

# Dictionary-guided Scene Text Recognition: Supplementary material

Nguyen Nguyen<sup>1</sup>, Thu Nguyen<sup>1,2,4</sup>, Vinh Tran<sup>6</sup>,  
Minh-Triet Tran<sup>3,4</sup>, Thanh Duc Ngo<sup>2,4</sup>, Thien Huu Nguyen<sup>1,5</sup>, Minh Hoai<sup>1,6</sup>

<sup>1</sup>VinAI Research, Hanoi, Vietnam; <sup>2</sup>University of Information Technology, VNU-HCM, Vietnam;

<sup>3</sup>University of Science, VNU-HCM, Vietnam;

<sup>4</sup>Vietnam National University, Ho Chi Minh City, Vietnam;

<sup>5</sup>University of Oregon, Eugene, OR, USA; <sup>6</sup>Stony Brook University, Stony Brook, NY, USA

{v.nguyennm, v.thunm15, v.thiennh4, v.hoainm}@vinai.io

thanhnhd@uit.edu.vn, tmtriet@fit.hcmus.edu.vn, tquangvinh@cs.stonybrook.edu

This supplementary material contains additional qualitative results that could be included in the main paper due to page limit.

As explained in the main paper, using the naive way of using a dictionary, where the dictionary word with the smallest edit distance to the intermediate result is taken as the final output, does not always work. Fig. 1 shows two cases where this naive approach fails, while the proposed method correctly recognize the text instances in the photos.

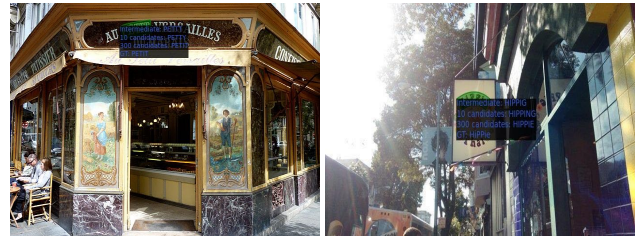


(a) Intermediate: CAFO	(b) Intermediate: Ylower/Tower
Closest: CFO	Closest: lower/lower
Proposed: CAFE	Proposed: Flower/Power
Ground-truth: CAFE	Ground-truth: Flower/Power

Figure 1: Several cases where our proposed approach of using the dictionary leads to correct recognition performance, while the naive approach of casting the intermediate result to a dictionary word fails.

As reported in the main paper, the performance the proposed method ABCNet+D increases as the number of candidates considered during inference increases. This can be attributed to its compatibility scoring model to calculate the compatibility score between the visual feature and the candidates. Fig. 2 shows several cases where having the candidate set of size 10 is not enough, but the text instances are correctly recognized when we increase the size of the

candidate sets to 300.



(a) Intermediate: PETTT	(b) Intermediate: HIPPIG
10 candidates: PETTY	10 candidates: HIPPIING
300 candidates: PETIT	300 candidates: HIPPIE
Ground-truth: PETIT	Ground-truth: HIPPIE



(c) Intermediate: HARBORT	(d) Intermediate: TOWAV
10 candidates: HARBORS	10 candidates: TWAT
300 candidates: HARBOR	300 candidates: COWAY
Ground-truth: HARBOR	Ground-truth: COWAY

Figure 2: The accuracy of the proposed method correlates with the size of the candidate set considered during inference. As shown in the main paper, using the candidate set size 10 improves the recognition performance significantly. However, sometimes, 10 might not be enough; the model can correctly recognize the text instances with 300 candidates but not 10.

The remaining set of figures display representative detection and recognition results of ABCNet+D on several randomly chosen images from the datasets considered in this paper.









