# Scene Text Recognition with Image-Text Matching-guided Dictionary

Jiajun Wei[1], Hongjian Zhan[1,2], Xiao Tu[1], Yue Lu[1]⋆, and Umapada Pal[3]

[1] Shanghai Key Laboratory of Multidimensional Information Processing, East China Normal University, Shanghai, China
[2] Chongqing Institute of East China Normal University. Chongqing. 401120. China
[3] CVPR Unit, Indian Statistical Institute, Kolkata, India
jjwei@stu.ecnu.edu.cn, ecnuhjzhan@foxmail.com, xtu@cee.ecnu.edu.cn,
ylu@cs.ecnu.edu.cn, umapada@isical.ac.in

**Abstract.** Employing a dictionary can efficiently rectify the deviation between the visual prediction and the ground truth in scene text recognition methods. However, the independence of the dictionary on the visual features may lead to incorrect rectification of accurate visual predictions. In this paper, we propose a new dictionary language model leveraging the **S**cene **I**mage-**T**ext **M**atching(SITM) network, which avoids the drawbacks of the explicit dictionary language model: 1) the independence of the visual features; 2) noisy choice in candidates etc. The SITM network accomplishes this by using Image-Text Contrastive (ITC) Learning to match an image with its corresponding text among candidates in the inference stage. ITC is widely used in vision-language learning to pull the positive image-text pair closer in feature space. Inspired by ITC, the SITM network combines the visual features and the text features of all candidates to identify the candidate with the minimum distance in the feature space. Our lexicon method achieves better results(93.8% accuracy) than the ordinary method results(92.1% accuracy) on six mainstream benchmarks. Additionally, we integrate our method with ABINet and establish new state-of-the-art results on several benchmarks.

**Keywords:** Dictionary Language Model · Scene Image-Text Matching · Image-Text Contrastive Learning · Scene Text Recognition.

## 1 Introduction

Deep learning-based scene text recognition has been developed for years. The accuracy of scene text recognition has vastly increased as the appropriate design of model architecture and the expansion of model size. Previous methods[4,37,5,50] can address a variety of recognition issues, but the inherent ambiguities, such as complicated background or diversity of font, etc, render the recognized results inaccurate.

---

⋆ Corresponding author

Due to the unique characteristics of text recognition, it is feasible to employ human language priors to rectify the output of a vision recognition model. Utilizing a pre-trained language model is one of the common methods. Fang et al[7] pre-train a language model using WikiText-103[28]. The pre-trained language model rectifies the visual prediction through learning grammar and the construction of words in the human language system. Another popular approach is to search for a word that has minimum edit distance(Levenshtein distance[19]) with the visual prediction in a dictionary. Nguyen et al[31] present a method for incorporating a dictionary into the training pipeline. They use the dictionary to generate a certain number of candidates and then output the most compatible one with the highest compatibility scores in a probability matrix $\mathbf{P}$, which is generated by the visual feature $\boldsymbol{F}_v$. But they still disregard the interaction between visual features and text features in the inference stage.

The aforementioned methods utilizing explicit language models have several problems. First, regardless of the pre-trained language model or dictionary language model, the independence of the language model from the visual feature may erroneously rectify the correct prediction results. Second, it is illogical to utilize human language priors to rectify texts that appear infrequently or have no linguistic information(*e.g.* ngee, tsc), since neither a pre-trained language model nor a dictionary can rectify texts without human language logic.
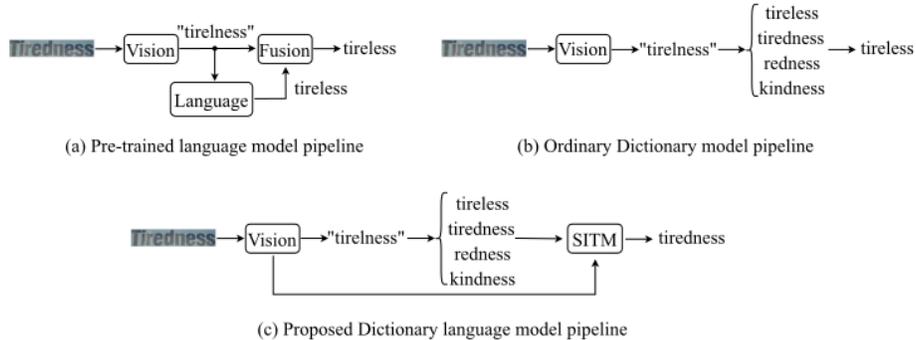


**Fig. 1.** Comparison with different pipelines.

In this paper, we present an effective method for incorporating a dictionary into scene text recognition that possesses two advantages: 1) taking visual prediction into account. When generating candidates for the visual prediction, the visual prediction is also included in the candidate set. In this case, the initial visual prediction has a probability to be the ultimate outcome. 2) Integrating visual features into the inference stage. Image-Text Contrastive(ITC) Learning is an unsupervised learning method that aims to make positive image-text pairs have higher similarity scores. Inspired by ITC learning, we additionally inte-

grate a **S**cene **I**mage-**T**ext **M**atching(SITM) network to match visual features and text features by ITC Learning. Nevertheless, when we merely employ other label texts in the same batch as negatives, the image-text matching accuracy in the inference stage is not as good as the training stage. We address the problem by generating hard negatives that resemble the shape of label texts. The difference between three methods is depicted in Fig. 1

The main contributions of this paper are summarized as follows: 1) we propose a novel method to integrate a dictionary into scene text recognition that avoids the drawbacks of an ordinary dictionary language model. 2) We also offer a new strategy that employs labels to generate resemblant words as hard negatives in the SITM training stage. 3) A Scene Image-Text Matching Module is introduced, which matches positive image-text pairs in the inference stage.

