

# Visual-Semantic Transformer for Scene Text Recognition

Xin Tang\*, Yongquan Lai\*, Ying Liu, Yuanyuan Fu, Rui Fang

Visual Computing Group, Ping An Property & Casualty Insurance Company, Shenzhen, China  
tangxin051@pingan.com.cn, tangxint@gmail.com

## Abstract

Modeling semantic information is helpful for scene text recognition. In this work, we propose to model semantic and visual information jointly with a Visual-Semantic Transformer (VST). The VST first explicitly extracts primary semantic information from visual feature maps with a transformer module and a primary visual-semantic alignment module. The semantic information is then joined with the visual feature maps (viewed as a sequence) to form a pseudo multi-domain sequence combining visual and semantic information, which is subsequently fed into a transformer-based interaction module to enable learning of interactions between visual and semantic features. In this way, the visual features can be enhanced by the semantic information and vice versa. The enhanced version of visual features are further decoded by a secondary visual-semantic alignment module which shares weights with the primary one. Finally, the decoded visual features and the enhanced semantic features are jointly processed by the third transformer module obtaining the final text prediction. Experiments on seven public benchmarks including regular/irregular text recognition datasets verifies the effectiveness our proposed model, reaching state of the art on four of the seven benchmarks.

## Introduction

Scene text recognition (STR) is the task of recognizing text from images taken in complex scenes such as street view. It is an inherently difficult task due to the variance of shape, color, scale and appearance of the embedded text. Clutter background, large perspective distortion, lighting condition, degraded image quality due to motion/out-of-focus blur also impose severe challenges to the successfully solving the task, resulting severe miss-predictions.

Despite its difficulty, STR has many real-world applications ranging from self-driving cars (Yu et al. 2019), street image understanding to applications such as instant translation and intelligent text reading in smart-phones (Wu et al. 2019). For decades, STR has been an active research direction (Baek et al. 2019; Li et al. 2019; Shi et al. 2018; Wan et al. 2020; Long, He, and Yao 2021), attracting many efforts in designing new models and creating new datasets in order to solve the problem.

\*These authors contributed equally.

A complete approach to recognizing text from scene images usually involves text detection and text recognition. In this paper, we assume that text detection is done and only focus on the recognition part. That is, we assume the input to our model is a cropped image with regular or irregular characters lying in it. A large portion of the previous work and many open-source datasets (Karatzas et al. 2013, 2015) follow from this assumption. In this work, we will refer to STR as just recognizing text from cropped images.

The approaches to solving STR problem can be roughly divided into two categories: linguistic-based and linguistic-free. Linguistic-based methods refer to those which incorporate vocabulary (dictionary) or lexicon (parts of words), while linguistic-free methods use only the images themselves without relying on explicit language modeling, whether from internal or external, pretrained or from-scratch.

The **motivation** of this work is multi-fold. Firstly, to deal with the cases when visual information alone is inadequate, semantic or linguistic features has been introduced in various effects, among which (Qiao et al. 2020) propose a semantic enhanced encoder-decoder framework to recognize low-quality scene texts. A semantic module is designed to directly produce semantic features that are consistent with the word embedding learned from a pretrained language model. Inspired by their efforts in exploiting semantic features, we also explicitly model semantic information. But unlike their approach, we achieve semantic modeling without relying on external language models but instead using an alignment module.

On the other hand, our work is also partially inspired from wav2vec (Baevski et al. 2020) in audio community. The wav2vec model first extract primary audio features from waveform using 1D convolution. The features are processed by a succeeding transformer module, resulting in secondary contextualized audio features. The secondary audio features interact with the primary features by predicting them back at the next few time steps. Similarly, we extract semantic and visual features at different stages, making them interact with each other so that the overall recognition performance is improved.

Apart from the aforementioned work, we have also noted that many efforts have been devoted to modeling visual-semantic relation. We argue that extracting semantic fea-

