bai, and Serge J. Belongie, "Detecting oriented text in natural images by linking segments," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3482–3490, 2017.

[18] Dan Deng, Haifeng Liu, Xuelong Li, and Deng Cai, "Pixellink: Detecting scene text via instance segmentation," *CoRR*, vol. abs/1801.01315, 2018.

[19] Wenhao He, Xu-Yao Zhang, Fei Yin, and Cheng-Lin Liu, "Deep direct regression for multi-oriented scene text detection," *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 745–753, 2017.