

Supplementary Material: MIST: Multilingual Incidental Dataset for Scene Text Detection

Anonymous WACV Algorithms Track submission

Paper ID 1327

001 **Contents**

002 **A Annotation Guideline**

003 **B Dataset Comparison Complete**

004 **C Dataset Split and Curation**

005 **D Evaluation Metrics**

006 **E Training Configurations**

007 **F. Insights and Takeaways**

008 **G VLM performance on MIST**

009 **H Visual Samples and Results**

010 **A. Annotation Guideline**

011 The key annotation guidelines provided to annotators are as
012 follows:

- 013 Bounding box annotations follow the COCO-Text approach,
014 where a word is defined as an uninterrupted sequence of characters separated by spaces.
- 015 Word-level annotations use polygons, similar to Total-
016 Text and CTW1500. Given the dataset's multilingual
017 and multi-oriented text, polygonal bounding boxes ensure
018 precise annotations.
- 019 All human-legible text must be tightly annotated, regard-
020 less of challenging conditions such as poor illumination
021 or motion blur. If a text instance appears illegible at first
022 glance, annotators should adjust contrast or brightness to
023 determine its readability.
- 024 Occluded words should be annotated as a whole rather
025 than in segments.
- 026 We also annotate the illegible text and mark it as *do not*
027 *care* during evaluation.
- 028 Personally Identifiable Information (PII) must be anno-
029 tated and labeled as such PII. During post-processing,
030 these marked areas are blurred to protect sensitive data.

032 **B. Dataset Comparison Complete**

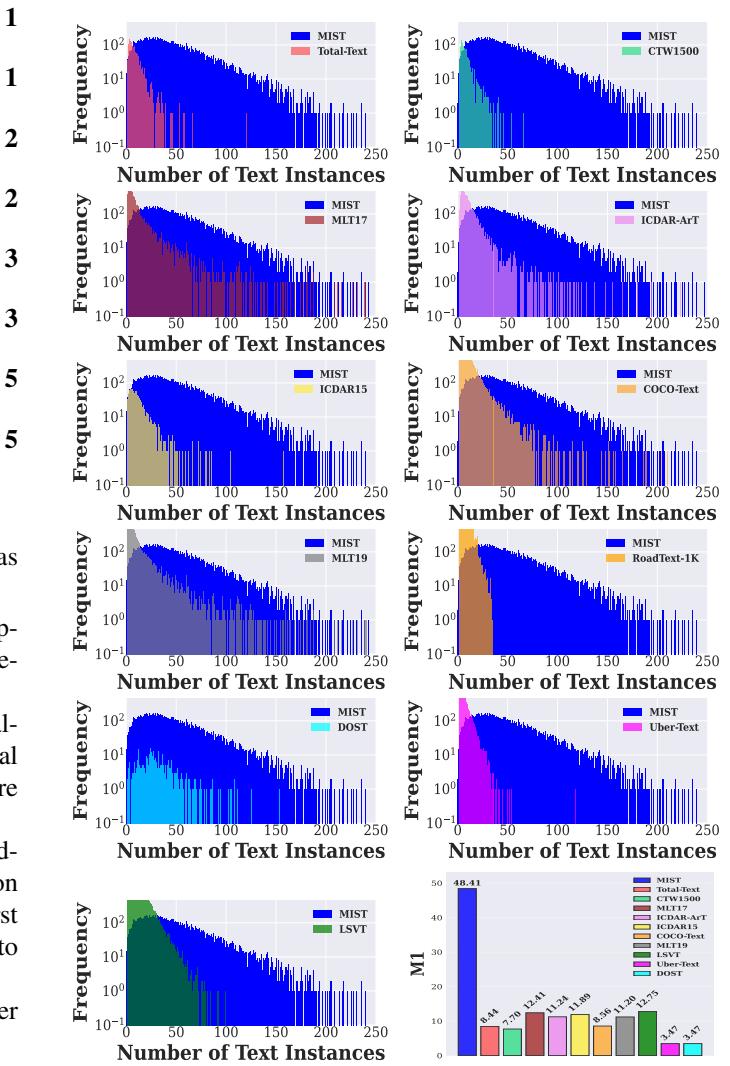


Figure 1. Compares the **distribution of text instances in scene images** (M_1) of against existing datasets: Total-Text, CTW1500, MLT17, ICDAR-ArT, ICDAR15, COCO-Text, and RoadText-1K.

