

---

# MMVText: A Large-Scale, Multidimensional Multilingual Dataset for Video Text Spotting

---

Anonymous Author(s)

Affiliation

email

## Abstract

1 Video text spotting is crucial for numerous real application scenarios, but most  
2 existing video text reading benchmarks are challenging to evaluate the performance  
3 of advanced deep learning algorithms due to the limited amount of training data  
4 and tedious scenarios. To address this issue, we introduce a new large-scale bench-  
5 mark dataset named **Multidimensional Multilingual Video Text (MMVText)**, the  
6 first large-scale and multilingual benchmark for video text spotting in a variety of  
7 scenarios. There are mainly three features for MMVText. Firstly, we provide **510**  
8 videos with more than **1,000,000** frame images, four times larger than the existing  
9 largest dataset for text in videos. Secondly, our dataset covers 30 open categories  
10 with a wide selection of various scenarios, *e.g.*, *life vlog*, *sports news*, *automatic*  
11 *drive*, *cartoon*, *etc.* Besides, caption text and scene text are separately tagged for the  
12 two different representational meanings in the video. The former represents more  
13 theme information, and the latter is the scene information. Thirdly, the MMVText  
14 provides multilingual text annotation to promote multiple cultures live and commu-  
15 nication. In the end, a comprehensive experimental result and analysis concerning  
16 text detection, recognition, tracking, and end-to-end spotting on MMVText are pro-  
17 vided. We also discuss the potentials of using MMVText for other video-and-text  
18 research. The dataset and code can be found at [github.com/wei jiawu/MMVText](https://github.com/wei jiawu/MMVText).

## 19 1 Introduction

20 Text reading [18, 12] has received increasing attention due to its numerous applications in computer  
21 vision, *e.g.*, document analysis, image-based translation, image retrieval [29, 23], *etc.* With the advent  
22 of deep learning and abundance in digital data, reading text from images has made extraordinary  
23 progress in recent years with a lot of great public datasets [8, 13, 5] and algorithms [35, 44, 19, 17]. By  
24 contrast, video text spotting almost remains at a standstill for the lack of large-scale multidimensional  
25 practical datasets, which limited numerous applications of video text, *e.g.*, video understanding [32],  
26 video retrieval [7], video text translation, and license plate recognition [1], *etc.*

27 Most existing algorithms [44, 35, 16] in text detection and recognition deal with only static frames.  
28 Therefore, one intuitive drawback of these approaches is that they do not necessarily work well  
29 in the video domain, while at the same time they do not take advantage of the extra information  
30 present in the video (*e.g.*, tracking already detected regions). Moreover, the quality of the image  
31 is generally worse than static images, due to motion blur and out of focus issues, while video  
32 compression might create further artefacts. Due to these interferences, methods designed for still  
33 images, may fail to obtain reliable detection and recognition results when applied to a video frame.  
34 Most importantly, these methods based on image-level can not obtain text tracking information in  
35 video. However, spatio-temporal information in video is vital for a number of real-world applications.  
36 For example, video understanding and video caption translation all require temporal text information  
37 in sequential frames. There have been a few previous works [40, 38] in the community for attempting

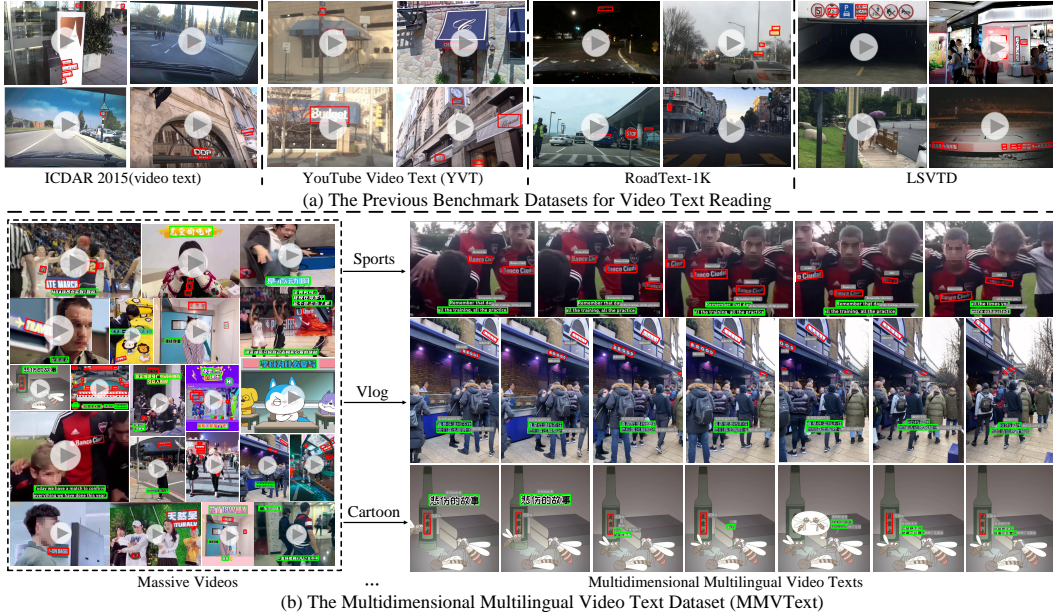


Figure 1: **Example Sequences and Annotations.** Unlike the previous benchmarks, our MMVText contains a wide variety of scenarios and multi-languages. The caption text and scene text are separately tagged for the two different representational meanings.

38 to develop text reading in videos, and there is a handful of datasets [25, 14] that support the research.  
 39 ICDAR2015 (Text in Videos) [13], as one of the common datasets, was introduced during the ICDAR  
 40 Robust Reading Competition in 2015 and mainly includes a training set of 25 videos (13,450 frames  
 41 in total) and a test set of 24 videos (14,374 frames in total). The videos were categorized into  
 42 seven scenarios: walking outdoors, searching for a shop in a shopping street, browsing products in  
 43 a supermarket, etc. YouTube Video Text (YVT) [25] dataset harvested from YouTube, contains 30  
 44 videos (13,500 frames in total), 15 for training, and 15 for testing. The text content in the dataset  
 45 can be divided into two categories, overlay text (*e.g.*, captions, songs title, logos) and scene text (*e.g.*,  
 46 street signs, business signs, words on shirt). RoadText-1K [26] are sampled from BDD100K [42],  
 47 includes 700 videos (210,000 frames) for training and 300 videos for testing. The texts in the  
 48 dataset are all obtained from driving videos and match for driver assistance and self-driving systems.  
 49 LSVTD [4] includes 100 text videos, 13 indoor (*e.g.*, bookstore, shopping mall) and 9 outdoor (*e.g.*,  
 50 highway, city road) scenarios. The existing video text benchmarks are limited by the amount of  
 51 training data (less than 300k frames) and tedious data scenarios, as shown in Figure. 1 (a). There  
 52 are only a few outdoor scene text videos with 13k frames in ICDAR2015 (video text). Similar situation  
 53 for YVT, RoadText-1k and LSVTD, the training set is limited and the dataset scenarios are tedious.  
 54 This makes it difficult to evaluate the effectiveness of more advanced deep learning models.