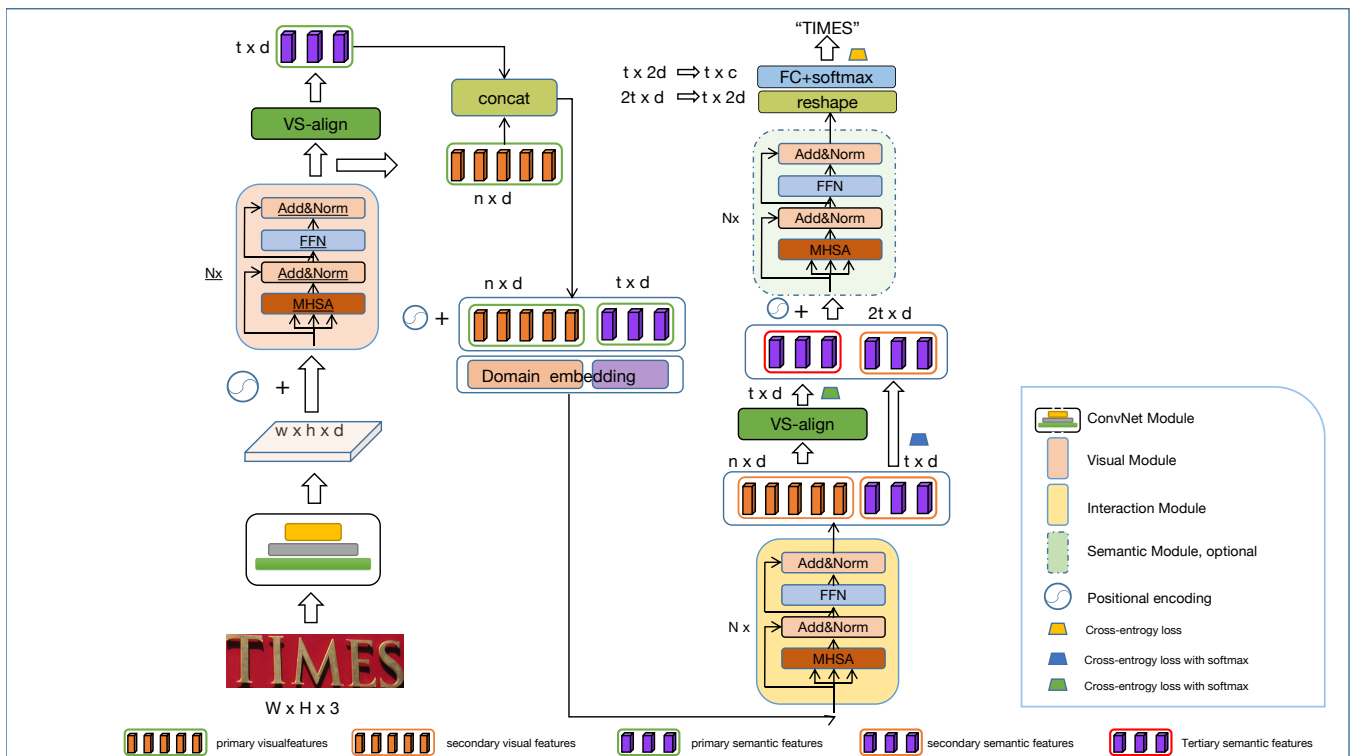


Figure 1: The architecture of Visual-Semantic Transformer (VST). The VST consists of several key modules, namely ConvNet (C), Visual module (V), Interaction module (I), Semantic module (S) and two weight-sharing Visual-Semantic Alignment (A, or vs align) modules. Best view in color.

tures from appearance and then enforce interaction with visual features is one key to successfully solving STR. To this end, we propose a novel transformer model for scene text recognition. Our model is a unified end-to-end text recognizer which converts image patches into text in parallel (i.e., non-autoregressively), without using complex decoder such as CTC (Graves et al. 2006). The model, which we coin *Visual-Semantic Transformer (VST)*, is able to learn semantic features from appearance and combine them back with visual features to enable visual-semantic interaction using multi-head self-attentions (MHSA).

The overall architecture of VST is shown in Fig. 1. The VST first explicitly extracts primary semantic information from visual feature maps with a transformer module and a primary visual-semantic alignment module (vs-align). The semantic information is then joined with the visual feature maps (viewed as a sequence) to form a pseudo multi-domain sequence combining visual and semantic information, which is subsequently fed into the second transformer module to enable learning of interactions between visual and semantic features. In this way, the visual features can be enhanced by the semantic information and vice versa. The visual features are further decoded with a secondary vs-align module which shares weights with the primary one. Finally, the decoded visual features and the semantic features can be jointly processed by the third transform module and the final softmax layer to obtain the resulting text. Overall, our ap-

proach is end-to-end and conceptually simple. Experiments on public benchmarks on regular/ irregular text recognition tasks demonstrate the effectiveness of our proposed model. The main **contributions** of this paper are as follows,

- We propose a novel visual-semantic transformer to effectively solve STR problem, surpassing or on par with state-of-the-art in most STR datasets.
- We design weight-sharing visual-semantic alignment modules to explicitly enforce the learning of semantic information without external language models.
- We introduce an interaction module that allows visual and semantic information to globally interact with and enhance from each other.

In this work, **semantic information** refers to the information that connects visual appearance and the underlying linguistic information. In other words, it is the information extracted from visual features which are very closely related to the text represented by the scene text image. Semantic information is distinguished from the term language or linguistic information, because the later usually refers to the information extracted directly from text solely. Semantic feature extraction is not language modeling or word embedding because features are not computed directed from real character sequence, but rather learned from image appearance and forced to be consistent with language. The semantic feature is ready to convert into text characters using a simple linear

probe. Under this definition, our model can be categorized as linguistic-free, as we do not require an external language model. Instead, we explicitly model semantic information which can be seen as pseudo-linguistic information.

In the following sections, we will briefly review the related work and then introduce our VST model in details, followed by extensive experiments and conclusion.

## Related Work

Text recognition has been an active research area for decades. See (Long, He, and Yao 2021; Chen et al. 2021) for comprehensive reviews. Due to page limitation, we can only list a portion of recent work here.

Recently, (Nguyen et al. 2021) incorporate a dictionary in training and inference stage to help selecting the most compatible outcome for STR. (Feng et al. 2021) introduce character center segmentation branch to extract semantic features which encode the category and position of characters for improving video text detection performance. (Patel et al. 2016) propose to generate contextualized lexicons for scene images with only visual information. (Sabir, Moreno-Noguer, and Padró 2018) use language model to build the semantic correlation between scene and text in order to re-rank the recognition results. (Zheng, Wang, and Betke 2019) also propose to use pretrained language models for correcting predictions. Very recently, (Wang et al. 2021b) propose to learn the linguistic rules in the visual space by randomly mask out some characters from the image and predict them back.

Similar to speech recognition, scene text recognition can be treated as a sequence-to-sequence (seq2seq) mapping problem (Jaderberg, Vedaldi, and Zisserman 2014; Jaderberg et al. 2015; Qiao et al. 2020; Yue et al. 2020). (Li et al. 2019) combine convolution and LSTM as an encoder, then use another LSTM as decoder to predict text attentively, quite similar to the *show, attend and tell* work on image captioning (Xu et al. 2015). ASTER (Shi et al. 2018) takes a two-stage approach to first rectify curved text images and then perform recognition using seq2seq model with attention. CRNN (Shi, Bai, and Yao 2016) adapts CNN to obtain visual features, which are then fed into LSTM module with CTC loss (Graves et al. 2006) for text prediction. STAR-Net (Liu et al. 2016) uses spatial transformer (Jaderberg et al. 2015) to tackle challenges brought by image distortion. Attention can be added to the seq2seq models (Li et al. 2019; Bhunia et al. 2021) in a straightforward way to alleviate bottleneck effects brought by seq2seq models. (Litman et al. 2020) propose stacked block architecture with intermediate supervision to train a deep BiLSTM encoder, while attention is used in decoding stage to exploit contextualized visual features. (Aberdam et al. 2021) propose seq2seq contrastive learning of visual representations which can be applied to text recognition. DAN (Wang et al. 2020) decouples alignment from the decoding stage into the early conventional encoder network. (Wang et al. 2021a) propose an alignment module enabling text recognizer to recognize document level image. (Wang et al. 2019b) use adversarial loss to handle low-resolution image.

