

On Vocabulary Reliance in Scene Text Recognition

Zhaoyi Wan^{1*}, Jielei Zhang^{1*}, Liang Zhang², Jiebo Luo^{3,†}, Cong Yao^{1†}

¹Megvii, ²China University of Mining and Technology, ³University of Rochester
i@wanzy.me, {yctmzjl,yaocong2010}@gmail.com, zhangliang04@hotmail.com, jluo@cs.rochester.edu

Abstract

The pursuit of high performance on public benchmarks has been the driving force for research in scene text recognition, and notable progress has been achieved. However, a close investigation reveals a startling fact that the state-of-the-art methods perform well on images with words within vocabulary but generalize poorly to images with words outside vocabulary. We call this phenomenon “vocabulary reliance”. In this paper, we establish an analytical framework to conduct an in-depth study on the problem of vocabulary reliance in scene text recognition. Key findings include: (1) Vocabulary reliance is ubiquitous, i.e., all existing algorithms more or less exhibit such characteristic; (2) Attention-based decoders prove weak in generalizing to words outside vocabulary and segmentation-based decoders perform well in utilizing visual features; (3) Context modeling is highly coupled with the prediction layers. These findings provide new insights and can benefit future research in scene text recognition. Furthermore, we propose a simple yet effective mutual learning strategy to allow models of two families (attention-based and segmentation-based) to learn collaboratively. This remedy alleviates the problem of vocabulary reliance and improves the overall scene text recognition performance.

1. Introduction

As a pivotal task in many visual recognition and comprehension systems [42, 25, 17, 35, 22, 21], scene text recognition has been an active research field in computer vision for decades [24, 45, 43, 44, 32, 39, 36]. Recently, the pursuit of high performance on benchmarks has drawn much attention from the community. Driven by deep learning [50, 31, 2, 33, 12] and large volume of synthetic data [13, 29, 46], the recognition accuracy on standard benchmarks has escalated rapidly. For instance, the accuracy on IIIT-5k [27] without lexicon has increased from 78.2% [31] to 96.0% [12] in a very short period.

* Authors contribute equally

† Corresponding author



Figure 1: The recurrent memory mechanism in RNN-attention based methods [33] is actually a double-edged sword. The positive aspect is that for text images with words in the vocabulary (Left), even though image qualities are degraded (blur or partial occlusion), the content can be still correctly recognized. The negative aspect, which is previously neglected, lies in that for text images with words outside the vocabulary (Right), mistakes (marked in red) might easily occur.

However, an important issue has been overlooked for a long time: Even though achieving high accuracy on various benchmarks, state-of-the-art algorithms actually demonstrate obviously higher performance on images with words in the vocabulary¹ than on those with words outside it. The gap is not caused by the image quality. As shown in Fig. 1, a top-performing text recognizer [33] can correctly read the content even for images with poor quality but might make mistakes for images with better quality. The secret lies in the vocabulary: state-of-the-art methods seem inclined to memorize words that have been seen in the training phase. We call this phenomenon “vocabulary reliance”.

To further verify whether *vocabulary reliance* is common in scene text recognition, we reproduce a number of representative methods for scene text recognition, including CRNN [31], FAN [2], CA-FCN [23] and ASTER [33]. The same backbone network (ResNet-50 [8]) and training data (SynthText [7]) are used for these methods, in order to rule out interference factors. As can be observed from Tab. 1, the performance gaps between test images with words in and outside the vocabulary are significant for all evaluated methods. It reveals that *vocabulary reliance is ubiquitous*.

In this paper, we systematically investigate the problem

¹To be more specific, vocabulary in this work consists of all the words that appear in the training set.

Table 1: Accuracy gap between test images with words in and outside vocabulary on IIIT-5k. “InVoc.” and “OutVoc.” stand for in and outside the vocabulary, respectively.

Methods	All	InVoc.	OutVoc.	Gap
CRNN [31]	86.8	91.1	68.7	22.5
FAN [2]	89.9	93.1	75.3	17.8
CA-FCN [23]	89.3	91.6	76.3	15.3
ASTER [33]	89.2	92.9	74.6	18.4

of vocabulary reliance in scene text recognition. An evaluation framework is established, in which training datasets with controlled vocabularies and targeted metrics are devised to assess and compare different module combinations.

Using training data with controlled vocabularies, we are able to inspect the impact of vocabulary on algorithm performance and abilities of different algorithms in learning language prior. Meanwhile, targeted metrics allows for the evaluation of the strengths and weaknesses of different module combinations in a quantitative and precise manner. Through experiments, we obtain a series of valuable observations and findings and accordingly give a few guidelines for choosing module combinations and suggestions for developing scene text recognition algorithms in the future.

Furthermore, in order to alleviate vocabulary reliance in existing methods, we propose a novel mutual learning strategy, which allows models with different PRED layers, i.e., attention-based decoder and segmentation-based decoder, to complement each other during training. Experimental results demonstrate its effectiveness in improving the accuracy and generalization ability of both attention decoders and segmentation-based methods.

The contributions of this work are as follows:

- We raise the problem of vocabulary reliance and propose an analytical framework for investigating it.
- We discovered through experiments the advantages and limitations of current PRED layers. Attention-based decoders generalize poorly from the learned vocabulary but perform well when trained on data with a random corpus. Segmentation-based methods can accurately extract visual features while the CTC- family generally has weaker visual observation ability.
- We found that the effect of CNTX modules, which perform context modeling, is highly coupled with PRED layers. We thus provide guidelines for choosing the CNTX modules according to PRED layers.
- Moreover, we present a simple yet effective mutual learning approach to allow models of different families to optimize collaboratively, which can alleviate the problem of vocabulary reliance.

2. Proposed Analytical Framework

In this section, we describe our analytical framework, including data, modules, and metrics, in detail.

2.1. Test Data

To conduct experiments, we adopt various evaluation benchmarks, among which some are commonly used in prior works. We first briefly introduce public test datasets with real-word images, whose details are referred to [1].

