

# Learning Generative Structure Prior for Blind Text Image Super-resolution (Supplementary Material)

Xiaoming Li<sup>1</sup>    Wangmeng Zuo<sup>2</sup>    Chen Change Loy<sup>1</sup>(✉)  
<sup>1</sup>S-Lab, Nanyang Technological University    <sup>2</sup>Harbin Institute of Technology  
 csxmli@gmail.com    wmzuo@hit.edu.cn    ccloy@ntu.edu.sg

This supplemental material mainly contains:

- Analyses of model complexity in Section I
- Architecture details in Section II
- Broader impact in Section III
- Structure images of 182 font families that are used in synthesizing the HR text images in Figure A
- More visual comparison on synthetic and real-world LR text images in Figures B, C, D and E
- Analyses of recognition prior in competing methods and our style vector  $w$  in Figure F
- Analyses of failure cases in Figure G.

## I. Model Complexity

The whole framework of our MARCONet consists of three parts, 1) learnable codebook and StyleGAN with 27.97 M parameters, 2) Transformer encoder with 59.66 M parameters, and 3) SR network with 16.86 M parameters. It takes 0.21s on average to super-resolve a  $32 \times 512$  LR input to a  $128 \times 2048$  SR result with one Tesla V100 GPU.

## II. Architecture Details

The Transformer encoder follows the default settings of the vanilla Transformer [5]. Here we mainly show the architectural details of the structure prior transform module and final reconstruction module. The StyleGAN used for text images in this work has  $i \in \{8, 16, 32, 64, 128\}$  scales in total and each scale  $i$  represents the size of intermediate features. We adopt scales  $i \in \{32, 64\}$  on the multi-scale structure prior transform modules. Details are presented in Table A.  $\text{SNConv}(d, k, s)$  represents a convolutional layer followed by a spectral normalization [3], where  $d, k$  and  $s$  are output dimension, kernel size and stride, respectively.  $\text{LReLU}(c)$  is leaky ReLU with a negative slope  $c$ . Swish is a self-gated activation [4] and is defined as  $f(x) = x \cdot \text{sigmoid}(x)$ .

## III. Broader Impact

Our text super-resolution framework is of great value in real-world scenarios, e.g., the restoration of old printed media for digitalization and preservation, and the enhancement

of captions in digital media, license plates, and invoices. It can also be employed in image acquisition devices to enhance the visual quality of text regions and further improve optical character recognition (OCR). As with other computational imaging techniques, the text SR technology may be misused, causing information leakage or other issues in privacy. We note that the accuracy of text SR will be compromised under severe degradation, causing errors in the final character output. Such an error may be unacceptable for forensic document analysis. It is worth noting that the positive impact of text SR outweighs the potential problems. We call on people to use this technology without harming personal information and take the SR result as a reference.

Table A. Architecture details of our multi-scale structure prior transform modules and final reconstruction module.

Structure Prior Transform on $32 \times 32$		
Input	LR Feature	Character Structure Prior
<b>RoIAlign</b>		
<b>For Each Character</b>	$F_{LR}^c$	$\mathcal{P}^c$
	AdaIN ( $F_{LR}^c, \mathcal{P}^c$ ) Concatenate output with $F_{LR}^c$	
<b>ResBlock</b>	GroupNorm, Swish, SNConv (256, 3, 1) GroupNorm, Swish, SNConv (256, 3, 1) Element-wise Addition	
<b>Spatial Feature Transform</b>	SNConv (256, 3, 1) LReLU (0.2)	SNConv (256, 3, 1) LReLU (0.2)
	SNConv (256, 3, 1) Scale $\alpha$	SNConv (256, 3, 1) Shift $\beta$
$\alpha \cdot F_{LR}^c + \beta$		
<b>Reverse RoIAlign</b>		
<b>UpSample</b> (32→64)	Bilinear UpSample SNConv (256, 3, 1), LReLU (0.2) ResBlock, SNConv (256, 3, 1)	
<b>Structure Prior Transform on <math>64 \times 64</math></b> (Kindly refer to the former steps)		
<b>Recon- struction</b> (64→128)	SNConv (128, 3, 1), LReLU (0.2) Bilinear Upsample SNConv (64, 3, 1), LReLU (0.2) ResBlock SNConv (3, 3, 1) Tanh	
<b>Output</b>	SR Result	

