

An End-to-end Video Text Detector with Online Tracking

Hongyuan Yu^{13*†}, Chengquan Zhang^{2*}, Xuan Li², Junyu Han², Errui Ding², and Liang Wang¹³⁴

¹University of Chinese Academy of Sciences (UCAS)

²Department of Computer Vision Technology(VIS), Baidu Inc.

³Center for Research on Intelligent Perception and Computing (CRIPAC),
National Laboratory of Pattern Recognition (NLPR)

⁴Chinese Academy of Sciences Artificial Intelligence Research (CAS-AIR)

Email: hongyuan.yu@cripac.ia.ac.cn, {zhangchengquan,lixuan12,hanjunyu,dingerrui}@baidu.com, wangliang@nlpr.ia.ac.cn

Abstract—Video text detection is considered as one of the most difficult tasks in document analysis due to the following two challenges: 1) the difficulties caused by video scenes, i.e., motion blur, illumination changes, and occlusion; 2) the properties of text including variants of fonts, languages, orientations, and shapes. Most existing methods attempt to enhance the performance of video text detection by cooperating with video text tracking, but treat these two tasks separately. In this work, we propose an end-to-end video text detection model with online tracking to address these two challenges. Specifically, in the detection branch, we adopt ConvLSTM to capture spatial structure information and motion memory. In the tracking branch, we convert the tracking problem to text instance association, and an appearance-geometry descriptor with memory mechanism is proposed to generate robust representation of text instances. By integrating these two branches into one trainable framework, they can promote each other and the computational cost is significantly reduced. Experiments on existing video text benchmarks including ICDAR2013 Video, Minetto and YVT demonstrate that the proposed method significantly outperforms state-of-the-art methods. Our method improves F-score by about 2% on all datasets and it can run realtime with 24.36 fps on TITAN Xp.

I. INTRODUCTION

With the rapid development of mobile internet, video-related applications become more and more popular in our daily life. The analysis of videos therefore becomes an important task for practical applications. Among various types of objects appearing in videos, text usually contains abundant semantic information and plays an important role in many applications, such as video annotation, multimedia retrieval and industrial automation [25], [27].

In the last few years, we have witnessed significant efforts in tackling video text detection. Previous works on this problem [16], [19], [24], [11] are generally carried out in two steps: text in individual frame is detected first, the data association is then performed. However, those two-step methods suffer from the following problems: (1) single-frame detection in video does not make full use of the temporal context in the video; (2) tracking after detection needs additional networks to extract tracking features, which leads to additional computational cost, and most of these

*Equal contribution. †This work is done when Hongyuan Yu is intern at VIS, Baidu Inc.

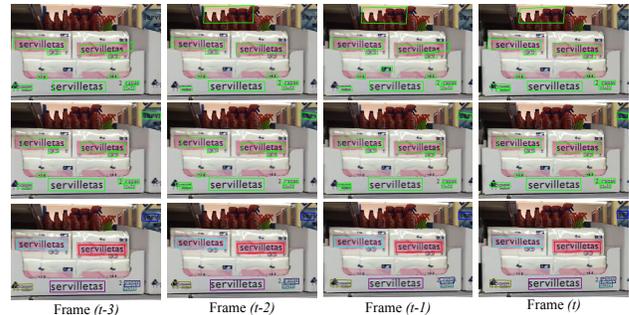


Figure 1. Video text detection and tracking. The first row is the detection results of EAST [30]. The second row is the results from our video text detection branch with ConvLSTM. The final row is our detection results with online tracking, where boxes with the same color in different frames belong to the same trajectory. (Best view in color.)

tracking methods are offline, which is limited in practical applications; (3) these two parts are separately trained, and they cannot fully utilize the supervision information of each other. Video text detection and tracking tasks are closely related

In this work, we present a novel end-to-end video text detector with online tracking. By integrating detection and tracking together, they can exploit the supervision information of each other. The whole pipeline of our proposed approach is depicted in Fig. 2. Specifically, to take full advantage of the temporal domain information and texture properties of scene text while preserving its structure, ConvLSTM [23] layer is introduced to the video text detection branch. In addition, the appearance-geometry descriptor (AGD) and corresponding estimated appearance-geometry descriptor (EAGD) are proposed to model short-term target association for online tracking. Simultaneously, these descriptors are updated in time steps to capture the information of the long-term multiple targets changing process, i.e., new target entry and old target departure. As shown in Fig. 1, our method can significantly improve the performance.