ICDAR2013 (IC13) [15] is a dataset of ICDAR 2013 Robust Reading Competition for camera-captured scene text. **ICDAR2015 (IC15)** [14] comes from scene text images collected by Google glasses, where cropped text images are blurred, oriented and with low-resolution. **Street View Text (SVT)** [37] is an outdoor street images collection from Google Street View, including noisy, blurry or low-resolution images. **SVT Perspective (SVTP)** [28] focuses on curved text images. The dataset contains 645 evaluation images, which are severely distorted by non-frontal perspectives. **CUTE80 (CT)** [30] consists of 80 natural scene images, from which 288 cropped word images are generated for scene text recognition.

Basically, as shown in Fig. 1, the recognition of text images with difficulty in visual features, such as blur, stain, and irregular fonts, relies more on speculation according to the vocabulary. Thus, we group 5 datasets mentioned above into a set Ω . The ground truths of Ω are collected as our corpus for synthetic training data. Therefore, Ω and its complement Ω^c stand for the set of text images in and outside vocabulary, respectively.

Another evaluation dataset, **IIIT-5k (IIIT)** [27], is excluded from corpus collecting, which generally contains regular text and is of clear appearance. We choose IIIT as the stand-alone set to conduct the Ω^c due to its relatively large amount of images and visual clearance. By the collected vocabulary, 1354 images in vocabulary are divided into Ω and the left 1646 images make Ω^c . They are named as **IIIT-I** and **IIIT-O**, respectively.

The size of the datasets and the number of their vocabularies are shown in Tab. 2. Besides, there are 3172 distinct words in the vocabulary of Ω .

2.2. Training Data

Recent works for scene text recognition use synthetic data [7, 13] for training. **SynthText (ST)** is a dataset generated by a synthetic engine proposed in [7], whose background images are extracted from Google Image Search. It contains 80k images, from which researchers cropped about 7 million text instances for training.

As shown in Tab. 2, ST is generated from a large corpus from Newgroup20 [16] dataset, which has tens of thousands of words in the vocabulary. The large vocabulary of ST obfuscates the impact and cause of vocabulary reliance on such training data. Therefore we generate new training data for study by constraining the vocabulary.

Specifically, as stated in Sec. 2.1, our corpus is collected from test datasets. Using the synthetic engine of ST, three

Table 2: The number of words and images in training and evaluation data. ‘‘Voc.’’ is the vocabulary of datasets. ‘‘Test’’ is the vocabulary collected from test images except IIIT.

Dataset	Voc.	Images		Words	
		InVoc.	OutVoc.	InVoc.	OutVoc.
ST	ST	7266715	-	76222	-
IC13	ST	857	158	549	142
IC15	ST	1369	442	669	348
SVT	ST	530	117	333	94
SVTP	ST	536	109	300	80
CT	ST	218	70	171	63
IIIT	ST	2429	571	1277	495
IIIT	Test	1354	1646	502	1270

datasets with a similar appearance and diverse corpus are conducted for thorough and controlled comparison. Examples are illustrated in Fig. 2.

LexiconSynth (LS) From collected ground truth words, we build the corpus for LS by uniformly sampling from instances. As the vocabulary of Ω is covered by LS, models trained with LS data acquire the facilitation of vocabulary learning when evaluated on Ω . However, this purified corpus also exacerbates the over-fitting to words in vocabulary. In observation of the performance gap, properties about vocabulary learning of models can be dogged out.

RandomSynth (RS) In contrast to LS, the corpus of RS data is generated from characters in a random permutation. The lengths of the pseudowords are of the same distribution with those in LS, but the distribution of character classes is uniform. That is, the accuracy of models trained on RS is achieved without the assistance of vocabulary prior.

MixedSynth (MS) An intuitive solution for preventing algorithms from vocabulary reliance is to mix RS data into LS data. In our experiments, MS data is the union of LS and RS. Instances are sampled from RS and LS with ratio $r : (1 - r), r \in [0, 1]$. The training steps are fixed in all experiments. In comparison with datasets with a large vocabulary, the mixture of RS and LS is more practicable in real-world situations where the vocabulary is seldom completely given in advance.

Synthesis Details As the annotation of evaluation datasets serves in different manners on how to treat the case and punctuation of words, we collect the corpus as case-insensitive words without punctuation. During the rendering of LS data, each gathered word generates three instances with different variants: Uppercase, lowercase, and first-letter-capitalized case. Besides, words are inserted with a randomly chosen punctuation by a chance of 10%.

For the corpus of RS data, the proportion of letters, digits, and punctuation is about 6:3:1. Each word is rendered in the same three cases as LS data. Following the scale of ST, about 7 million cropped images are generated for RS and LS data respectively. When without special statements, the ratio r of MS data is set as 0.5 empirically.



Figure 2: Samples of generated training data. From top to bottom: all uppercase, all lowercase, and the first-letter-capitalized case. The left 2 columns are images picked up from LS, while the right 2 columns are ones from RS.

2.3. Module Combinations

According to [1], a typical scene text recognition method can be divided into four stages, *transformation (TRAN)*, *feature extraction (FEAT)*, *context modeling (CNTX)*, and *prediction (PRED)*. The CNTX stage is similar to sequence modeling (Seq.) in [1]. We extend to modeling context as we also take segmentation-based methods into consideration, for the sake of discussing the problem of vocabulary reliance in a broader perspective. The pipeline of scene text recognition is shown in Fig. 3.

In our experiments and analyses, we focus on CNTX and PRED stages, as these two stages are highly relevant to vocabulary reliance. TRAN and FEAT stages are fixed to control variables: No transformation layer is adopted and ResNet50 backbone is used in all combinations. Below, we will introduce three PRED layers and three choices for the CNTX stage.

Prediction Layers CTC [6] and attention-based decoders [3, 40] are two dominating approaches in the choices of prediction layers. As illustrated in Fig. 3d, CTC aligns the frame-wise predictions into the target string. Frames with the same characters without ‘‘BLANK’’, which is introduced to stand for no characters, are removed in final outputs. CTC is widely used in many real-world applications [20] and academic researches [4, 9], due to its superior inference speed [1].

Attention-based (Atten. for short) decoders [2, 33] are state-of-the-art methods in the field of scene text recognition. A glimpse vector is generalized from the feature sequence, upon which an RNN is adopted to produce attention vectors over the feature sequence and produce character classification each in order (see Fig. 3c).

Recently, MaskTextSpotter [26] introduces instance segmentation to localize and classify each character separately and inspires following works [5, 23, 41]. Although segmentation-based (Seg. for short) methods directly extract characters by finding connected components in the segmentation map, the large receptive field of deep convolutional networks might bring vocabulary reliance.