Transformers have also been successfully applied to

STR. ABINet (Fang et al. 2021) enforces the bidirectional language-model (LM) to only learn linguistic rules by gradient-stopping in training. The decoding is in an iterative way allowing the predictions to be refined progressively. Their conv+transformer visual module and transformer-LM can be separately pretrained to improve performance. HRGAT (Yang et al. 2020) connects CNN feature maps to a transformer-based autoregressive decoder, where the cross-attention is guided by holistic representation obtained by average-pooling of 2d feature maps.

Perhaps more related, SRN (Yu et al. 2020) incorporates visual-to-semantic embedding block and cross-entropy loss to align with ground-truth text, but they use argmax embedding, which is different from our direct use of probability vector that enables smooth gradient flow in training. Our work is also innovative in many ways such as visual-semantic alignment, multiple-stage semantic processing and transformer-based visual-semantic interaction. Instead of using argmax, (Bhunia et al. 2021) use Gumbel-softmax (Jang, Gu, and Poole 2016) for extracting semantic information, which is then fed into succeeding transformer-based visual-semantic reasoning module. The decoding involves complex multiple-stage attentional LSTM that couples with feature pyramid networks. Our approach is not only conceptually much simpler and computationally more efficient, but also more effective in solving STR.

## Visual-Semantic transformer

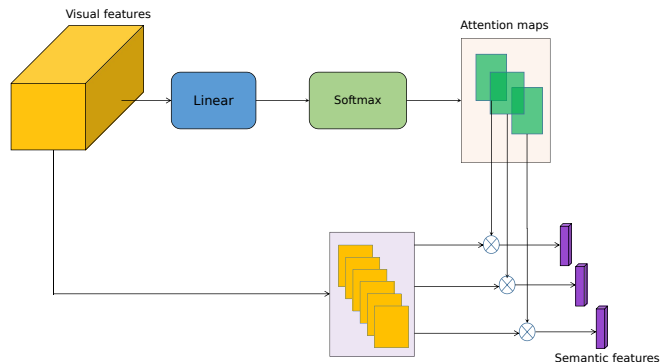


Figure 2: The architecture of vs-align module. Visual features (viewed as feature maps) are projected using a linear layer and further normalized using softmax operator, obtaining  $t$  attention maps, each of which has the same spatial dimension as the origin feature map. The  $t$  attention maps and the original  $d$  visual feature maps together will reduce (by multiplication) to  $t$  semantic features in  $R^d$ .

We will introduce the visual-semantic transformer (VST) in this section. The VST consists of several key modules, namely ConvNet (C), Visual module (V), Interaction module (I), Semantic module (S) and two weight-sharing Visual-Semantic Alignment (A, or vs-align) module. Module C can be any convnet, Module V, I, S are basically transformer blocks, while vs-align is an attention-based alignment block. Since a large portion of our model consists of transformer

blocks, we name it *visual-semantic transformer*, implying that it is a transformer that explicitly models visual and semantic information.

The overview of the architecture is shown Fig. 1. The convnet module extracts local visual features from cropped images. The visual module transforms the local visual features into contextualized global visual features, which we call primary visual features. The first vs-align module converts the primary visual features into primary semantic features. The interaction module then transforms primary visual and primary semantic features by allowing them to interact with each other through MHSA, producing secondary visual and secondary semantic features respectively. The secondary visual features are fed into the second vs-align module (which shares weights with the first one), obtaining tertiary (third) semantic features. Secondary and tertiary semantic features together will be processed by the final semantic module to obtain the resulting text prediction.

To enforce the learning of semantic information, the two vs-align modules **share weights**. In other words, the improvement of the second vs-align module is immediately helpful for improving the accuracy in aligning the primary semantic features, which in return will help improve the overall recognition accuracy. We found that this design is very useful for extracting primary semantic information, as seen in Fig. 3.

The VST has two variants, namely VST-F (full) or VST-B (basic), depending on whether the module S is used or not. In VST-F, there are three losses in training and the output of final softmax is used for text decoding. For VST-B, there are two losses (blue and green in Fig. 1) and the semantic stream of module I’s output are passed to softmax layer, while the visual stream will go through additional vs-align module before softmax and CE loss. The decoding for VST-B can be evaluated separately or jointly (by probability voting).

### The ConvNet (C) Module

We use resnet-like architecture in extracting local features. The input images are first resized to have the same height, with aspect ratio kept unchanged. During training and testing, replication-padding (padding using values from the image border) is used for batch processing.

Assuming that input image is of size  $W \times H \times 3$ , the convnet module will produce feature maps of size  $w \times h \times d$ . The features will be fed into the visual module.

### The Visual (V) Module

Convnets are good at learning local features but hard to learn global correlation between features that locate far apart. Inspired from the success of transformers when applied to vision task, we insert a transformer module here in order to enhance the feature maps by allowing them to interact with each other. The visual module takes the convolution feature maps as input, producing primary visual features. The primary visual features together encode completely the appearance information of the input image. They are then converted to primary semantic features by the vs-align module.

