

MENTOR: Multilingual tExt detectionN TOward leaRning by analogy



Hsin-Ju Lin, Tsu-Chun Chung, Ching-Chun Hsiao, Pin-Yu Chen*, Wei-Chen Chiu, and Ching-Chun Huang

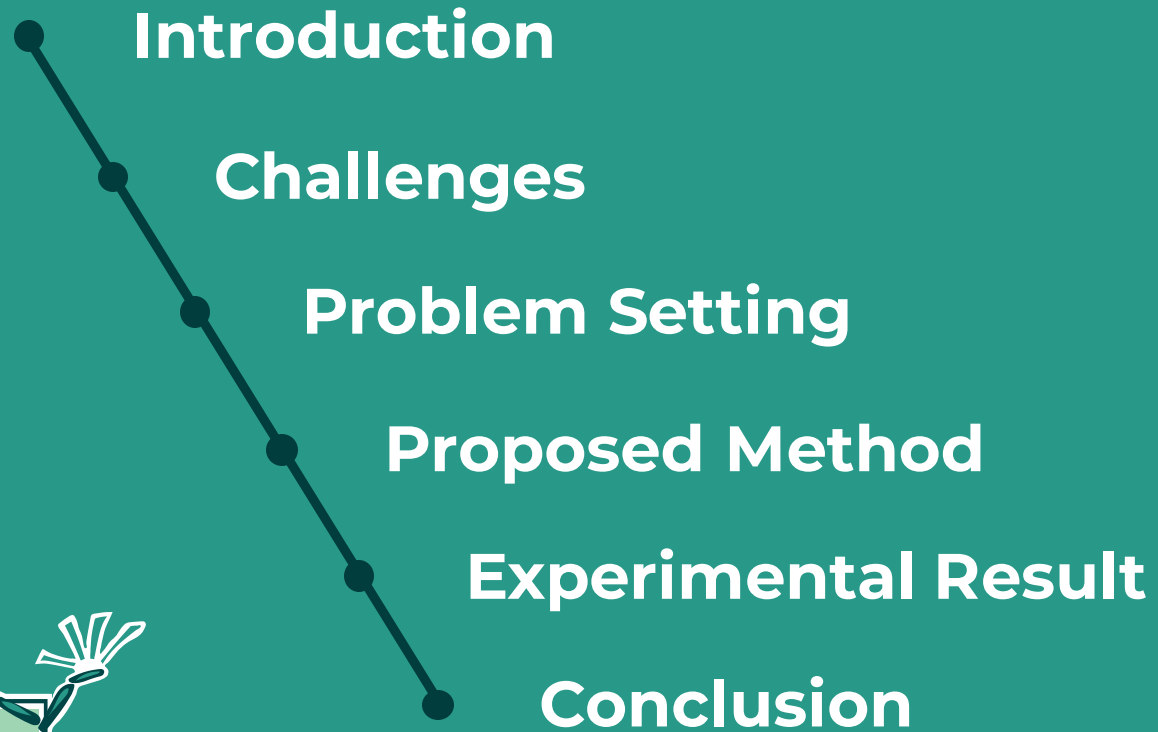
Organization: IBM research*

National Yang Ming Chiao Tung University

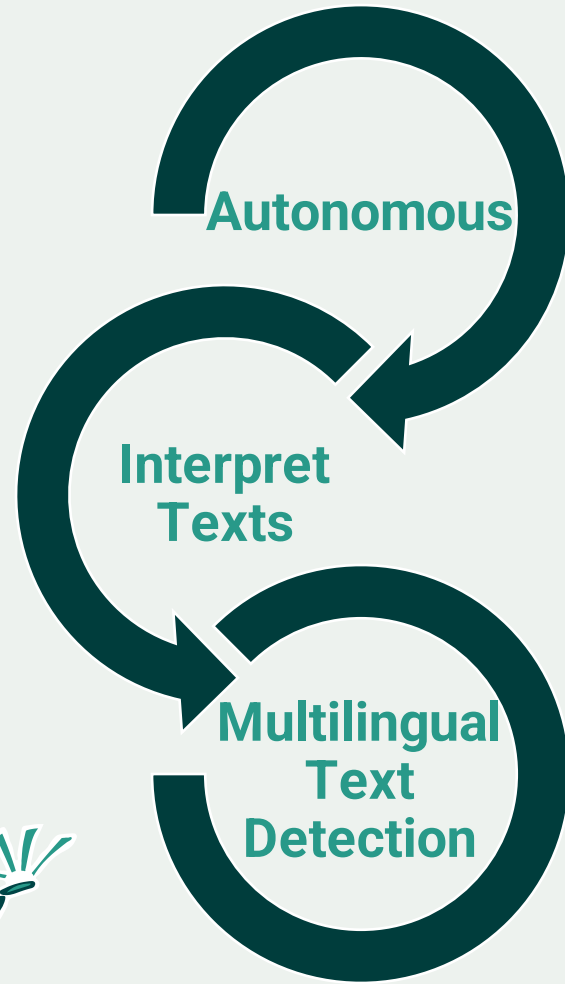
Reporter: Lucy Lin (Hsin-Ju Lin)