Context Modules Bi-directional LSTM (BLSTM) [11] is employed for context modeling on top of feature maps ex-

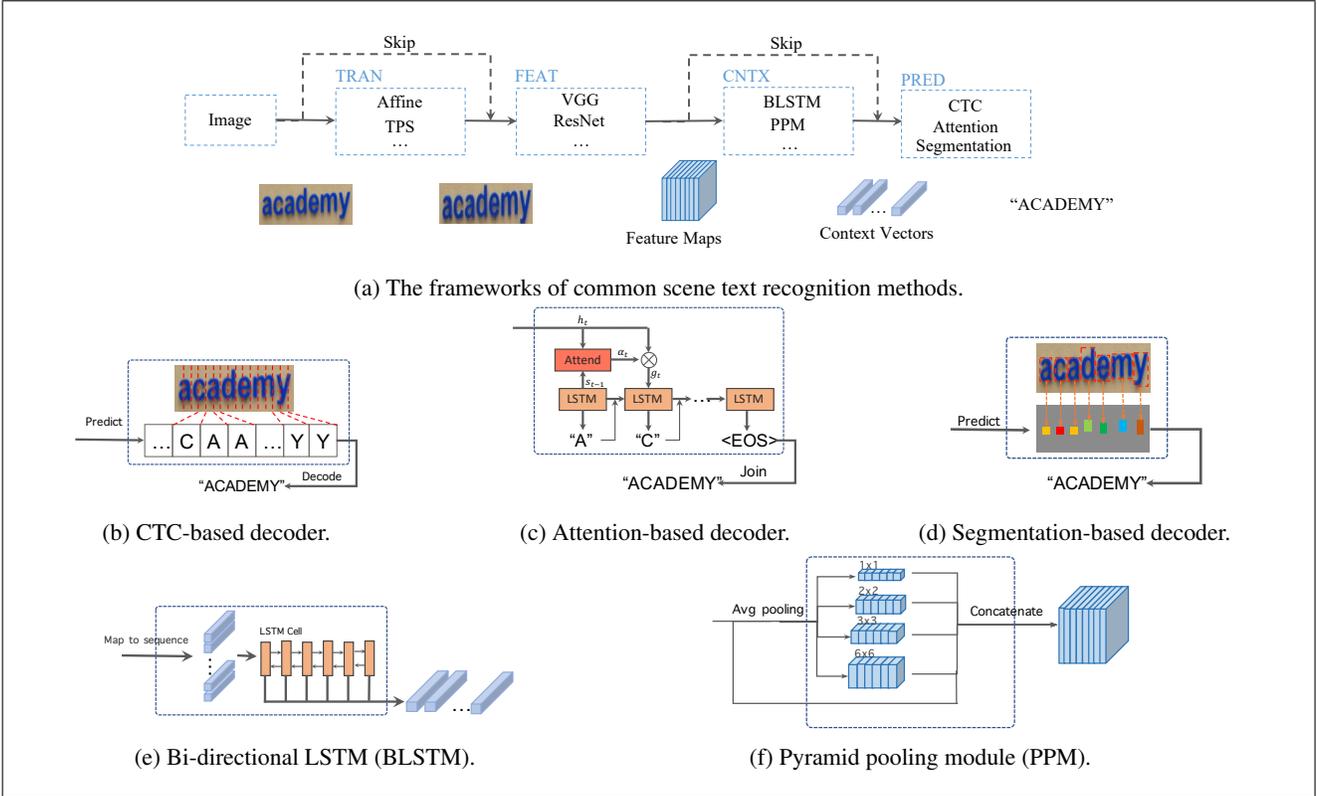


Figure 3: Pipeline and typical modules applied in scene text recognition. ‘‘Skip’’ indicates stages which are not requisite thus can be omitted in specific recognition methods.

tracted by CNNs in recent works [33, 19].

As illustrated in Fig. 3e, the BLSTM module takes feature sequences as input, which are transformed from feature maps by pooling or convolution with strides. It is a common practice in many scene text recognition methods [34, 40] for context modeling, as the BLSTM scans and maps features in the bi-directional order.

Pyramid pooling module (PPM) [49] shown in Fig. 3f is another choice for context modeling, which is proved effective on segmentation-based methods [18]. It utilizes adaptive average pooling to pool feature maps into different square resolutions (1, 3, 4, 6 in our experiments). Pooled features are then resized to the input resolution by bi-linear interpolation and concatenated with original features to gain global context information in different scales. Since segmentation-based methods are incompatible with BLSTM, PPM is a practical module for context modeling. Our experiments also validate its effectiveness in enhancing the vocabulary learning of models.

Besides, the explicit contextual modeling is not requisite for robust text recognition, as deep convolutional networks usually have large receptive fields [38, 47]. Though, in our experiments, context modeling modules do bring diversity in vocabulary learning and reliance.

The raw results are shown in Tab. 3, in which module combinations are named with circled numbers.

2.4. Metrics

Using our re-designed training data, we can evaluate the performance of algorithms on several training data. Several metrics are proposed for benchmarking the properties of models in aspects.

Firstly, we introduce a conventional metric for performance evaluation, **General Accuracy** ($\mathcal{G}\mathcal{A}$). The current practice for evaluating algorithms of scene text recognition is to evaluate models on public benchmarks with real-world images. We define the recognition accuracy on all test images of the mentioned evaluation datasets as $\mathcal{G}\mathcal{A}$, corresponding to the common evaluation in previous works.

In addition to the general metric, we further propose three specific metrics and their harmonic mean to fully reflect particular properties of different methods. For clarity, let’s define two functions. $Acc(X_{train}, X_{test})$ is the accuracy of models trained on dataset X_{train} and tested on dataset X_{test} . $Gap(\cdot)$ is defined as the performance gap on IIIT-I and IIIT-O with the same training data X_{train} :

$$Gap(X_{train}) = Acc(X_{train}, IIIT-I) - Acc(X_{train}, IIIT-O). \quad (1)$$

Observation Ability ($\mathcal{O}\mathcal{A}$) Accurate visual feature extraction and recognition is the fundamental ability of scene text recognition methods. We define $\mathcal{O}\mathcal{A}$ as how accurately an algorithm recognizes words without any vocabulary given

Table 3: The raw accuracy of models, which are numbered with circled numbers. ‘‘Gap’’ is the accuracy gap between IIIT-I and IIIT-O. ‘‘NGap’’ is normalized by recognition accuracy on IIIT.

