

# Robust Reference-based Super-Resolution via $C^2$ -Matching

Yuming Jiang<sup>1</sup> Kelvin C.K. Chan<sup>1</sup> Xintao Wang<sup>2</sup> Chen Change Loy<sup>1</sup> Ziwei Liu<sup>1✉</sup>  
<sup>1</sup>S-Lab, Nanyang Technological University <sup>2</sup>Applied Research Center, Tencent PCG  
 {yuming002, chan0899, ccloy, ziwei.liu}@ntu.edu.sg xintao.wang@outlook.com

## Abstract

Reference-based Super-Resolution (Ref-SR) has recently emerged as a promising paradigm to enhance a low-resolution (LR) input image by introducing an additional high-resolution (HR) reference image. Existing Ref-SR methods mostly rely on **implicit correspondence matching** to borrow HR textures from reference images to compensate for the information loss in input images. However, performing local transfer is difficult because of two gaps between input and reference images: the transformation gap (e.g. scale and rotation) and the resolution gap (e.g. HR and LR). To tackle these challenges, we propose  $C^2$ -Matching in this work, which produces **explicit robust matching cross-invariant transformation and resolution**. 1) For the transformation gap, we propose a contrastive correspondence network, which learns transformation-robust correspondences using augmented views of the input image. 2) For the resolution gap, we adopt a teacher-student correlation distillation, which distills knowledge from the easier HR-HR matching to guide the more ambiguous LR-HR matching. 3) Finally, we design a dynamic aggregation module to address the potential misalignment issue. In addition, to faithfully evaluate the performance of Ref-SR under a realistic setting, we contribute the Webly-Referenced SR (WR-SR) dataset, mimicking the practical usage scenario. Extensive experiments demonstrate that our proposed  $C^2$ -Matching significantly outperforms state of the arts by over 1dB on the standard CUFED5 benchmark. Notably, it also shows great generalizability on WR-SR dataset as well as robustness across large scale and rotation transformations<sup>1</sup>.

## 1. Introduction

Reference-based Super-Resolution (Ref-SR) [40, 39, 34, 26] has attracted substantial attention in recent years. Compared to Single-Image-Super-Resolution (SISR) [6, 13, 14, 17, 25, 4], where the only input is a single low-resolution (LR) image, Ref-SR super-resolves the LR image with the guidance of an additional high-resolution (HR) reference

<sup>1</sup>Codes and datasets are available at <https://github.com/yumingj/C2-Matching>.

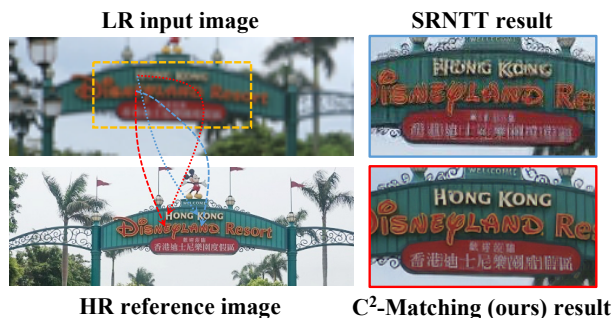


Figure 1. **Cross transformation and Cross resolution matching are performed in our  $C^2$ -Matching.** Our proposed  $C^2$ -Matching successfully transfers the HR details of the reference image by finding more accurate correspondences. The correspondences found by our method are marked in red and the correspondences found by SRNTT [39] are marked in blue.

image. Textures of the HR reference image are transferred to provide more fine details for the LR image.

The key step in texture transfer for Ref-SR is to find correspondences between the input image and the reference image. Existing methods [39, 34, 32] perform correspondence matching implicitly. Their correspondences are computed based on the content and appearance similarities, which are then embedded into the main framework. However, it is a difficult task to accurately compute the correspondences under real-world variations due to two major challenges: **1)** the underlying transformation gap between input images and reference images; **2)** the resolution gap between input images and reference images. In Ref-SR, same objects or similar texture patterns are often present in both input images and reference images, but their appearances vary due to scale and rotation transformations. In this case, correspondences computed purely by appearance are inaccurate, leading to an unsatisfactory texture transfer. For the resolution gap, due to the imbalance in the amount of information contained in an LR input image and an HR reference image, the latter is often downsampled (to an LR image) to match the former (in resolution). The downsampling operation inevitably results in information loss, hampering the search for accurate correspondences, especially for the fine-texture regions.