55 To address this issue, our work intends to contribute a large-scale, multidimensional multilingual  
 56 benchmark dataset (MMVText) to the community for developing and testing video text reading  
 57 systems that can fare in a realistic setting. Our dataset has several advantages. Firstly, the large  
 58 training set (*i.e.*, 1,010,848 video frames) enables the development of deep design specific for video  
 59 text spotting. Secondly, MMVText is a multilingual multidimensional dataset. Abundant videos  
 60 in various scenarios (*e.g.*, driving, street view, news reports, cartoon) are provided for representing  
 61 real-world scenarios, as shown in Figure. 1 (b). Thirdly, caption and scene text are separately tagged  
 62 for the two different representational meanings in the video. This is in favor of other tasks, such as  
 63 video understanding and video retrieval. The main contributions of this work are three folds:

- 64 • We propose a large-scale, multidimensional, and multilingual video text reading benchmark  
 65 named MMVText. The proposed dataset span various video scenarios, text types, multi-stage  
 66 tasks and is four times the existing largest dataset.
- 67 • Caption text and scene text are separately tagged for the two different representational  
 68 meanings in the video. This favors other tasks, such as video understanding, video retrieval,  
 69 and video text translation.

70 • We evaluate the current state-of-the-art techniques for scene text detection, recognition, text  
71 tracking, and end-to-end video text spotting. Besides, a thorough analysis of performance  
72 on this dataset is provided.

## 73 2 Related Work

### 74 2.1 End-to-End Text Reading

75 For image-level text reading, various methods [15, 9, 19] based on deep learning have been proposed  
76 and have improved the performance considerably. Li et al. [15] proposed the first end-to-end trainable  
77 scene text spotting method. The method successfully uses a RoI Pooling [27] to joint detection  
78 and recognition features via a two-stage framework. Liao et al. [19] propose a Mask TextSpotter  
79 which subtly refines Mask R-CNN and uses character-level supervision to detect and recognize  
80 characters simultaneously. However, these methods based on the static image can not obtain temporal  
81 information in the video, which is essential for some downstream tasks such as video understanding.

82 Compared to text reading in a static image, video text spotting methods are rare. Yin et al. [41]  
83 provides a detailed survey, summarizes text detection, tracking and recognition methods in video  
84 and their challenges. Wang et al. [36] introduced an end-to-end text recognition method to detect  
85 and recognize text in each frame of the input video. Multi-frame text tracking is employed through  
86 associations of texts in the current frame and several previous frames to obtain final results. Cheng  
87 et al. [4] propose a video text spotting framework by only recognizing the localized text one-  
88 time. To promote text reading in the video, we attempt to establish a standardized evaluation and  
89 benchmark (MMVText), covering various open scenarios and multilingual text annotation.

### 90 2.2 Text Reading Datasets for Static Images

91 The various and practical benchmark datasets [13, 33, 14, 5] contribute to the huge success of  
92 scene text detection and recognition at the image level. ICDAR2015 [13] was provided from the  
93 ICDAR2015 Robust Reading Competition, which is commonly used for oriented scene text detection  
94 and spotting. Google glasses capture these images without taking care of position, so text in the  
95 scene can be in arbitrary orientations. ICDAR2017MLT [24] is a large-scale multilingual text dataset,  
96 which is composed of complete scene images which come from 9 languages, and text regions in this  
97 dataset can be in arbitrary orientations, so it is more diverse and challenging. ICDAR2013 [14] is  
98 a dataset proposed in the ICDAR 2013 Robust Reading Competition, which focuses on horizontal  
99 text detection and recognition in natural images. The COCO-Text dataset [33] is currently the largest  
100 dataset for scene text detection and recognition. It contains 50,000+ images for training and testing.  
101 The COCO-Text dataset is very challenging since the text in this dataset is in arbitrary orientation.

### 102 2.3 Text Reading Datasets for Videos

103 The development of video text spotting is limited in recent years due to the lack of efficient data  
104 sets. ICDAR 2015 Video [14] consists of 28 videos lasting from 10 seconds to 1 minute in indoors  
105 or outdoors scenarios. Limited videos (*i.e.*, 13 videos) used for training and 15 for testing. Minetto  
106 Dataset [22] consists of 5 videos in outdoor scenes. The frame size is 640 x 480 and all videos  
107 are used for testing. YVT [25] contains 30 videos, 15 for training and 15 for testing. Different  
108 from the above two datasets, it contains web videos except for scene videos. USTB-VidTEXT [40]  
109 with only five videos mostly contain born-digital text (captions and subtitles) sourced from Youtube.  
110 RoadText-1K provides a driving videos dataset with 1000 videos. The 10-second long video clips in  
111 the dataset are sampled from BDD100K [42]. As shown in Table. 1, the existing datasets contain a  
112 limited training set and tedium video scenarios. To promote the development of video text reading  
113 and extension of application based on video text, we create a large scale, multidimensional and  
114 multilingual dataset, and attempt to provide a more reasonable metric.

## 115 3 MMVText Benchmark

116 This section firstly introduces the collection and annotation of MMVText and provides a comprehen-  
117 sive analysis and comparison. And then, the related tasks and corresponding metrics are described.  
118 Finally, we discuss the link to application scenarios and potential impacts.



Figure 2: **Distributions of MMVText.** (a) Chinese caption and English scene text. (b) Only Chinese caption. (c) Multilingual caption and English scene text. (d) The benchmark dataset covers a wide and open range of life scenes (30 categories) with multilingual texts. Caption text (blue box) and scene text (red box) are distinguished in MMVText, which is favorable for downstream tasks.