PRED	CNTX	No.	Data.	$\subseteq \Omega$							$\subseteq \Omega^c$	Gap/NGap
				AVG	IC13	IC15	SVT	SVTP	CUTE	IIIT-I	IIIT-O	
Atten.	None	①	RS	68.5	82.2	55.1	71.7	57.0	54.2	83.2	73.3	9.8/12.6
			MS	81.8	89.9	72.2	86.4	75.2	65.6	93.0	80.1	12.9/15.0
			LS	85.7	92.7	77.4	90.5	82.3	71.5	93.7	61.0	32.7/43.2
	PPM	②	RS	70.3	84.6	57.1	74.1	58.2	55.2	84.7	77.5	7.3/9.0
			MS	81.6	88.6	71.8	85.0	75.6	71.9	92.8	80.7	12.2/14.2
			LS	85.5	92.1	77.0	89.4	81.8	74.0	94.2	69.5	24.7/30.7
	BLSTM	③	RS	68.6	82.4	55.4	70.9	57.0	53.5	82.9	73.8	9.4/12.0
			MS	82.7	89.3	74.5	86.6	77.8	67.0	92.7	81.0	11.7/13.6
			LS	87.0	92.7	79.8	92.0	84.2	73.3	94.2	63.9	30.3/39.1
CTC	None	④	RS	64.1	80.4	47.8	66.1	49.1	55.2	81.8	71.5	10.3/13.5
			MS	69.8	81.0	56.5	72.7	57.6	57.6	86.7	74.3	12.4/15.5
			LS	77.8	87.0	65.8	81.9	68.8	66.0	91.6	73.6	18.0/22.0
	PPM	⑤	RS	62.5	76.5	48.0	62.8	47.2	49.0	81.6	68.0	13.6/18.5
			MS	75.9	86.2	64.2	79.2	64.5	62.1	90.6	77.0	13.6/16.3
			LS	84.8	90.9	76.0	89.8	79.2	76.0	94.2	70.1	24.1/29.8
	BLSTM	⑥	RS	66.1	81.2	52.3	67.9	51.9	51.4	82.4	72.6	9.8/12.7
			MS	74.9	85.9	62.0	77.5	64.5	62.5	90.0	78.3	11.8/14.1
			LS	80.0	88.1	69.3	82.7	71.6	68.8	93.1	73.5	19.6/23.8
Seg.	None	⑦	RS	68.9	80.4	56.1	71.6	57.9	55.2	84.2	73.3	10.9/13.9
			MS	76.9	85.4	65.7	81.5	66.4	64.2	91.2	80.6	10.6/12.4
			LS	79.7	88.4	68.7	85.7	72.1	62.2	92.3	78.8	13.5/15.9
	PPM	⑧	RS	69.3	82.4	56.5	70.5	56.8	59.0	84.5	74.4	10.1/12.8
			MS	77.6	87.3	66.8	81.5	67.1	64.2	90.9	79.9	11.0/13.0
			LS	81.6	89.3	72.3	85.8	75.2	64.6	92.9	76.8	16.1/19.2
Atten.+Mut.	None	⑨	RS	70.4	82.8	57.0	72.7	58.8	56.9	86.3	75.8	10.5/13.1
			MS	82.0	89.9	72.3	86.4	75.2	68.1	93.1	80.7	12.4/14.3
			LS	85.8	91.9	77.2	90.8	83.1	72.7	94.5	77.6	16.9/19.9
Seg.+Mut.	None	⑩	RS	70.0	82.4	56.1	70.8	57.4	59.0	84.3	74.7	10.0/12.1
			MS	78.3	87.8	66.7	82.1	67.7	68.0	91.2	79.3	12.4/14.4
			LS	82.3	89.4	71.3	86.4	78.6	72.5	93.6	80.0	13.6/15.7

Table 4: The computation of proposed metrics. Therein, $Acc(\cdot)$ and $Gap(\cdot)$ are defined in Sec. 2.4.

Metrics.	Computation
$\mathcal{G}\mathcal{A}$	$Acc(X_{train}, \Omega \cup \Omega^c)$
$\mathcal{O}\mathcal{A}$	$Acc(RS, \Omega \cup \Omega^c)$
$\mathcal{V}\mathcal{A}$	$Acc(LS, \Omega)$
$\mathcal{V}\mathcal{G}$	$1 - (Gap(LS) - Gap(RS))$
$\mathcal{H}\mathcal{M}$	$3(\frac{1}{\mathcal{O}\mathcal{A}} + \frac{1}{\mathcal{V}\mathcal{A}} + \frac{1}{\mathcal{V}\mathcal{G}})^{-1}$

in training data. In the context of our framework, $\mathcal{O}\mathcal{A}$ is measured by evaluating models trained on RS data with test images from all benchmarks (7406 images in total). As the recognition accuracy purely comes from the observation of visual features without learning any vocabulary, it indicates the ability of models to utilize visual observation.

Vocabulary Learning Ability ($\mathcal{V}\mathcal{A}$) As stated in Sec. 1, it is likely for algorithms to employ learned vocabulary to refine or constrain recognition results of text images. Similar to $\mathcal{O}\mathcal{A}$, $\mathcal{V}\mathcal{A}$ is suggested for evaluating the recognition accu-

racy on limited vocabularies. In our experiments, measuring of $\mathcal{V}\mathcal{A}$ is to train models with LS data and evaluate the recognition accuracy on all images in Ω . $\mathcal{V}\mathcal{A}$ is meaningful for choosing models in text recognition tasks where lexicon is provided in advance.

Vocabulary Generalization ($\mathcal{V}\mathcal{G}$)

Human beings can easily generalize things from what they learnt, which inspires us to evaluate the vocabulary generalization ($\mathcal{V}\mathcal{G}$) of an algorithm by measuring the performance of models trained with LS data on words out of vocabulary. In fact, we witness the vocabulary generalization of current recognition methods in our experiments. To fairly evaluate $\mathcal{V}\mathcal{G}$, the influence of image visual feature on the dataset, which brings an intrinsic gap between two image sets, is supposed to be eliminated. Therefore $\mathcal{V}\mathcal{G}$ is indicated by

$$\mathcal{V}\mathcal{G} = 1 - (Gap(LS) - Gap(RS)) \quad (2)$$

where the score is subtracted from 1 in order to unify the monotonicity.