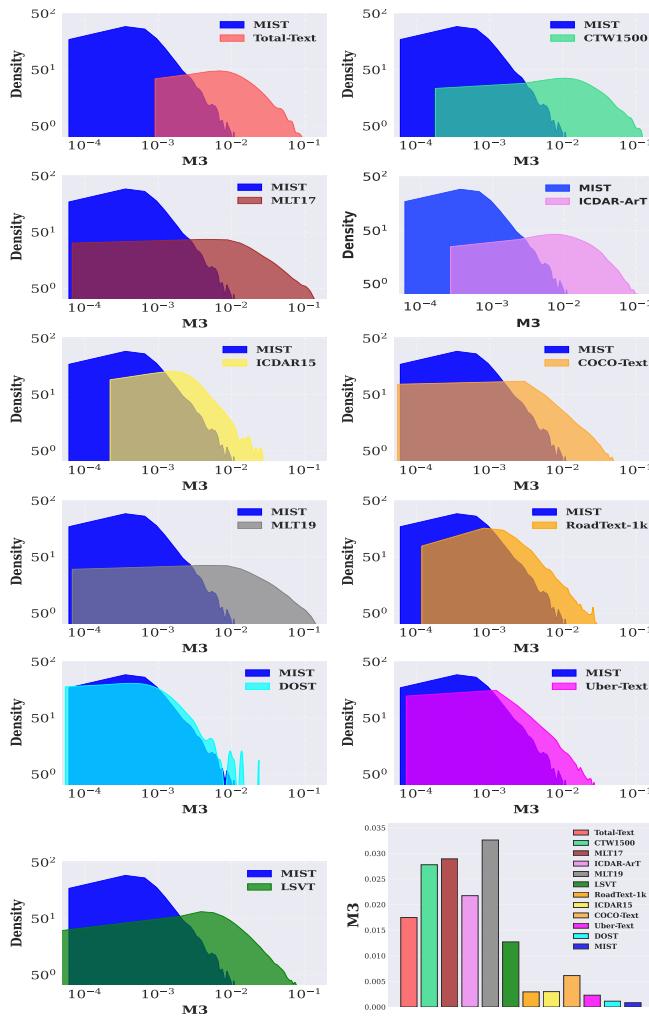


Figure 2. Compares the **average area of a text instance about the scene image** (M_3) of MIST against: Total-Text, CTW1500, MLT17, ICDAR-ArT, ICDAR15, COCO-Text, and RoadText-1K. The analysis employs kernel density estimation to create continuous distribution curves, with both axes displayed on logarithmic scales – base 50 for the y-axis and base 2 for the x-axis.

In terms of the occlusion/coverage metric M_3 , DOST and Uber-Text approach MIST, but both have practical limitations. **DOST**, the closest to MIST in its incidental capture, offers only a static-image split with **338 training images**. Additionally, its frames are subsampled at every 10th interval, introducing redundancy and limiting diversity. **Uber-Text**, on the other hand, has a very low average of ≈ 3 text instances per image, so it is used more often for scene text *recognition* than for detection.

C. Dataset Split and Curation

1. Dataset layout. We collect 150 h across 10 regions (15 h/region), with $M=15$ one-hour sequences per re-

gion.

2. Per-sequence quotas. Target 12,000 frames overall (1,200 per region). For each region, compute

$$b = \lfloor 1200/M \rfloor = 80, \quad \rho = 1200 \bmod M = 0,$$

and assign sequence-wise quotas

$$B_i = b + 1\{i \leq \rho\}, \quad i = 1, \dots, M.$$

(Here $\rho=0 \Rightarrow B_i=80$ for all i .)

3. Stratified temporal sampling. For each sequence i with quota B_i :

- (a) Partition the timeline into B_i equal temporal bins.
- (b) Select one frame per bin (either the bin midpoint or uniformly at random within the bin) to obtain B_i frames with near-uniform temporal coverage.

4. Region-level shuffling and split. Within each region, pool the $\sum_{i=1}^M B_i = 1200$ frames, apply a random permutation, and slice into train/val/test according to the region’s target counts (summing to 1,200).

5. Disjointness guarantee. Train/val/test splits are strictly disjoint at the *frame* level (no duplicate or overlapping frames across splits), with stratification ensuring near-uniform temporal coverage per sequence.

D. Evaluation Metrics

We adopt the DetEval protocol for evaluating scene text detection models, as used in ICDAR15 and Total-Text. This protocol employs three matching strategies — One-to-One, One-to-Many, and Many-to-One.

One-to-One: This matching strategy is used when a detected instance and a ground truth instance uniquely correspond to each other within the specified thresholds, with no additional overlapping candidates.