## 2 Related Work

### 2.1 Scene Text Recognition

**Language-free Methods.** Language-free approaches typically provide a prediction based on visual features, regardless of context information. CTC-based methods[8] utilize CNN to extract visual features, RNN to model sequence features, and CTC loss to train the entire recognition network end-to-end [10,39,11]. Segmentation-based methods[22] segment each character region before classifying and recognizing. The recognition results of all character areas compose the entire text sequence[27,48,42]. However, due to the absence of context information interaction, these approaches cannot attain exceptional performance.

**Language-based Methods.** In previous works, [13,14] use explicit language models to improve model recognition accuracy. CNN is employed to extract visual features to predict bags of N-grams of text strings. Recently, [7] regards the explicit language model as a spell checker to rectify visual prediction results. Some implicit language-based approaches connect visual features with context information by utilizing RNN[18,38] or attention mechanisms [45,36]. First, an image encoder is employed to extract features from word images, follower by an attention-based method for integrating visual features and context information. [37,5,6,4] focus on relevant information from 1D image features, and [50,47,23,20] from 2D image features. Some performance-enhancing approaches focus on learning new feature representations. [1,26,49] train their models by sequence contrastive learning, masked image modeling, and a mix of the two, respectively.

### 2.2 Vision-Language Learning

There are two categories of visual-language representation learning. In the first category, text features and image features are fused using a multi-mode encoder [40,24,25,21]. This type of approaches has achieved outperformance in

downstream tasks such as NLVR[41] and VQA[2]. The Second category focus on learning separate texts and images encoders[33,15]. CLIP[34] employs contrastive loss to train the image encoder and text encoder on a massive quantity of network image-text pairs. We opt for the second category to reuse the visual encoder trained in the recognition stage instead of training a new visual encoder. Then, the text encoder is trained from scratch in the matching stage.

## 3   Method

We propose a new method to incorporate a dictionary into scene text recognition. The dictionary is used to generate the certain number of candidates, which will subsequently be matched with the visual features by SITM to output the candidate with the highest similarity score. In this section, the details of the overall architecture are presented. We will also describe the Resemblant word generation strategy and the SITM network. The objective training function is finally introduced.
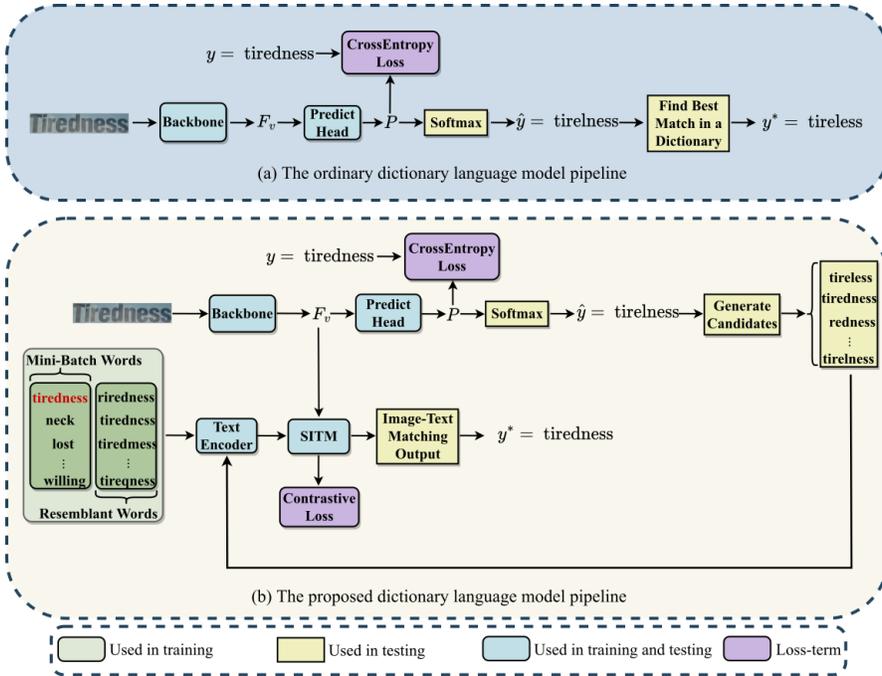


**Fig. 2.** Ordinary dictionary language model pipeline (a) and proposed dictionary language model pipeline (b). In the ordinary pipeline, the prediction is forced to be one with the smallest edit distance in the dictionary. In the proposed pipeline, the ultimate prediction is determined by SITM

### 3.1   Overall Architecture

As can be seen in Fig. 3, a general scene text recognition framework usually consists of a feature extraction module, a sequence modeling module, and a prediction module, which was proposed by Baek et al[3]. Our proposed dictionary method can combine any scene text recognition method with the above framework. We utilize the output of the sequence modeling as the visual features $F_v$. Specifically in this paper, we employ the vision module of the ABINet[7] as our baseline network.
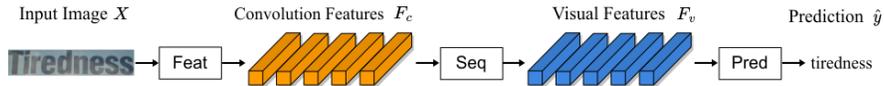


**Fig. 3.** General Vision Architecture of Scene Text Recognition.

Fig. 2b describes the procedure of our proposed recognition pipeline, which employs a forward-forward method. For the initial forward, we input the image, generate visual prediction $\hat{y}$ with the text recognition network, and then utilize $\hat{y}$ to find candidates with the top N smallest edit distance in the dictionary. The candidate set is comprised of the N candidates and the visual prediction $\hat{y}$. For example, if $\hat{y} = tirelness$, the candidates will be: *tireless, tiredness, redness, ..., kindness, tirelness*. For the second forward pass, the inputs are the image and the candidates obtained from the first forward. Then, utilizing the SITM network to match the text in the candidates with the image, and output the text with the highest similarity. The inference procedure is depicted in Algorithm 1.