Table 5: Metrics of models. The circled number corresponds to different combination of different module. No. is referred to Tab. 3.

No.	PRED	$\mathcal{G}\mathcal{A}$	$\mathcal{V}\mathcal{A}$	$\mathcal{V}\mathcal{G}$	$\mathcal{O}\mathcal{A}$	$\mathcal{H}\mathcal{M}$
①	Atten.	81.0	85.7	77.1	69.6	76.9
②	Atten.	81.3	85.5	82.6	71.9	79.5
③	Atten.	83.1	87.0	79.1	69.8	78.0
④	CTC	75.8	77.8	92.4	65.8	77.1
⑤	CTC	80.1	84.8	89.5	63.5	77.5
⑥	CTC	78.4	79.9	90.2	67.6	78.1
⑦	Seg.	80.8	79.7	97.3	69.9	80.8
⑧	Seg.	81.3	81.6	94.0	70.5	80.9

Harmonic Mean ($\mathcal{H}\mathcal{M}$) For a overall metric, the harmonic mean of $\mathcal{O}\mathcal{A}$, $\mathcal{V}\mathcal{A}$, and $\mathcal{V}\mathcal{G}$ is adopted as the summary score:

$$\mathcal{H}\mathcal{M} = 3\left(\frac{1}{\mathcal{O}\mathcal{A}} + \frac{1}{\mathcal{V}\mathcal{A}} + \frac{1}{\mathcal{V}\mathcal{G}}\right)^{-1}. \quad (3)$$

$\mathcal{H}\mathcal{M}$ can be taken as a standard for general comparison of different models.

Besides, evaluation on random string can be a metric, however, there is no standard benchmark that contains pure random labels with real-world complexity. Thus, it will not be discussed in this paper.

3. Comparisons and Analyses

Using our proposed framework in Sec. 2, we provide comparisons and analyses on various module combinations. Metrics of models are shown in Fig 5. Based on the specific evaluation, we assess and analyze module combinations in different aspects.

3.1. Effect of Training Data

Fundamentally, we should first validate the effectiveness of the proposed dataset and explore the relevance of vocabulary reliance on training data. Experiments are conducted by gradually adjusting the ratio r in MS data from 0 to 1. Three models, ①, ④ and ⑦ in Tab. 3, are adopted for comparison. Besides the recognition accuracy on IIIT, we observe the probability of predicted words falling into the vocabulary, as shown in Fig. 4.

With RS data mixed into the LS data, recognition accuracy on IIIT is improved as models trained with the mixed data are less prone to be misled by vocabulary reliance. Especially for model ①, the recognition accuracy on IIIT increases from 77.8% to 84.4%, benefiting from the mixed RS data with a ratio of 25%.

The improvement in accuracy ceases when r reaches around 0.5. On one hand, the reduction of the probability to produce word prediction in vocabulary proves it effective to countervail vocabulary reliance with RS data. On the other hand, it requires a sufficient ratio of LS data to learn vocabulary from training data.

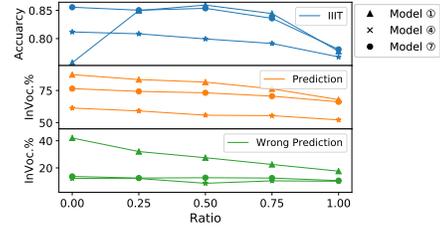


Figure 4: Probability of model ①, ④ and ⑦ on making prediction inside vocabulary. “Ratio” is the ratio of RS in MS data.

3.2. Comparison of Prediction Layers

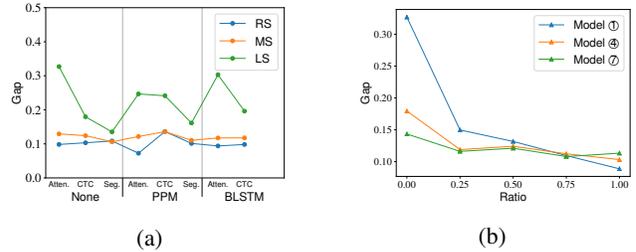
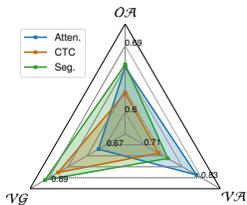


Figure 5: The accuracy gap between IIIT-I and IIIT-O. (a) Performance gap on IIIT-I and IIIT-O of module combinations. (b) The gap changes with adjusted ratio of RS data.

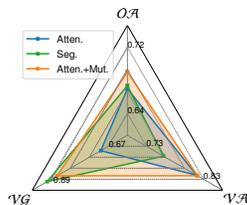
From Fig. 5a, we perceive the consistent performance gap between models trained with RS, MS, and LS data, despite PRED layers nor CNTX modules. It shows that all the combinations suffer from the problem of vocabulary reliance, but the severity differs.

Moreover, we illustrate the performance gap on IIIT of model ①, ④ and ⑦ trained with different training data. The models are built without CNTX modules, using the Atten., CTC, and Seg. PRED layers, respectively. The attention-based decoder starts with the highest gap on the point where $r = 0$ (LS data), as shown in Fig. 5b. With more RS data mixed into the training set, the gap of attention-based decoder decreases. The trend verifies the advantage of attention-based decoders on vocabulary learning and inferiority on vocabulary reliance.

In addition to vocabulary reliance, a thorough comparison of our proposed metrics of the PRED layers is illustrated in Fig. 6a. The performance of CTC is generally covered by the other two prediction layers, on metrics including both accuracy and generalization. Attention-based and segmentation-based decoders gain advantages in $\mathcal{V}\mathcal{A}$ and $\mathcal{V}\mathcal{G}$ respectively. They also perform similarly well in $\mathcal{O}\mathcal{A}$, indicating the ability to accurate recognition according to visual features only.



(a)



(b)

Figure 6: Performance of PRED layers on our metrics. All models are built without CNTX module. (a) The comparison of PRED layers. (b) OA and VA improvement of mutual learning.