One-to-Many: Once all One-to-One matches have been identified, this matching algorithm assigns a single ground truth instance to multiple detections. This approach ensures that multiple valid segmented detections related to a single ground truth are not penalized, which is especially important for long or multilingual text instances. Each segmented detection and its corresponding ground truth receive a partial score. This scoring acknowledges the partial correctness of the detection and the degree of information captured.

Many-to-One: After completing the previous matches, this algorithm assigns a single detection to multiple ground truths. This issue commonly arises in crowded environments, where several text instances may be inaccurately identified as just one detection. Similar to one-to-many matching, both the ground truths and the single detection receive partial scores.

091 **E. Training Configurations**

092 The training code for the compared methods can be
 093 found at the following repositories: DBNet++¹, MixNet²,
 094 TextBPN++³, and DPText-DETR⁴.

095 In Table 1, the training configurations for the
 096 TextBPN++ models specific to other datasets have been pro-
 097 vided. We used these models in Sec 3.5 and Sec 5.2 in the
 098 main paper.

Fine-tuning dataset	Pretraining used
Total-Text [1]	MLT17 [4]
CTW1500 [3]	MLT17 [4]
17 [4]	SynthText [2]

Table 1. Training configurations for TBPN models for Total-text, CTW1500 and MLT17

099 **Learning rate and schedule for TextBPN++ and DPText-
 100 DETR** We optimize the network with an initial learning
 101 rate $\eta_0 = 1 \times 10^{-4}$. A step-decay scheduler is used: at
 102 fixed intervals of s epochs, the learning rate is multiplied
 103 by a decay factor $\gamma = 0.9$, i.e.,

$$\eta_t = \eta_0 \cdot \gamma^{\lfloor t/s \rfloor},$$

105 where t counts epochs and s is the step size.

Hyperparameter	Value
Initial learning rate	1×10^{-4}
Learning rate scheduler	Step decay
Decay factor (per step)	0.9

Table 2. Optimization hyperparameters.

106 **F. Insights and Takeaways**107 **Generalization Capability of MIST**

- 108 We manually re-annotated CTW1500 at the word level
 109 for a fair comparison, as its official annotations are at
 110 the text-line/phrase level. We exclude CTW1500 and
 111 ICDAR-ArT from training to avoid train-test leakage
 112 arising from their line/phrase-level annotations and over-
 113 laps with other benchmarks.
- 114 Because MLT17 is a competition dataset with a non-
 115 public test set, we use the provided validation set to ana-
 116 lyze metrics and fine-tuning performance.
- 117 • We do not report the performance of a model trained
 118 on ICDAR-ART when evaluated on Total-Text or

¹<https://github.com/MhLiao/DB>

²<https://github.com/D641593/MixNet>

³<https://github.com/GXYM/TextBPN-Plus-Plus/>

⁴<https://github.com/ymy-k/DPText-DETR>

119 CTW1500, since the ICDAR-ART training set contains
 120 samples overlapping the test sets of those datasets, which
 121 would confound evaluation.

122 **Transfer Learning Capability Score**

123 Table 3 summarizes the performance of DPText-DETR and
 124 TextBPN++ pre-trained on MIST and fine-tuned on Total-
 125 Text, CTW1500, and MLT17.

Test Set	Model	P	R	F	F^α
Total-Text	H1	92.06	87.05	89.48	89.00
	H2	92.60	88.85	90.69	90.13
CTW1500	H1	91.58	86.72	89.08	88.80
	H2	88.56	86.83	87.69	86.49
MLT	H2	91.57	74.29	82.02	81.19

126 Table 3. Shows performance of DPText-DETR [5] and
 127 TextBPN++ [6] pre-trained on MIST and fine-tuned on Total-
 128 Text and CTW1500. H1 and H2 indicate DPText-DETR and
 129 TextBPN++, respectively. F^α is reported F-Measure for
 130 Total-Text, CTW1500, and MLT17 using these models without
 131 pre-trained on MIST. Bold and underscore indicate the best and
 second-best value, respectively. Since DPText-DETR is not eval-
 uated on MLT17, we only use TextBPN++ for MLT17.

132 **Few Shot Transfer Learning Capability**

Sample Size	F-Measure (Total-Text)	F-Measure (CTW1500)
100	85.10	82.24
500	86.90	86.43
1000	87.63	87.47
2000	88.98	88.06
4000	89.04	88.67
8000	89.20	89.02
Full	89.48	89.08

133 Table 4. F-Measure across different sample sizes for two evalua-
 134 tion settings. The final row indicates performance when using the
 135 entire dataset (17,600 and 13,000 samples respectively).

136 We demonstrate the few-shot transfer learning capa-
 137 bilities of TextBPN++ and DPText-DETR, pre-trained on
 138 MIST, when applied to existing datasets such as CTW1500
 139 and Total-Text. The training configurations remain the same
 140 as in the main paper.

