

A Benchmark for Chinese-English Scene Text Image Super-resolution

Jianqi Ma^{1,2}, Zhetong Liang², Wangmeng Xiang¹, Xi Yang^{1,2}, Lei Zhang^{1,2}

¹The Hong Kong Polytechnic University; ²OPPO Research

{csjma, cswmxiang, cslzhang, csxyang}@comp.polyu.edu.hk, zhetongliang@163.com

Abstract

Scene Text Image Super-resolution (STISR) aims to recover high-resolution (HR) scene text images with visually pleasant and readable text content from the given low-resolution (LR) input. Most existing works focus on recovering English texts, which have relatively simple character structures, while little work has been done on the more challenging Chinese texts with diverse and complex character structures. In this paper, we propose a real-world Chinese-English benchmark dataset, namely Real-CE, for the task of STISR with the emphasis on restoring structurally complex Chinese characters. The benchmark provides 1,935/783 real-world LR-HR text image pairs (contains 33,789 text lines in total) for training/testing in $2\times$ and $4\times$ zooming modes, complemented by detailed annotations, including detection boxes and text transcripts. Moreover, we design an edge-aware learning method, which provides structural supervision in image and feature domains, to effectively reconstruct the dense structures of Chinese characters. We conduct experiments on the proposed Real-CE benchmark and evaluate the existing STISR models with and without our edge-aware loss. The benchmark, including data and source code, is available at <https://github.com/mjq11302010044/Real-CE>.

1. Introduction

Text images are different from natural images in that the main contents are composed of words and characters to express different meanings and ideas. Due to limited sensor resolution and long photographing distance, the captured text images often have degraded quality with blurry and noisy contents, impairing the readability of the text. Therefore, scene text image super-resolution (STISR) is demanded to reconstruct clear and legible text contents.

STISR has long been studied in the computer vision community [30, 40, 36, 31]. The traditional STISR methods investigate various priors on text restoration and hand-craft the text super-resolution process [1, 8]. Since the manually designed priors cannot represent the complex text struc-

tures and degradation process, the traditional methods have limited performance. Deep learning based STISR methods train convolutional neural networks (CNNs) on datasets with low-resolution (LR) and high-resolution (HR) text image pairs, which can learn the complex text priors through data and reconstruct high-quality text images.

In deep learning based STISR [30, 40, 36, 31], datasets play an important role in model training and evaluation, because the image pairs encode the text transformation from low to high resolution. In the early stage, synthetic datasets are widely used [30, 40, 31], in which high-quality text images are collected as HR ground truths, and the LR images are generated by imposing synthetic degradations (e.g., bicubic downsampling or blurring) on the HR images. Since the real-world degradations are quite different from the synthetic ones, the STISR models trained on the synthetic datasets have limited performance on real-world LR text images. To alleviate this problem, Wang *et al.* [35] built a real-world text image dataset called TextZoom. The LR and HR text images in TextZoom are captured with different camera focal lengths and undergo the real-world degradation process. TextZoom provides a benchmark for the STISR task, which allows standardized evaluation of STISR methods in terms of text recognition precision.

Though the TextZoom dataset has largely facilitated the research of real-world STISR [35, 4, 25, 6, 46, 26, 47], it has some limitations. First, TextZoom only contains English texts composed of limited number of characters (*i.e.*, 26 letters) with simple stroke structures. As a result, models trained on TextZoom will produce inferior results on structurally complex characters like Chinese. Examples are shown in Figure 1(b). One can see that the model trained on TextZoom produces visually unpleasant artifacts on the reconstructed Chinese texts. This is because Chinese texts have a much larger number of characters, and many of them have complex structures. Thus, it is a more challenging task for performing STISR on Chinese texts. Moreover, TextZoom focuses on small and fixed-size text images (*i.e.*, 32×128), and thus the models trained on TextZoom cannot generalize to texts with various resolutions. Therefore, new dataset and benchmark are highly demanded for the re-

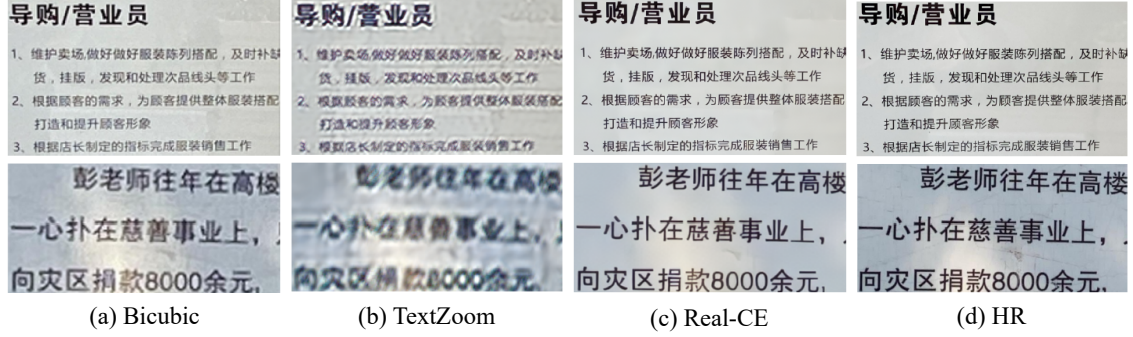


Figure 1. Comparison of STISR results on Chinese text images by methods trained on TextZoom and Real-CE datasets. From left to right are (a) bicubic LR images, STISR outputs by RRDB model [38] trained on (b) TextZoom [35] and (c) our Real-CE, and (d) the ground-truth HR text images. Please zoom in for more details.

search of STISR on Chinese text images.

To tackle the above-mentioned problems, in this work we develop a novel real-world Chinese-English benchmark dataset, termed Real-CE, for the training and evaluation of STISR models on both Chinese and English texts. The benchmark provides 1,935 real-world LR-HR image pairs for training, and 783 for testing (261 and 522 pairs for $4\times$ and $2\times$ zooming modes, respectively). It contains 24,666 Chinese text lines and 9,123 English text lines in total with different sizes. Detailed annotations on the image pairs, including detection boxes and text transcripts, are also provided to assist the training and evaluation. We also design the evaluation process to adapt to different sizes of text lines, aiming to preserve the visual quality of SR text images from resizing. Furthermore, we propose an edge-aware learning method for the reconstruction of Chinese texts with complex stroke structures. The text edge map is introduced as the network input as well as a structural loss in the training process, enhancing the learning on text structural regions. Experimental results show that models trained on our Real-CE data achieve superior performance over TextZoom on Chinese text super-resolution (as shown in Figure 1(c)) and the edge-aware learning can further promote the reconstruction quality on text regions.