The visual module consists of MHSA, layer-norm, and feed-forward network, as described in (Vaswani et al. 2017),

but with minor modification that puts layer-norm before MHSA (Dosovitskiy et al. 2020; Wang et al. 2019a). Note that Fig. 1 does NOT reflect this modification. The feature maps are viewed/reshaped as a sequence of dimension  $d$  and length  $n = w \times h$ . the module V, takes the feature sequence as input and generate a sequence of the same dimension and length. The interaction between visual features themselves mainly takes place in the MHSA layers.

### The Visual-Semantic Alignment (A) Module

Parallel attentions has been used to map visual features into semantics(Wang et al. 2021b; Lyu et al. 2019). Compared with previous work, the vs-align module in this work is much simpler but still effective. Note that we do not use the name semantic decoder here because the output of the module is not actually the character index, but rather the probability distribution of the characters. In other words, the semantic feature encodes the probability distribution at each candidate location, which can be readily converted into character index with *argmax* operation in decoding stage (not in training). This allows the gradients to flow smoothly inside the VST.

The architecture of vs-align is depicted in Fig. 2. As shown in the figure, the visual features (viewed as feature maps) are projected using a linear layer and normalize using softmax operator, obtaining  $t$  attention maps each of which has the same spatial dimension as the origin feature map. The alignment module can also be formulated mathematically as follows,

$$S = \text{softmax}(QV^T)V \quad (1)$$

where  $S \in R^{t \times d}$  is the semantic sequence,  $V \in R^{n \times d}$  is the visual sequence and  $Q \in R^{t \times d}$  is the trainable projection matrix.

Our VS-align module is able to learn the relation between the semantic and visual features. This is different from CTC (Graves et al. 2006), which is good at aligning two sequences with beam search, but not capable of handling information along the height direction.

As shown in Fig. 1, there are two vs-align modules, which share weights during training and inference. This weight-sharing scheme is the key to successfully learning 1st semantic information. It enables the weights learned in aligning 2nd visual features to transfer to aligning 1st visual features, making the early learning of 1st semantics possible. The 1st semantic in return enhances the 1st visual features through interaction module, further improving the training of the 2nd vs-align module. The early learning of semantics provides more chances to correct the semantic features through interaction in later stages. If semantic had been learned in the last stage, we would have not change to correct it, unless incorporating additional lexicon or dictionary information (Nguyen et al. 2021)

### The Interaction (I) Module

The interaction module plays the key role of mixing primary semantic features and visual features. The modules takes



Method	Regular test datasets				Irregular test datasets		
	IIIT	SVT	IC03	IC13	IC15	SVTP	CUTE
AON (Cheng et al. 2018)	87.0	82.8	91.5	-	68.2	73.0	76.8
ASTER (Shi et al. 2018)	93.4	89.5	94.5	91.8	76.1	78.5	79.5
NRTR (Sheng, Chen, and Xu 2019)	86.5	88.3	95.4	94.7	-	-	-
SAR (Li et al. 2019)	91.5	84.5	-	91.0	69.2	76.4	83.3
DAN (Wang et al. 2020)	94.3	89.2	95.0	93.9	74.5	80.0	84.4
HRGAT (Yang et al. 2020)	94.7	88.9	-	93.2	79.5	80.9	85.4
SRN (Yu et al. 2020)	94.8	91.5	-	95.5	82.7	85.1	87.8
SCATTER (Litman et al. 2020)	93.7	92.7	96.3	93.9	82.2	86.9	87.5
GTC (Hu et al. 2020)	95.5	92.9	95.2	94.3	82.5	86.2	92.3
RobustScanner (Yue et al. 2020)	95.3	88.1	-	94.8	77.1	79.5	90.3
(Bhunja et al. 2021)	95.2	92.2	-	95.5	84.0	85.7	89.7
PREN (Yan et al. 2021)	95.6	94.0	95.8	96.4	83.0	87.6	91.7
ABINet-SV (Fang et al. 2021)	95.4	93.2	-	96.8	84.0	87.0	88.9
ABINet-LV (Fang et al. 2021)	96.2	93.5	-	<b>97.4</b>	<b>86.0</b>	<b>89.3</b>	89.2
VST-B	<b>96.3</b>	<b>93.8</b>	96.4	96.4	85.4	88.7	<b>95.1</b>
VST-F	96.1	93.1	<b>97.1</b>	96.4	85.4	89.1	94.8

Table 1: When compared with previous work, our approach achieves very competitive results. VST-B denotes the C+V+A+I basic model and VST-F the C+V+A+I+S full model. '-' denotes data not available or config not the same.

two streams, namely the primary semantic features and primary visual features as input, producing the secondary visual and semantic features. Fixed positional encoding are added into the visual stream. For the semantic stream, learnable position embedding is used. Other positional encoding schemes probably help as long as it breaks permutation invariant property of transformers.

The design of the interaction module follows from module V, comprising layers of transformer blocks. The module enables the feature interaction between a) the visual stream and semantic will interact with each other, and b) the features themselves inside the same stream.

The interpretation of module can be seen from eq. 2. Following from the notations in (Vaswani et al. 2017), the interaction in the MHSA layer is obvious by judging from the following equation,

$$S = \text{softmax}\left(\frac{[Q_s; Q_v][K_s; K_v]^T}{\sqrt{d_k}}\right)[V_s; V_v], \quad (2)$$

where the subscripts s and v denote semantic and visual part respectively, ';' is column representation.

In this way, the primary semantic features will look for support from the spatial visual features. The interaction is useful for enhancing semantic features because they now have attentive access to all spatial location of the visual features. On the other hand, visual features are also enhanced not only because then can interact with all other spatial locations regardless of distance, but they can also learn from semantic features, which we hypothesize is useful for dealing with appearance degradation.

The semantic features can be seen as pseudo-linguistic features, distinguished from visual domain. Hence we add domain embedding before feeding two combined streams into the interaction module.