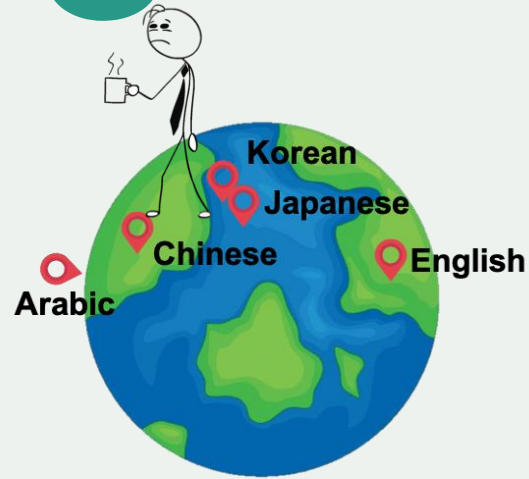
Outline



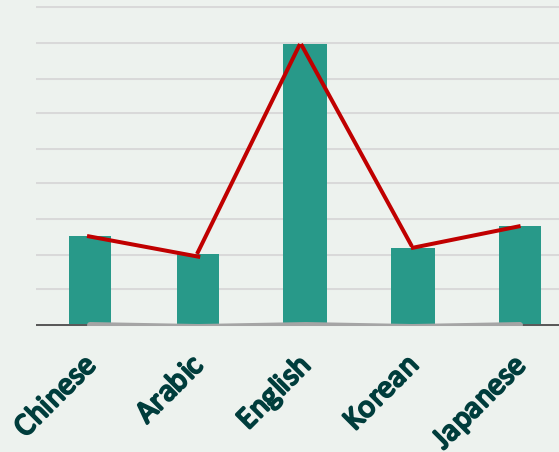
Introduction



3 Challenges



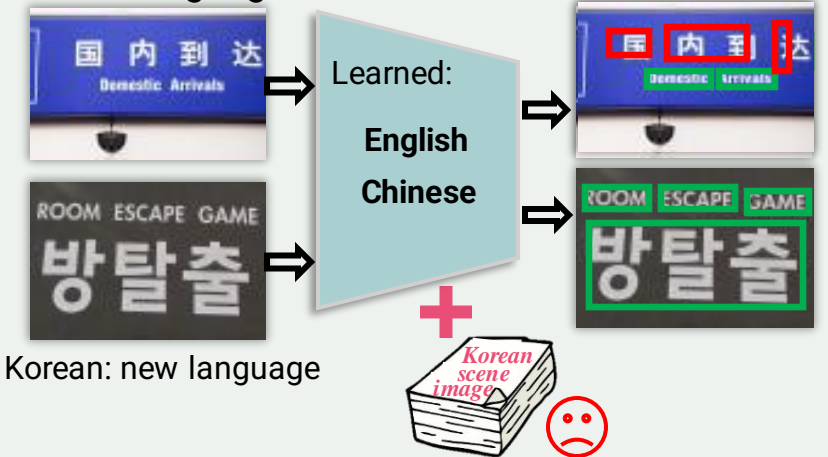
A **universal text dataset** containing all languages for supervised learning is **not available**.



Text training datasets that are **imbalanced between languages**.

Fine-tune Detector

trained languages



Korean: new language



The detection model must be **re-trained or fine-tune** to **detect new languages**.