141 **CTW1500:** We conduct transfer learning on TextBPN++
 142 using sample of 100, 500, and 1000 of CTW1500. As
 143 shown in Figure 3(a), the model consistently surpasses the

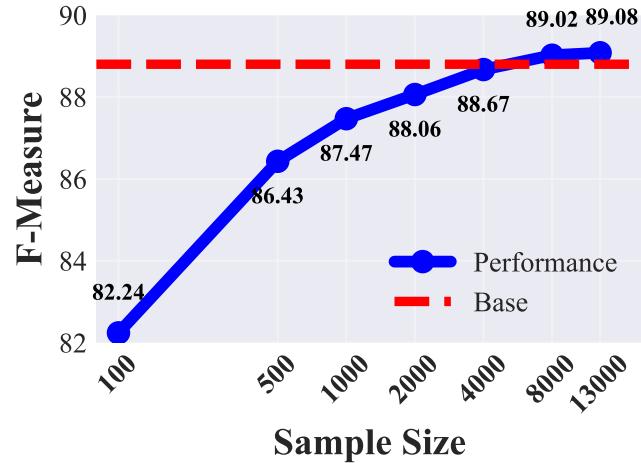
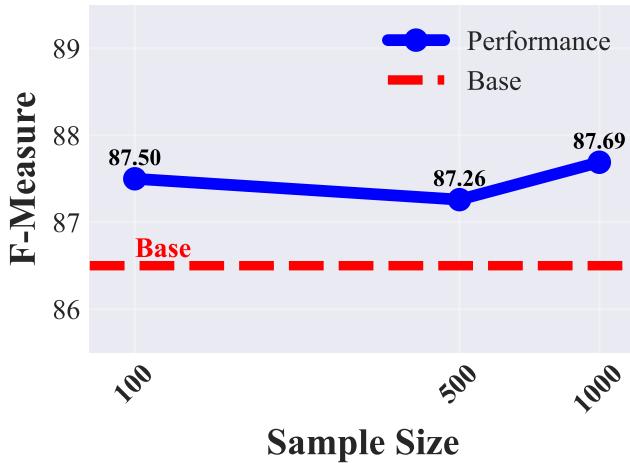


Figure 3. (a) shows the few-shot transfer learning performance of TextBPN++ pre-trained on MIST for CTW500. (b) illustrates the few-shot transfer learning performance of DPTText-DETR pre-trained on MIST for CTW1500.

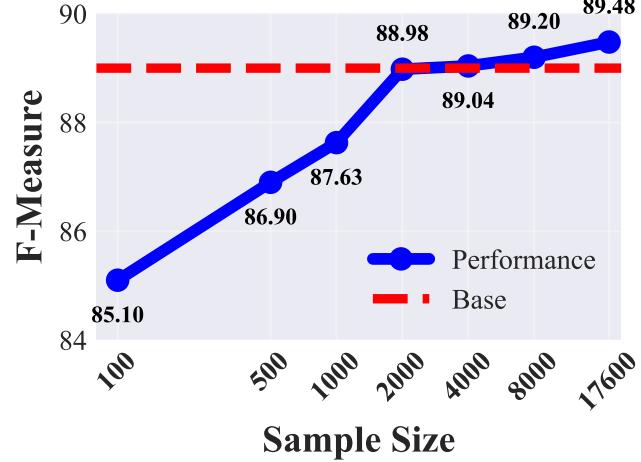
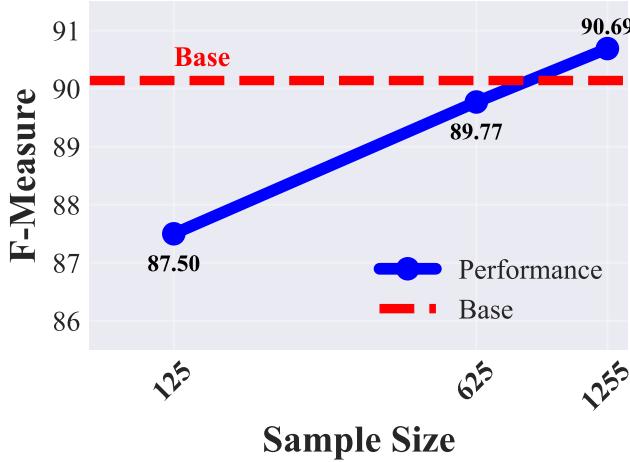


Figure 4. (a) shows the few-shot transfer learning performance of TextBPN++ pre-trained on MIST for Total-Text. (b) illustrates the few-shot transfer learning performance of DPTText-DETR pre-trained on MIST for Total-Text.