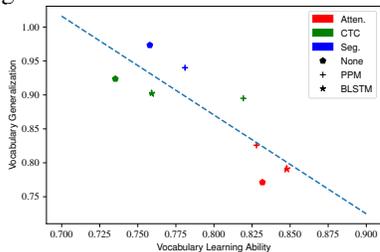


Figure 7: VA and corresponding VG of module combinations.

3.3. Comparison of Context Modules

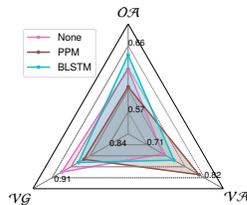
Basically, the adoption of context modules improves the vocabulary learning of models, as validated by the VA of module combinations. For example, PPM, which is not widely used in prior scene text recognition methods, brings boost on VA in combination with PRED layers: 3.9% for Seg. and 10.5% for CTC. On the other hand, as shown in Fig. 7, the strength in VA usually carries a decrease in VG .

Similar to PRED layers, the evaluation results of CNTX modules are illustrated in Fig. 8a and Fig. 8b. We find that the effect of CNTX modules in detail is highly coupled with prediction layers.

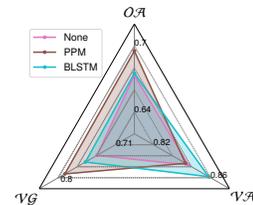
As stated in Sec. 3.2, attention-based decoders are more powerful in learning vocabulary from training data. Consequently, it brings less change in VA and VG to add more context modules to attention-based PRED layers. Besides, context modules, which perform as contextual information extractor, in fact, facilitates visual observation of attention-based and segmentation-based decoders.

As for CTC-family models, the situation is different. PPM and BLSTM significantly improve their VA and impair the VG , as the CTC decoder itself lacks of proper context modeling. The performance change in the three metrics brought by context modules on CTC-family models is shown in Fig. 8a.

In summary, it is effective to strengthen the vocabulary learning of models with proper context modules: BLSTM for attention-based, PPM for CTC and segmentation-based decoder. After all, it is a trade-off between VA and VG .



(a) CTC PRED



(b) Attention-based PRED

Figure 8: Comparison of CNTX modules.

3.4. Combination Recommendation

Based on Tab. 5 and the previous analyses, we recommend two combinations for different situations, depending on whether the vocabulary of target images are given.

Model ③, attention-based with BLSTM, achieves the best VA benefiting from the powerful CNTX module and PRED layer. This merit of model ③ in vocabulary learning also leads to the best GA , corresponding to the performance on conventional benchmarks. It is evidenced by the high score in VA and GA that ③ can perform well in applications where the vocabulary of test images are mostly a restricted subset of training data. Accordingly, model ③, similar to [40] in network design, is our first recommended combination for strong vocabulary learning ability.

As for many applications in the industry, algorithms trained with data in limited vocabulary are supposed to generalize well to more general words. Model ⑦ maintains good vocabulary generalization ability as it gets the best VG . Therefore, we recommend the combination ⑦, which is a CA-FCN-like [23] structure, for scenarios where the generalization of vocabulary is concerned.

4. Remedy by Mutual Learning

Previous sections demonstrate the trade-off between VA and VG and the diverse advantages of models. In this section, we propose a simple yet effective training strategy for combining advantages of models in different prediction layers, i.e., attention-based and segmentation-based decoders.

The idea is basically inspired by knowledge distillation [10] and deep mutual learning [48]. Similar to knowledge distillation, mutual learning of two models is a training strategy where models learn collaboratively. Knowledge distillation strategy transfers knowledge from a pre-trained powerful teacher network to student networks, while our approach optimizes two models simultaneously from scratch.

We choose the ensemble of the segmentation-based decoder and attention-based decoder as base models due to their advantages revealed in Fig. 6a. We suppose the generalization of segmentation-based decoders supervises attention-based decoders to learn to alleviate vocabulary reliance, and the accurate attention of attention-based decoders improves segmentation-based decoders in return.

4.1. Optimization

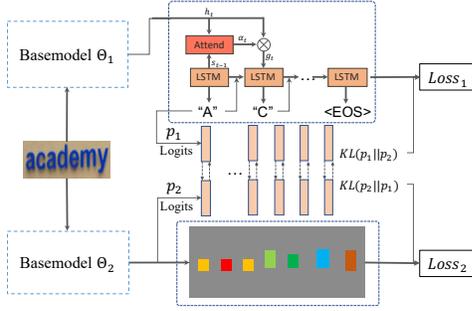


Figure 9: The mutual learning of attention-based decoder(top) and segmentation-based decoder(bottom). The KL divergence of logits are computed as auxiliary supervision, which makes the models learn collaboratively.

Let Θ_1 and Θ_2 be the network applying attention-based PRED layer and segmentation-based PRED layer, respectively. In addition to the original loss of the network L_{Θ_1} and L_{Θ_2} , the Kullback Leibler (KL) Divergence is computed as an auxiliary loss. Then the ultimate loss for Θ_1 and Θ_2 are:

$$L_1 = \sum_i^Y D_{KL}(p_2^i \| p_1^i) + L_{\Theta_1} \quad (4)$$

$$L_2 = \sum_i^Y D_{KL}(p_1^i \| p_2^i) + L_{\Theta_2}$$

where p_1, p_2 are the sequence of logits produced by Θ_1 and Θ_2 , respectively. D_{KL} is the KL Divergence and Y is the sequential label. Note that for segmentation-based decoders, the logits are “voted” scores [23] inside the shrunken region of characters.

From the Eq. 4, we can optimize the networks mutually supervised. The optimization is described in Alg. 1.

Input : Training data X and label Y .

Initialize: Θ_1 and Θ_2 separately.

while not converged **do**

$p_1 \leftarrow \text{Forward}(\Theta_1, X)$;

$p_2 \leftarrow \text{Forward}(\Theta_2, X)$;

 Compute L_1 from Eq. 4;

 Backward(Θ_1, L_1);

$p_1 \leftarrow \text{Forward}(\Theta_1, X)$;

$p_2 \leftarrow \text{Forward}(\Theta_2, X)$;

 Compute L_2 from Eq. 4;

 Backward(Θ_2, L_2);

end

Algorithm 1: Optimization of mutual learning.