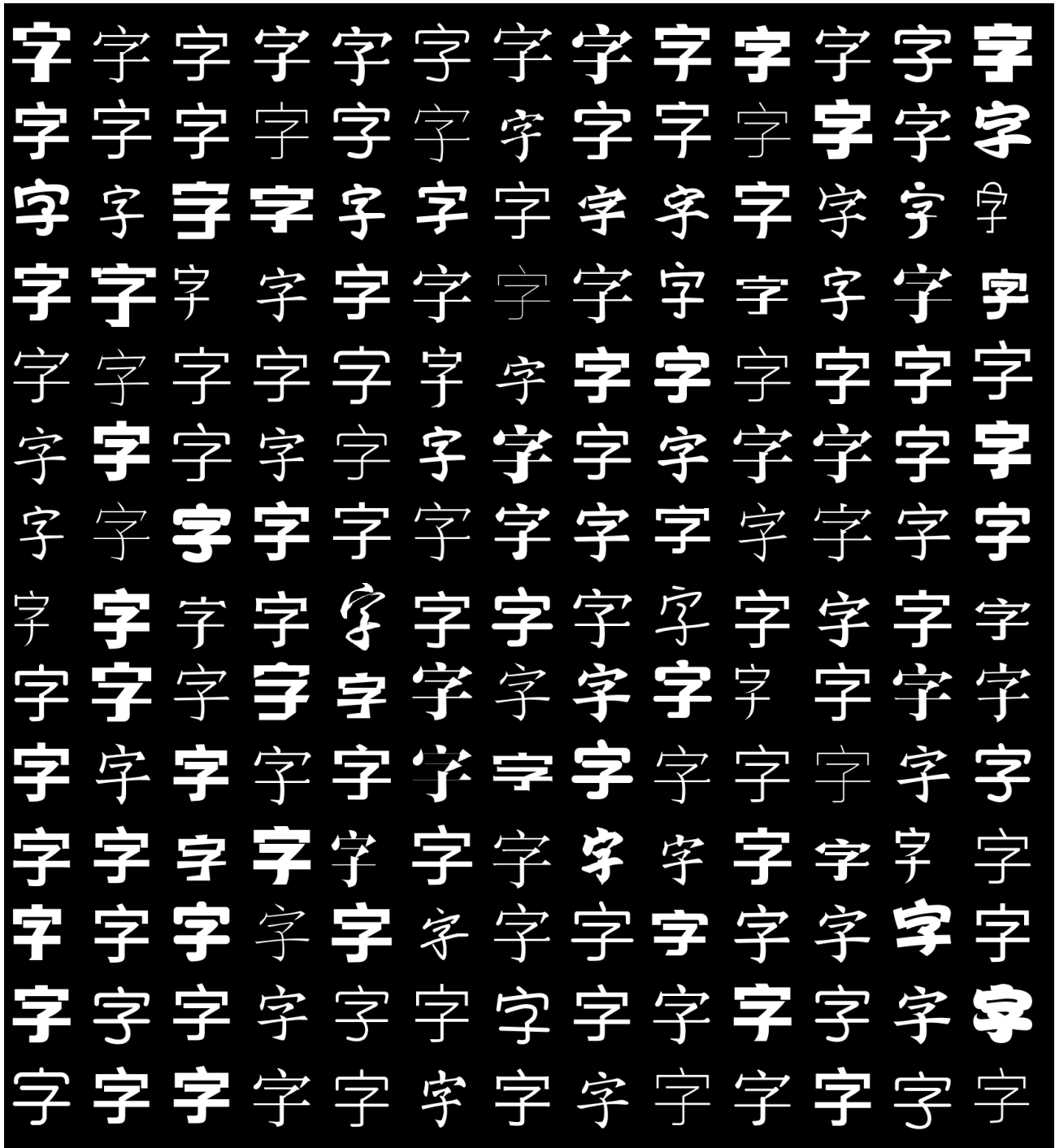


Figure A. Structure images of the 182 font families used in synthesizing our training data.