Problem Setting



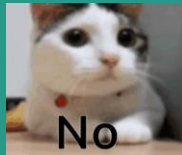
A generalizable multilingual text detector



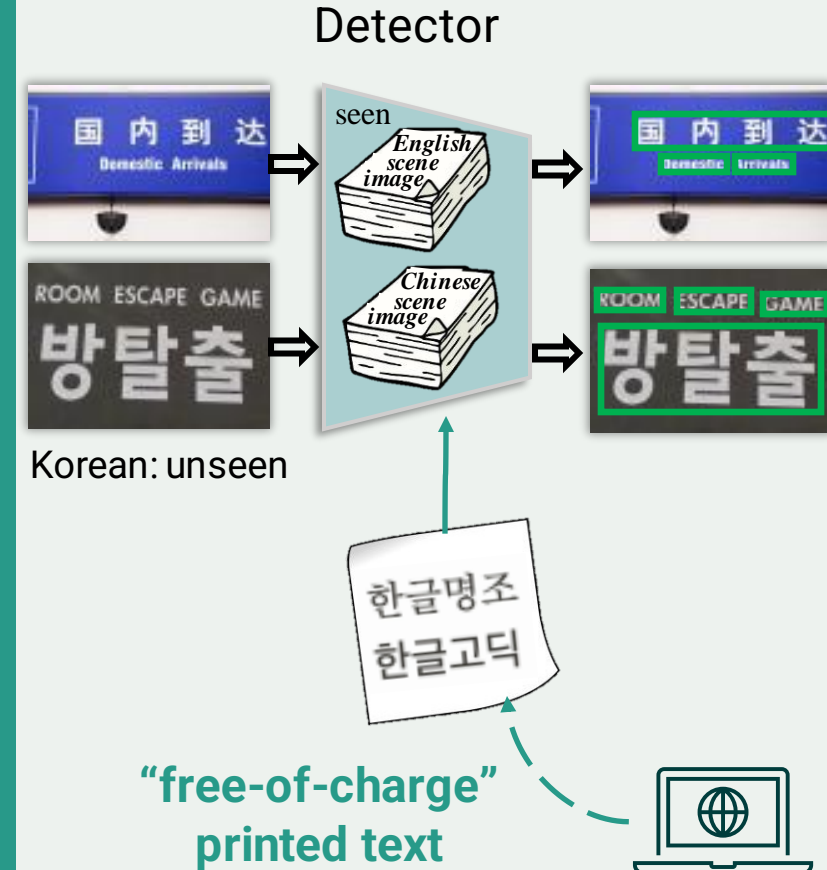
SHOULD detect both seen and unseen language regions



NO need to collect supervised training data for unseen language



NO model re-training required for unseen language



printed text
images
generation



11

[illegible]