One of the outputs of interaction module is secondary visual sequence in  $R^{n \times d}$ , which will be converted to tertiary

Module	Regular test datasets				Irregular test datasets		
	IIIT	SVT	IC03	IC13	IC15	SVTP	CUTE
CV	95.10	91.94	95.72	95.52	81.83	86.23	91.33
CVA	95.63	91.94	96.41	95.55	82.31	86.95	91.76
CVAI / $s_3$	96.40	93.51	96.54	96.35	85.09	88.84	94.79
CVAI / $s_2$	96.26	<b>93.81</b>	96.30	96.35	85.14	88.68	95.13
CVAI / $p$	<b>96.36</b>	<b>93.81</b>	96.42	<b>96.45</b>	85.36	88.68	<b>95.13</b>
CVAIS	96.06	93.50	<b>97.11</b>	96.40	<b>85.37</b>	<b>89.15</b>	94.79

Table 2: Recognition accuracy increases when extra module is added. / $s_2$ , / $s_3$  and / $p$  denote decoding from 2nd/3rd semantics ( $s_3$ ) and by probability voting from both semantics respectively, CVAIS means modules C+V+A+I+S and so on.

semantic features using vs-align module following by cross-entropy loss with softmax, against the ground-truth text. The other is semantic features in  $R^{n \times d}$ , which will be directly compared against ground-truth text with the same activation ans loss type. Note that at inference time, the two loss branches are unneeded.

### The Semantic (S) Module

The semantic module is inserted optionally to play the role of further fusing the two semantic streams. When inserted, our model is named VST-F (full); otherwise, it is named VST-B (basic). The two streams are the secondary semantic features generated from the intersection module and the tertiary (third) semantic features generated from the second vs-align module respectively.

The layer configuration of this module also follows from the design of module V and I. Fixed positional encoding is added, but there is no domain or segment embedding because the two streams are both semantic. Since the transformer blocks keep the input shape, we concatenate the two streams along channel direction to  $t \times 2d$  and feed it to linear and softmax layer, obtaining the final text predictions.



Figure 3: Visualization of attention maps for decoding each character.

## Mathematical interpretation

As shown in Fig. 1, there are three losses indicated by the ladders in different colors. The green one measures the discrepancy between the true text labels and semantic features output from second vs-align module using softmax activation and cross-entropy loss. The blue one measures how different are the secondary semantic features from the true text. The yellow one is the main loss for decoding final text for VST-F, while the other two are auxiliary for improving convergence performance. For VST-B, since there is no module S, only blue and green loss are used. For all loss branches, the discrepancy are measured by cross-entropy loss with softmax activation.

From the optimization perspective, the VST-F solves the following optimization problem,

$$\begin{aligned} \min_{\theta_v, \theta_a, \theta_i, \theta_s} \quad & L(s_2, y) + L(s_3, y) + L(f_{\theta_s}(s_2, s_3), y) \\ \text{subject to} \quad & v_1 = f_{\theta_v}(x), s_1 = f_{\theta_a}(v_1) \\ & s_2, v_2 = f_{\theta_i}(s_1, v_1), s_3 = f_{\theta_s}(v_2) \end{aligned} \quad (3)$$

where  $\theta_v, \theta_a, \theta_i, \theta_s$  denote the parameters of module C+V, A, I, and S respectively,  $L$  is cross-entropy loss,  $s_1, s_2, s_3, v_1, v_2$  are the primary/secondary/tertiary semantic/visual features resp.,  $x, y$  are the input image and text label resp. For VST-B, the third term and the corresponding constraint in the above objective function is discarded.

## Experiments

### Datasets

Our model is trained on three datasets: SynthText (ST) (Gupta, Vedaldi, and Zisserman 2016), MJSynth (MJ) (Jaderberg et al. 2014, 2016) and SynthAdd (SA) (Li et al. 2019). The training dataset contains totally 15.7 million synthetic images. For evaluation, we use four regular text datasets: IIIT 5K-words (IIIT) (Mishra, Alahari, and Jawahar 2012), Street View Text (SVT) (Wang, Babenko, and Belongie 2011), ICDAR 2003 (IC03) (Lucas et al. 2003), ICDAR 2013 (IC13) (Karatzas et al. 2013), and three irregular text datasets: ICDAR 2015 (IC15) (Karatzas et al. 2015), Street View Text Perspective (SVTP) (Phan et al. 2013) and CUTE (Risnumawan et al. 2014).

IIIT contains 3000 cropped images collected from Google image search. SVT is collected from Google Street View

containing 647 testing images. Following from (Wang, Babenko, and Belongie 2011), we select 867 cropped images from IC03 for testing. IC13 contains 1095 testing images and we discard images containing non-alphanumeric characters or contains fewer than three characters, resulting in 1015 images for testing, following from previous work (Yu et al. 2020).

For irregular datasets, IC15 contains 2077 cropped images by Google Glasses. We use 1811 images after discarding extremely distorted images. SVTP and CUTE contains 639 and 288 images respectively.

### Implementation Details

Module C is implemented as a four-layers resnet, with each layer having 1,2,5 and 3 blocks respectively. Module V, I and S each consists of 3 transformer layers. The number of parameters are about 63.99M for VST-F, and 63.65M for VST-B.

The image patches are resized with aspect-ratio unchanged so that  $H = 48$ , and the width is trimmed or padded to  $W = 160$  pixels. The feature map size is  $6 \times 40 \times 512$ ,  $w = 40, h = 6, n = 240$ . The number of class is set to 38, including 0-9, a-z, [unk], [eos]. We assume that the maximum character length is 25, i.e.,  $t = 25$ . We use online data augmentations including distortion, stretch, perspective transform, blurring, colour jitter etc. In training, we set the batch size to 256 and sample from ST, MJ, SA datasets with weight 0.4, 0.4 and 0.2 to balance the datasets (Baek et al. 2019). We use Adam optimizer with the initial learning rate  $1e-4$  and decreased lr to  $1e-5$  when the loss plateaus. We use for and the training takes roughly 3 days to converge on four Tesla V100 GPUs.