### 119 3.1 Data Collection and Annotation

120 **Data Collection.** To obtain abundant and various text videos, we first start by acquiring a large list  
 121 of text videos class using *KuaiShou*<sup>1</sup> - an online resource that contains billions of videos with various  
 122 scene text from cartoon movies to human relation. Then, we choose 30 live video categories, *i.e.*, ,  
 123 *E-commerce, Game, Home, Fashion, and Technology*, as shown in Figure. 2 (d). With each raw video  
 124 category, we first choose the video clips with text, then make two rounds of screening to remove  
 125 the ordinary videos. As a result, we obtain 512 videos with 1, 010, 848 video frames, as shown in  
 126 Table 1. Finally, to fair evaluation, we divide the dataset into two parts: the training set with 641, 049  
 127 frames from 331 videos, and the testing set with 369, 799 frames from 179 videos. As shown in  
 128 Figure 2 (a), different from the existing data sets, which only focus on one type of video text and the  
 129 video scene is limited, our dataset not only care about scene text reading in the real world, but also  
 130 focus on caption texts in the video. For the most part, caption text represents more global information  
 131 than scene text, which is quite favorable for some downstream tasks, *e.g.*, *video understanding, video*  
 132 *caption translation*. Therefore, the MMVText can cover a wider and open range of life scenes, and  
 133 contains various text with a more comprehensive description of the video.

134 **Data Annotation.** We invite a professional annotation team to label each video text with four kinds of  
 135 description information: the bounding box describing the location information, judging the tracking  
 136 identification (ID) of the same text instance, identifying the content of the text information, and  
 137 distinguishing the category label of the caption or scene text. To save the annotation cost, we first  
 138 sample the videos, annotate each sampled video frame at an instance level, and then transform the  
 139 annotation information from the sampled video frame to the unlabeled video frame by interpolation.  
 140 *For video sampling*, we use uniform sampling with a sampling frequency of 7 to sample all the videos  
 141 in the dataset, and obtain the sampled video frame set. *For sampling video frame annotation*, each text  
 142 instance is labeled in the same quadrilateral way as in the ICDAR 2015 incidental text dataset [45].  
 143 In addition, the text instance also will be marked with two description information: the category of  
 144 the caption or scene, and the recognition content. After the spatial location, content, and category of  
 145 the video text are determined, the annotator will determine the tracking ID by browsing the length  
 146 of the same video text in the continuous sampling video frames. We also invited other text-related  
 147 people to conduct two rounds of cross-checking to ensure the annotation quality. *For video frame*  
 148 *recovery*, each text instance is marked with tracking ID and recognition content, so we can judge  
 149 whether different texts in adjacent sampling frames are the same text. After determining the same text  
 150 instance, we first determine whether the text annotation of the sampled video frame is the starting and  
 151 end frame of the text instance. If not, we look forward and backward for the starting and end position  
 152 of the text instance and label it. Then we use the linear interpolation way to calculate the position of  
 153 the text object in the middle of the unmarked video frame, and give tracking ID, recognition content,

<sup>1</sup><https://www.kuaishou.com/en>

Table 1: **Statistical Comparison.** Comparisons between MMVText and existing datasets for caption and scene text in videos. *D*, *T*, and *S* denotes the Detection, Tracking, and Spotting respectively.

Dataset	Category	MLingual	Scenario	Videos	Frames	Texts	Task
AcTiV-D [43]	Caption	-	News video	8	1,843	5,133	D
UCAS-STLData [3]	Caption	-	Teleplay video	3	57,070	41,195	D
USTB-VidTEXT [40]	Caption	-	Web video	5	27,670	41,932	D&S
YVT [25]	Scene	-	Incidental	30	13,500	16,620	D&T&S
ICDAR 2015 VT [45]	Scene	-	Incidental	51	27,824	143,588	D&T&S
LSVTD [4]	Scene	✓	Incidental	100	66,700	569,300	D&T&S
RoadText-1K [26]	Scene	-	Driving	1000	300,000	1,280,613	D&T&S
MMVText (ours)	Both	✓	Open	510	<b>1,010,848</b>	<b>4,513,525</b>	D&T&S

154 and category. After all the video annotations are restored, we carry out another round of double  
 155 detection correction. As a labor-intensive job, the labeling process takes 30 men in two months, *i.e.*,  
 156 20,160 man-hours, to complete about 200,000 sampled video frame annotations.

### 157 3.2 Dataset Analysis

158 **Statistic Comparison.** The qualitative and statistic comparison between the established MMVText  
 159 and other datasets are visualized in Figure. 1, and summarized in Table. 1. *Category* denotes the  
 160 category of the text type in the corresponding dataset. *MLingual* denotes whether the dataset contains  
 161 multiple language texts. *Scenario* denotes the scene range of the video. *Videos*, *Frames*, *Texts*  
 162 represents the number of videos, video frames, video texts in the dataset, respectively. *Task* denotes  
 163 which tasks the dataset supports. **Caption Text and Scene Text.** For comprehensive evaluation and  
 164 research, we not only expand the scale of the dataset (*i.e.*, , the number of videos, video frame, and  
 165 video text), and label the spatial quadrilateral position, recognition content, and tracking ID, but also  
 166 additionally collect and annotate the category of caption or scene for each text instance. As shown  
 167 in Figure. 2 (a), in a video, different types of text instances may exist simultaneously, and they are  
 168 helpful to understand videos synergistically. Concretely, caption text can directly show the dialogue  
 169 between people in video scenes and represent the time or topic of the video scenes, scene text can  
 170 unambiguously define the object and can identify important localization and road paths in video  
 171 scenes. Besides, nowadays, caption text frequently exists in all kinds of life scenarios video. Even  
 172 for some videos, without any scene texts, there is a lot of caption text, as shown in Figure. 2 (b). To  
 173 favor downstream tasks (*e.g.*, video text translation, video understanding, and video retrieval), we  
 174 also provide multilingual text annotations, as shown in Figure. 2 (c).