The contributions of our method are summarized as follows: (1) to the best of our knowledge, this is the first end-to-end video text detection and online tracking framework; (2) we introduce ConvLSTM to our detection branch, which is very useful for capturing spatial-temporal information; (3) the proposed appearance-geometry descriptor has been proven to be robust and effective for multiple text instance

association; (4) extensive experiments have shown the effectiveness of our method, and we have obtained state-of-the-art results in multiple public benchmarks.

II. RELATED WORK

Single frame text detection has made great progress recently. However, for video text detection, how to make full use of the context information in video is still not addressed well. In this section, the development of single image detection and video text detection will be reviewed.

Single frame text detection It is known that numerous methods for text detection have been successfully proposed in recent years. Specifically, component based methods, segmentation based methods and detection based methods are the main kinds of single frame detection. *Component based methods* [2], [13] usually detect parts or components of text first, after then a set of complex procedures including component grouping, filtering and word partition are followed to obtain final detection results in word-level. *Segmentation based methods* [28], [22] regard all the pixels in one word / text-line as a whole instance. It can handle arbitrary shape of text, but rely too much on fine grained segmentation and also need some discontinuous post-processing operations. *Detection based methods* [4], [30] draw inspiration from general object detection, and output text detection results in word / text-line level directly.

Video text detection Standing on the shoulder of single frame text detection, many video mode text detection methods have also been proposed. Researchers try to enhance the detection result through tracking, namely tracking based text detection methods. These methods utilize some specific tracking techniques, such as multi-strategy tracking methods [31], dynamic programming [16] and network flow based methods [24], to track text and then heuristically combine detected results in passed frames. But they are essentially based on single-frame detection methods, and the training of detector is independent of the tracking. Moreover, they do not make full use of the abundant temporal information of video. Wang et al. [17] notice that the cues of background regions can promote video text detection. However, they only consider the short-term dependencies. There is no doubt that some structures like optical flow, Conv3D [5] and ConvLSTM [23] are efficient to catch spatial-temporal information, which have been explored in general object tracking. Inspired by that, in this paper, we employ ConvLSTM in video text detection branch in our framework. With the help of long-term spatial-temporal memory and online tracking, our end-to-end video text detector achieves better performance.

III. APPROACH

A. Overall Architecture

The proposed method, as shown in Fig. 2, integrates video text detection and tracking in an unified framework

through the descriptor generation module. Given a video, all frames should pass a backbone network (ResNet50 [3] + U-Net [12]) to extract common features for detection and tracking. For video text detection, we adopt the anchor-free regression manner [30] to detect the quadrangles of words in a per-pixel manner. Notice that a ConvLSTM block is followed with common features to extract spatial-temporal information. The benefits from ConvLSTM block are shown in Tab. I. For video text tracking, detected proposals and common feature maps of current frame are fed into the descriptor generation module first, and then corresponding appearance-geometry descriptors AGD_t are output. In order to associate text instances in sequential frames, the descriptors of frame $(t - 1)$ namely AGD_{t-1} are passed through GRU units to generate $EAGD_{t-1}$, meaning the estimated states of appearance-geometry descriptors at time t . After then, a similarity matrix is build on $EAGD_{t-1}$ and AGD_t , where two text proposals belonging to the same trajectory should have a small metric distance value in the similarity matrix. With the help of online text tracking, our methods are able to improve the performance of video text detection.

B. Text Detection Branch

Text in videos always appears in sequential frames with abundant temporal information. In most of existing methods, video text detection is usually performed in individual frame or integrated temporal information with short term dependencies. To address this problem, we incorporate a ConvLSTM block into our text detection branch to propagate frame-level information across time while maintaining the structural properties. Accordingly, the formulation of inferring feature maps F_t at the t -th frame is:

$$(F_t, s_t) = ConvLSTM(M(I_t), s_{t-1}) \quad (1)$$

In this equation, $M(I_t)$ is the common feature maps of the t -th frame (I_t) obtained by backbone network. s_{t-1} and s_t mean the hidden sates of ConvLSTM at time $t - 1$ and time t respectively. By this way, features can be directly modulated by their previous frames and recursively depend on other frames in a long time range.