The paper is organized as follows. Section 2 reviews the works on STISR research. Section 3 introduces the Real-CE benchmark in detail. Section 4 describes the edge-aware learning method. Section 5 shows the experimental results on the benchmark and Section 6 concludes the paper.

2. Related Work

Our work is related to single image super-resolution (SISR), scene text image super-resolution (STISR), and English and Chinese text recognition, as reviewed below.

SISR. SISR estimates a high-resolution (HR) output by intaking the low-resolution (LR) image as input. Traditional approaches apply manually-designed priors for this task in

terms of statistical information [14], self-similarity [28] and sparsity [41]. Recent deep-learning methods employ convolutional neural networks (CNNs) for SISR and achieve significantly better performance. As a pioneer work, SRCNN [11] adopts a three-layer CNN to perform HR estimation. Later on, more elaborate designs on network architecture further upgrade the SISR performance, including residual connection [22], Laplacian pyramid [19], dense connection block [38] and the Transformer architecture [21, 44]. Adversarial learning techniques have also been applied for more photo-realistic results [20, 37].

STISR benchmarks and methods. STISR focuses on scene text images. It aims to reconstruct the text shape by upgrading the image resolution in order to benefit the downstream recognition task. The early methods of STISR directly adopt the CNN architectures used in general SISR tasks. In [12], Dong *et al.* adopted SRCNN [11] to text images, and achieved state-of-the-art performance in ICDAR 2015 competition [30]. PlugNet [29] employs a pluggable super-resolution unit to learn the semantics in LR images in feature domain. TextSR [36] utilizes the text recognition loss to supervise SR recovery learning and improve the text recognition. Aiming to learn text image deblurring and super-resolution, Xu *et al.* [40] and Quan *et al.* [31] collected high-quality document text data to evaluate synthetic image deblurring and super-resolution.

To address real-world STISR problems, Wang *et al.* [35] built an STISR benchmark, namely TextZoom, which provides the LR and HR text image pairs extracted from real-world SISR datasets [43, 2]. They also proposed TSRN [35] by applying the sequential residual block to model the sequential semantics in image features. SC-GAN [40] adopts GAN loss to supervise the STISR model for more realistic text images. Quan *et al.* [31] proposed a cascading network for reconstructing high-quality text images in both high-frequency domain and image domain. Chen *et al.* [4, 6], Zhao *et al.* [46] and Ma *et al.* [26] upgraded the network block structures to enhance the STISR

performance with transformer-based networks or text prior.

However, current STISR methods are designed only for English-based text line images (*i.e.*, the TextZoom) with many limitations. We therefore make an attempt to build a bilingual benchmark to fill in the blank.

Scene text recognition. Scene text recognition (STR) aims to recognize the semantic meaning in the text image by predicting the characters or the whole word [18, 15, 17, 23]. It can be considered as an image-to-sequence problem. CRNN [32] uses recurrent neural networks to model semantic information. Recently, attention-based methods have achieved great success due to their robustness against shape variations of text images [9, 10, 33]. However, most methods are proposed for English text, and Chinese scene text recognition receives less attention [13, 16, 7]. To promote research along this line, Chen *et al.* [5] attempted to benchmark the Chinese scene text recognition with unified input and evaluation metrics. In this paper, text recognizers in both languages are adopted for evaluating text recognition after STISR.

3. Real-CE Benchmark

The proposed benchmark includes a dataset with Chinese and English LR-HR text image pairs and an evaluation protocol with five metrics.

3.1. Dataset Construction

The dataset is constructed by several steps, including data collection, registration, text cropping and text labeling, which are illustrated in Figure 2.

Data collection. We adopt iPhone 11 pro and iPhone 12 pro for text image collection. Both of them are equipped with camera modules of three fixed focal lengths (13 mm, 26 mm and 52 mm), which allow us to capture the same scene with different focal lengths simultaneously. Image pairs collected by these devices enable the training of STISR models in 2 \times (from 13 mm to 26 mm, and from 26 mm to 52 mm) and 4 \times (from 13 mm to 52 mm) zooming modes. We capture the Chinese and English text images from various scenes and resources, including band and curve outdoor street signboards, subway notifications, deformed books and hospital billboards, so that the diversity of text contents, presentation and lighting conditions can be ensured. Since the three cameras may have different image processing pipelines, we use CameraPixels app¹ to align the colors and brightness of the three captured images. Figure 3 shows some typical scenes in the collected dataset.

Image registration. We adopt the image registration method proposed in Cai *et al.* [2] for the alignment of LR and HR text image pairs. Specifically, we take the images captured by 52mm lens as the ground-truth HR images since

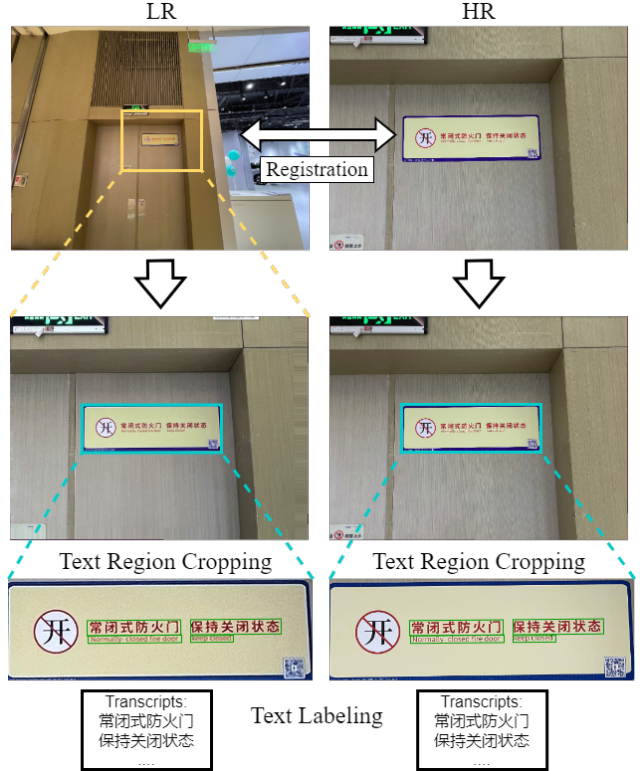


Figure 2. The pipeline of data processing. From top to bottom: the center area of the LR image is first registered to the HR image, then the corresponding text regions in the HR and LR images are cropped and manually aligned; finally, the text lines are annotated and the transcripts are labeled.

they have the best quality, and iteratively register shorter-length LR counterparts to it. The algorithm also enables finer-grained adjustment, which reduces the color and brightness differences between the LR-HR image pairs.