175 To provide the community with unified text-level quantitative descriptions, and facilitate controlled  
 176 evaluation for different approaches, we will compare our dataset with caption or scene text datasets  
 177 from four aspects, *i.e.*, text description, video scene, dataset size, and supported tasks. *For text*  
 178 *description attribute (i.e., Category, MLingual)*, our MMVText contains both types (caption and  
 179 scene) of video text and has multi-language features, which obviously has more extensive description  
 180 ability than caption or scene text dataset. *For video scene attribute (i.e., Scenario)*, the caption  
 181 text datasets choose videos with certain professional purposes (*e.g.*, news reports, TV dramas, and  
 182 documentaries), which shows that the scenes they cover are relatively limited. And the existing  
 183 scene text datasets often choose some video scenes captured by mobile shooting, and the number of  
 184 collectors is small, the range of captured scenes is also limited. However, the videos in our dataset  
 185 are from videos uploaded voluntarily by all kinds of users. Therefore, the proposed MMVText  
 186 covers various scenarios, but it also brings significant challenges to researchers. *For the size of*  
 187 *the dataset (i.e., Videos, Frames, Texts)*, we can find that our MMVText has advantages over the  
 188 superimposed caption text dataset and the scene text dataset in the indicators of videos, frames, and  
 189 texts. The number of videos in RoadText-1K is more than ours (1,000 vs. 510), but the number of  
 190 frames in RoadText-1K is far less than ours (300,000 vs. 1,010,848), which imply that the average  
 191 video length of RoadText-1K is much shorter than ours (300 vs. 1,982). *For the supported tasks*,  
 192 the proposed MMVText supports four common video text tasks: detection, recognition, video text  
 193 tracking, end to end video text spotting. The focus and application scenarios of each task is entirely  
 194 different. For example, detection task used in the static image focus on localization performance,  
 195 paving the way for recognition task, which apply to license plate recognition. End to end video text  
 196 spotting task focuses on recognition and tracking performance, which apply to video understanding

197 and video retrieval. In conclusion, the high efficiency of MMVText for evaluating advanced deep  
 198 learning methods is very favorable for promoting various text reading applications in real life.

### 199 3.3 MMVText Tasks and Metrics

200 Standardized benchmark metrics are crucial as same as the dataset for the majority of computer vision  
 201 applications, and we attempt to provide a reasonable evaluation for video text reading methods. The  
 202 proposed MMVText mainly includes two tasks: (1) Video Text Tracking, aimed at describing text  
 203 location information in continuous frames. (2) End to End Text Spotting in Videos, to understand  
 204 text and track multiple frames. For the detection and recognition task, we also provide corresponding  
 205 experimental results and analysis in the experiments.

206 Most tracking tasks all use the *MOT* metrics [2], which was launched to establish a standardized  
 207 evaluation of multiple object tracking methods. The same case for video text tracking, the ICDAR2013  
 208 Robust Reading Challenge [14] for video text reading adopts *MOTP* (Multiple Object Tracking  
 209 Precision) and *MOTA* (Multiple Object Tracking Accuracy) as the metrics. Following the previous  
 210 works [14, 26], MMVText evaluates text tracking methods in video and compares their performance  
 211 with the *MOTA* and *MOTP*. Besides,  $ID_{F1}$  as the new metrics for tracking is presented from some  
 212 tracking works [6, 28] in recent year.  $ID_{F1}$  is the ratio of correctly identified detections over the  
 213 average number of ground-truth and computed detections. And the metric is more reasonable to  
 214 evaluate ID switches in some cases. We also evaluate the metrics in MMVText by:

$$ID_{F1} = \frac{2ID_{tp}}{2ID_{tp} + ID_{fp} + ID_{fn}}, \quad (1)$$

215 where  $ID_{tp}$ ,  $ID_{fp}$  and  $ID_{fn}$  refer to true positive, false positive and false negative of matching ID.  
 216 Besides, the ID metric [6] also includes *MT* (Mostly Tracked) Number of objects tracked for at least  
 217 80 percent of lifespan, *ML* (Mostly Lost) Number of objects tracked less than 20 percent of lifespan.

218 In Task2 (End to End Text Spotting in Videos), the objective of this task is to recognize words in the  
 219 video as well as localize them in terms of time and space. And we argue that the final recognition  
 220 result is more important than text localization in videos. Thus, we modify the  $ID_{F1}$  to  $TID_{F1}$ , which  
 221 focuses on text instance ID tracking and recognition results that be required by many downstream  
 222 tasks. More specifically,

$$TID_{tp} = \sum_h \sum_t m(h, o, \Delta_t, \Delta_s, \Delta_r), \quad (2)$$

$$TID_{F1} = \frac{2TID_{tp}}{2TID_{tp} + TID_{fp} + TID_{fn}}, \quad (3)$$

223 where  $\Delta_t$ ,  $\Delta_s$  and  $\Delta_r$  refer to time matching, space location matching and recognition result  
 224 matching. And  $h$  and  $o$  denote hypothesis and true text trajectory with recognition result. The match  
 225 of  $h$  and  $o$  is a true positives of text ID (*i.e.*,  $TID_{tp}$ ) when these conditions (*i.e.*,  $\Delta_t$ ,  $\Delta_s$  and  $\Delta_r$ )  
 226 are met. Similarly, false positive (*i.e.*,  $TID_{fp}$ ) and false negative (*i.e.*,  $TID_{fn}$ ) of text ID can be  
 227 obtained for  $TID_{F1}$  calculation. More details concerning metrics in supparmentary material.

### 228 3.4 Methods

229 Text detection and recognition in the static image have made tremendous progress, and abundant  
 230 great work [35, 44, 30] be proposed. By contrast, the counterparts in video text reading are rare and  
 231 lack quality open-source algorithms. Therefore, we adopt various mature techniques in the static  
 232 image to better evaluate the efficiency of MMVText.

233 **Detection.** The deep learning-based text detection methods can be roughly divided into two cate-  
 234 gories: regression-based method and segmentation-based method. EAST [44] as one of the popular  
 235 regression-based methods is used to test our MMVtext. The method adopts FCNs to predict shrink-  
 236 able text score maps, rotation angles and perform per-pixel regression, followed by a post-processing  
 237 NMS. For segmentation based methods, we adopt PSENet [35] and DB [16] to evaluate our MMVtext.  
 238 PSENet [35] generates various scales of shrinked text segmentation maps, then gradually expands  
 239 kernels to generate the final instance segmentation map. Similarly, DB [16] utilizes the shrinked

240 text segmentation maps and differentiable binarization to detect text instances. **Recognition.** Recent  
241 methods mainly use two techniques to train the scene text recognition model, namely Connectionist  
242 Temporal Classification (CTC) and attention mechanism. In CTC-based methods, CRNN [30] as  
243 the representation, which introduced CTC decoder into scene text recognition with a Bidirectional  
244 Long Short-Term Memory (BiLSTM) to model the feature sequence. In Attention-based methods,  
245 RARE [31] firstly normalizes the input text image using the Spatial Transformer Network (STN [11]),  
246 then utilizes CNN to extract feature and captures the contextual information within a sequence of  
247 characters. Finally, it estimates the output character sequence from the identified features with the  
248 attention module.