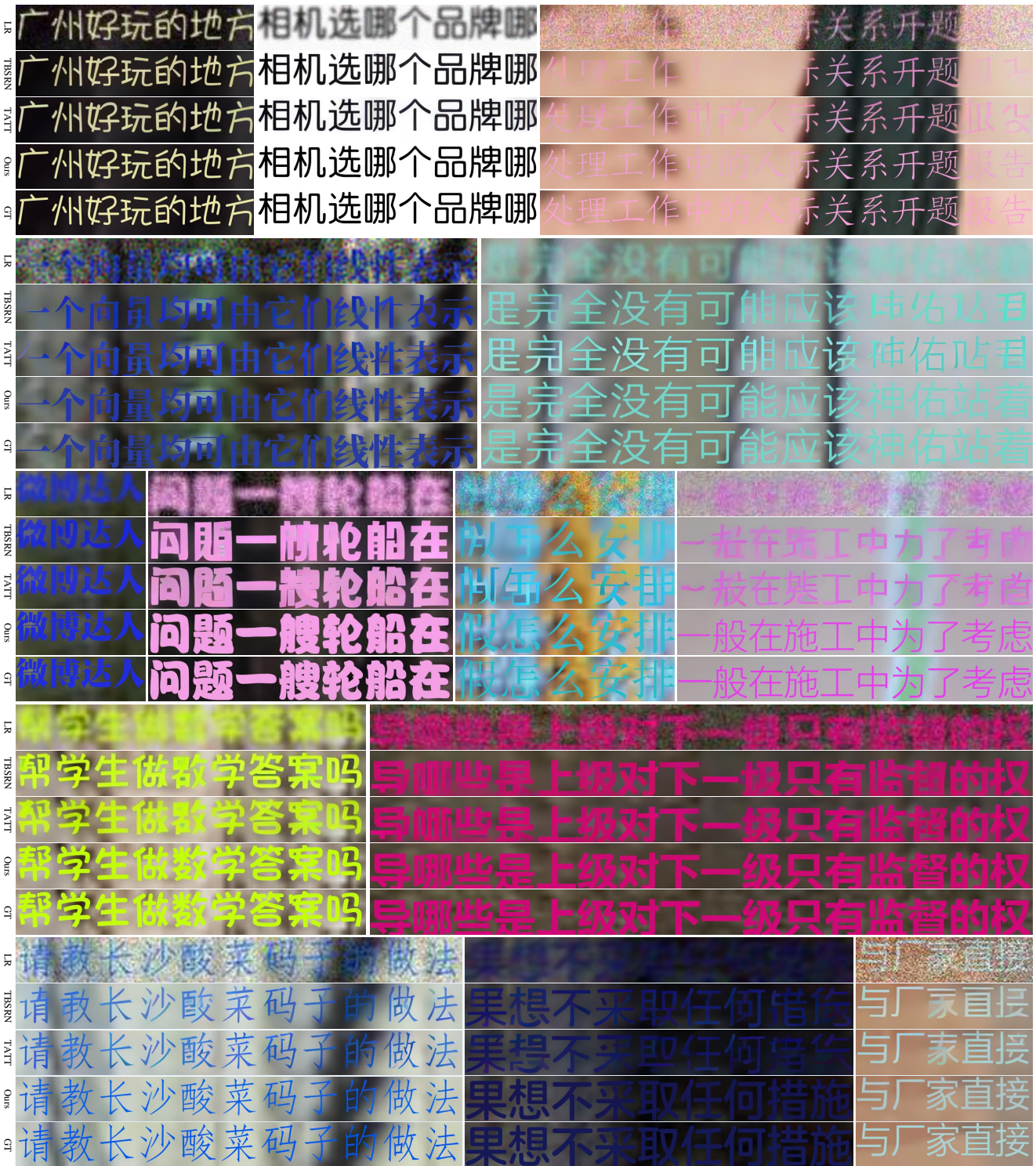


Figure B. More visual comparison with competing methods (*i.e.*, retained TBSRN [1] and TATT [2]) on our synthetic LR text images.



Figure C. Results obtained with retrained TBSRN [1], TATT [2] and our approach on several segments of real-world LR text images.