135 baseline performance of TextBPN++ on CTW1500, as re-
136 ported in the original paper. Notably, even with just 100
137 samples, the model performs exceptionally well. However,
138 its best performance is achieved with the largest sample
139 size, surpassing the base performance by 1.20%

140 We perform transfer learning on DPTText-DETR us-
141 ing sample sizes of 100, 500, 1000, 2000, 4000, 8000,
142 and 13,000 as shown in Figure 3(b). The samples were
143 sourced from data provided by the authors of TextBPN++
144 and DPTText-DETR, resulting in different sample sizes.
145 The model surpasses DPTText-DETR’s best-reported per-
146 formance on CTW1500 by 0.28%, a notable achievement
147 given that DPTText-DETR is the state-of-the-art model for
148 this dataset. Remarkably, it even surpasses the baseline

149 performance with just 8000 training samples and achieves
150 near baseline performance with 4000 samples. This further
151 demonstrates the ability of MIST trained models to adapt to
152 specific domains.

Total-Text: For Total-Text, we conduct transfer learning
153 on TextBPN++ using sample sizes of 125, 625, and 1255.
154 Our model surpasses the base performance of TextBPN++,
155 as reported in its original paper (see Figure 4 (a)). MIST
156 achieves an F-measure close to the base performance even
157 with a small sample size, demonstrating its strong adapt-
158 ability to specific domains.

159 For DPTText-DETR, transfer learning is performed with
160 sample sizes of 100, 500, 1000, 2000, 4000, 8000, and
161

162 17,600. The model reaches performance close to its base-
 163 line with just 2,000 our of the 17600 samples, falling short
 164 by only 0.02%. Beyond this point, increasing the amount of
 165 training samples to 4000 and 8000 surpasses the base per-
 166 formance. This reinforces the adaptability of MIST-trained
 167 models on different data distributions (see Figure 4 (b)).

168 In addition to its value as an incidental scene text dataset
 169 and benchmark, MIST stands out as the most generaliz-
 170 able dataset in the literature, well-suited for real-world chal-
 171 lenges and highly effective for pretraining. Its strong pre-
 172 training performance stems from the diverse and complex
 173 text instances it captures, providing the model with samples
 174 to tackle macro-level disturbances. Additionally, MIST can
 175 achieve near-baseline performance on specific datasets with
 176 minimal training data, highlighting its exceptional value in
 177 scene text research.

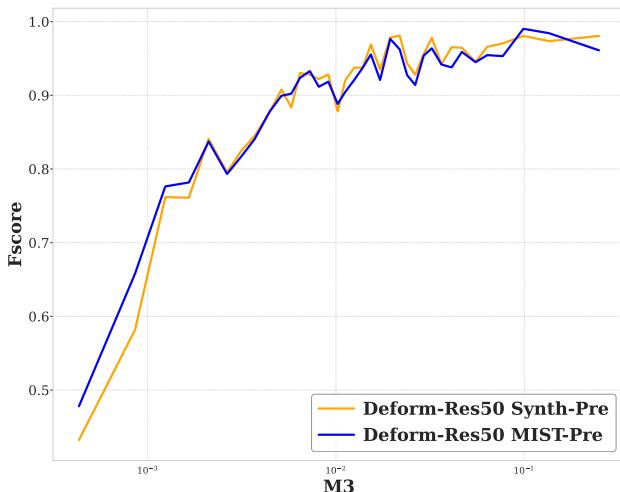


Figure 5. Compares the F-measure and M_3 of TextBPN++ pre-trained with SynthText and MIST, trained on MLT17, while evaluating on MLT17.

178 Issue Mitigation

179 In the main paper, we raised the topic of issue mitigation
 180 by fine-tuning a MIST pre-trained model. This sec-
 181 tion explores issue mitigation on MLT17 and the metric
 182 M_2 . Through Figure 5, the effect of the MIST-pre-trained
 183 model was seen, improving TextBPN++’s performance on
 184 relatively lower M_3 . We perform the same experiment
 185 on MLT17, grouping the MLT17’s samples into batches
 186 of 50 and calculating the F-measure for each batch. In
 187 the case of MLT17, the mitigation of the issue is am-
 188 biguous. The MIST pre-trained model initiates with im-
 189 proved performance on lower M_3 ; however, after the first
 190 4 batches, the performance becomes ambiguous, with rel-
 191 atively fewer batches corresponding to improved perfor-
 192 mances than Total-Text. Out of the first 10 batches, MIST

193 pre-trained models outperform the base SynthText pre-
 194 trained model 6 times. When coupled with the model out-
 195 performing the base TextBPN++ model on the validation
 196 set, this observation suggests a positive result of issue miti-
 197 gation at smaller M_3 .