249 **Text Tracking Trajectory Generation.** With text detection and recognition in a static image, we  
250 only obtain text localization and recognition information without temporal information, which are  
251 insufficient for video spotting evaluation (*e.g.*,  $TID_{F1}$ ,  $MOTA$  and  $MOTP$ ). The work [36] based  
252 on multi-frame tracking provides a method to track text instances temporally based on attributes of  
253 the text objects in multiple frames. Following the work [36], we link and match text objects in the  
254 current frame and several frames by IOU and edit distance of text.

### 255 3.5 Link to Real Applications

256 Text understanding in static images has numerous application scenarios: (1) Automatic data entry.  
257 SF-Express<sup>2</sup> utilizes OCR techniques to accelerate the data entry process. NEBO<sup>3</sup> performs instant  
258 transcription as the user writes down notes. (2) Autonomous vehicle [21, 20]. Text-embedded  
259 panels carry important information, *e.g.*, geo-location, current traffic condition, navigation, and etc.  
260 Similarly, there are many application demands for video text understanding across various industries  
261 and in our daily lives. We list the most outstanding ones that significantly impact, improving our  
262 productivity and life quality. Firstly, automatically describing video with natural language [39, 37]  
263 can bridge video and language. Secondly, video text automatic translation<sup>4</sup> can be extremely helpful  
264 as people travel, and help video-sharing websites<sup>5</sup> to cut down language barriers. More details and  
265 analyses for application scenarios concerning MMVText in the supplementary material.

## 266 4 Experimental

267 In this section, we conduct experiments on our MMVText to demonstrate the effectiveness of the  
268 proposed benchmark. Note that we denote Ground Truth of ID tracking in all the experiments, Mostly  
269 Tracked and Mostly Lost as ‘GT’, ‘MT’ and ‘ML’, respectively.

### 270 4.1 Implementation Details

271 All of the experiments use the same strategy: (1) Training detector and recognizer with MMVText.  
272 (2) Matching text objects with corresponding text tracking trajectory id. *Detection:* without pretrained  
273 model, we train detectors directly with training set (*i.e.*, 641,049 frame images) of MMVText.  
274 *Recognition:* the network is pre-trained on the *chinese ocr*<sup>6</sup> and MJSynth [10], and further fine-tuned  
275 on our MMVText. All of our experiments are conducted on 8 V100 GPUs. PSENet [35], EAST [44]  
276 and DB [16] are adopted as the base detectors because of their popularity. CRNN [30] and RARE [31]  
277 as the base text recognizers to evaluate our MMVText. In the PSENet, EAST, DB, CRNN and RARE  
278 experiments, all settings follow the original reports.

### 279 4.2 Attribute Experiments Analysis

280 **Text Tracking in Different Scenarios.** Figure. 3 (a) gives the tracking performance  $ID_{F1}$  of  
281 EAST [44] in different scenarios of MMVText. The model achieves the best performance with a  
282  $ID_{F1}$  of 57% in cartoon videos, since the conspicuous text instances and simple background are

<sup>2</sup><https://www.sf-express.com/cn/sc/>

<sup>3</sup><https://www.myscript.com/nebo/>

<sup>4</sup><https://translate.google.com/intl/en/about/>

<sup>5</sup><https://www.youtube.com/>

<sup>6</sup>[https://github.com/YCG09/chinese\\_ocr](https://github.com/YCG09/chinese_ocr)

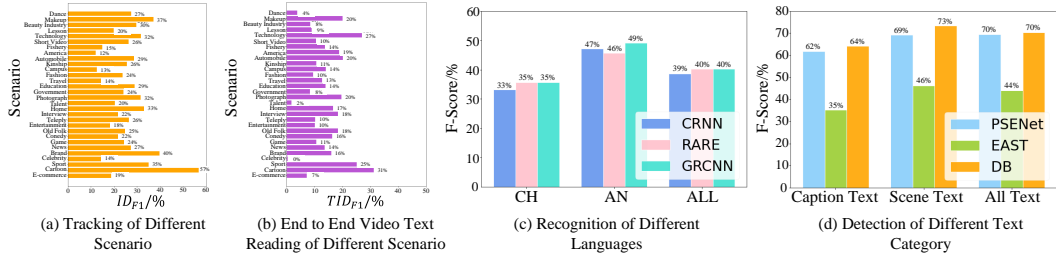


Figure 3: **Attribute Experiments of MMVText.** (a) Tracking performance (*i.e.*,  $ID_{F1}$ ) with EAST [44] in different scenarios. (b) End to end video text spotting performance (*i.e.*,  $TID_{F1}$ ) with PSENet [35] and CRNN [30] in different scenarios. (c) Recognition accuracy of different models in different languages. (d) Detection performance of different model in caption or scene text. ‘CH’, ‘AN’ and ‘ALL’ refer to ‘Chinese Characters’, ‘Alphanumeric Characters’ and ‘All Characters’.

Table 2: **Detection and Recognition Performance on MMVText.** Frame level text Detection and Recognition performance of existing models on MMVText. ‘CH’, ‘AN’ and ‘ALL’ refer to ‘Chinese Characters’, ‘Alphanumeric Characters’ and ‘All Characters’.

Detection Performance/%				Recognition Performance/%						
Method	Precision	Recall	F-score	Method	Pretrained			Fine tuned		
					CH	AN	ALL	CH	AN	ALL
EAST [44]	52.2	38.1	44.1	CRNN [30]	26.0	32.1	23.2	33.2	47.1	38.6
PSENet [35]	74.3	65.2	69.5	RARE [31]	25.2	34.2	23.5	35.6	45.7	40.2
DB [16]	77.2	64.5	70.3	GRCNN [34]	23.1	39.8	26.7	35.6	49.2	40.3

283 designed in cartoon videos. By comparison, several scene categories obtain extremely dissatisfied  
 284 performance due to complex background and various text appearance, such as *Campus* and *Travel*.

285 **End to End Text Spotting in Different Scenarios.** Figure. 3 (b) gives the end-to-end performance  
 286  $TID_{F1}$  using PSENet [35] and CRNN [30] in different scenarios of MMVText. Similar to tracking  
 287 performance using EAST [44], the end-to-end video spotting performance shows the best performance  
 288 with a  $TID_{F1}$  of 31% in scenario of *Cartoon*.

289 **Text Recognition for Different Language.** As shown in Figure. 3 (c), the text recognition results  
 290 for different languages are provided. In summary, the alphanumeric recognition result (about 47%)  
 291 is better than the Chinese recognition result (about 35%), regardless of the models. The final  
 292 results (about 40%) for all characters are satisfactory, can not meet the requirement of the application.