---

**Algorithm 1** Inference procedure

---

**Input:** $x$: Input Image; $n$: Forward State
**Output:** Prediction $y^*$

1: initial $n = 1$
2: **for** $n = 1, 2$ **do**
3:     **if** $n = 1$ **then**
4:         Input $x$ to get $\hat{y} = \mathbf{V}(x)$, where $\mathbf{V}$ is the vision module
5:         Construct candidate set $C$ of $\hat{y}$ using dictionary
6:     **else**
7:         Input $x$ and $C$ to calculate the similarity scores $S$
8:         Get the $c \in C$ with the highest score in $S$ as the final prediction $y^*$
9: **return** $y^*$

---

During training, we generate text candidates for Image-Text Contrastive Learning by using labels in the same mini-batch. In addition, we create a certain

number of resemblant words as hard negatives. A contrastive loss function, which is defined based on the similarity cross-entropy function depicted in Section 3.3, and the recognition loss are then employed for training.

## 3.2  Resemblant Words Generation

| Image | Label | Hard Dictionary Negatives |
|---|---|---|
| Tiredness | tiredness | 'redness' 'tireless' 'kindness' 'likeness' 'redress' 'timeless' 'nakedness' |
| short | short | 'sort' 'shot' 'shorts' 'shirt' 'sport' 'snort' 'shout' |
| break | break | 'bleak' 'bread' 'creak' 'freak' 'wreak' 'beak' 'beau' |
| could | could | 'cold' 'would' 'bold' 'gourd' 'cod' 'should' 'court' |

**Fig. 4.** Qualitative hard negatives in the inference stage.

There is a gap between the SITM training stage and the inference stage if we utilize the normal contrastive learning method. Specifically, in the training stage, the negatives are the other labels in the same mini-batch for a single image. For example, if a mini-batch contains: *tiredness, kills, short, break, could, save, **your**, life*, the text negatives are *tiredness, kills, short, break, could, save, life* and the text positive is ***your*** for the image ***your***. However, in the inference stage, the negatives are candidates from the dictionary with the top N smallest edit distance, which is similar to the ground truth. For example, for the image ***your***, the negatives in the inference stage are  *pour, you, tour, hour, dour, sour, four*. Fig. 4 exhibits some hard negatives for the labels in the test set. We find the gap would cause some mismatches between image-text pairs in the inference stage and degrade the performance of the dictionary.

We address the problem with our proposed Resemblant Words Generation strategy. When we train the SITM network, in addition to using the text labels corresponding to other images in the same batch as negatives, we present a strategy for constructing hard negatives using labels. Specifically, we initially establish a similar character lookup table containing five similar characters for each English character. We observe the difference between visual prediction and ground truth and record the wrong predicted characters as the composition of the lookup table. For character $a$, we select $d, e, o, q$ and $u$ as the similar characters. Then we randomly replace a character in a label having a similar appearance. For example, if y $= tiredness$ and the number of the resemblant words is 4, the hard negatives will be *riredness*, *tiredncss*, *tiredmess*, *tireqness*.

### 3.3 Scene Image-Text Matching Module



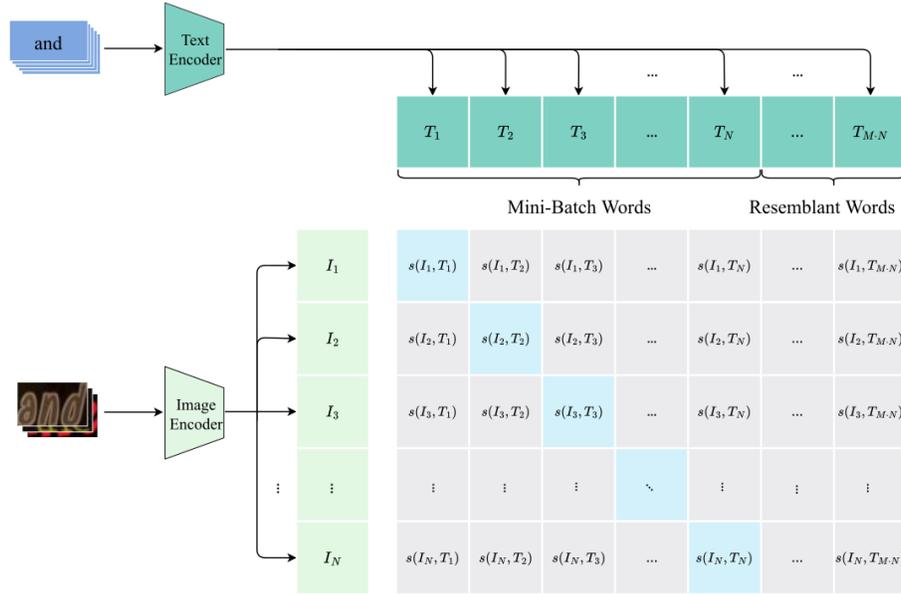| | | Mini-Batch Words | | | | Resemblant Words | |
|---|---|---|---|---|---|---|---|
| | $T_1$ | $T_2$ | $T_3$ | ... | $T_N$ | ... | $T_{M \cdot N}$ |
| $I_1$ | $s(I_1,T_1)$ | $s(I_1,T_2)$ | $s(I_1,T_3)$ | ... | $s(I_1,T_N)$ | ... | $s(I_1,T_{M \cdot N})$ |
| $I_2$ | $s(I_2,T_1)$ | $s(I_2,T_2)$ | $s(I_2,T_3)$ | ... | $s(I_2,T_N)$ | ... | $s(I_2,T_{M \cdot N})$ |
| $I_3$ | $s(I_3,T_1)$ | $s(I_3,T_2)$ | $s(I_3,T_3)$ | ... | $s(I_3,T_N)$ | ... | $s(I_3,T_{M \cdot N})$ |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\vdots$ |
| $I_N$ | $s(I_N,T_1)$ | $s(I_N,T_2)$ | $s(I_N,T_3)$ | ... | $s(I_N,T_N)$ | ... | $s(I_N,T_{M \cdot N})$ |