### Comparison with State-of-the-Art

We compare VST-B and VST-F with previous work over seven public benchmarks in table 1. The proposed VST-B and VST-F both achieve on par with state-of-the-art approaches. The large version of the very recent work by (Fang et al. 2021) performs slightly better than ours on three of the seven datasets, but they used extra cumbersome pre-trained language models. Moreover, our VST-B and VST-F are about three times faster than theirs in inference. On the rest four datasets, we improve SOTA by large margins (0.2% IIIT, 0.3% SVT, 0.8% IC03, 2.8% CUTE).

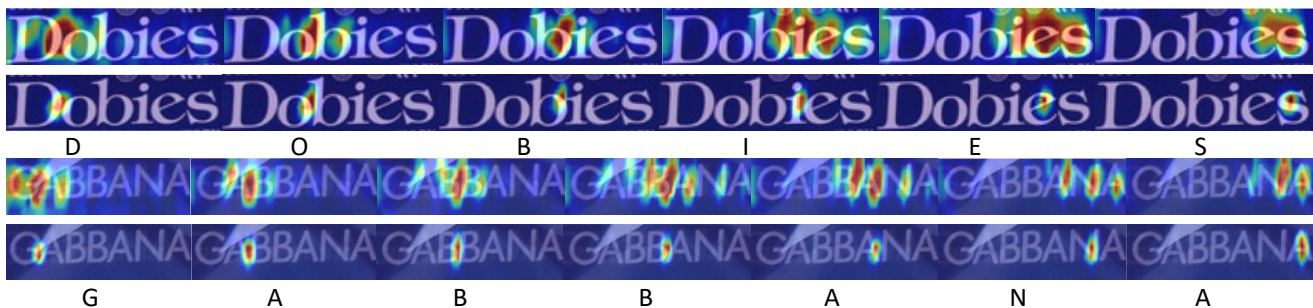


Figure 4: Visualization of attention heatmaps of the primary and secondary vs-align modules. For each of the two examples, the top (bottom) row shows the heatmaps of the primary (secondary) vs-align module.

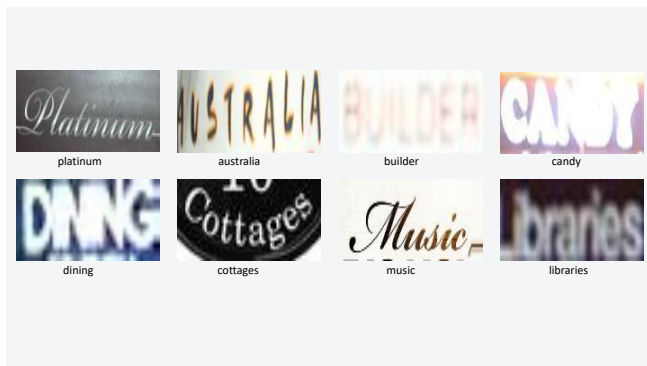


Figure 5: Successful cases for fairly hard examples.

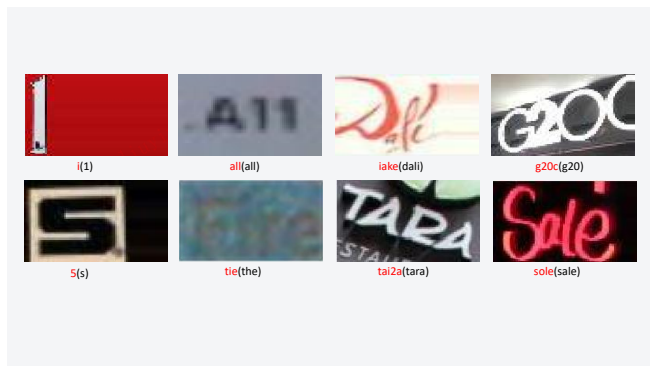


Figure 6: Failure cases for some confusing images.

## Ablation Study

To verify whether each module is critical for the final performance, ablation study is conducted by subsequently adding one module starting from the most basic C+V configuration. The results are shown in table 2, from which we observe a consistent performance gain when a new module is added. Also note that C+V+A also performs better than C+V, which verifies the effectiveness of vs-align module. For the C+V+A+I configuration, there are three ways of decoding: decoding from secondary semantics ( $s_2$ ), tertiary semantics ( $s_3$ ) or by probability voting from both semantics.

## Visualization

Inside the module I, each semantic vector will attend to visual features corresponding to some spatial locations. We can visualize which locations are attended to using heatmaps. Fig. 3 shows three examples. The heatmaps are computed by averaging all 8 attention heads and overlapping onto images after upsampling. Note that inside the module each head focuses on different aspects of the images, but on average the attention is mainly at the right spatial location for each decoded character. For some characters, the focus is slightly shifted, probably due to the translation-invariant property of convnet, which means that the visual features cannot be mapped back uniformly to the corresponding spatial locations. Overall, the visualization proves that the interaction module works as expected.

Likewise, we also visualize the attention heatmaps of the primary and secondary vs-align modules. Even though the two modules share the same weights, their attention maps should be different since the visual and semantic features are different. We expect that the attention maps of 2nd vs-align be more precise than the first one. This is verified in Fig. 4.

Beside visualization, we select some successful and failure cases for illustration purpose, shown in Fig. 5 and Fig. 6 respectively.

## Conclusion

In this paper, we propose visual-semantic transformer to solve scene text recognition problem. The VST consists of three model explicitly extract semantic features in the early stage by using weight-sharing visual-semantic alignment module, making the visual and semantic fusion and interaction possible in the following transforms modules. Our model is unified and end-to-end trainable, and it does not require autoregressive decoding. Extensive experiments on regular and irregular scene text recognition datasets have verified the effectiveness of our model.