293 **Text Detection for Different Text Category.** As shown in Figure. 3 (d), we provide the detection  
 294 performance comparison for different models in different text categories (*i.e.*, caption text or scene  
 295 text) of MMVText. It is obvious that the performance for scene text is better than the counterpart of  
 296 caption text, regardless of which detection model. The prime reason is that caption texts are all long  
 297 text, a different case to detect without any model refinement.

### 298 4.3 Text Detection and Recognition in Images

299 Although text detection and recognition in static images are not the focus in this work, we provide  
 300 the corresponding performance for comparison, as shown in Table. 2. For text detection, we adopts  
 301 EAST [44], PSENet [35] and DB [16] to evaluate the proposed MMVText. We observe that frame-  
 302 level text detection and recognition results on MMVText are not unsatisfactory, with lower results than  
 303 these methods report on existing scene text datasets. For example, EAST only obtains an f-score of  
 304 44.1% compared to the F-score of 80.7% on icdar2015 [45]. For text recognition, CRNN [30] based  
 305 on CTC loss, RARE [31] with attention mechanism and GRCNN [34] as the base text recognizers to  
 306 test our MMVText. The text annotation in our MMVText covers two languages (*i.e.*, English and  
 307 Chinese), thus we conduct several experiments for each language. ‘CH’ and ‘AN’ refer to Chinese  
 308 text instances and alphanumeric characters. ‘ALL’ denotes all characters regardless of which language.  
 309 Similar to the detection task, the recognition model only yields about 40% accuracy on our dataset,  
 310 but the same model reports  $> 90+$  on most benchmark datasets [14] for scene text recognition. The  
 311 main reasons have two points: (1) The proposed MMVText is multilingual, and the category number



Table 3: **Text Tracking Performance on MMVText.** Text tracking trajectory id generation use a method proposed in [36].

Method	MOTP	MOTA	ID <sub>P</sub> /%	ID <sub>R</sub> /%	ID <sub>F1</sub> /%	GT	MT	ML
EAST [44]	0.275	-0.301	23.5	22.9	23.2	48321	9680	35802
PSENet [35]	0.112	0.334	34.7	26.7	29.9	48321	12755	33410
DB [16]	0.102	0.438	33.7	29.9	<b>31.7</b>	48321	14958	31444

Table 4: **End to End Video Text Spotting Performance on MMVText.** Text tracking trajectory id generation use a method proposed in [36].  $TID_P$ ,  $TID_R$ ,  $TID_{F1}$ ,  $MOTA_T$  and  $MOTP_T$  refer to the corresponding metrics with recognition results in Table. 3.

Method		TID <sub>P</sub> /%	TID <sub>R</sub> /%	TID <sub>F1</sub> /%	MOTA <sub>T</sub>	MOTP <sub>T</sub>	MT	ML
Detection	Recognition							
EAST [44]	CRNN [30]	5.3	5.1	5.2	-0.835	0.173	1564	45963
	RARE [31]	3.0	3.6	3.2	-1.130	0.173	1265	47104
PSENet [35]	CRNN [30]	14.7	9.8	11.8	-0.300	0.197	3790	42957
	RARE [31]	15.2	10.4	12.4	-0.280	0.201	3821	42417
DB [16]	CRNN [30]	15.6	9.6	11.9	-0.284	0.230	3356	43246
	RARE [31]	20.1	15.2	<b>17.3</b>	-0.293	0.150	4230	39650

312 of Chinese characters in real-world images is much larger than those of Latin languages. (2) The  
 313 video texts are quite blurred, out-of-focus, and the distribution of characters is relatively smaller than  
 314 the static image counterparts.

#### 315 4.4 Text Tracking and Spotting in Videos

316 **Video Text Tracking.** Table. 3 shows the comparing results of text tracking on MMVText. We  
 317 observe that the overall performances of the used detectors are dissatisfactory on MMVText. Besides,  
 318 the ID<sub>F1</sub> of EAST [44] is lower with 6.7% gap than that of PSENet [35]. The main reason is that  
 319 MMVText contains a mass of long text instances, but regression-based EAST can not deal with  
 320 the long text cases well. The performance of DB is similar to that of PSENet for both all are the  
 321 segmentation-based methods. According to Table. 3,  $MOTP$  shows a better performance than  
 322  $MOTA$ . We argue that detectors such as PSENet or DB provide strong detecting capacity, but the  
 323 tracking ability is relatively weak. By comparison,  $IDF_1$  is a comprehensive metric for object ID  
 324 tracking. ID<sub>F1</sub> (31.7%) of DB achieves the best performance of the three detectors, and EAST shows  
 325 the worst performance with a ID<sub>F1</sub> of 23.2%.

326 **End to End Text Spotting in Video.** Detection or text tracking tasks are paving the way for the  
 327 recognition task. Table. 4 shows the performance of text spotting in the video. And  $TID_{F1}$  in  
 328 Equation. 3 as an integrated metric to evaluate algorithms in spatial location, content, and temporal  
 329 information three dimensions. Similar to the text tracking performance of EAST, the corresponding  
 330 performance  $TID_{F1}$  using CRNN [30] as the recognizer in video text spotting is still not satisfied  
 331 with a 5.2%  $TID_{F1}$ . The combination of DB [16] and RARE [31] achieves the best performance  
 332 with a 17.3%  $TID_{F1}$  among all the cases, but the performance still is inadequate to meet application  
 333 requirements. MT (Mostly Tracked) and ML (Mostly Lost) as the metrics concerning statistical  
 334 number can be used to evaluate from another aspect. For the combination of DB [16] and RARE [31],  
 335 39650 text tracking trajectories are lost, less than 20 percent of lifespan. By comparison, only 4230  
 336 tracking trajectories are satisfactory, more than 80 percent of lifespan tracked.

## 337 5 Conclusion and Future Work

338 In this paper, we establish a large-scale multidimensional and multilingual dataset for video text  
 339 tracking and spotting, termed as MMVText, with four description information, *i.e.*, bounding box,  
 340 tracking ID, recognition content, and text category label. Compare with the existing benchmarks, the  
 341 proposed MMVText mainly contains three advantages: large-scale, multidimensional, multilingual.  
 342 MMVText spans various video scenarios, text types, and multi-stage tasks, promoting video text  
 343 research. We also conduct several experiments on this dataset and shed light on what attributes are  
 344 especially difficult for the current task, which cast new insight into the video text tracking, spotting  
 345 field. In general, we hope the MMVText would facilitate the advance of video-and-text research.