**start
position**

synthetic text generation



2

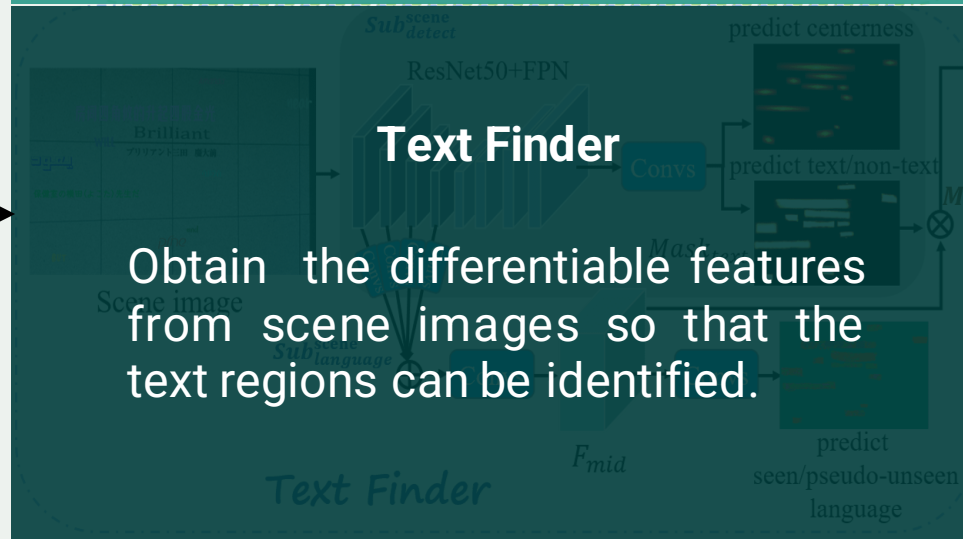
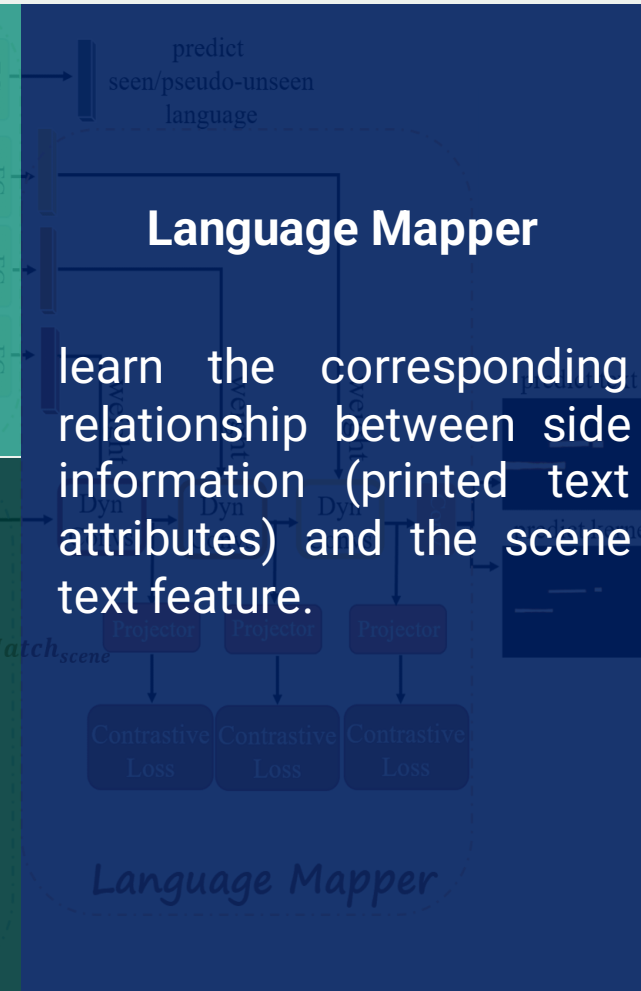
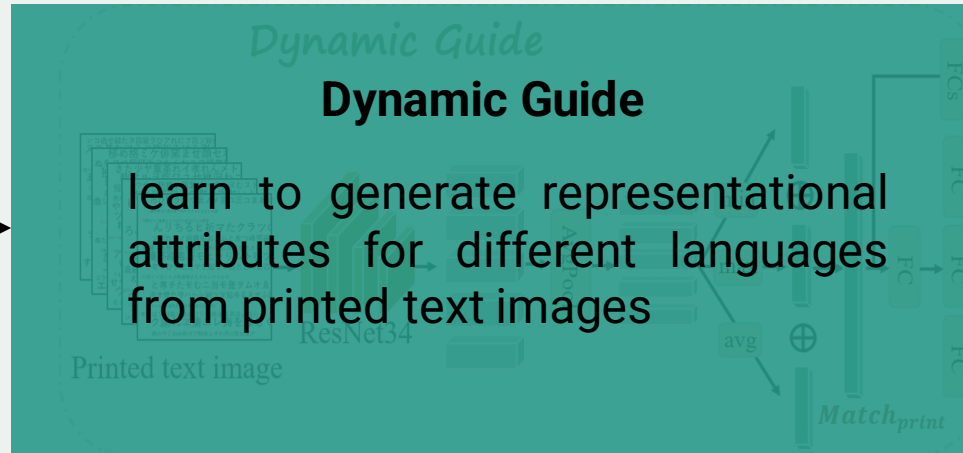
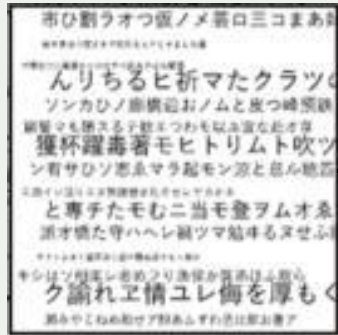
[illegible]

color

position



Proposed Method



Experimental Result

Quantitative evaluation (F1-score) on MLT17 dataset and Malayalam in IIIT-ILST dataset

Method	English	Arabic	Bangla	Chinese	Japanese	Korean	Malayalam
E2E-MLT [1]	55.43	55.431	3.027	50.594	12.9	32.715	x
MultiplexedOCR [2]	83.284	80.074	78.104	56.251	70.986	67.862	11.932
(a) Ours	84.031	80.952	81.76	83.585	76.896	72.479	65.837
(b) Ours	82.527	82.092	80.909	69.046	83.51	84.013	55.895
(c) Ours	81.778	54.646	82.134	82.554	80.737	83.143	43.992



[1] M.Buřta, Y.Patel, and J.Matas, “E2e-mlt-an unconstrained end-to-end method for multi- language scene text,” ArXiv:1801.09919, 2018.