After integrating temporal information, convolution operation is applied to make dense per-pixel predictions in word-level. Similar to EAST [30], pixels within the quadrangle annotation of text instance are considered as positive. For each positive sample, the offsets to the 4 vertexes of the quadrangle are predicted at the following 8 channels. Therefore, the loss of detection branch is composed of two terms: text/non-text classification term and quadrangle offset regression term. The detailed definition of detection loss in t -th frame is illustrated as follows:

$$L_{det}(t) = L_{cls}(t) + \alpha L_{off}(t) \quad (2)$$

where $L_{cls}(t)$ measures the text/non-text classification by dice loss and $L_{off}(t)$ is the smooth-L1 loss to measure the

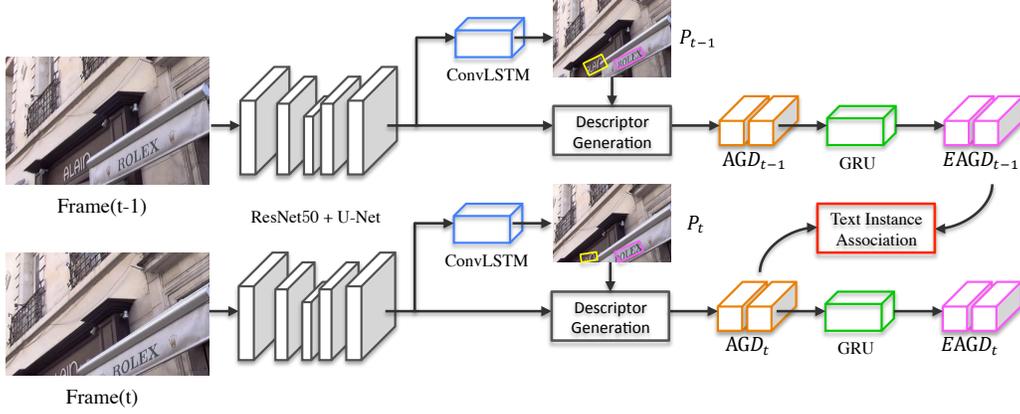


Figure 2. The proposed architecture of video text detection with online tracking.

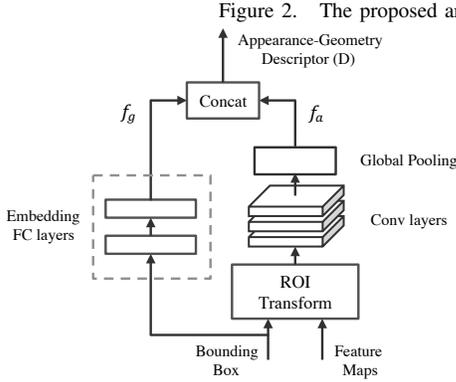


Figure 3. Descriptor Generation.

quality of regression offset. α is a hyper-parameter, which is set to 5 in our experiments. In addition, we use NMS (Non-maximum suppression) to get preliminary detection results and feed top-K proposals into the next tracking branch.

C. Text Tracking Branch

In order to improve the robustness of text instance representation in some difficult circumstances such as occlusion, motion blur, and etc, this section proposes an effective and efficient descriptor generation module to generate the descriptors of text candidates, which contains not only geometry features but also appearance features. And the descriptor of the same text candidate appears in the next frame is estimated through a GRU unit which can make use of the trajectory history information and capture the information of long-term multiple targets changing process.

We define a novel descriptor, namely *appearance-geometry descriptor* (AGD), for each text candidate. The descriptors of the kept K proposals from detection branch, denoted as AGD_t , contain two parts: the first part is the *appearance feature*, which is extracted by ROI Transform layer [15] from the valid regions of text candidates in the common feature maps $M(I_t)$. And the second part is the *geometry feature* that is composed of the embedding values of quadrilateral coordinates.

As shown in Fig. 3, we firstly use ROI Transform layer to

extract initial text feature block of K text candidates from the shared feature map. Then three convolution layers with $3 * 3$ kernel and one global pooling layer are followed to generate the final appearance feature, denoted as f_t^a for time t . Next, we feed all normalized coordinate vectors $\{g_n | n = 0, \dots, 7\}$ of detected K proposals into the geometry embedding layers that are composed of two fully connection layers to get final geometry feature at time t , namely f_t^g . Finally, the appearance features f_t^a and geometry features f_t^g are concatenated to generate the appearance-geometry descriptor AGD_t , which can be formulated as follows:

$$AGD_t = Concat([f_t^a, f_t^g]) \quad (3)$$

Then we feed the descriptor into a GRU unit to estimate the descriptor of the same instance appeared in the next frame, namely *estimated appearance-geometry descriptor* ($EAGD$). As we know, GRU is an efficient structure to capture the temporal changing information. Therefore, instead of building similarity matrix on appearance-geometry descriptors between two adjacent frames, we choose to match descriptors of the current frame with the estimated descriptors of the previous frame. The estimated descriptor $EAGD_t$ in current frame can be expressed as:

$$(EAGD_t, h_t) = GRU(AGD_t, mask_t * h_{t-1}) \quad (4)$$

where AGD_t is the appearance-geometry feature of text candidates obtained in current frame, h_{t-1} is the hidden state value of GRU in previous frame. Specially, $mask_t$ is the hidden state mask to control whether we need to reset GRU hidden states or not. It will be set to zero when the instance does not exist in the previous frame; otherwise, it will be set to one.

Video text tracking tries to match text instances belonging to the same object in adjacent frames while maintaining the identities of text instances. To simplify this problem, we convert the text instance association to pairwise matching by defining an association objective function, where descriptor representations should be close for the positive pairs and far for the negative pairs. Therefore, the contrastive loss is

suitable for this task, then our tracking loss L_{track} at time t can be represented by:

$$L_{track}(t) = \frac{1}{K^2} \sum_{i=1}^K \sum_{j=1}^K y d^2 + (1-y) \max(m-d, 0)^2 \quad (5)$$

where d denotes euclidean distance of text instances between adjacent frames and y is the pairs label $L_{i,j}^t$ whose value is 1 for positive pairs and 0 for negative pairs. And m is the margin value that is set to 1.0 in our experiments.

Finally, combined with detection loss $L_{det}(t)$ in Eq. 2, the full multi-task loss function is:

$$L_{d\&t} = \frac{1}{N} \sum_{t=1}^N L_{det}(t) + \beta L_{track}(t) \quad (6)$$

where N is length of video frames, and β is a hyperparameter to control the trade-off between detection and tracking loss. β is set to 0.1 in our experiments, where different β values have little effect on the final result.

D. Inference

In the inference phase, we propose an efficient and robust online trajectory generation method to improve the performance of video text detection, which is present in algorithm 1. Besides, the inference speed of our method on TITAN Xp can reach 24.36fps.

Algorithm 1 Online Trajectory Generation

Input: A frame of current time step I_t , the previous detection results D_{t-1} , corresponding estimated appearance-geometry descriptors ($EAGD_{t-1}$), and tracklet set T_{t-1} at time $t-1$.

Output: Current frame detection results D_t , the corresponding estimated appearance-geometry descriptors ($EAGD_t$) for next time step, and the tracklet set T_t .

- 1: Feed the I_t into the network to get the primary detection results D_t^* by θ_l at time t , and obtain the corresponding appearance-geometry descriptors AGD_t .
 - 2: Calculate the similarity matrix $S_t \leftarrow \text{similarity}(EAGD_{t-1}, AGD_t)$.
 - 3: Use Kuhn-Munkres algorithm with threshold value θ_m to find the matching pairs M .
 - 4: Update the part of tracklet set which find matching text instances in current frame, namely T_{update} .
 - 5: Reward the matched instance confidence score by $\tau * \ln(\text{length}(\text{tracklet}))$. If the scores of no-matching candidates are higher than θ_h , new trajectories T_{new} are built for them.
 - 6: Obtain the full tracklet set of current time: $T_t \leftarrow T_{update} + T_{new}$.
 - 7: Mark the corresponded detected boxes of T_t at current time step as D_t , and feed the corresponding AGD_t into the GRU to get the $EAGD_t$.
-

IV. EXPERIMENTS

A. Datasets

- **ICDAR 2013 Video** [6] This dataset consists of 28 videos lasting from 10 seconds to 1 minute in indoors or outdoors scenarios. 13 videos used for training and 15 for testing. Its frame size ranges from 720 x 480 to 1280 x 960.
- **Minetto Dataset** [8] Minetto Dataset consists of 5 videos in outdoor scenes. The frame size is 640 x 480 and all videos are used for test.
- **YVT** [10] This dataset contains 30 videos, 15 for training and 15 for testing. Different from the above 2 datasets, it contains web videos except for scene videos. The frame size is 1280 x 720.