## References

Aberdam, A.; Litman, R.; Tsiper, S.; Anshel, O.; Slossberg, R.; Mazor, S.; Manmatha, R.; and Perona, P. 2021.

- Sequence-to-sequence contrastive learning for text recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 15302–15312.
- Baek, J.; Kim, G.; Lee, J.; Park, S.; Han, D.; Yun, S.; Oh, S. J.; and Lee, H. 2019. What Is Wrong With Scene Text Recognition Model Comparisons? Dataset and Model Analysis. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 4715–4723.
- Baevski, A.; Zhou, H.; Mohamed, A.; and Auli, M. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *arXiv preprint arXiv:2006.11477*.
- Bhunia, A.; Sain, A.; Kumar, A.; Ghose, S.; Chowdhury, P. N.; and Song, Y.-Z. 2021. Joint Visual Semantic Reasoning: Multi-Stage Decoder for Text Recognition. *ArXiv*, abs/2107.12090.
- Chen, X.; Jin, L.; Zhu, Y.; Luo, C.; and Wang, T. 2021. Text recognition in the wild: A survey. *ACM Computing Surveys (CSUR)*, 54(2): 1–35.
- Cheng, Z.; Xu, Y.; Bai, F.; Niu, Y.; Pu, S.; and Zhou, S. 2018. AON: Towards Arbitrarily-Oriented Text Recognition. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5571–5579.
- Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Zhai, X.; Unterthiner, T.; Dehghani, M.; Minderer, M.; Heigold, G.; Gelly, S.; et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Fang, S.; Xie, H.; Wang, Y.; Mao, Z.; and Zhang, Y. 2021. Read Like Humans: Autonomous, Bidirectional and Iterative Language Modeling for Scene Text Recognition.
- Feng, W.; Yin, F.; Zhang, X.-Y.; and Liu, C.-L. 2021. Semantic-Aware Video Text Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1695–1705.
- Graves, A.; Fernández, S.; Gomez, F.; and Schmidhuber, J. 2006. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, 369–376.
- Gupta, A.; Vedaldi, A.; and Zisserman, A. 2016. Synthetic Data for Text Localisation in Natural Images. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- Hu, W.; Cai, X.; Hou, J.; Yi, S.; and Lin, Z. 2020. GTC: Guided Training of CTC Towards Efficient and Accurate Scene Text Recognition. In *AAAI*.
- Jaderberg, M.; Simonyan, K.; Vedaldi, A.; and Zisserman, A. 2014. Synthetic Data and Artificial Neural Networks for Natural Scene Text Recognition. *arXiv preprint arXiv:1406.2227*.
- Jaderberg, M.; Simonyan, K.; Vedaldi, A.; and Zisserman, A. 2015. Deep Structured Output Learning for Unconstrained Text Recognition. In *ICLR 2015 : International Conference on Learning Representations 2015*.
- Jaderberg, M.; Simonyan, K.; Vedaldi, A.; and Zisserman, A. 2016. Reading Text in the Wild with Convolutional Neural Networks. *International Journal of Computer Vision*, 116(1): 1–20.
- Jaderberg, M.; Simonyan, K.; Zisserman, A.; et al. 2015. Spatial transformer networks. In *Advances in neural information processing systems*, 2017–2025.
- Jaderberg, M.; Vedaldi, A.; and Zisserman, A. 2014. Deep Features for Text Spotting. In *European Conference on Computer Vision*, 512–528.
- Jang, E.; Gu, S.; and Poole, B. 2016. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*.
- Karatzas, D.; Gomez-Bigorda, L.; Nicolaou, A.; Ghosh, S.; Bagdanov, A.; Iwamura, M.; Matas, J.; Neumann, L.; Chandrasekhar, V. R.; Lu, S.; Shafait, F.; Uchida, S.; and Valveny, E. 2015. ICDAR 2015 competition on Robust Reading. In *2015 13th International Conference on Document Analysis and Recognition (ICDAR)*, 1156–1160.
- Karatzas, D.; Shafait, F.; Uchida, S.; Iwamura, M.; Bigorda, L. G.; Mestre, S. R.; Mas, J.; Mota, D. F.; Almazan, J. A.; and de las Heras, L.-P. 2013. ICDAR 2013 Robust Reading Competition. In *2013 12th International Conference on Document Analysis and Recognition*, 1484–1493.
- Li, H.; Wang, P.; Shen, C.; and Zhang, G. 2019. Show, attend and read: A simple and strong baseline for irregular text recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 8610–8617.
- Litman, R.; Anshel, O.; Tsiper, S.; Litman, R.; Mazor, S.; and Manmatha, R. 2020. Scatter: selective context attentional scene text recognizer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11962–11972.
- Litman, R.; Anshel, O.; Tsiper, S.; Litman, R.; Mazor, S.; and Manmatha, R. 2020. SCATTER: Selective Context Attentional Scene Text Recognizer. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 11962–11972.
- Liu, W.; Chen, C.; Wong, K.-Y. K.; Su, Z.; and Han, J. 2016. STAR-Net: A SpaTial Attention Residue Network for Scene Text Recognition. In *BMVC*, volume 2, 7.
- Long, S.; He, X.; and Yao, C. 2021. Scene Text Detection and Recognition: The Deep Learning Era. *International Journal of Computer Vision*, 129(1): 161–184.
- Lucas, S.; Panaretos, A.; Sosa, L.; Tang, A.; Wong, S.; and Young, R. 2003. ICDAR 2003 robust reading competitions. In *Seventh International Conference on Document Analysis and Recognition, 2003. Proceedings.*, volume 3, 682–687.
- Lyu, P.; Yang, Z.; Leng, X.; Wu, X.; Li, R.; and Shen, X. 2019. 2d attentional irregular scene text recognizer. *arXiv preprint arXiv:1906.05708*.
- Mishra, A.; Alahari, K.; and Jawahar, C. 2012. Scene Text Recognition using Higher Order Language Priors. In *Proceedings of the British Machine Vision Conference*, 127.1–127.11. BMVA Press. ISBN 1-901725-46-4.
- Nguyen, N.; Nguyen, T.; Tran, V.; Tran, M.-T.; Ngo, T. D.; Nguyen, T. H.; and Hoai, M. 2021. Dictionary-Guided Scene Text Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7383–7392.