## References

- 346
- 347 [1] Christos-Nikolaos E Anagnostopoulos, Ioannis E Anagnostopoulos, Ioannis D Psoroulas, Vassili  
348 Loumos, and Eleftherios Kayafas. License plate recognition from still images and video  
349 sequences: A survey. *IEEE Transactions on intelligent transportation systems*, 9(3):377–391,  
350 2008.
- 351 [2] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: the  
352 clear mot metrics. *EURASIP Journal on Image and Video Processing*, 2008:1–10, 2008.
- 353 [3] Yuanqiang Cai, Weiqiang Wang, Shao Huang, Jin Ma, and Ke Lu. Spatiotemporal text localiza-  
354 tion for videos. *Multimedia Tools and Applications*, 77(22):29323–29345, 2018.
- 355 [4] Zhazhan Cheng, Jing Lu, Yi Niu, Shiliang Pu, Fei Wu, and Shuigeng Zhou. You only recognize  
356 once: Towards fast video text spotting. In *ACM International Conference on Multimedia*, pages  
357 855–863, 2019.
- 358 [5] Chee Kheng Ch’ng and Chee Seng Chan. Total-text: A comprehensive dataset for scene  
359 text detection and recognition. In *IEEE International Conference on Document Analysis and  
360 Recognition*, pages 935–942, 2017.
- 361 [6] Patrick Dendorfer, Hamid Rezaatofghi, Anton Milan, Javen Shi, Daniel Cremers, Ian Reid,  
362 Stefan Roth, Konrad Schindler, and Laura Leal-Taixe. Cvpr19 tracking and detection challenge:  
363 How crowded can it get? *arXiv preprint arXiv:1906.04567*, 2019.
- 364 [7] Jianfeng Dong, Xirong Li, Chaoxi Xu, Xun Yang, Gang Yang, Xun Wang, and Meng Wang.  
365 Dual encoding for video retrieval by text. *IEEE Transactions on Pattern Analysis and Machine  
366 Intelligence*, 2021.
- 367 [8] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation  
368 in natural images. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages  
369 2315–2324, 2016.
- 370 [9] Tong He, Zhi Tian, Weilin Huang, Chunhua Shen, Yu Qiao, and Changming Sun. An end-to-end  
371 textspotter with explicit alignment and attention. In *IEEE conference on computer vision and  
372 pattern recognition*, pages 5020–5029, 2018.
- 373 [10] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Synthetic data and  
374 artificial neural networks for natural scene text recognition. *arXiv preprint arXiv:1406.2227*,  
375 2014.
- 376 [11] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial trans-  
377 former networks. In *Neural Information Processing Systems*, pages 2017–2025, 2015.
- 378 [12] Keechul Jung, Kwang In Kim, and Anil K Jain. Text information extraction in images and  
379 video: a survey. *Pattern recognition*, 37(5):977–997, 2004.
- 380 [13] Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew  
381 Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar,  
382 Shijian Lu, et al. Icdar 2015 competition on robust reading. In *IEEE International Conference  
383 on Document Analysis and Recognition*, pages 1156–1160, 2015.
- 384 [14] Dimosthenis Karatzas, Faisal Shafait, Seiichi Uchida, Masakazu Iwamura, Lluís Gomez i Big-  
385 orda, Sergi Robles Mestre, Joan Mas, David Fernandez Mota, Jon Almazan Almazan, and  
386 Lluís Pere De Las Heras. Icdar 2013 robust reading competition. In *IEEE International  
387 Conference on Document Analysis and Recognition*, pages 1484–1493, 2013.
- 388 [15] Hui Li, Peng Wang, and Chunhua Shen. Towards end-to-end text spotting with convolutional  
389 recurrent neural networks. In *IEEE International Conference on Computer Vision*, pages  
390 5238–5246, 2017.
- 391 [16] Minghui Liao, Zhaoyi Wan, Cong Yao, Kai Chen, and Xiang Bai. Real-time scene text  
392 detection with differentiable binarization. In *AAAI Conference on Artificial Intelligence*, pages  
393 11474–11481, 2020.

- 394 [17] Xuebo Liu, Ding Liang, Shi Yan, Dagui Chen, Yu Qiao, and Junjie Yan. Fots: Fast oriented text  
395 spotting with a unified network. In *Proceedings of the IEEE conference on computer vision and*  
396 *pattern recognition*, pages 5676–5685, 2018.
- 397 [18] Shangbang Long, Xin He, and Cong Yao. Scene text detection and recognition: The deep  
398 learning era. *International Journal of Computer Vision*, 129(1):161–184, 2021.
- 399 [19] Pengyuan Lyu, Minghui Liao, Cong Yao, Wenhao Wu, and Xiang Bai. Mask textspotter:  
400 An end-to-end trainable neural network for spotting text with arbitrary shapes. In *European*  
401 *Conference on Computer Vision*, pages 67–83, 2018.
- 402 [20] Abdelhamid Mammeri, Azzedine Boukerche, et al. Mser-based text detection and commu-  
403 nication algorithm for autonomous vehicles. In *2016 IEEE symposium on computers and*  
404 *communication (ISCC)*, pages 1218–1223. IEEE, 2016.
- 405 [21] Abdelhamid Mammeri, El-Hebri Khiari, and Azzedine Boukerche. Road-sign text recognition  
406 architecture for intelligent transportation systems. In *2014 IEEE 80th Vehicular Technology*  
407 *Conference (VTC2014-Fall)*, pages 1–5. IEEE, 2014.
- 408 [22] Rodrigo Minetto, Nicolas Thome, Matthieu Cord, Neucimar J Leite, and Jorge Stolfi. Snooper-  
409 track: Text detection and tracking for outdoor videos. In *IEEE International Conference on*  
410 *Image Processing*, pages 505–508, 2011.
- 411 [23] Anand Mishra, Karteek Alahari, and CV Jawahar. Image retrieval using textual cues. In  
412 *Proceedings of the IEEE International Conference on Computer Vision*, pages 3040–3047,  
413 2013.
- 414 [24] Nibal Nayef, Fei Yin, Imen Bizid, Hyunsoo Choi, Yuan Feng, Dimosthenis Karatzas, Zhenbo  
415 Luo, Umapada Pal, Christophe Rigaud, Joseph Chazalon, et al. Icdar2017 robust reading  
416 challenge on multi-lingual scene text detection and script identification-rrc-mlt. In *IEEE*  
417 *International Conference on Document Analysis and Recognition*, volume 1, pages 1454–1459,  
418 2017.
- 419 [25] Phuc Xuan Nguyen, Kai Wang, and Serge Belongie. Video text detection and recognition:  
420 Dataset and benchmark. In *IEEE winter conference on applications of computer vision*, pages  
421 776–783, 2014.
- 422 [26] Sangeeth Reddy, Minesh Mathew, Lluís Gomez, Marçal Rusinol, Dimosthenis Karatzas, and  
423 CV Jawahar. Roadtext-1k: Text detection & recognition dataset for driving videos. In *IEEE*  
424 *International Conference on Robotics and Automation*, pages 11074–11080, 2020.
- 425 [27] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time  
426 object detection with region proposal networks. pages 91–99, 2015.
- 427 [28] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance  
428 measures and a data set for multi-target, multi-camera tracking. In *Workshops of European*  
429 *conference on computer vision*, pages 17–35, 2016.
- 430 [29] Georg Schroth, Sebastian Hilsenbeck, Robert Huitl, Florian Schweiger, and Eckehard Steinbach.  
431 Exploiting text-related features for content-based image retrieval. In *2011 IEEE international*  
432 *symposium on multimedia*, pages 77–84. IEEE, 2011.
- 433 [30] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-  
434 based sequence recognition and its application to scene text recognition. *IEEE transactions on*  
435 *pattern analysis and machine intelligence*, 39(11):2298–2304, 2016.
- 436 [31] Baoguang Shi, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Robust scene text  
437 recognition with automatic rectification. In *IEEE conference on computer vision and pattern*  
438 *recognition*, pages 4168–4176, 2016.
- 439 [32] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video  
440 representations using lstms. In *International Conference on Machine Learning*, pages 843–852,  
441 2015.