Text region cropping. Though the text images are centered on text content, they still contain a large proportion of background area. We therefore crop the central text region from the LR and HR images to exclude background areas, followed by a manual adjustment to ensure accurate LR-HR image pair alignment.

Text labeling. Besides the HR ground truths, we provide two extra text labels, including detection boxes and text transcripts. The detection boxes provide the location of the text areas, while the text transcripts record the semantics of the texts. For the detection box annotation, we first apply some text detection methods (*e.g.*, RRPN [27, 24]) to provide a coarse detection result, then we refine the detection results manually to provide a precise boundary of each text line in the cropped text region. For the text transcripts annotation, we employ a text recognizer pretrained in Chinese and English [32, 5] to obtain the initial transcripts, followed by a manual refinement. With precise text labeling, STISR models can be evaluated from the aspect of text recognition

¹<https://apps.apple.com/us/app/camerapixels-pro/id1148178499>



Figure 3. Typical scenes in our collected Real-CE dataset.

	Region		Text line	
SR factor	4×	2×	4×	2×
train	645	1, 290	7, 849	15, 698
test	261	522	3, 414	6, 828
Max resolution	4, 032 × 3, 024		1, 156 × 2, 883	
Min resolution	228 × 396		16 × 22	
Chinese	-		8, 222	16, 444
English	-		3, 041	6, 082

Table 1. Statistics of the constructed Real-CE dataset.

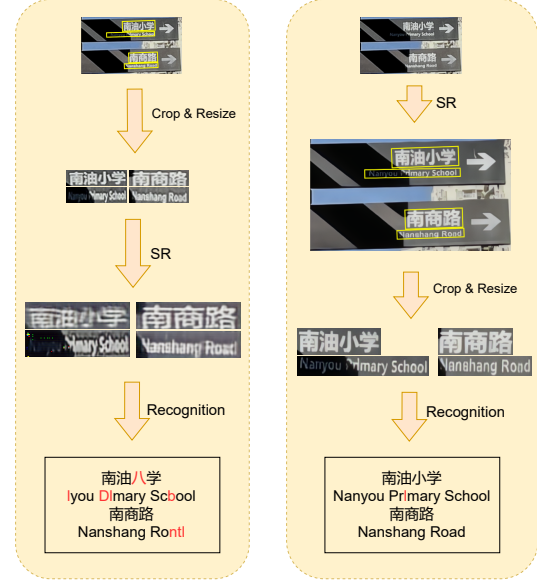
on the test set of our Real-CE benchmark.

3.2. Dataset Statistics

Our dataset contains both 2× and 4× zooming modes for training and testing. The detailed statistics of our dataset are shown in Table 1. Our Real-CE dataset contains 33, 789 text line pairs. In particular, 24, 666 of them are Chinese texts while the rest are English texts.

Text region pairs. Our dataset contains 2, 718 text region pairs, 1, 935 of which are training pairs and the rest are testing pairs. Among the testing pairs, there are 261 pairs for 4× (13mm to 52mm) zooming and another 522 pairs for 2× (26mm to 52mm and 13mm to 26mm) zooming evaluation. All the cropped HR text regions are ranged from size of 228 × 396 to 4, 032 × 3, 024. Each text region contains one or more text lines.

Text lines. The text semantics and language are distinguished with text lines. Text boxes and recognition are also annotated by lines. There are 23, 547 text lines for training, 3, 414 text lines for 4× zooming evaluation, and 6, 828 text lines for 2× zooming evaluation. The size of text lines ranges from 16 × 22 to 1, 156 × 2, 883. The category of the characters in Real-CE is 3, 755 in total.



(a) TextZoom Evaluation

(b) Real-CE Evaluation

Figure 4. Comparison of the TextZoom evaluation and Real-CE evaluation protocols. Wrong recognition results are in red. Please zoom in for more details.

3.3. Evaluation Protocol

To evaluate the performance of STISR models on Real-CE, we employ 5 metrics, including structural similarity index measure (SSIM) [39], peak signal-to-noise ratio (PSNR), learned perceptual image patch similarity (LPIPS) [42], normalized edit distance (NED) and word accuracy (ACC). Among them, PSNR, SSIM and LPIPS measure the errors between the reconstructed HR images and the ground truths. In particular, PSNR and SSIM are evaluated in image space while LPIPS is evaluated in feature space. ACC and NED employ text recognition models to evaluate the recognition accuracy of the reconstructed HR images. Here we adopt pre-trained CRNN [32, 5] as the text recognition model for evaluation. Particularly, ACC computes the word-level accuracy of the predicted sequence. NED between the predicted text sequence P and the ground truth text image label G are computed as follows:

$$NED(P, G) = 1 - \frac{ED(P, G)}{\max(|P|, |G|)}, \quad (1)$$

where $ED(\cdot)$ stands for the edit distance calculation, $|P|$ and $|G|$ refer to the length of the prediction and the ground-truth label. Therefore, the predicted sequence is more accurate and closer to the ground-truth label when the NED is larger. When we measure long texts, the ACC index may not fully reflect the recognition correctness at character level, while NED can measure it in a finer-grained manner.

In the testing process, the trained STISR models are performed on the original LR text region image to obtain the

reconstructed HR images first. Then the text lines in reconstructed HR images are cropped and kept in their original ratio for recognition evaluation in terms of ACC and NED. The evaluation process is illustrated in Figure 4 (b). Compared with the evaluation protocol of TextZoom [35] (see Figure 4 (a)), which trains and evaluates text lines with fixed sizes and shapes, our protocol can avoid the text deformation brought by the resizing operation. As shown in Figure 4 (a), such an arrangement is unfriendly for Chinese long text in Real-CE (often presented as sentences), resulting in low reconstruction quality and recognition accuracy.