B. Implementation Details

The ResNet50 [3] pretrained on ImageNet is employed as our initialized model. Then Adam is employed to train our model with the initial learning rate being 10^{-4} which decays by 0.94 times every 10 thousand iterations. All training videos are harvested from the training set of ICDAR2013 and YVT with data augmentation. Random crop and resize operations are applied for the first frame with the scale chosen from [0.5, 1.0, 2.0, 3.0]. Other frames in the video clip perform the same operation as the first frame. The frame interval is selected randomly from 1 to 5 to improve robustness. In our experiment, each video clip has 24 frames and every frame is resized and padded to 512 x 512 during the training process. Each frame contains 10 detection boxes, including positive samples and negative samples. The shape of detection box extracted by ROI Transform layer is set to 8 x 64. The size of the appearance descriptor and the geometry descriptor for one instance is 128 and 8, so the size of AGD in our experiment is 136. All experiments are conducted on 8 P40 GPUs and each GPU has 1 batch.

C. Evaluation of Video Text Detection

In this section, we evaluate the effect of the short-term and long-term memory for video text detection. As shown in Tab. I, the optical flow is adopted and there is about 0.2% drop in f-measure, indicating that the common object detection tracking method is not applicable to video text. Then, we exploit the ConvLSTM block in our detection branch and compare it with Conv3D [5]. As can be seen in Tab. I, both of these methods outperform single frame detection and the optical flow, and ConvLSTM gives a 0.66 points F-measure gain over Conv3D. The improvement in video text detection is mainly due to the fact that the temporal information in sequential frames is beneficial to video text detection since the text in video generally does not change as sharply as a general object. And as a valid long-term memory extractor, ConvLSTM can take advantage of more sequential frames information than Conv3D, resulting in an improvement in performance.

Method	Precision	Recall	F-measure
EAST [30]	64.13	53.22	56.44
Ours detection	75.08	52.28	61.64
Ours with optical flow	68.49	55.69	61.43
Ours with Conv3D	78.94	54.76	64.66
Ours with ConvLSTM	79.97	55.21	65.32

Table I
COMPARISON OF VIDEO TEXT DETECTION PERFORMANCES ON ICDAR 2013 DATASET [6].

Method	MOTP	MOTA
Zuo et al. [31]	73.07	56.37
Pei et al. [11]	73.07	57.71
Ours with appearance descriptor	75.90	72.79
Ours with geometry descriptor	76.66	74.04
matching AGD with AGD	74.70	75.62
matching AGD with EAGD	75.72	81.31

Table II
COMPARISON OF VIDEO TEXT TRACKING PERFORMANCES ON MINETTO DATASET [8].

D. Evaluation of Video Text Tracking

Tab. II shows the influence of different types of tracking descriptors. We adopt the widely-used CLEAR MOT metrics [1], including MOTP (Multi-Object Tracking Precision) and MOTA (Multi-Object Tracking Accuracy) as text tracking evaluation metrics. The MOTP is the mean error of estimated positions for matched pairs of all frames, while the MOTA accounts for errors made by trackers, i.e., false positives, misses and mismatches. At first, the appearance and geometry descriptors are studied respectively. Compared with appearance feature, the performance of using geometry features is 1.25% ahead. However, there is a significant gain about 7% when they are concatenated. As the appearance and geometry features are able to capture different local information in the tracking process, combining them together is more robust. Then, in order to evaluate the effectiveness of the GRU in our proposed estimated descriptor, we try to associate *AGD* of the adjacent frames directly, instead of matching the current frame *AGD* with *EAGD* obtained from the previous frame. Consequently, this matching method results in nearly 6% loss on MOTA. This highlights that the temporal changing information captured by the GRU also plays an important role in the video text tracking.

E. Comparison with State-of-the-Art Video Text Detection Methods

In this section, we compare our method with state-of-the-art methods on three public video text datasets. As shown in Tab. III, we have summarized various video text detection methods, and our method outperforms all the other methods by more than 2% on f-measure.

The precision of video text detection is improved by a large margin along with high recall ratio with our method.

This is mainly due to our end-to-end joint training and the introduction of long-term memory mechanisms. Specifically, in our end-to-end inference process, the detection results can be refined along with the updating of trajectories, and the long time memory can also inhibit some negative samples, resulting in more accurate results denoted by “Our end-to-end detection with online tracking” in Tab. III.