4.2. Experimental Validation

We evaluate the mutual learning strategy using the proposed evaluation framework and exhibit the raw accuracy and performance on our metrics in Tab. 3 and Tab. 6, respectively. Experimental results demonstrate the significant

Table 6: Performance comparison of mutual learning strategy on our metrics. “Mut.” indicates using mutual learning or not. The raw accuracy is shown in Tab. 3.

No.	PRED	Mut.	\mathcal{VA}	\mathcal{VG}	\mathcal{OA}	\mathcal{FM}
①	Atten.	✗	83.2	77.1	69.6	76.9
⑨	Atten.	✓	85.8	93.6	71.5	82.6
⑦	Seg.	✗	75.8	97.3	69.9	80.8
⑩	Seg.	✓	82.3	96.0	70.7	81.7

improvement of base models brought by the mutual learning strategy.

These two models united by the mutual learning strategy maintains diverse properties and distinguishable advantage. The joint training procedure combines their inclination to visual features and vocabularies by harmonizing their estimation with the KL divergence. As evidence indicates, the \mathcal{OA} and \mathcal{VA} of both models are improved, which verifies the effectiveness of the mutual learning strategy.

Moreover, the vocabulary reliance of attention-based decoder is neutralized by the segmentation-based decoder. In the training of attention-based decoder, the prediction of the segmentation-based model, which relies more on visual features, acts as an extra visual regularization. In addition to minimizing L_{Θ_1} , Θ_1 is driven to fit the observation probability of Θ_2 . Quantitatively, the \mathcal{VG} of Θ_1 is boosted from 77.1% to 93.6%. In raw accuracy, the performance gap between images with words in and out of the vocabulary on LS data is almost halved (32.7% to 16.9%).

The qualitative comparison of the proposed mutual learning strategy is shown in Fig. 6b. Notable improvement on benchmarks demonstrates the effectiveness of the proposed mutual learning strategy, thus validating it reasonable to integrate the advantages of different PRED layers.

5. Conclusion

In this paper, we investigate an important but long-neglected problem: vocabulary reliance in scene text recognition methods. A comprehensive framework is built for comparing and analyzing individual text recognition modules and their combinations. Based on this framework, a series of key observations and findings have been acquired, as well as valuable recommendations, which could be conducive to the future research of scene text recognition. Moreover, we have analyzed current contextual and prediction modules and proposed a mutual learning strategy for enhancing their vocabulary learning ability or generalization ability to words out of vocabulary.

Acknowledgement

This research was supported by National Key R&D Program of China (No. 2017YFA0700800). The authors would like to thank Minghui Liao for helpful discussions.

References

- [1] Jeonghun Baek, Geewook Kim, Junyeop Lee, Sungrae Park, Dongyoon Han, Sangdoo Yun, Seong Joon Oh, and Hwal-suk Lee. What is wrong with scene text recognition model comparisons? dataset and model analysis. *2019 IEEE International Conference on Computer Vision*, 2019.
- [2] Z. Cheng, F. Bai, Y. Xu, G. Zheng, S. Pu, and S. Zhou. Focusing attention: Towards accurate text recognition in natural images. In *ICCV 2017*, pages 5086–5094, Oct 2017.
- [3] Zhanzhan Cheng, Yangliu Xu, Fan Bai, Yi Niu, Shiliang Pu, and Shuigeng Zhou. AON: towards arbitrarily-oriented text recognition. In *CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 5571–5579, 2018.
- [4] Yunze Gao, Yingying Chen, Jinqiao Wang, Hanqing Lu, INational Lab of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing, and China. Reading scene text with attention convolutional sequence modeling. 1709.
- [5] Yunze Gao, Yingying Chen, Jinqiao Wang, Ming Tang, and Hanqing Lu. Reading scene text with fully convolutional sequence modeling. *Neurocomputing*, 339:161–170, 2019.
- [6] Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. Connectionist temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd International Conference on Machine learning*, pages 369–376, Pittsburgh, Pennsylvania, USA, 2006. IMLS.
- [7] Ankush Gupta, Andrea Vedaldi, and Andrew Zisserman. Synthetic data for text localisation in natural images. In *CVPR*, pages 2315–2324, 2016.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *PAMI*, pages 770–778, 2016.
- [9] Pan He, Weilin Huang, Yu Qiao, Chen Change Loy, and Xiaoou Tang. Reading scene text in deep convolutional sequences. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA.*, pages 3501–3508, 2016.
- [10] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [11] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [12] Wenyang Hu, Xiaocong Cai, Jun Hou, Shuai Yi, and Zhiping Lin. Gtc: Guided training of ctc towards efficient and accurate scene text recognition. *arXiv preprint arXiv:2002.01276*, 2020.
- [13] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Synthetic data and artificial neural networks for natural scene text recognition. *NIPS Deep Learning Workshop*, 2014.
- [14] Dimosthenis Karatzas, Lluís Gomez-Bigorda, Angelos Nicolaou, Suman Ghosh, Andrew Bagdanov, Masakazu Iwamura, Jiri Matas, Lukas Neumann, Vijay Ramaseshan Chandrasekhar, Shijian Lu, et al. Icdar 2015 competition on robust reading. In *2015 13th ICDAR*, pages 1156–1160. IEEE, 2015.
- [15] D. Karatzas, F. Shafait, S. Uchida, M. Iwamura, L. G. i. Bigorda, S. R. Mestre, J. Mas, D. F. Mota, J. A. Almazàn, and L. P. de las Heras. Icdar 2013 robust reading competition. In *2013 12th ICDAR*, pages 1484–1493, Aug 2013.
- [16] Tom Mitchell Ken Lang. Newsgroup 20 dataset. 1999.
- [17] Vijeta Khare, Palaiahnakote Shivakumara, Chee Seng Chan, Tong Lu, Liang Kim Meng, Hon Hock Woon, and Michael Blumenstein. A novel character segmentation-reconstruction approach for license plate recognition. *Expert Systems with Applications*, 131:219–239, 2019.
- [18] Praveen Krishnan, Kartik Dutta, and CV Jawahar. Word spotting and recognition using deep embedding. In *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*, pages 1–6. IEEE, 2018.
- [19] Zhengchao Lei, Sanyuan Zhao, Hongmei Song, and Jianbing Shen. Scene text recognition using residual convolutional recurrent neural network. *Machine Vision and Applications*, 29(5):861–871, 2018.
- [20] Hui Li, Peng Wang, and Chunhua Shen. Toward end-to-end car license plate detection and recognition with deep neural networks. *IEEE Transactions on Intelligent Transportation Systems*, 20(3):1126–1136, 2018.
- [21] Minghui Liao, Pengyuan Lyu, Minghang He, Cong Yao, Wenhao Wu, and Xiang Bai. Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [22] Minghui Liao, Zhaoyi Wan, Cong Yao, Kai Chen, and Xiang Bai. Real-time scene text detection with differentiable binarization. *ArXiv*, abs/1911.08947, 2019.
- [23] Minghui Liao, Jian Zhang, Zhaoyi Wan, Fengming Xie, Jiajun Liang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Scene text recognition from two-dimensional perspective. In *AAAI*, 2019.
- [24] Shangbang Long, Xin He, and Cong Yao. Scene text detection and recognition: The deep learning era. *arXiv preprint arXiv:1811.04256*, 2018.
- [25] Shangbang Long, Jiaqiang Ruan, Wenjie Zhang, Xin He, Wenhao Wu, and Cong Yao. Textsnake: A flexible representation for detecting text of arbitrary shapes. In *ECCV*, 2018.
- [26] Pengyuan Lyu, Minghui Liao, Cong Yao, Wenhao Wu, and Xiang Bai. Mask textspotter: An end-to-end trainable neural network for spotting text with arbitrary shapes. In *ECCV*, pages 67–83, 2018.
- [27] Anand Mishra, Karteek Alahari, and CV Jawahar. Scene text recognition using higher order language priors. In *BMVC-British Machine Vision Conference*. BMVA, 2012.
- [28] T. Q. Phan, P. Shivakumara, S. Tian, and C. L. Tan. Recognizing text with perspective distortion in natural scenes. In *2013 IEEE International Conference on Computer Vision*, pages 569–576, Dec 2013.
- [29] Xiaohang Ren, Kai Chen, and Jun Sun. A cnn based scene chinese text recognition algorithm with synthetic data engine. *arXiv preprint arXiv:1604.01891*, 2016.
- [30] Anhar Risnumawan, Palaiahankote Shivakumara, Chee Seng Chan, and Chew Lim Tan. A robust arbitrary text detection