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

(a) Real-world LR segment from an old newspaper

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

建屋发展局也已停止接受市区组屋的申请。李玉胜次长解释说，以目前的环境而言，要在市区建造1千个单位的组屋也非常困难，这是由于土地的极端缺乏。他说，在市区一年只能建造750个单位的组屋，申请人数如有6,000人，这批申请者便要等待10年才能全部分配到组屋。这主要是因为市区里可供兴建组屋的地皮近乎用完。

(b) Restoration results of our MARCONet

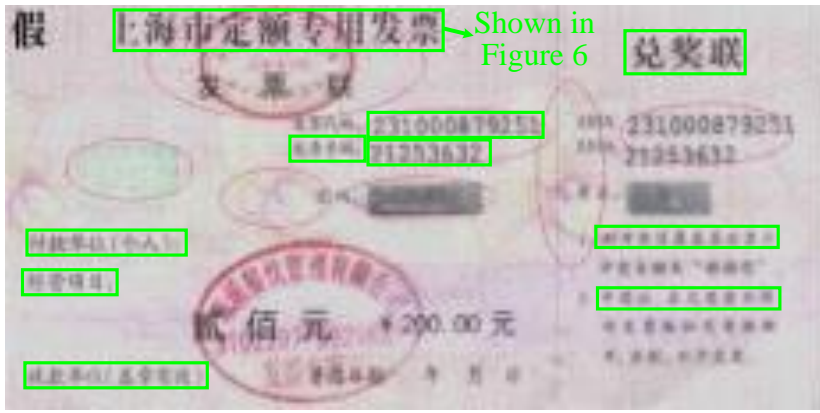
Figure D. Results of our approach on super-resolving the segments from a real-world LR newspaper.



老重庆	烂摊摊	火锅
老重庆	烂摊摊	火锅
订餐电话:		15822870771
全球最难吃的		火锅
全球最难吃的		火锅



学府路	学府路
江滨路	江滨路
宗泽路	宗泽路
象山路	象山路
谷阳路	谷阳路

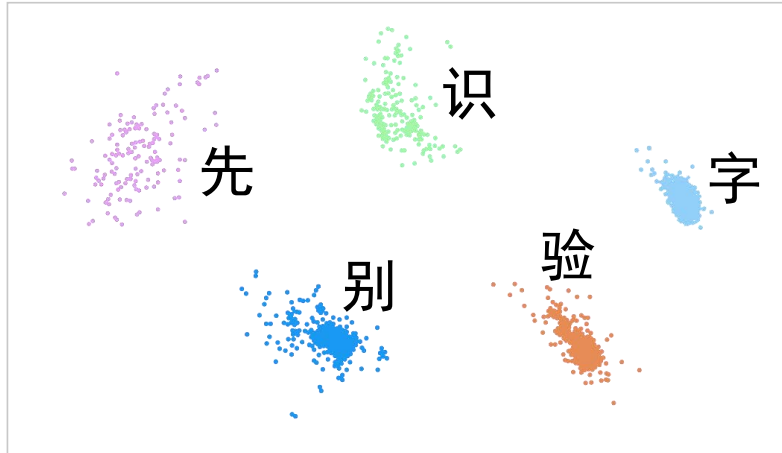


发票号码:	免奖联
发票号码:	免奖联
231000879251	
231000879251	
付款单位(个人):	
付款单位(个人):	