**Fig. 5.** Scene Image-Text Matching Module Architecture.

The Scene Image-Text Matching module contains image encoder, text encoder and Scene Image-Text Contrastive Learning module. The image encoder consists of a backbone network that shares parameters with the backbone network of the recognition module and a parallel attention layer that is used to convert visual features to sequence features. The text encoder consists of two layers of transformer encoder. Scene Image-Text Contrastive Learning module consists of two liner layers. Image sequence features $I \in \mathbb{R}^{L \times C}$ and text sequence features $T \in \mathbb{R}^{L \times C}$ are obtained by image encoder and text encoder, respectively. The image features $I$ and text features $T$ pass through the linear projection layer before Image-Text Contrastive Learning. Fig. 5 shows the details of our SITM module.

During the training stage, we employ the image-text contrastive learning task in vision-language learning to complete the scene image-text matching task. Image-Text Contrastive Learning aims to learn a representation of distinct modal features. It learns the cosine similarity function $s(I, T) = \frac{l_v(I) \cdot l_t(T)}{||l_v(I)|| \cdot ||l_t(T)||}$, where $l$ represents linear layer and $I \in \mathbb{R}^{L \times C}$, $T \in \mathbb{R}^{L \times C}$ represent image features

and text features, respectively. The matched image-text pair will have a higher similarity score. We calculate image-to-text(i2t) and text-to-image(t2i) similarity and normalize the results using softmax. The formulas are as below:

$$p_m^{i2t}(I) = \frac{\exp(s(I, T_m)/\tau)}{\sum_{n=1}^{MN} \exp(s(I, T_n)/\tau}, \quad p_m^{t2i}(T) = \frac{\exp(s(T, I_m)/\tau)}{\sum_{n=1}^{N} \exp(s(T, I_n)/\tau}, \quad (1)$$

where $\tau$ is a temperature parameter, $m$ is the order indication of the image or text, $s$ is the cosine similarity function and N is the batch size and M-1 is the number of resemblant words of one label. As can be seen in Fig. 6, the parallel attention layer focuses on the main character features in the image to guide the matching procedure.
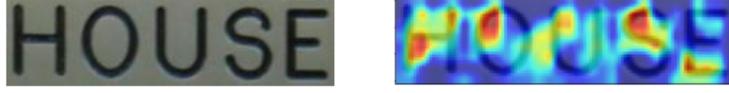


**Fig. 6.** An example of the parallel attention layer Gradient-weighted Class Activation Mapping. The left image is the input and the right image is the activation mapping.

During the inference stage, we only calculate the image-to-text(i2t) similarity to find the highest score among candidates from the candidate set:

$$p_m^{i2t}(I) = \frac{\exp(s(I, T_m)/\tau)}{\sum_{n=1}^{T} \exp(s(I, T_n)/\tau}, \quad (2)$$

where T is the number of candidates in the candidates set which is depicted in Section 3.1.

### 3.4   Overall Objective Function

A supervised cross-entropy loss function is utilized for training and optimizing the scene text recognizer. By minimizing the negative log-likelihood sequence probability loss function, the difference between the prediction and the ground truth is quantified. The specific formula is given below:

$$\mathcal{L}_{recog} = \boldsymbol{E}_{(x,y)\sim(X,Y)}\{-\mathrm{log}p(y|F_v(x))\}. \quad (3)$$

The full loss function of the Scene Image-Text Matching module consists of $\mathcal{L}_{itc}$

$$\mathcal{L}_{SITM} = \mathcal{L}_{itc}, \quad (4)$$

$$\mathcal{L}_{itc} = \frac{1}{2}\boldsymbol{E}_{(I,T)\sim D}\big[\mathrm{H}(\boldsymbol{y}^{i2t}(I), \boldsymbol{p}^{i2t}(I)) + \mathrm{H}(\boldsymbol{y}^{t2i}(T), \boldsymbol{p}^{t2i}(T))\big], \quad (5)$$

where $\boldsymbol{y}^{i2t}(I)$ and $\boldsymbol{y}^{t2i}(T)$ denote the ground-truth one-hot similarity, in which the negative pair probability is 0 and the positive pair probability is 1. The H is defined as the cross-entropy loss.