[2] M.Buřta, Y.Patel, and J.Matas, “J. Huang, G. Pang, R. Kovvuri, M. Toh, K. J. Liang, P. Krishnan, X. Yin, and T. Hassner, “A multiplexed network for end-to-end, multilingual ocr,” In IEEE Conference on Computer Vision and Pattern Recognition (CVPR),, 2021.

Experimental Result

Quantitative evaluation (F1-score) on **synthetic** MLT17 dataset and **synthetic** Malayalam in IIIT-ILST dataset

Method	English	Arabic	Bangla	Chinese	Japanese	Korean	Malayalam
E2E-MLT [1]	50.19	54.67	4.027	55.138	22.346	34.317	x
MultiplexedOCR [2]	67.392	75.511	78.874	54.914	62.024	71.21	11.159
(a) Ours	71.76	74.949	79.044	74.351	71.556	50.534	55.607
(b) Ours	63.119	71.811	80.393	48.905	63.472	75.385	42.497
(c) Ours	63.424	46.061	80.328	80.328	59.254	75.385	36.586



[1] M.Buřta, Y.Patel, and J.Matas, “E2e-mlt-an unconstrained end-to-end method for multi- language scene text,” ArXiv:1801.09919, 2018.

[2] M.Buřta, Y.Patel, and J.Matas, “J. Huang, G. Pang, R. Kovvuri, M. Toh, K. J. Liang, P. Krishnan, X. Yin, and T. Hassner, “A multiplexed network for end-to-end, multilingual ocr,” In IEEE Conference on Computer Vision and Pattern Recognition (CVPR),, 2021.

Experimental Result

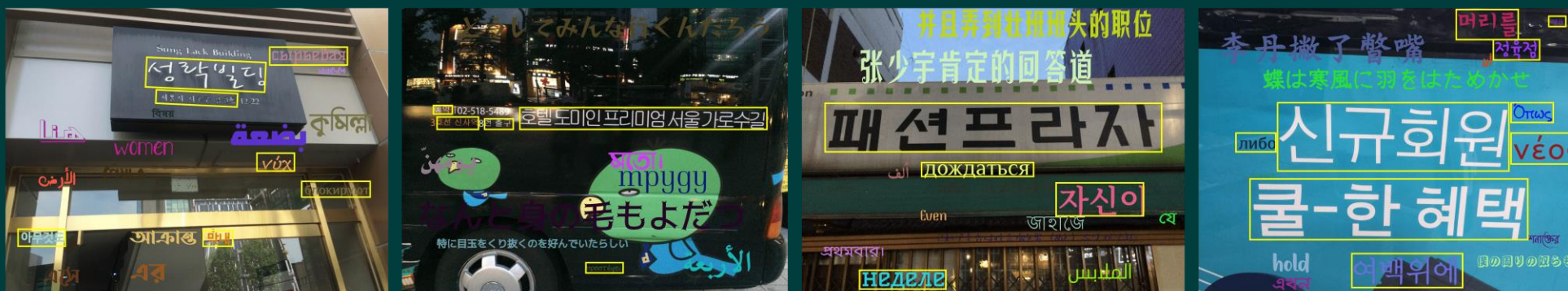
The visualization results of detecting language-agnostic text regions

- yellow bounding box means the text region belongs to unseen languages.

masking result



original image



Experimental Result

The detection results of synthetic MLT17 from different target languages



(a) Japanese



(b) Arabic



(c) Bangla



(d) Korean



(e) English



(f) Chinese (unseen)

Conclusion

- 1) We introduced a **new problem setting** for multilingual scene text detection and proposed a novel method “MENTOR” to **deal with text detection for both seen and unseen languages**.
- 2) Our “**Dynamic Guide**”, a dynamic and learnable module, and “**Text Finder**” module can **extract language specific features for seen and unseen languages** from printed text image and scene image, respectively.
- 3) We can identify the text regions of unseen languages via our “**LM**” module, by **comparing pixel-wise scene text features with** language-specific printed text feature, that is, **its mentor**.
- 4) Experiments show that our “MENTOR” can **achieve comparable results** with supervised methods **for seen languages** and **outperforms** other methods **in detecting unseen languages**.



Thank you for listening

Refer to the github for more detail



SCAN ME