4. Text Edge-aware STISR

Different from English characters, Chinese characters are composed of more basic radical-level parts (one can refer to [5] for more details) and have more complicated internal structures. Therefore, elaborated designs are needed to enhance the model capacity for Chinese text reconstruction. In this section, we propose an edge-aware learning method, which uses the text edge map as input and an edge-aware loss for supervision.

4.1. Text Edge Map

The text information in an image is inevitably blended with complex background. This will weaken the saliency of the text structures and somehow impairs the text reconstruction process. Text edge information is helpful to tackle this problem because it can effectively guide an STISR model to be better aware of the text structures and strokes.

We adopt the Canny edge detector [3] to compute a text edge map, denoted as \mathcal{C} , in the training process. The text edge map assigns value 1 to the text contour area and 0 to the background. Thus, the text edge map contains text structures and excludes the background information. From Fig. 5, one can see that the character shape and structure may be unclear in the LR-HR image pairs, while in their Canny edge maps, the text shapes and structures are enhanced. We compute edge maps for both LR and HR images in the dataset. The LR edge map \mathcal{C}_{LR} is concatenated with the LR image in channel dimension as the network input, which is shown in Fig. 6. With this extra input, the STISR model can learn a stronger feature representation of the finer-grained text structure.

4.2. Edge-aware Loss

We propose an edge-aware loss based on the computed edge map. First, the STISR model is modified to output both the reconstructed HR text image $\hat{\mathcal{I}}_H$ and an estimated HR text edge map $\hat{\mathcal{C}}_H$. This estimated text edge map is used in the training stage to gain extra supervision, but is discarded in the testing stage. The EA loss is computed between the estimated text edge map and the ground truth edge map at pixel level and feature level.

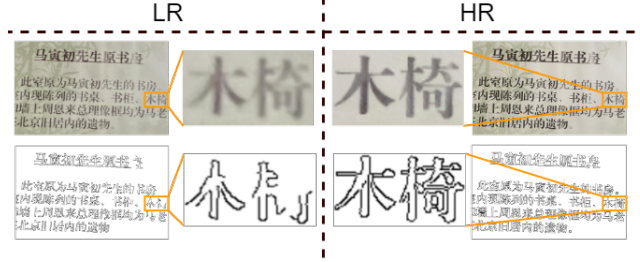


Figure 5. LR-HR RGB images (top) and their Canny edge maps (bottom). Foreground edges are drawn in black, and background in white for better visualization.

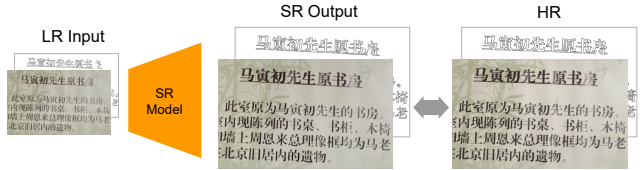


Figure 6. Illustration of the edge-aware STISR model learning. The edge map of the LR image is extracted and input to the network, and the edge map of the HR image is used to supervise the network training.

At pixel level, we adopt the \mathcal{L}_1 loss in image domain between the estimated HR edge map $\hat{\mathcal{C}}_H$ and the ground truth HR text map \mathcal{C}_H . Therefore, the EA loss at the pixel level \mathcal{L}_{EA}^P is calculated as:

$$\mathcal{L}_{EA}^P = |\mathcal{C}_H - \hat{\mathcal{C}}_H|. \quad (2)$$

Besides the pixel-level supervision, we compute the feature level EA loss \mathcal{L}_{EA}^F as follows:

$$\mathcal{L}_{EA}^F = |\mathcal{F}(\hat{\mathcal{I}}_H) \cdot \mathcal{F}(\hat{\mathcal{C}}_H) - \mathcal{F}(\mathcal{I}_H) \cdot \mathcal{F}(\mathcal{C}_H)|, \quad (3)$$

where \mathcal{F} denotes a pretrained feature extractor network (VGG19 [34] is used in this paper). $\mathcal{F}(\hat{\mathcal{I}}_H)$ and $\mathcal{F}(\mathcal{I}_H)$ denote the feature representation of the estimated and the ground truth HR text images, respectively. $\mathcal{F}(\hat{\mathcal{C}}_H)$ and $\mathcal{F}(\mathcal{C}_H)$ denote the feature representation of the estimated and the ground truth HR text edge maps, respectively. The image features are weighted by the edge features via element-wise multiplication (e.g., $\mathcal{F}(\hat{\mathcal{I}}_H) \cdot \mathcal{F}(\hat{\mathcal{C}}_H)$) to strengthen the structural areas. Finally, an \mathcal{L}_1 loss is imposed on the strengthened features between the estimated and ground truth ones. One can view the **supplementary file** for more detailed analysis.

Finally, together with the \mathcal{L}_1 loss on RGB images and EA loss terms, the overall loss function \mathcal{L} is formulated as:

$$\mathcal{L} = \mathcal{L}_1 + \alpha \mathcal{L}_{EA}^P + \beta \mathcal{L}_{EA}^F, \quad (4)$$

where α and β are balancing parameters.