The overall objective function $\mathcal{L}_{overall}$ is defined as:

$$\mathcal{L}_{overall} = \lambda_1 \mathcal{L}_{recog} + \lambda_2 \mathcal{L}_{SITM}, \tag{6}$$

where $\lambda_1$ and $\lambda_2$ are the hyper-parameters used to control the training stages. We respectively set $\lambda_1 = 1, \lambda_2 = 0$ and $\lambda_1 = 0, \lambda_2 = 1$ when we train the recognition module and Scene Image Text Matching module.

## 4  Experimental Results

### 4.1  Datasets

The common synthetic datasets SynthText[9] and MJSynth[12] are utilized to train our proposed model. We employ six widely used benchmarks to evaluate the performance of the model, including three regular text datasets ICDAR2013, SVT, IIIT5K and three irregular text datasets ICDAR2015, SVTP and CUTE80. Following are the specifics of the datasets:

ICDAR2013(IC13)[17] has 1015 test images. The dataset contains only horizontal text instances.

Street View Text (SVT)[44] contains 647 images collected from Google Street View. This dataset contains fuzzy, blurry, and low-resolution text images.

IIIT5K[29] contains 3000 test images crawled from Google image searches with query words. Most text instances are rules for horizontal layout.

ICDAR2015(IC15)[16] contains 1811 test images created for the ICDAR 2015 Robust Reading competitions. Most instances of text are irregular (noisy, blurry, perspective or curved).

Street View Text Perspective (SVTP)[32] contains 645 cropped images from Google Street View. Many of the images have a distorted perspective.

CUTE80[35] is collected from nature scenes and contains 288 cropped images for verification. Most of them are curved text.

### 4.2  Training Setting

PyTorch is applied to implement the model proposed in this paper. All the experiments are conducted on a 24GB-memory NVIDIA3090. All input images are scaled to $32 \times 128$ while maintaining their aspect ratio.The character set includes 37 classes, which contains 10 digits, 26 lowercase letters, and an EOS token. The maximum sequence length is 25. Adam is selected as the optimizer, and the batch size is set to 320.

The training procedure consists of two stages: the recognition stage and the matching stage. In the recognition stage, the text recognition network is merely trained to minimize the text recognition loss function. We trained the recognition network 8 epochs on SynthText and MJSynth from scratch. During the matching stage, the SITM network is unfrozen. Image-Text contrastive loss is applied to train the text encoder.

### 4.3   Comparison with the Ordinary Dictionary Method and the State-of-the-art

In this part, we compare the accuracy of the baseline, ordinary dictionary-guided baseline and the proposed dictionary-guided baseline on the six benchmarks. The baseline described in Section 3.1 serves as a comparison standard in our experiments. We utilize the same lexicon, which comprises approximately 20K words and is composed of numbers, common English words and common English trademarks. *Full Lexicon* is not utilized to construct the dictionary, which means that some words in the test set may not be in the lexicon. As we consider this would be a more realistic dictionary composition with some words included and some were excluded. For a fair comparison, all the methods are trained on the SynthText and MJSynth datasets.

**Table 1.** Comparison with Ordinary Dictionary-guided Baseline.

| Method | Regular Text | | | Irregular Text | | | Average |
|---|---|---|---|---|---|---|---|
| | IC13 | SVT | IIIT5K | IC15 | SVTP | CUTE80 | |
| Baseline | 94.9 | 90.4 | 94.6 | 81.7 | 84.2 | 86.5 | 89.8 |
| Baseline+Dict Guided | 95.8 | 92.1 | 96.2 | 85.6 | 87.4 | 90.8 | 92.1 |
| Baseline+Our Method | **97.8** | **94.1** | **97.1** | **88.0** | **89.3** | **93.4** | **93.8** |
| Improvement | **+2.0** | **+2.0** | **+0.9** | **+2.4** | **+1.9** | **+2.6** | **+1.7** |

As can be seen from Table 1, utilizing ordinary dictionary guidance would enhance performance, but the improvement on some benchmarks is insignificant. On six benchmarks, our proposed dictionary-guided method outperforms the ordinary method with 2.0%, 2.0%, 0.9%, 2.4%, 1.9% and 2.6% on IC13, SVT, IIIT5K, IC15, SVTP and CUTE80 datasets, respectively. We also discover that our method has superiority on irregular datasets IC15, SVTP and CUTE80 as they contain low-quality images such as curved and blurred texts. As the visual prediction of the irregular datasets often have more severe deviation from the ground truth, the candidate with smallest edit distance may not be the correct answer.

The weaker performance of the ordinary dictionary method stems from two aspects. 1) It disregards visual prediction. The ordinary dictionary pipeline takes the word in the dictionary as output with the smallest edit distance for the visual prediction that is not a component of the dictionary, which makes the correct prediction incorrect. 2) For words in the dictionary with the same edit distance, the traditional dictionary pipeline is unable to determine which output is right. The random selection will fail to output the correct outcome among the candidates.

Our proposed pipeline effectively avoids the aforementioned issues. In addition to the text recognition network, we also train a SITM network. When a dictionary is employed, we combine the prediction and the top N smallest edit distance dictionary words as candidates set, and the SITM network is used to

| Input |  |  |  |  |
|---|---|---|---|---|
| Label | ronaldo | ebizu | bud | finest |
| Baseline | ronaldo | ebizu | bod | vinest |
| Ordinary Dict | renal | biz | bold | vines |
| Proposed Dict | ronaldo | ebizu | bud | finest |

**Fig. 7.** Qualitative results for the ordinary dictionary method and our proposed method.

determine which one is correct. Fig. 7 illustrates instances successfully recognized by our method while ordinary dictionary method could not. The second and third columns represent that the visual prediction of the scene text recognition network is correct, but there is no corresponding in the dictionary. In this case, the ordinary method generates the wrong answer. However, our proposed method can find the visual prediction output in candidates. The fourth and fifth columns represents the deviation between the visual predictions and the ground truths. When facing candidates with the same edit distance, the ordinary method can only randomly output, while our proposed method can find the correct candidate word according to the SITM network.