For training, except for the configuration introduced in Sec. IV-B, the YVT model is finetuned on its own training set as it contains many web videos. For testing, the longer sides of input images in ICDAR 2013 are resized to 1280. While, input images not resized for Minetto and YVT during evaluation. Besides, multi-scale testing is not conducted since it is too slow and unpractical for video especially.

V. CONCLUSIONS AND FUTURE WORK

In this work, we present an end-to-end framework for video text detection with online tracking according to the characteristics of video scene text. The proposed *AGD* and *EAGD* are employed to convert the long-term multiple targets changing process to a trainable model. By sharing convolutional features, the text tracking branch is nearly cost-free. In the inference phase, the text detection results are obtained along with the trajectory online generation. Experiments on video text datasets show that our method significantly outperforms previous methods in both detection and tracking. However, there are still some drawbacks: the trajectory generation is not incorporated into the training process and semantic information has not been exploited. In the future, we plan to add video text detection, tracking, and recognition to an end-to-end framework.

ACKNOWLEDGMENTS

This work is jointly supported by National Key Research and Development Program of China (2016YFB1001000), National Natural Science Foundation of China (61525306, 61633021, 61721004, 61420106015, 61806194), Capital Science and Technology Leading Talent Training Project (Z181100006318030), Beijing Science and Technology Project (Z181100008918010) and CAS-AIR.

REFERENCES

- [1] K. Bernardin and R. Stiefelhagen. Evaluating multiple object tracking performance: the clear mot metrics. *Journal on Image and Video Processing*, 2008:1, 2008.
- [2] B. Epshtein, E. Ofek, and Y. Wexler. Detecting text in natural scenes with stroke width transform. In *CVPR*, pages 2963–2970, 2010.
- [3] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [4] W. He, X.-Y. Zhang, F. Yin, and C.-L. Liu. Deep direct regression for multi-oriented scene text detection. In *ICCV*, pages 745–753, 2017.
- [5] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE TPAMI*, 35(1):221–231, 2013.

Method	ICDAR 2013			Minetto dataset			YVT		
	P	R	F	P	R	F	P	R	F
Epshtein et al. [2]	39.80	32.53	35.94	—	—	—	68.00	76.00	72.00
Zhao et al. [29]	47.02	46.30	46.65	—	—	—	34.00	41.00	37.00
Minetto et al. [8]	—	—	—	61.00	69.00	63.00	—	—	—
Yin et al. [26]	48.62	54.73	51.56	—	—	—	—	—	—
Moslen et al. [9]	—	—	—	—	—	—	79.00	72.00	75.00
Wu et al. [21]	63.00	68.00	65.00	—	—	—	81.00	73.00	77.00
Zuo et al. [31]	—	—	—	84.00	68.00	75.00	—	—	—
Khare et al. [7]	57.91	55.90	51.70	—	—	—	—	—	—
Shivakumara et al. [14]	61.00	57.00	59.00	—	—	—	79.00	73.00	76.00
Pei et al. [11]	—	—	—	89.00	84.00	86.00	—	—	—
Wang et al. [18]	58.34	51.74	54.45	88.80	87.53	88.14	—	—	—
Wang et al. [20]	71.90	58.67	62.65	83.03	84.22	83.30	—	—	—
Our two-stage detection	79.97	55.21	65.32	82.90	80.66	81.77	70.15	64.52	67.22
Our two-stage detection with online tracking	81.36	55.41	65.93	83.92	79.89	81.85	71.64	63.98	67.59
Our end-to-end detection	80.10	57.07	66.65	86.69	91.17	88.87	78.51	71.54	74.86
Our end-to-end detection with online tracking	82.36	56.36	66.92	91.27	89.38	90.32	89.12	71.03	79.05

Table III
COMPARISON OF VIDEO TEXT DETECTION PERFORMANCES ON SEVERAL PUBIC DATASETS.