- system for natural scene images. *Expert Systems with Applications*, 41(18):8027 – 8048, 2014.
- [31] Baoguang Shi, Xiang Bai, and Cong Yao. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *PAMI*, 39(11):2298–2304, 2017.
- [32] Baoguang Shi, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Robust scene text recognition with automatic rectification. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4168–4176, 2016.
- [33] Baoguang Shi, Mingkun Yang, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. Aster: An and attentional scene and text recognizer and with flexible and rectification. In *PAMI*, pages 1–1. IEEE, 2018.
- [34] Bolan Su and Shijian Lu. Accurate scene text recognition based on recurrent neural network. In Daniel Cremers, Ian Reid, Hideo Saito, and Ming-Hsuan Yang, editors, *Computer Vision – ACCV 2014*, pages 35–48, Cham, 2015. Springer International Publishing.
- [35] Shu Tian, Xu-Cheng Yin, Ya Su, and Hong-Wei Hao. A unified framework for tracking based text detection and recognition from web videos. *IEEE transactions on pattern analysis and machine intelligence*, 40(3):542–554, 2017.
- [36] Zhaoyi Wan, Mingling He, Haoran Chen, Xiang Bai, and Cong Yao. Textscanner: Reading characters in order for robust scene text recognition. *ArXiv*, abs/1912.12422, 2019.
- [37] Kai Wang, Boris Babenko, and Serge Belongie. End-to-end scene text recognition. In *ICCV, ICCV '11*, pages 1457–1464, Washington, DC, USA, Nov 2011. IEEE Computer Society.
- [38] Zecheng Xie, Zenghui Sun, Lianwen Jin, Ziyong Feng, and Shuye Zhang. Fully convolutional recurrent network for handwritten chinese text recognition. In *2016 23rd International Conference on Pattern Recognition (ICPR)*, pages 4011–4016. IEEE, 2016.
- [39] Mingkun Yang, Yushuo Guan, Minghui Liao, Xin He, Kaigui Bian, Song Bai, Cong Yao, and Xiang Bai. Symmetry-constrained rectification network for scene text recognition. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9146–9155, 2019.
- [40] MingKun Yang, Yushuo Guan, Minghui Liao, Xin He, Kaigui Bian, Song Bai, Cong Yao, and Xiang Bai. Symmetry-constrained rectification network for scene text recognition. *2019 IEEE International Conference on Computer Vision*, 2019.
- [41] Xiao Yang, Dafang He, Zihan Zhou, Daniel Kifer, and C. Lee Giles. Learning to read irregular text with attention mechanisms. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence, IJCAI'17*, pages 3280–3286. AAAI Press, 2017.
- [42] Cong Yao, Xiang Bai, and Wenyu Liu. A unified framework for multi-oriented text detection and recognition. 2014.
- [43] Cong Yao, Xiang Bai, Wenyu Liu, Yi Ma, and Zhuowen Tu. Detecting texts of arbitrary orientations in natural images. *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1083–1090, 2012.
- [44] Cong Yao, Xiang Bai, Baoguang Shi, and Wenyu Liu. Strokelets: A learned multi-scale representation for scene text recognition. *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 4042–4049, 2014.
- [45] Qixiang Ye and David Doermann. Text detection and recognition in imagery: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 37(7):1480–1500, 2014.
- [46] Fangneng Zhan, Shijian Lu, and Chuhui Xue. Verisimilar image synthesis for accurate detection and recognition of texts in scenes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 249–266, 2018.
- [47] Yaping Zhang, Shuai Nie, Wenju Liu, Xing Xu, Dongxiang Zhang, and Heng Tao Shen. Sequence-to-sequence domain adaptation network for robust text image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2740–2749, 2019.
- [48] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4320–4328, 2018.
- [49] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6230–6239, Honolulu, HI, USA, July 2017.
- [50] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. East: An efficient and accurate scene text detector. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2642–2651, 2017.