- 442 [33] Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. Coco-text:  
443 Dataset and benchmark for text detection and recognition in natural images. *arXiv preprint*  
444 *arXiv:1601.07140*, 2016.
- 445 [34] Jianfeng Wang and Xiaolin Hu. Gated recurrent convolution neural network for ocr. In *Neural*  
446 *Information Processing Systems*, pages 334–343, 2017.
- 447 [35] Wenhai Wang, Enze Xie, Xiang Li, Wenbo Hou, Tong Lu, Gang Yu, and Shuai Shao. Shape  
448 robust text detection with progressive scale expansion network. In *IEEE conference on computer*  
449 *vision and pattern recognition*, pages 9336–9345, 2019.
- 450 [36] Xiaobing Wang, Yingying Jiang, Shuli Yang, Xiangyu Zhu, Wei Li, Pei Fu, Hua Wang, and  
451 Zhenbo Luo. End-to-end scene text recognition in videos based on multi frame tracking. In  
452 *IEEE International Conference on Document Analysis and Recognition*, pages 1255–1260,  
453 2017.
- 454 [37] Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuan-Fang Wang, and William Yang Wang. Vatex:  
455 A large-scale, high-quality multilingual dataset for video-and-language research. In *Proceedings*  
456 *of the IEEE/CVF International Conference on Computer Vision*, pages 4581–4591, 2019.
- 457 [38] Liang Wu, Palaiahnakote Shivakumara, Tong Lu, and Chew Lim Tan. A new technique for  
458 multi-oriented scene text line detection and tracking in video. *IEEE Transactions on multimedia*,  
459 17(8):1137–1152, 2015.
- 460 [39] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for  
461 bridging video and language. In *Proceedings of the IEEE conference on computer vision and*  
462 *pattern recognition*, pages 5288–5296, 2016.
- 463 [40] Chun Yang, Xu-Cheng Yin, Wei-Yi Pei, Shu Tian, Ze-Yu Zuo, Chao Zhu, and Junchi Yan.  
464 Tracking based multi-orientation scene text detection: A unified framework with dynamic  
465 programming. *IEEE Transactions on Image Processing*, 26(7):3235–3248, 2017.
- 466 [41] Xu-Cheng Yin, Ze-Yu Zuo, Shu Tian, and Cheng-Lin Liu. Text detection, tracking and  
467 recognition in video: a comprehensive survey. *IEEE Transactions on Image Processing*,  
468 25(6):2752–2773, 2016.
- 469 [42] Fisher Yu, Wenqi Xian, Yingying Chen, Fangchen Liu, Mike Liao, Vashisht Madhavan, and  
470 Trevor Darrell. Bdd100k: A diverse driving video database with scalable annotation tooling.  
471 *arXiv preprint arXiv:1805.04687*, 2018.
- 472 [43] Oussama Zayene, Mathias Seuret, Sameh Masmoudi Touj, Jean Hennebert, Rolf Ingold, and  
473 Najoua Essoukri Ben Amara. Text detection in arabic news video based on SWT operator and  
474 convolutional auto-encoders. In *Workshop on Document Analysis Systems*, pages 13–18, 2016.
- 475 [44] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang.  
476 East: an efficient and accurate scene text detector. In *IEEE conference on computer vision and*  
477 *pattern recognition*, pages 5551–5560, 2017.
- 478 [45] Xinyu Zhou, Shuchang Zhou, Cong Yao, Zhimin Cao, and Qi Yin. Icdar 2015 text reading in  
479 the wild competition. *arXiv preprint arXiv:1506.03184*, 2015.

## 480 Checklist

- 481 1. For all authors...
- 482 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
483 contributions and scope? [Yes]
- 484 (b) Did you describe the limitations of your work? [Yes] We describe the limitations in  
485 supplementary material.
- 486 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 487 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
488 them? [Yes]

- 489 2. If you are including theoretical results...
- 490 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 491 (b) Did you include complete proofs of all theoretical results? [Yes]
- 492 3. If you ran experiments (e.g. for benchmarks)...
- 493 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
- 494 perimental results (either in the supplemental material or as a URL)? [Yes] We have
- 495 provided the URL concerning the coding and the data to promote further research.
- 496 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 497 were chosen)? [Yes]
- 498 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 499 ments multiple times)? [Yes]
- 500 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 501 of GPUs, internal cluster, or cloud provider)? [Yes]
- 502 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 503 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 504 (b) Did you mention the license of the assets? [Yes]
- 505 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 506 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 507 using/curating? [N/A]
- 508 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 509 information or offensive content? [Yes] We have blurred identifiable information or
- 510 offensive content.
- 511 5. If you used crowdsourcing or conducted research with human subjects...
- 512 (a) Did you include the full text of instructions given to participants and screenshots, if
- 513 applicable? [Yes]
- 514 (b) Did you describe any potential participant risks, with links to Institutional Review
- 515 Board (IRB) approvals, if applicable? [N/A]
- 516 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 517 spent on participant compensation? [Yes] We have paid salary to the related participants.