198 We perform a similar analysis on M_2 by labeling text in-
 199 stances larger than size thresholds as ‘do not care’ regions
 200 and calculating the F-measure for the filtered test set. We
 201 plot the F-Measure against M_2 , as seen in Figure 6. Un-
 202 like the performance on M_3 , where MLT17 was ambigu-
 203 ous, here, the performance on MLT17 surpasses the base
 204 model’s performance on all thresholds. In contrast, the
 205 model fine-tuned on Total-Text is extremely ambiguous at
 206 all thresholds.

207 MIST’s issue mitigation abilities at various scales are
 208 highlighted through these observations. This property can
 209 be attributed to its collection of varied text instances and
 210 many scenes in harsh conditions.

211 G. VLM performance on MIST

212 We qualitatively evaluated general-purpose vision language
 213 models (VLMs) ChatGPT-5 and Gemini 2.5 Pro on 5 MIST
 214 test images to probe their ability to perform scene text de-
 215 tection directly from prompts. While these VLMs often
 216 demonstrate strong high-level image understanding, their
 217 detection outputs on MIST were underwhelming and ex-
 218 hibited several recurrent issues:

- **Sparse and low-recall detections.** Models produced a significantly low recall rate of 15%, missing the majority of text instances.
- **Resolution/scale mismatch.** The coordinates returned by the models appeared to be defined in the (downsampled/compressed) internal inference resolution rather than the original image size; without explicit post-processing. This led to systematic misalignment of boxes when overlaid on the 1920×1080 input.

219 Visual results provided in Fig. 7

220 H. Visual Samples and Results

221 A Few Samples from MIST

222 Figure 8 shows a few sample images from MIST with di-
 223 verse text instances.

224 Comparison of Visual Result

225 Figure. 9 visually compares text detection results in MIST
 226 by (a) TextBPN++ trained on MIST, (b) state-of-the-art
 227 model for Total-Text, (c) state-of-the-art model for ICDAR-
 228 ArT, (d) state-of-the-art model for MLT.

229 Visual Result of Cross-Domain Scene Images

230 Figure 10 shows detected text in Total-Text and MLT by
 231 TextBPN++ trained on MIST, state-of-the-art model for

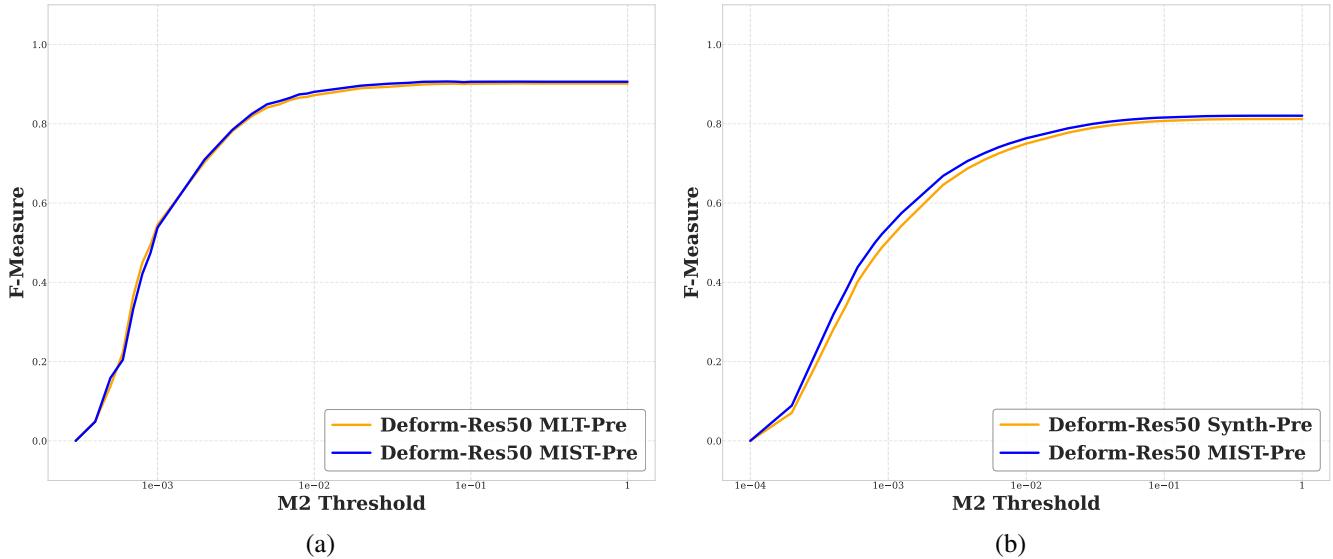


Figure 6. (a) Compares the F-measure and M_2 of TextBPN++ pre-trained with MLT17 and MIST, trained on Total-Text, while evaluating on Total-Text. (b) Compares the F-measure and M_2 of TextBPN++ pre-trained with SynthText and MIST, trained on MLT17, while evaluating on MLT17.