**Table 2.** Comparison with State-of-the-art Methods and Ordinary Dictionary-guided State-of-the-art Methods.

| Methods | Ordinary Dict Guide | Proposed Dict Guide | Regular Text | | | Irregular Text | | |
|---|---|---|---|---|---|---|---|---|
| | | | IC13 | SVT | IIIT5K | IC15 | SVTP | CUTE80 |
| PlugNet[30] | - | - | 95.0 | 92.3 | 94.4 | 82.2 | 84.3 | 85.0 |
| SRN[51] | - | - | 95.5 | 91.5 | 94.8 | 82.7 | 85.1 | 87.8 |
| RobustScanner[52] | - | - | 94.1 | 89.3 | 95.4 | 79.2 | 82.9 | 92.4 |
| TextScanner[43] | - | - | 94.9 | 92.7 | 95.7 | 83.5 | 84.8 | 91.6 |
| AutoSTR[53] | - | - | 94.2 | 90.9 | 94.7 | 81.8 | 81.7 | - |
| VisionLAN[46] | - | - | 95.7 | 91.7 | 95.8 | 83.7 | 86.0 | 88.5 |
| CRNN[3] | - | - | 88.8 | 78.9 | 84.3 | 61.5 | 64.8 | 61.3 |
| ABINet[7] | - | - | **97.4** | 93.5 | 96.2 | 86.0 | **89.3** | 89.2 |
| PARSeq[4] | - | - | 97.0 | **93.6** | **97.0** | **86.5** | 88.9 | **92.2** |
| CRNN[3] | ✓ | | 95.2 | 90.8 | 91.5 | 83.0 | 84.0 | 78.5 |
| CRNN[3] | - | ✓ | **96.9** | **92.1** | **93.2** | **84.6** | **88.5** | **80.9** |
| ABINet[7] | ✓ | | 97.7 | 94.1 | 96.8 | 87.5 | 90.0 | 90.3 |
| ABINet[7] | - | ✓ | **98.4** | **95.8** | **98.0** | **88.6** | **90.1** | **91.3** |

To verify the effectiveness of our method, we combine two existing scene text recognition frameworks with our proposed dictionary method. We select

the CRNN and the state-of-the-art ABINet to validate our proposed approach.

Table 2 shows that our proposed method still outperforms the ordinary dictionary method. As can be seen from the comparison, in the CRNN[3], our proposed dictionary-guided method outperforms the ordinary method with 1.7%, 1.3%, 1.7%, 1.6%, 4.5% and 2.4% on IC13, SVT, IIIT5K, IC15, SVTP, CUTE80 datasets, respectively. In ABINet[7], the improvements on the six benchmarks are 0.7%, 1.7%, 1.2%, 1.1%, 0.1% and 1.0%, respectively. In the meanwhile, we find that the utilization of a dictionary to rectify the visual prediction is a highly effective way of enhancing performance. When employing a dictionary to rectify visual prediction, the CRNN[3] exceeds numerous state-of-the-art methods on some benchmarks.

## 4.4   Ablation Study

**Table 3.** Comparison of recognition accuracy on different numbers of candidates.

| candidates | 1 | 5 | 10 | 20 | 30 | 80 | 150 | 300 |
|---|---|---|---|---|---|---|---|---|
| Recognition Accuracy | 92.1 | 93.8 | 94.1 | 94.2 | 94.2 | 94.3 | 94.3 | 94.3 |

**The recognition accuracy of baseline as the numbers of candidates varies:** The quantity of candidates is one of the primary distinctions between our approach and the ordinary pipeline. Table 3 demonstrates how the amount of candidate words affects the accuracy of the pipeline. The second column, when the number of candidates is 1, corresponds to the ordinary dictionary-guided method. As can be observed, a substantial improvement of 2% in accuracy occurs when the candidate number increases from 1 to 10, which explains that the correct word is not necessarily the one with the smallest edit distance. The average accuracy marginally improves as the number of candidates increases from 10 to 80. The saturation appears when the number arrives at 150. Table 3 illustrates the primary benefit of the proposed method, which can select the correct output from a group of options.

**Table 4.** Comparison of recognition accuracy on different numbers of resemblant words.

| Number of resemblant words | 0 | 3 | 7 | 15 | 31 |
|---|---|---|---|---|---|
| Recognition Accuracy | 91.1 | 93.8 | 93.9 | 93.9 | 93.9 |

**The discussion of resemblant word function:** For image-text pairs to be successfully matched, a certain number of hard negatives are included in the

training process. To illustrate the efficacy of this strategy, we arrange a variety of resemblant words: 0, 3, 7, 15 and 31. The recognition accuracy of the entire pipeline is shown in Table 4. The second column 0 indicates that no hard negative is used. It can be seen that recognition accuracy improves as the number of hard negatives increases. However, it will not be improved until a certain number of hard negatives has been accumulated. In contrast, when the number of hard negatives is equal to 0, the SITM network cannot complete the image-text matching task, therefore some incorrect matching pairs are produced. The performance(91.1% accuracy) is significantly worse than the ordinary dictionary method performance(92.1% accuracy). Table 4 demonstrates that the model is capable of learning more fine-grained distinctions between different text features through resemblant word generation strategy.