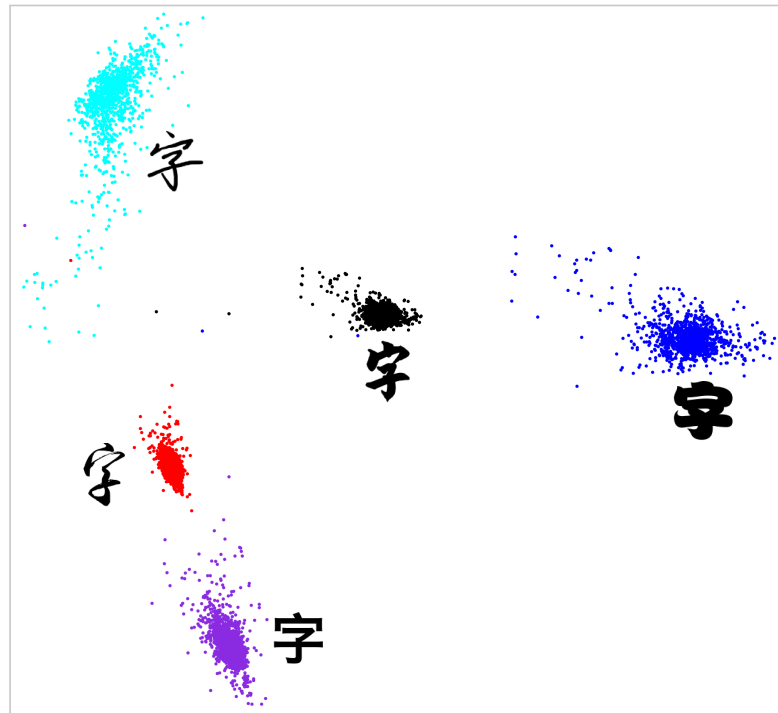
经营项目:	21253632	收款单位(盖章有效)
经营项目:	21253632	收款单位(盖章有效)

刮开奖区覆盖层后显示	中奖后,在兑奖前不得
刮开奖区覆盖层后显示	中奖后,在兑奖前不得

Figure E. More results of our approach on super-resolving the segments from different LR sources.

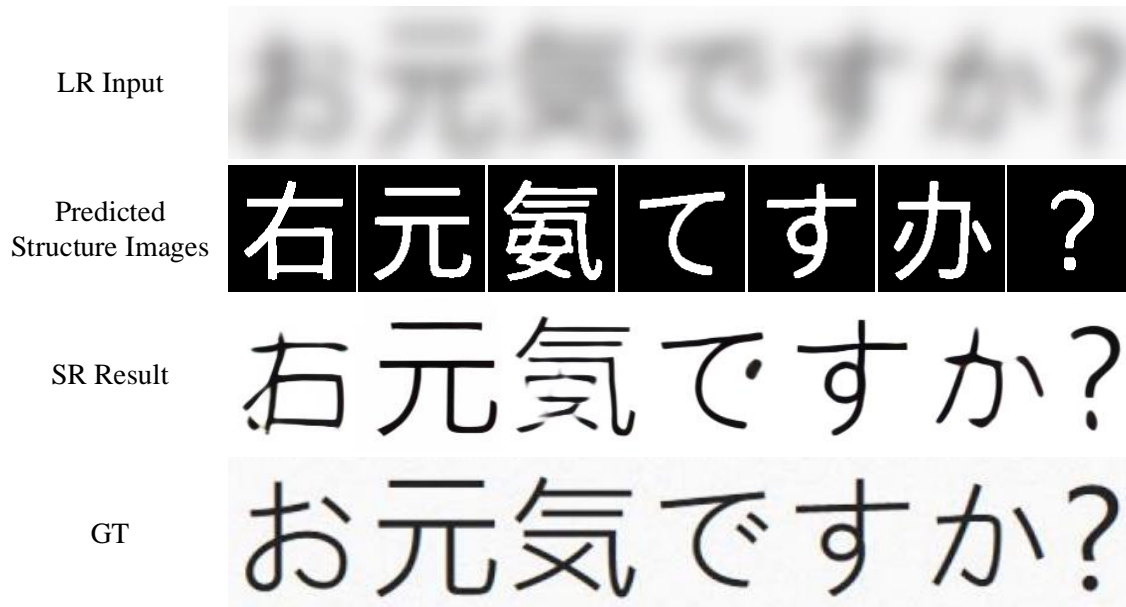


(a) The t-SNE results on recognition prior from the final layer of the character recognizer. We adopt 5 characters (see them in the plot) and synthesize their images with random degradation and 182 font families, respectively.