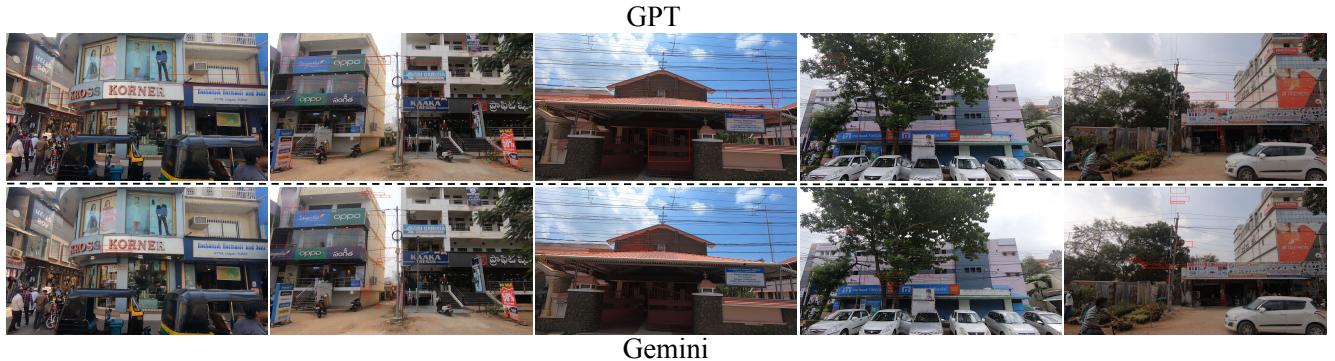


Figure 7. The predicted boxes are annotated in red. As we can see, the boxes are not aligned with the text. However, there seems to be some sign that the models understand the layout of the text and were constrained due to post-processing.

241 Total-Text, and state-of-the-art model for MLT.

242 Visual Results of Out-of-Domain Scene Images

243 We download a few scene text images from the Internet by searching. We provide text detection results of TextBPN++ trained on MIST, Total-Text, MLT, and ICDAR-ArT. Figure 11 shows detected text using different models. It is visually shown that TextBPN++ trained on MIST is able to detect the majority of text than other models.



Figure 8. Shows ample images contains text instances in poor light, complex background, crowed surrounding, 3D and smaller text, occlusion, perspective distortion.



Figure 9. Shows predicted text detection results of MIST. (a) detected text by TextBPN++ trained on MIST. (b) detected text by state-of-the-art model for Total-Text. (c) detected text by state-of-the-art model for ICDAR-ArT. (d) detected text by state-of-the-art model for MLT.