## 5    Conclusion

In this paper, we propose a new dictionary-guided scene text recognition method, which integrates the visual features into the inference stage and can effectively boost the performance of dictionary language model. In addition, the SITM is designed to indicate the correctness of explicit language model rectification. The resemblant words generation strategy, which utilizes labels to generate hard negatives in the training stage, is presented to improve the matching accuracy of SITM network. The experiments on six mainstream benchmarks demonstrate that our method outperforms the ordinary dictionary method and also show superiority in other state-of-the-art scene text recognition methods.

## Acknowledgement

## References

1. Aberdam, A., Litman, R., Tsiper, S., Anschel, O., Slossberg, R., Mazor, S., Manmatha, R., Perona, P.: Sequence-to-sequence contrastive learning for text recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15302–15312 (2021)
2. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: Vqa: Visual question answering. In: Proceedings of the IEEE international conference on computer vision. pp. 2425–2433 (2015)
3. Baek, J., Kim, G., Lee, J., Park, S., Han, D., Yun, S., Oh, S.J., Lee, H.: What is wrong with scene text recognition model comparisons? dataset and model analysis. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 4715–4723 (2019)
4. Bautista, D., Atienza, R.: Scene text recognition with permuted autoregressive sequence models. In: European Conference on Computer Vision. pp. 178–196. Springer (2022)

5. Cheng, Z., Bai, F., Xu, Y., Zheng, G., Pu, S., Zhou, S.: Focusing attention: To-wards accurate text recognition in natural images. In: Proceedings of the IEEE international conference on computer vision. pp. 5076–5084 (2017)
6. Cheng, Z., Xu, Y., Bai, F., Niu, Y., Pu, S., Zhou, S.: Aon: Towards arbitrarily-oriented text recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5571–5579 (2018)
7. Fang, S., Xie, H., Wang, Y., Mao, Z., Zhang, Y.: Read like humans: Autonomous, bidirectional and iterative language modeling for scene text recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7098–7107 (2021)
8. Graves, A., Fernández, S., Gomez, F., Schmidhuber, J.: Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In: Proceedings of the 23rd international conference on Machine learning. pp. 369–376 (2006)
9. Gupta, A., Vedaldi, A., Zisserman, A.: Synthetic data for text localisation in natural images. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2315–2324 (2016)
10. He, P., Huang, W., Qiao, Y., Loy, C.C., Tang, X.: Reading scene text in deep convolutional sequences. In: Thirtieth AAAI conference on artificial intelligence (2016)
11. Hu, W., Cai, X., Hou, J., Yi, S., Lin, Z.: Gtc: Guided training of ctc towards efficient and accurate scene text recognition. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 11005–11012 (2020)
12. Jaderberg, M., Simonyan, K., Vedaldi, A., Zisserman, A.: Synthetic data and artificial neural networks for natural scene text recognition. arXiv preprint arXiv:1406.2227 (2014)
13. Jaderberg, M., Simonyan, K., Vedaldi, A., Zisserman, A.: Deep structured output learning for unconstrained text recognition. ICLR (2015)
14. Jaderberg, M., Vedaldi, A., Zisserman, A.: Deep features for text spotting. In: European conference on computer vision. pp. 512–528. Springer (2014)
15. Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: International Conference on Machine Learning. pp. 4904–4916. PMLR (2021)
16. Karatzas, D., Gomez-Bigorda, L., Nicolaou, A., Ghosh, S., Bagdanov, A., Iwamura, M., Matas, J., Neumann, L., Chandrasekhar, V.R., Lu, S., et al.: Icdar 2015 competition on robust reading. In: 2015 13th international conference on document analysis and recognition (ICDAR). pp. 1156–1160. IEEE (2015)
17. Karatzas, D., Shafait, F., Uchida, S., Iwamura, M., i Bigorda, L.G., Mestre, S.R., Mas, J., Mota, D.F., Almazan, J.A., De Las Heras, L.P.: Icdar 2013 robust reading competition. In: 2013 12th international conference on document analysis and recognition. pp. 1484–1493. IEEE (2013)
18. Lee, C.Y., Osindero, S.: Recursive recurrent nets with attention modeling for ocr in the wild. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2231–2239 (2016)
19. Levenshtein, V.I., et al.: Binary codes capable of correcting deletions, insertions, and reversals. In: Soviet physics doklady. vol. 10, pp. 707–710. Soviet Union (1966)
20. Li, H., Wang, P., Shen, C., Zhang, G.: Show, attend and read: A simple and strong baseline for irregular text recognition. In: Proceedings of the AAAI conference on artificial intelligence. vol. 33, pp. 8610–8617 (2019)