		SR factor	4×					2×				
	Approach	train set	PSNR ↑	SSIM ↑	LPIPS ↓	ACC ↑	NED ↑	PSNR ↑	SSIM ↑	LPIPS ↓	ACC ↑	NED ↑
SISR Methods	Bicubic		19.65	0.6684	0.3987	0.2759	0.6173	20.82	0.7106	0.2100	0.3475	0.6982
	SRRes [20]	TZ [35]	19.72	0.6808	0.3872	0.2201	0.5992	20.28	0.6762	0.3467	0.2742	0.6401
		RS [2]	18.60	0.6576	0.3736	0.2642	0.6087	19.10	0.6872	0.3244	0.2977	0.6671
		RC	20.22	0.7224	0.2665	0.2879	0.6361	20.72	0.7360	0.2116	0.3499	0.6996
	RRDB [38]	TZ [35]	18.95	0.6575	0.4495	0.1463	0.3776	19.43	0.6899	0.3887	0.1962	0.4665
		RS [2]	19.59	0.6703	0.2765	0.2590	0.6201	20.34	0.7312	0.2267	0.3307	0.6772
		RC	20.23	0.7231	0.2626	0.2920	0.6421	21.10	0.7535	0.2065	0.3494	0.7003
	EDSR [22]	TZ [35]	18.88	0.6512	0.4860	0.0799	0.3414	19.32	0.6904	0.4012	0.1369	0.4077
		RS [2]	19.59	0.6728	0.3295	0.2702	0.6231	20.02	0.7291	0.2837	0.3189	0.6708
		RC	20.16	0.7195	0.2883	0.2882	0.6330	20.74	0.7448	0.2258	0.3468	0.6954
	RCAN [45]	TZ [35]	18.97	0.6277	0.4816	0.0810	0.3424	19.48	0.6488	0.4075	0.1507	0.5420
		RS [2]	19.55	0.6661	0.3475	0.2450	0.5989	20.05	0.7044	0.2806	0.2968	0.6776
		RC	20.33	0.7232	0.2878	0.2879	0.6321	20.98	0.7435	0.2173	0.3484	0.7006
	ELAN [44]	TZ [35]	19.21	0.6459	0.3796	0.1778	0.4764	20.10	0.6653	0.3241	0.2254	0.5467
		RS [2]	19.60	0.6660	0.3348	0.2674	0.6228	20.48	0.6907	0.2732	0.3104	0.6642
		RC	20.39	0.7299	0.2892	0.2953	0.6404	21.16	0.7480	0.2201	0.3508	0.6992
STISR Methods	TSRN [35]	TZ [35]	17.47	0.4853	0.1990	0.1796	0.3874	18.73	0.5676	0.1855	0.2471	0.4622
		RS [2]	17.83	0.4899	0.2154	0.1733	0.3759	19.06	0.5322	0.1892	0.2675	0.4526
		RC	18.11	0.4850	0.1981	0.2316	0.4159	18.99	0.5233	0.1677	0.2854	0.4809
	TPGSR [25]	TZ [35]	17.37	0.4913	0.1896	0.2076	0.3842	17.99	0.5312	0.1686	0.2655	0.4423
		RS [2]	17.65	0.4772	0.1947	0.2203	0.3930	18.56	0.5462	0.1754	0.2952	0.4658
		RC	18.07	0.4758	0.1843	0.2326	0.4123	18.83	0.5562	0.1661	0.3007	0.4913
	TBSRN [4]	TZ [35]	17.59	0.4919	0.1767	0.2246	0.4133	18.41	0.5456	0.1588	0.2905	0.4896
		RS [2]	17.69	0.4762	0.1849	0.2235	0.4021	18.69	0.5309	0.1666	0.2895	0.4644
		RC	18.33	0.4826	0.1715	0.2527	0.4444	19.01	0.5366	0.1652	0.3181	0.5294
	TATT [26]	TZ [35]	17.43	0.5010	0.2003	0.2100	0.3926	18.24	0.5667	0.1827	0.2755	0.4993
		RS [2]	17.66	0.4993	0.2256	0.2092	0.3916	18.47	0.5253	0.1930	0.2749	0.4702
		RC	17.96	0.4904	0.1804	0.2330	0.4342	19.06	0.5772	0.1590	0.3127	0.5240
HR		-	-	-	0.4807	0.8342	-	-	-	0.4514	0.8038	

Table 2. Experimental results on Real-CE test set with SISR and STISR models trained on different training sets. TZ, RS and RC refer to TextZoom [35], RealSR [2] and Real-CE datasets, respectively. It should be noted that the evaluated metric scores of SISR and STISR methods are very different because SISR models intake global images as input, while STISR models intake text lines as input.

5. Experimental Results

In this section, we first validate the effectiveness of our established Real-CE dataset by comparing STISR models trained on it and other text image datasets, and then validate the proposed EA loss in improving STISR model performance. All models are all trained with the Adam optimizer. When trained on our Real-CE training set, the number of epochs is set to 400. The learning rate is set to 2×10^{-4} . In the calculation of \mathcal{L}_{EA} , we adopt the *Conv5.4* features of pre-trained VGG19 [34]. The balancing parameters α and β in Eq. (4) are set to 1 and 5×10^{-4} , respectively (one can refer to **supplementary file** for the parameter selection details). When computing recognition-based metrics, we first crop the text lines from the global text image and then rescale the cropped SR text line image to fit the recognizer.

5.1. Effectiveness of Real-CE Dataset

In this section, we perform experiments to validate the advantages of the proposed Real-CE dataset over existing real-world SR datasets, including TextZoom [35] and RealSR [2]. TextZoom is built for real-world English text super-resolution, which lacks dense character structures in the dataset. RealSR is built for real-world natural image super-resolution. We evaluate five state-of-the-art SISR models and four state-of-the-art STISR models on the three datasets. The five SISR models are SRRes [20],

RRDB [38], EDSR [22], RCAN [45] and ELAN [44], where the first four are CNN-based models and the last one is a transformer-based model. The four STISR models are TSRN [35], TPGSR [25], TBSRN [4] and TATT [26], where the first two are CNN-based models, and the rest are transformer-based models.

All the STISR and SISR models are trained on Real-CE, TextZoom and RealSR, respectively, and tested on the testing set of Real-CE. Since SISR models generally support arbitrary input sizes, the original test images are set as the input, and the PSNR, SSIM and LPIPS metrics are computed on the original sizes. Note that this is the default evaluation protocol of our benchmark, as described in Section 3.3. However, most of the STISR models [35, 25, 4, 26] only support inputs with fixed sizes. Thus, we first crop and reshape the test images to a fixed size as the network input, and then the network outputs are compared with the resized ground truth images to compute PSNR, SSIM and LPIPS.

The quantitative results of compared SISR and STISR models are shown in Table 2. One can see that, the models trained on TextZoom obtain inferior performance in terms of image-based metrics and recognition-based metrics. This is because the training data in TextZoom lacks complex character structures, and hence the trained models cannot handle the complex Chinese texts in Real-CE test set. Moreover, by using TextZoom, the SR models can only be trained with data of fixed sizes, which are hard to be generalized to



Figure 7. STISR results of different models trained on different training datasets. Note that SISR models (EDSR [37], RRDB [38] and ELAN [44]) intake global images as input and follow the Real-CE inference protocol, while STISR models like TSRN [35], TBSRN [4] and TATT [26] can only take text line as input.