- Patel, Y.; Gomez, L.; Rusinol, M.; and Karatzas, D. 2016. Dynamic lexicon generation for natural scene images. In *European Conference on Computer Vision*, 395–410. Springer.
- Phan, T. Q.; Shivakumara, P.; Tian, S.; and Tan, C. L. 2013. Recognizing Text with Perspective Distortion in Natural Scenes. In *2013 IEEE International Conference on Computer Vision*, 569–576.
- Qiao, Z.; Zhou, Y.; Yang, D.; Zhou, Y.; and Wang, W. 2020. Seed: Semantics enhanced encoder-decoder framework for scene text recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 13528–13537.
- Qiao, Z.; Zhou, Y.; Yang, D.; Zhou, Y.; and Wang, W. 2020. SEED: Semantics Enhanced Encoder-Decoder Framework for Scene Text Recognition. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 13528–13537.
- Risnumawan, A.; Shivakumara, P.; Chan, C. S.; and Tan, C. L. 2014. A robust arbitrary text detection system for natural scene images. *Expert Systems With Applications*, 41(18): 8027–8048.
- Sabir, A.; Moreno-Noguer, F.; and Padró, L. 2018. Visual re-ranking with natural language understanding for text spotting. In *Asian Conference on Computer Vision*, 68–82. Springer.
- Sheng, F.; Chen, Z.; and Xu, B. 2019. NRTR: A No-Recurrence Sequence-to-Sequence Model for Scene Text Recognition. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, 781–786.
- Shi, B.; Bai, X.; and Yao, C. 2016. An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE transactions on pattern analysis and machine intelligence*, 39(11): 2298–2304.
- Shi, B.; Yang, M.; Wang, X.; Lyu, P.; Yao, C.; and Bai, X. 2018. Aster: An attentional scene text recognizer with flexible rectification. *IEEE transactions on pattern analysis and machine intelligence*, 41(9): 2035–2048.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. In *Advances in neural information processing systems*, 5998–6008.
- Wan, Z.; Zhang, J.; Zhang, L.; Luo, J.; and Yao, C. 2020. On Vocabulary Reliance in Scene Text Recognition. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 11422–11431.
- Wang, K.; Babenko, B.; and Belongie, S. 2011. End-to-end scene text recognition. In *2011 International Conference on Computer Vision*, 1457–1464.
- Wang, K.; Babenko, B.; and Belongie, S. 2011. End-to-end scene text recognition. In *2011 International Conference on Computer Vision*, 1457–1464. IEEE.
- Wang, Q.; Li, B.; Xiao, T.; Zhu, J.; Li, C.; Wong, D. F.; and Chao, L. S. 2019a. Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*.
- Wang, T.; Zhu, Y.; Jin, L.; Luo, C.; Chen, X.; Wu, Y.; Wang, Q.; and Cai, M. 2020. Decoupled Attention Network for Text Recognition. In *AAAI*, 12216–12224.
- Wang, T.; Zhu, Y.; Jin, L.; Peng, D.; Li, Z.; He, M.; Wang, Y.; and Luo, C. 2021a. Implicit Feature Alignment: Learn to Convert Text Recognizer to Text Spotter. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5973–5982.
- Wang, W.; Xie, E.; Sun, P.; Wang, W.; Tian, L.; Shen, C.; and Luo, P. 2019b. TextSR: Content-aware text super-resolution guided by recognition. *arXiv preprint arXiv:1909.07113*.
- Wang, Y.; Xie, H.; Fang, S.; Wang, J.; Zhu, S.; and Zhang, Y. 2021b. From Two to One: A New Scene Text Recognizer with Visual Language Modeling Network. *arXiv preprint arXiv:2108.09661*.
- Wu, L.; Zhang, C.; Liu, J.; Han, J.; Liu, J.; Ding, E.; and Bai, X. 2019. Editing Text in the Wild. In *Proceedings of the 27th ACM International Conference on Multimedia*, 1500–1508.
- Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A.; Salakhudinov, R.; Zemel, R.; and Bengio, Y. 2015. Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning*, 2048–2057. PMLR.
- Yan, R.; Peng, L.; Xiao, S.; and Yao, G. 2021. Primitive Representation Learning for Scene Text Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 284–293.
- Yang, L.; Wang, P.; Li, H.; Li, Z.; and Zhang, Y. 2020. A Holistic Representation Guided Attention Network for Scene Text Recognition. *Neurocomputing*.
- Yu, D.; Li, X.; Zhang, C.; Liu, T.; Han, J.; Liu, J.; and Ding, E. 2020. Towards accurate scene text recognition with semantic reasoning networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12113–12122.
- Yu, H.; Zhang, C.; Li, X.; Han, J.; Ding, E.; and Wang, L. 2019. An End-to-End Video Text Detector with Online Tracking. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*, 601–606.
- Yue, X.; Kuang, Z.; Lin, C.; Sun, H.; and Zhang, W. 2020. RobustScanner: Dynamically Enhancing Positional Clues for Robust Text Recognition. In *European Conference on Computer Vision*, 135–151. Springer.
- Zheng, Y.; Wang, Q.; and Betke, M. 2019. Deep neural network for semantic-based text recognition in images. *arXiv preprint arXiv:1908.01403*.