- [6] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazan, and L. P. De Las Heras. Icdar 2013 robust reading competition. In *ICDAR*, pages 1484–1493. IEEE, 2013.
- [7] V. Khare, P. Shivakumara, R. Paramesran, and M. Blumenstein. Arbitrarily-oriented multi-lingual text detection in video. *Multimedia Tools and Applications*, 76(15):16625–16655, 2017.
- [8] R. Minetto, N. Thome, M. Cord, N. J. Leite, and J. Stolfi. Snooptrack: Text detection and tracking for outdoor videos. In *ICIP*, pages 505–508. IEEE, 2011.
- [9] A. Mosleh, N. Bouguila, and A. B. Hamza. Automatic inpainting scheme for video text detection and removal. *IEEE TIP*, 22(11):4460–4472, 2013.
- [10] P. X. Nguyen, K. Wang, and S. Belongie. Video text detection and recognition: Dataset and benchmark. In *WACV*, pages 776–783. IEEE, 2014.
- [11] W.-Y. Pei, C. Yang, L.-Y. Meng, J.-B. Hou, S. Tian, and X.-C. Yin. Scene video text tracking with graph matching. *IEEE Access*, 6:19419–19426, 2018.
- [12] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCA*, pages 234–241. Springer, 2015.
- [13] B. Shi, X. Bai, and S. Belongie. Detecting oriented text in natural images by linking segments. In *CVPR*, pages 3482–3490. IEEE, 2017.
- [14] P. Shivakumara, L. Wu, T. Lu, C. L. Tan, M. Blumenstein, and B. S. Anami. Fractals based multi-oriented text detection system for recognition in mobile video images. *Pattern Recognition*, 68:158–174, 2017.
- [15] Y. Sun, C. Zhang, J. Liu, J. Han, and E. Ding. Textnet: Irregular text reading from images with an end-to-end trainable network. In *ACCV*, 2018.
- [16] S. Tian, W.-Y. Pei, Z.-Y. Zuo, and X.-C. Yin. Scene text detection in video by learning locally and globally. In *IJCAI*, 2016.
- [17] L. Wang, Y. Wang, S. Shan, and F. Su. Scene text detection and tracking in video with background cues. In *ICMR*, 2018.
- [18] L. Wang, Y. Wang, S. Shan, and F. Su. Scene text detection and tracking in video with background cues. In *ICMR*, pages 160–168. ACM, 2018.
- [19] X. Wang, Y. Jiang, S. Yang, X. Zhu, W. Li, P. Fu, H. Wang, and Z. Luo. End-to-end scene text recognition in videos based on multi frame tracking. In *ICDAR*, volume 1, pages 1255–1260. IEEE, 2017.
- [20] Y. Wang, L. Wang, and F. Su. A robust approach for scene text detection and tracking in video. In *PCM*, pages 303–314. Springer, 2018.
- [21] L. Wu, P. Shivakumara, T. Lu, and C. L. Tan. A new technique for multi-oriented scene text line detection and tracking in video. *IEEE Trans. Multimedia*, 17(8):1137–1152, 2015.
- [22] Y. Wu and P. Natarajan. Self-organized text detection with minimal post-processing via border learning. In *ICCV*, 2017.
- [23] S. Xingjian, Z. Chen, H. Wang, D.-Y. Yeung, W.-K. Wong, and W.-c. Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. In *NeurIPS*, pages 802–810, 2015.
- [24] X.-H. Yang, W. He, F. Yin, and C.-L. Liu. A unified video text detection method with network flow. In *ICDAR*, volume 1, pages 331–336. IEEE, 2017.
- [25] Q. Ye and D. S. Doermann. Text detection and recognition in imagery: A survey. *IEEE TPAMI*, 37:1480–1500, 2015.
- [26] X.-C. Yin, X. Yin, K. Huang, and H.-W. Hao. Robust text detection in natural scene images. *IEEE TPAMI*, (1):1, 2013.
- [27] X.-C. Yin, Z.-Y. Zuo, S. Tian, and C.-L. Liu. Text detection, tracking and recognition in video: A comprehensive survey. *IEEE TIP*, 25:2752–2773, 2016.
- [28] Z. Zhang, C. Zhang, W. Shen, C. Yao, W. Liu, and X. Bai. Multi-oriented text detection with fully convolutional networks. In *CVPR*, 2016.
- [29] X. Zhao, K.-H. Lin, Y. Fu, Y. Hu, Y. Liu, and T. S. Huang. Text from corners: a novel approach to detect text and caption in videos. *IEEE TIP*, 20(3):790–799, 2011.
- [30] X. Zhou, C. Yao, H. Wen, Y. Wang, S. Zhou, W. He, and J. Liang. East: an efficient and accurate scene text detector. In *CVPR*, pages 2642–2651, 2017.
- [31] Z.-Y. Zuo, S. Tian, W.-y. Pei, and X.-C. Yin. Multi-strategy tracking based text detection in scene videos. In *ICDAR*, pages 66–70. IEEE, 2015.