	SR factor				4×					2×				
	Approach	\mathcal{L}_1	\mathcal{L}_{EA}^P	\mathcal{L}_{EA}^F	PSNR ↑	SSIM ↑	LPIPS ↓	ACC ↑	NED ↑	PSNR ↑	SSIM ↑	LPIPS ↓	ACC ↑	NED ↑
SISR Methods	Bicubic	-	-	-	19.65	0.6684	0.3987	0.2759	0.6173	20.82	0.7106	0.2100	0.3475	0.6982
	SRRes [20]	✓	×	×	20.22	0.7224	0.2665	0.2879	0.6361	20.72	0.7360	0.2116	0.3499	0.6996
		✓	✓	×	20.30	0.7219	0.2722	0.2909	0.6373	21.23	0.7551	0.2022	0.3496	0.7013
		✓	✓	✓	20.18	0.7102	0.2041	0.2917	0.6454	21.09	0.7489	0.1875	0.3540	0.7080
	RRDB [38]	✓	×	×	20.23	0.7231	0.2626	0.2920	0.6421	21.10	0.7535	0.2065	0.3494	0.7003
		✓	✓	×	20.42	0.7303	0.2630	0.2914	0.6399	21.21	0.7559	0.2010	0.3575	0.7063
		✓	✓	✓	20.14	0.7210	0.2031	0.3093	0.6622	21.00	0.7517	0.1852	0.3549	0.7130
	ELAN [44]	✓	×	×	20.39	0.7299	0.2892	0.2953	0.6404	21.16	0.7480	0.2201	0.3508	0.6992
		✓	✓	×	20.47	0.7330	0.2767	0.2982	0.6441	21.29	0.7557	0.1989	0.3549	0.7047
		✓	✓	✓	20.21	0.7245	0.2071	0.3061	0.6567	21.10	0.7479	0.1835	0.3524	0.7073
STISR Methods	TBSRN [4]	✓	×	×	18.33	0.4826	0.1715	0.2527	0.4444	19.01	0.5366	0.1652	0.3181	0.5294
		✓	✓	×	18.46	0.4881	0.1699	0.2555	0.4483	19.10	0.5378	0.1662	0.3185	0.5364
		✓	✓	✓	18.20	0.4796	0.1431	0.2615	0.4531	18.97	0.5334	0.1307	0.3289	0.5450
	TATT [26]	✓	×	×	17.96	0.4904	0.1804	0.2330	0.4342	19.06	0.5772	0.1590	0.3127	0.5240
		✓	✓	×	18.12	0.4916	0.1786	0.2324	0.4306	19.17	0.5825	0.1604	0.3166	0.5360
		✓	✓	✓	17.89	0.4822	0.1546	0.2417	0.4549	19.02	0.5723	0.1475	0.3239	0.5491
	HR	-	-	-	-	-	-	0.4807	0.8342	-	-	-	0.4514	0.8038

Table 3. Comparison of SISR and STISR models trained on Real-CE with different losses.

other text image sizes. The models trained on RealSR [2] also obtain inferior results, since RealSR is basically established for natural image SISR. In contrast, the STISR and SISR models trained on our Real-CE dataset show much better text recovery performance on all evaluation metrics. In addition, it should be noted that the evaluated metric scores of SISR and STISR methods are very different because SISR models intake global images as input, while STISR models intake text line as input.

Figure 7 visualizes the SR results of some representative SISR and STISR models trained on the three datasets. For

convenience, we input different images to different models for more comprehensive evaluation. One can see that the text recovery results by models trained on TextZoom are blurry and contain visual artifacts. This is because TextZoom lacks training data with complex Chinese character structures. Besides, TextZoom only supports fixed-size training and the trained model cannot generalize to other sizes of testing data. The results of models trained on RealSR have fewer artifacts but are blurry with mixed strokes. In contrast, the results of models trained on Real-CE dataset have clear edges and are highly readable on both Chinese

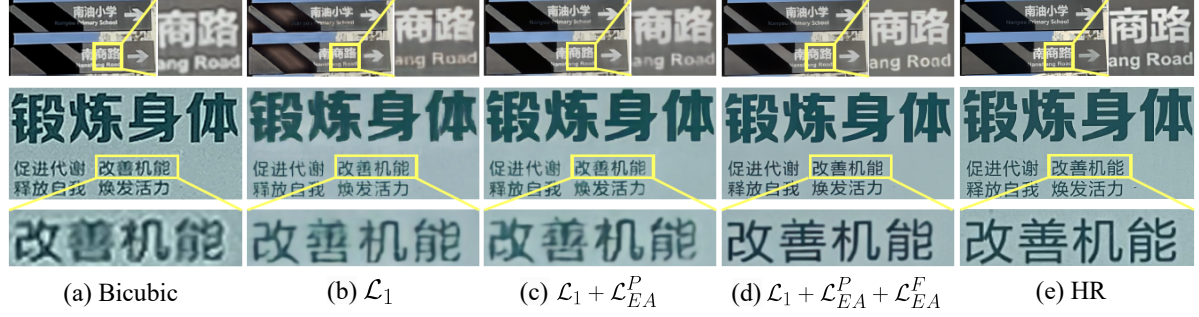


Figure 8. STISR results of RRDB models trained with different losses.

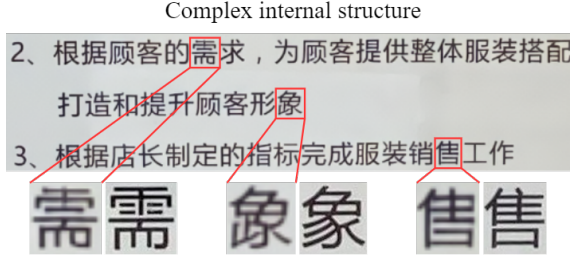


Figure 9. Examples of failure cases by our method.

and English characters. More visual results can be found in the **supplementary material**.