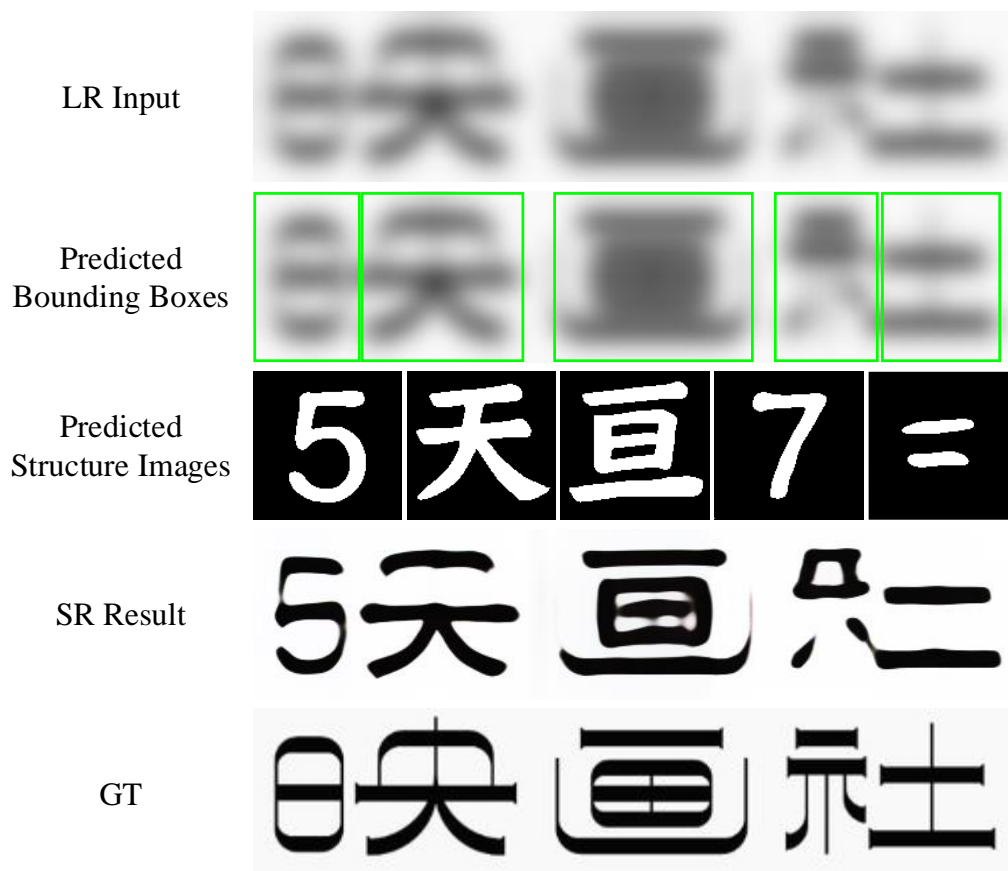


(b) The t-SNE results on our style space  $W$ . We build 5 groups of text images and each group has the same font family but is synthesized with random text and degradation (see their font style examples in the plot).

Figure F. Analyses of recognition prior used in competing methods and our style vector  $w$ . We can observe that (a) recognition prior can well distinguish each character but is not enough to aid SR task because it easily ignores the font styles which are crucial for super-resolving a faithful and accurate structure. So these recognition prior based methods (*i.e.*, TBSRN [1] and TATT [2]) have limited performance on this task; (b) in contrast, the style vector  $w$  in our framework can easily capture the styles of different font families, no matter what text and degradation they have. By combining style vector  $w$  with the codebook, our MARCONet can well benefit the SR task by providing accurate structure-level guidance.



(a) Results of our MARCONet which is trained on Chinese but is adopted on Kanji text image SR



(b) Results of our MARCONet on LR text image which has unusual font style

Figure G. Analyses of failure cases. We can observe that (a) when directly using our MARCONet on another language (*i.e.*, Kanji), the error in predicting character classification easily leads to poor guidance for the SR results, (b) when the LR input has an unusual font style, our MARCONet easily has an inaccurate prediction in character classification or bounding box regression, which easily results in the inconsistent structures on the SR results. For this case, we will add more font families in our training phase to improve the generalization performance.



## References

- [1] Jingye Chen, Bin Li, and Xiangyang Xue. Scene text telescope: Text-focused scene image super-resolution. In *CVPR*, 2021. 3, 4, 7
- [2] Jianqi Ma, Zhetong Liang, and Lei Zhang. A text attention network for spatial deformation robust scene text image super-resolution. In *CVPR*, 2022. 3, 4, 7
- [3] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. In *ICLR*, 2018. 1
- [4] Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017. 1
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017. 1