21. Li, L., Yatskar, M., Yin, D., Hsieh, C., Chang, K.: A simple and performant baseline for vision and language. arXiv preprint arXiv:1908.03557 (2019)
22. Li, Y., Qi, H., Dai, J., Ji, X., Wei, Y.: Fully convolutional instance-aware semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2359–2367 (2017)
23. Liao, M., Zhang, J., Wan, Z., Xie, F., Liang, J., Lyu, P., Yao, C., Bai, X.: Scene text recognition from two-dimensional perspective. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8714–8721 (2019)
24. Lu, J., Batra, D., Parikh, D., Lee, S.: Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. Advances in neural information processing systems **32** (2019)
25. Lu, J., Goswami, V., Rohrbach, M., Parikh, D., Lee, S.: 12-in-1: Multi-task vision and language representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10437–10446 (2020)
26. Luo, C., Jin, L., Chen, J.: Siman: Exploring self-supervised representation learning of scene text via similarity-aware normalization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1039–1048 (2022)
27. Lyu, P., Liao, M., Yao, C., Wu, W., Bai, X.: Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 67–83 (2018)
28. Merity, S., Xiong, C., Bradbury, J., Socher, R.: Pointer sentinel mixture models. arXiv preprint arXiv:1609.07843 (2016)
29. Mishra, A., Alahari, K., Jawahar, C.: Top-down and bottom-up cues for scene text recognition. In: 2012 IEEE conference on computer vision and pattern recognition. pp. 2687–2694. IEEE (2012)
30. Mou, Y., Tan, L., Yang, H., Chen, J., Liu, L., Yan, R., Huang, Y.: Plugnet: Degradation aware scene text recognition supervised by a pluggable super-resolution unit. In: European Conference on Computer Vision. pp. 158–174. Springer (2020)
31. Nguyen, N., Nguyen, T., Tran, V., Tran, M.T., Ngo, T.D., Nguyen, T.H., Hoai, M.: Dictionary-guided scene text recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7383–7392 (2021)
32. Phan, T.Q., Shivakumara, P., Tian, S., Tan, C.L.: Recognizing text with perspective distortion in natural scenes. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 569–576 (2013)
33. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning. pp. 8748–8763. PMLR (2021)
34. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International Conference on Machine Learning. pp. 8748–8763. PMLR (2021)
35. Risnumawan, A., Shivakumara, P., Chan, C.S., Tan, C.L.: A robust arbitrary text detection system for natural scene images. Expert Systems with Applications **41**(18), 8027–8048 (2014)
36. Sheng, F., Chen, Z., Xu, B.: Nrtr: A no-recurrence sequence-to-sequence model for scene text recognition. In: 2019 International conference on document analysis and recognition (ICDAR). pp. 781–786. IEEE (2019)
37. Shi, B., Wang, X., Lyu, P., Yao, C., Bai, X.: Robust scene text recognition with automatic rectification. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4168–4176 (2016)

38. Shi, B., Yang, M., Wang, X., Lyu, P., Yao, C., Bai, X.: Aster: An attentional scene text recognizer with flexible rectification. IEEE transactions on pattern analysis and machine intelligence **41**(9), 2035–2048 (2018)
39. Su, B., Lu, S.: Accurate recognition of words in scenes without character segmentation using recurrent neural network. Pattern Recognition **63**, 397–405 (2017)
40. Su, W., Zhu, X., Cao, Y., Li, B., Lu, L., Wei, F., Dai J, V.B.: Pre-training of generic visual-linguistic representations. In: Proceedings of the 8th International Conference on Learning Representations. pp. 1–14 (2020)
41. Suhr, A., Zhou, S., Zhang, A., Zhang, I., Bai, H., Artzi, Y.: A corpus for reasoning about natural language grounded in photographs. arXiv preprint arXiv:1811.00491 (2018)
42. Wan, Z., He, M., Chen, H., Bai, X., Yao, C.: Textscanner: Reading characters in order for robust scene text recognition. In: Proceedings of the AAAI conference on artificial intelligence. vol. 34, pp. 12120–12127 (2020)
43. Wan, Z., He, M., Chen, H., Bai, X., Yao, C.: Textscanner: Reading characters in order for robust scene text recognition. In: Proceedings of the AAAI conference on artificial intelligence. vol. 34, pp. 12120–12127 (2020)
44. Wang, K., Babenko, B., Belongie, S.: End-to-end scene text recognition. In: 2011 International conference on computer vision. pp. 1457–1464. IEEE (2011)
45. Wang, P., Yang, L., Li, H., Deng, Y., Shen, C., Zhang, Y.: A simple and robust convolutional-attention network for irregular text recognition. arXiv preprint arXiv:1904.01375 **6**(2),  1 (2019)
46. Wang, Y., Xie, H., Fang, S., Wang, J., Zhu, S., Zhang, Y.: From two to one: A new scene text recognizer with visual language modeling network. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14194–14203 (2021)
47. Wojna, Z., Gorban, A.N., Lee, D.S., Murphy, K., Yu, Q., Li, Y., Ibarz, J.: Attention-based extraction of structured information from street view imagery. In: 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR). vol. 1, pp. 844–850. IEEE (2017)
48. Xing, L., Tian, Z., Huang, W., Scott, M.R.: Convolutional character networks. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 9126–9136 (2019)
49. Yang, M., Liao, M., Lu, P., Wang, J., Zhu, S., Luo, H., Tian, Q., Bai, X.: Reading and writing: Discriminative and generative modeling for self-supervised text recognition. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 4214–4223 (2022)
50. Yang, X., He, D., Zhou, Z., Kifer, D., Giles, C.L.: Learning to read irregular text with attention mechanisms. In: IJCAI. vol. 1, p. 3 (2017)
51. Yu, D., Li, X., Zhang, C., Liu, T., Han, J., Liu, J., Ding, E.: Towards accurate scene text recognition with semantic reasoning networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12113–12122 (2020)
52. Yue, X., Kuang, Z., Lin, C., Sun, H., Zhang, W.: Robustscanner: Dynamically enhancing positional clues for robust text recognition. In: European Conference on Computer Vision. pp. 135–151. Springer (2020)
53. Zhang, H., Yao, Q., Yang, M., Xu, Y., Bai, X.: Autostr: efficient backbone search for scene text recognition. In: European Conference on Computer Vision. pp. 751–767. Springer (2020)