For comparisons on the STISR models trained on synthetic LR-HR data and our Real-CE data, please also refer to the **supplementary material**.

5.2. Effectiveness of the EA Loss

We then validate the effectiveness of our proposed EA losses by testing SISR and STISR models with different combinations of \mathcal{L}_1 , \mathcal{L}_{EA}^P and \mathcal{L}_{EA}^F losses. Here we employ three SISR models, including SRRes [20], RRDB [38] and ELAN [44], and two STISR models, including TBSRN [4] and TATT [26], in the experiments. The evaluation metrics are the same as that in Section 5.1.

Quantitative evaluation results of the losses are shown in Table 3. One can see that compared with the models trained with the \mathcal{L}_1 loss only, models trained with \mathcal{L}_1 and \mathcal{L}_{EA}^P demonstrate enhanced PSNR/SSIM scores. This is because \mathcal{L}_{EA}^P provides pixel-wise supervision on edge areas, resulting in improved pixel-wise metrics. However, the improvements on perceptual metrics (*i.e.*, LPIPS) and recognition metrics are still limited. By further adding \mathcal{L}_{EA}^F loss into training, all models demonstrate notable improvement on LPIPS and recognition accuracy, especially on $4\times$ results. This indicates that the character structural information is important for text legibility. Since the character structures can be well enhanced by using \mathcal{L}_{EA}^F in training, the text recognition is significantly improved.

By using the RRDB model, we visualize the STISR results by different losses in Figure 8. One can see that RRDB trained with only the \mathcal{L}_1 loss shows limited improve-

ment compared with the bicubic interpolation. By including \mathcal{L}_{EA}^P loss in training, the reconstructed text images are much enhanced with clearer character edges, as shown in Figure 8(c). By further incorporating the \mathcal{L}_{EA}^F loss, a significant enhancement in terms of edge clarity and local contrast can be observed, which greatly improves the legibility of Chinese text contents, as shown in Figure 8(d). More visualization results can be found in the **supplementary file**.

5.3. Failure Cases

Our proposed method may fail when the character has very low resolution and intricate structures, as shown in Figure 9. Though the output still has clear edges, some of the tiny strokes are wrong. This is because the tiny strokes are very obscure in the low-resolution input text image. In such cases, semantic information can be incorporated to assist the text restoration, which will be our future work.

6. Conclusions

In this paper, we established a Chinese-English benchmark, namely Real-CE, for scene text image super-resolution (STISR) model training. It contained 1,935 training and 783 testing images. The text region pairs contained 33,789 text lines, among which 24,666 were Chinese texts with complex structures. We further proposed an edge-aware (EA) learning method for the restoration of Chinese texts, which computed a text edge map from the given input image and employed an EA loss to guide the STISR model learning process. Experimental results demonstrated that the models trained on our Real-CE dataset can recover clearer and more readable Chinese texts than other STISR datasets, and the EA learning scheme can effectively improve text image quality. The Real-CE dataset provided a valuable benchmark for researcher to investigate the challenging Chinese text image recovery problems.

7. Acknowledgements

This work is supported by the Hong Kong RGC RIF grant (R5001-18). We thank Dr. Xindong Zhang for his help on the project.