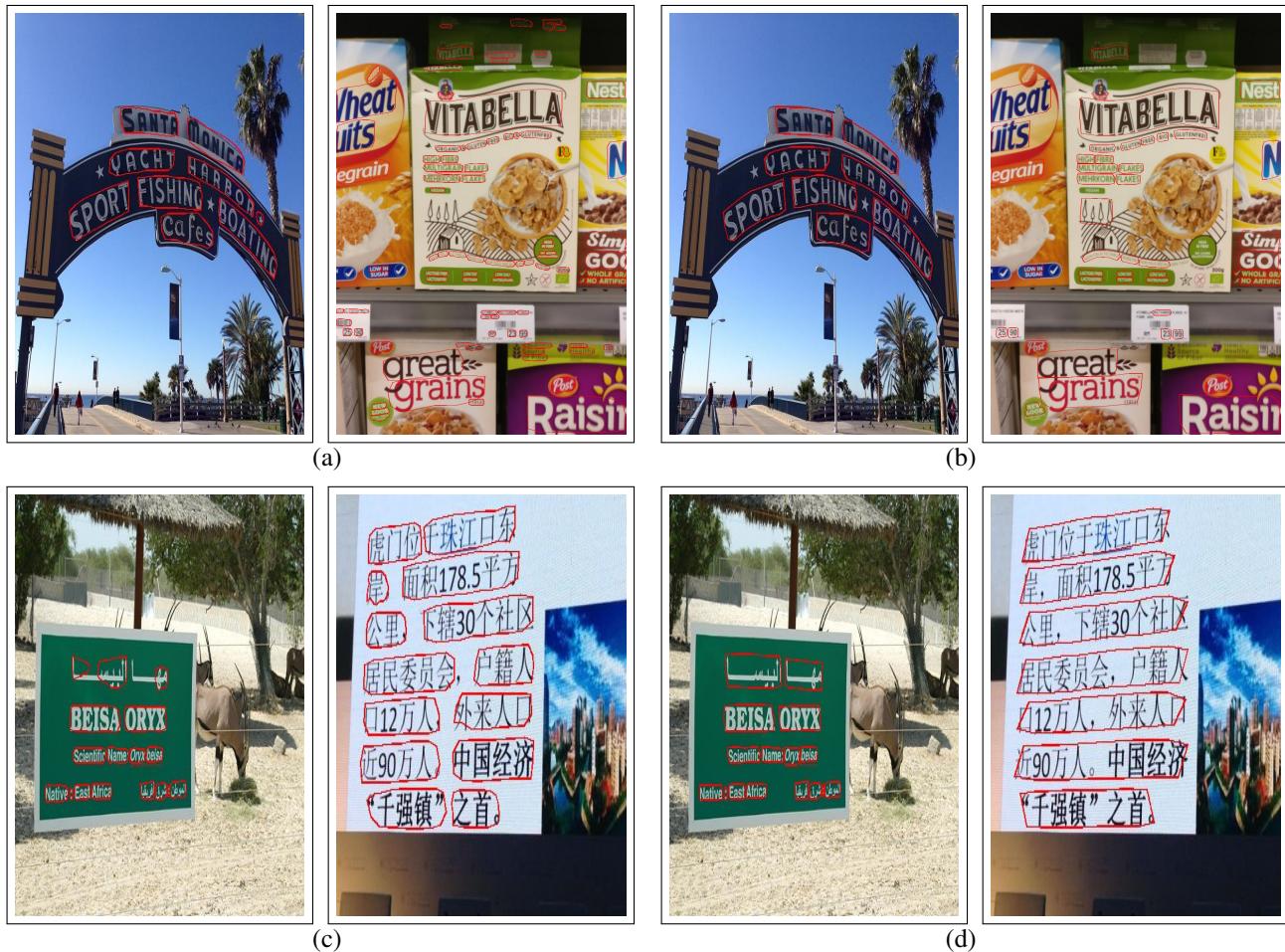


Figure 10. Shows predicted text detection results of Total-Text and MLT. (a) detected text in Total-Text by TextBPN++ trained on MIST. (b) detected text in Total-Text by state-of-the-art model for Total-Text. (c) detected text in MLT by TextBPN++ trained on MIST. (d) detected text in MLT by state-of-the-art model for MLT.



Figure 11. Shows predicted text detection results on out-of-distribution scene images. (a) detected text by TextBPN++ trained on MIST. (b) detected text by state-of-art model for Total-Text. (c) detected text by state-of-the-art model for MLT. (d) detected text by state-of-the-art model for ICDAR-ArT.

249 References

- 250 [1] Chee Kheng Ch'ng and Chee Seng Chan. Total-Text: A com-
251 prehensive dataset for scene text detection and recognition. In
252 *ICDAR*, pages 935–942, 2017.
- 253 [2] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Syn-
254 thetic data for text localisation in natural images. In *CVPR*,
255 pages 2315–2324, 2016.
- 256 [3] Yuliang Liu, Lianwen Jin, Shuitao Zhang, Canjie Luo, and
257 Sheng Zhang. Curved scene text detection via transverse and
258 longitudinal sequence connection. *PR*, 90:337–345, 2019.
- 259 [4] Nibal Nayef, Fei Yin, Imen Bizid, Hyunsoo Choi, Yuan
260 Feng, Dimosthenis Karatzas, Zhenbo Luo, Umapada Pal,
261 Christophe Rigaud, Joseph Chazalon, et al. ICDAR2017 ro-
262 bust reading challenge on multi-lingual scene text detection
263 and script identification-RRC-MLT. In *ICDAR*, pages 1454–
264 1459, 2017.
- 265 [5] Maoyuan Ye, Jing Zhang, Shanshan Zhao, Juhua Liu, Bo Du,
266 and Dacheng Tao. DPText-DETR: Towards better scene text
267 detection with dynamic points in transformer. In *AAAI*, pages
268 3241–3249, 2023.
- 269 [6] Shi-Xue Zhang, Chun Yang, Xiaobin Zhu, and Xu-Cheng
270 Yin. Arbitrary shape text detection via boundary transformer.
271 *TMM*, 26:1747–1760, 2023.