References

- [1] Michael S Brown and W Brent Seales. Image restoration of arbitrarily warped documents. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(10):1295–1306, 2004. [1](#)
- [2] Jianrui Cai, Hui Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. In *Int. Conf. Comput. Vis.*, pages 3086–3095, 2019. [2](#), [3](#), [6](#), [7](#)
- [3] John Canny. A computational approach to edge detection. *IEEE Trans. Pattern Anal. Mach. Intell.*, (6):679–698, 1986. [5](#)
- [4] Jingye Chen, Bin Li, and Xiangyang Xue. Scene text telescope: Text-focused scene image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 12026–12035, 2021. [1](#), [2](#), [6](#), [7](#), [8](#)
- [5] Jingye Chen, Haiyang Yu, Jianqi Ma, Mengnan Guan, Xixi Xu, Xiaocong Wang, Shaobo Qu, Bin Li, and Xiangyang Xue. Benchmarking chinese text recognition: Datasets, baselines, and an empirical study. *arXiv preprint arXiv:2112.15093*, 2021. [3](#), [4](#), [5](#)
- [6] Jingye Chen, Haiyang Yu, Jianqi Ma, Bin Li, and Xiangyang Xue. Text gestalt: Stroke-aware scene text image super-resolution. In *AAAI*, volume 36, pages 285–293, 2022. [1](#), [2](#)
- [7] Kai Chen, Li Tian, Haisong Ding, Meng Cai, Lei Sun, Sen Liang, and Qiang Huo. A compact cnn-dblstm based character model for online handwritten chinese text recognition. In *Int. Conf. Doc. Anal. Recog.*, 2017. [3](#)
- [8] Xiaogang Chen, Xiangjian He, Jie Yang, and Qiang Wu. An effective document image deblurring algorithm. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 369–376. IEEE, 2011. [1](#)
- [9] Zhazhan Cheng, Fan Bai, Yunlu Xu, Gang Zheng, Shiliang Pu, and Shuigeng Zhou. Focusing attention: Towards accurate text recognition in natural images. In *Int. Conf. Comput. Vis.*, pages 5076–5084, 2017. [3](#)
- [10] Zhazhan Cheng, Yangliu Xu, Fan Bai, Yi Niu, Shiliang Pu, and Shuigeng Zhou. AON: Towards arbitrarily-oriented text recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5571–5579, 2018. [3](#)
- [11] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(2):295–307, 2015. [2](#)
- [12] Chao Dong, Ximei Zhu, Yubin Deng, Chen Change Loy, and Yu Qiao. Boosting optical character recognition: A super-resolution approach. *arXiv preprint arXiv:1506.02211*, 2015. [2](#)
- [13] Jun Du, Zi-Rui Wang, Jian-Fang Zhai, and Jin-Shui Hu. Deep neural network based hidden markov model for offline handwritten chinese text recognition. In *ICPR*, 2016. [3](#)
- [14] Bahadır K. Gunturk, Yucel Altunbasak, and Russell M. Mersereau. Super-resolution reconstruction of compressed video using transform-domain statistics. *IEEE Trans. Image Process.*, 13(1):33–43, 2004. [2](#)
- [15] Pan He, Weilin Huang, Yu Qiao, Chen Change Loy, and Xiaoou Tang. Reading scene text in deep convolutional sequences. In *AAAI*, 2016. [3](#)
- [16] Jie Hu, Tszhang Guo, Ji Cao, and Changshui Zhang. End-to-end chinese text recognition. In *GlobalSIP*, 2017. [3](#)
- [17] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Reading text in the wild with convolutional neural networks. *Int. J. Comput. Vis.*, 116(1):1–20, 2016. [3](#)
- [18] Max Jaderberg, Andrea Vedaldi, and Andrew Zisserman. Deep features for text spotting. In *Eur. Conf. Comput. Vis.*, pages 512–528. Springer, 2014. [3](#)
- [19] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep Laplacian pyramid networks for fast and accurate super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 624–632, 2017. [2](#)
- [20] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4681–4690, 2017. [2](#), [6](#), [7](#), [8](#)
- [21] Jingyun Liang, Jie Zhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Int. Conf. Comput. Vis.*, pages 1833–1844, 2021. [2](#)
- [22] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 136–144, 2017. [2](#), [6](#)
- [23] Wei Liu, Chaofeng Chen, Kwan-Yee K Wong, Zhizhong Su, and Junyu Han. STAR-Net: A spatial attention residue network for scene text recognition. In *Brit. Mach. Vis. Conf.*, volume 2, page 7, 2016. [3](#)
- [24] Jianqi Ma. Rrpn++: Guidance towards more accurate scene text detection. *arXiv preprint arXiv:2009.13118*, 2020. [3](#)
- [25] Jianqi Ma, Shi Guo, and Lei Zhang. Text prior guided scene text image super-resolution. *IEEE Trans. Image Process.*, 32:1341–1353, 2023. [1](#), [6](#)
- [26] Jianqi Ma, Zhetong Liang, and Lei Zhang. A text attention network for spatial deformation robust scene text image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5911–5920, 2022. [1](#), [2](#), [6](#), [7](#), [8](#)
- [27] Jianqi Ma, Weiyuan Shao, Hao Ye, Li Wang, Hong Wang, Yingbin Zheng, and Xiangyang Xue. Arbitrary-oriented scene text detection via rotation proposals. *IEEE Trans. Multimedia*, 20(11):3111–3122, 2018. [3](#)
- [28] Julien Mairal, Francis R. Bach, Jean Ponce, Guillermo Sapiro, and Andrew Zisserman. Non-local sparse models for image restoration. In *Int. Conf. Comput. Vis.*, pages 2272–2279. IEEE Computer Society, 2009. [2](#)
- [29] Yongqiang Mou, Lei Tan, Hui Yang, Jingying Chen, Leyuan Liu, Rui Yan, and Yaohong Huang. PlugNet: Degradation aware scene text recognition supervised by a pluggable super-resolution unit. In *Eur. Conf. Comput. Vis.*, 2020. [2](#)
- [30] Clément Peyrard, Moez Baccouche, Franck Mamalet, and Christophe Garcia. ICDAR2015 competition on text image super-resolution. In *Int. Conf. Doc. Anal. Recog.*, pages 1201–1205. IEEE, 2015. [1](#), [2](#)

- [31] Yuhui Quan, Jietao Yang, Yixin Chen, Yong Xu, and Hui Ji. Collaborative deep learning for super-resolving blurry text images. *IEEE Trans. Comput. Imaging*, 6:778–790, 2020. 1, 2
- [32] 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. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(11):2298–2304, 2016. 3, 4
- [33] Baoguang Shi, Mingkun Yang, Xinggang Wang, Pengyuan Lyu, Cong Yao, and Xiang Bai. ASTER: An attentional scene text recognizer with flexible rectification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(9):2035–2048, 2018. 3
- [34] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2014. 5, 6
- [35] Wenjia Wang, Enze Xie, Xuebo Liu, Wenhai Wang, Ding Liang, Chunhua Shen, and Xiang Bai. Scene text image super-resolution in the wild. In *Eur. Conf. Comput. Vis.*, 2020. 1, 2, 5, 6, 7
- [36] Wenjia Wang, Enze Xie, Peize Sun, Wenhai Wang, Lixun Tian, Chunhua Shen, and Ping Luo. Textsr: Content-aware text super-resolution guided by recognition. *arXiv preprint arXiv:1909.07113*, 2019. 1, 2
- [37] Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 606–615, 2018. 2, 7
- [38] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *Eur. Conf. Comput. Vis. Worksh.*, pages 0–0, 2018. 2, 6, 7, 8
- [39] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Trans. Image Process.*, 13(4):600–612, 2004. 4
- [40] Xiangyu Xu, Deqing Sun, Jinshan Pan, Yujin Zhang, Hanspeter Pfister, and Ming-Hsuan Yang. Learning to super-resolve blurry face and text images. In *Int. Conf. Comput. Vis.*, pages 251–260, 2017. 1, 2
- [41] Jianchao Yang, John Wright, Thomas S. Huang, and Yi Ma. Image super-resolution via sparse representation. *IEEE Trans. Image Process.*, 19(11):2861–2873, 2010. 2
- [42] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 586–595, 2018. 4
- [43] Xuaner Zhang, Qifeng Chen, Ren Ng, and Vladlen Koltun. Zoom to learn, learn to zoom. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3762–3770, 2019. 2
- [44] Xindong Zhang, Hui Zeng, Shi Guo, and Lei Zhang. Efficient long-range attention network for image super-resolution. *Eur. Conf. Comput. Vis.*, 2022. 2, 6, 7, 8
- [45] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *Eur. Conf. Comput. Vis.*, pages 286–301, 2018. 6
- [46] Cairong Zhao, Shuyang Feng, Brian Nlong Zhao, Zhijun Ding, Jun Wu, Fumin Shen, and Heng Tao Shen. Scene text image super-resolution via parallelly contextual attention network. In *ACM Int. Conf. Multimedia*, pages 2908–2917, 2021. 1, 2
- [47] Minyi Zhao, Miao Wang, Fan Bai, Bingjia Li, Jie Wang, and Shuigeng Zhou. C3-stir: Scene text image super-resolution with triple clues. *IJCAI*, 2022. 1