To address the aforementioned challenges, we propose  $C^2$ -matching for Robust Reference-based Super-Resolution, where *Cross transformation* and *Cross resolution* matching are explicitly performed. To handle the transformation gap, a contrastive correspondence network is proposed to learn transformation-robust correspondences between input images and reference images. Specifically, we employ an additional triplet margin loss to minimize the distance of point-wise features before and after transformations while maximizing the distance of irrelevant features. Thus, the extracted feature descriptors are more robust to scale and rotation transformations, and can be used to compute more accurate correspondences.

As for the resolution gap, inspired by knowledge distillation, we propose a teacher-student correlation distillation. We train the teacher contrastive correspondence network for HR-HR matching. Since the teacher network takes two HR images as input, it is better at matching the regions with complicated textures. Thus, the knowledge of the teacher model can be distilled to guide the more ambiguous LR-HR matching. The teacher-student correlation distillation enables the contrastive correspondence network to compute correspondences more accurately for texture regions.

After obtaining correspondences, we then fuse the information of reference images through a dynamic aggregation module to transfer the HR textures. With  $C^2$ -Matching, we achieve over 1dB improvement on the standard CUFED5 dataset. As shown in Fig. 1, compared to SRNTT [39], our  $C^2$ -Matching finds more accurate correspondences (marked as red dotted lines) and thus has a superior restoration performance.

To facilitate the evaluation of Ref-SR tasks in a more realistic setting, we contribute a new dataset named Webly-Reference SR (WR-SR) dataset. In real-world applications, given an LR image, users may find its similar HR reference images through some web search engines. Motivated by this, for every input image in WR-SR, we search for its reference image through Google Image. The collected WR-SR can serve as a benchmark for real-world scenarios.

To summarize, our main contributions are: **1)** To mitigate the transformation gap, we propose the contrastive correspondence network to compute correspondences more robust to scale and rotation transformations. **2)** To bridge the resolution gap, a teacher-student correlation distillation is employed to further boost the performance of student LR-HR matching model with the guidance of HR-HR matching, especially for fine texture regions. **3)** We contribute a new benchmark dataset named Webly-Referenced SR (WR-SR) to encourage a more practical application in real scenarios.

## 2. Related Work

**Single Image Super-Resolution.** Single Image Super-Resolution (SISR) aims to recover the HR details of LR

images. The only input to the SISR task is the LR image. Dong *et al.* [6] introduced deep learning into SISR tasks by formulating the SISR task as an image-to-image translation problem. Later, SR networks had gone deeper with the help of residual blocks and attention mechanisms [25, 17, 14, 36, 37, 15, 4, 41]. However, the visual quality of the output SR images did not improve. The problem was the mean square error (MSE) loss function. In order to improve the perceptual quality, perceptual loss [12, 23], generative loss and adversarial loss [15] were introduced into the SR network [30, 35]. Knowledge distillation framework is also explored to improve the SR performance in [9, 16].

**Reference-based Image Super-Resolution.** Different from SISR, where no additional information is provided, the Reference-based Image Super-Resolution (Ref-SR) task [33, 38, 40] super-resolves input images by transferring HR details of reference images. Patch-Match method [1] was employed in [39, 34] to align input images and reference images. SRNTT [39] performed correspondence matching based on the off-the-shelf VGG features [27]. Recently, [34, 32] used learnable feature extractors, which were trained end-to-end accompanied with the main SR network. Even with the learnable feature extractors, the correspondences were computed purely based on contents and appearances. Recent work [26] introduced Deformable Convolution Network (DCN) [3, 42] to align input images and reference images. Inspired by previous work [29, 2], we also propose a similar module to handle the potential misalignment issue.

**Image Matching.** Scale Invariant Feature Transform (SIFT) [20] extracted local features to perform matching. With the advance of convolution neural networks (CNN), feature descriptors extracted by CNN were utilized to compute correspondences [8, 31, 7]. Recently, SuperPoint [5] was proposed to perform image matching in a self-supervised manner, and graph neural network was introduced to learn feature matching [24]. Needle-Match [19] performed image matching in a more challenging setting, where two images used for matching are degraded. Different from the aforementioned methods dealing with two images of the same degradation, in our task, we focus on cross resolution matching, *i.e.* matching between one LR image and one HR image.

## 3. Our Approach

The overview of our proposed  $C^2$ -Matching is shown in Fig. 2(a). The proposed  $C^2$ -Matching consists of two major parts: 1) Contrastive Correspondence Network and 2) Teacher-Student Correlation Distillation. The contrastive correspondence network learns transformation-robust correspondence matching; the teacher-student correlation distillation transfers HR-HR matching knowledge to LR-HR matching for a more accurate correspondence matching on

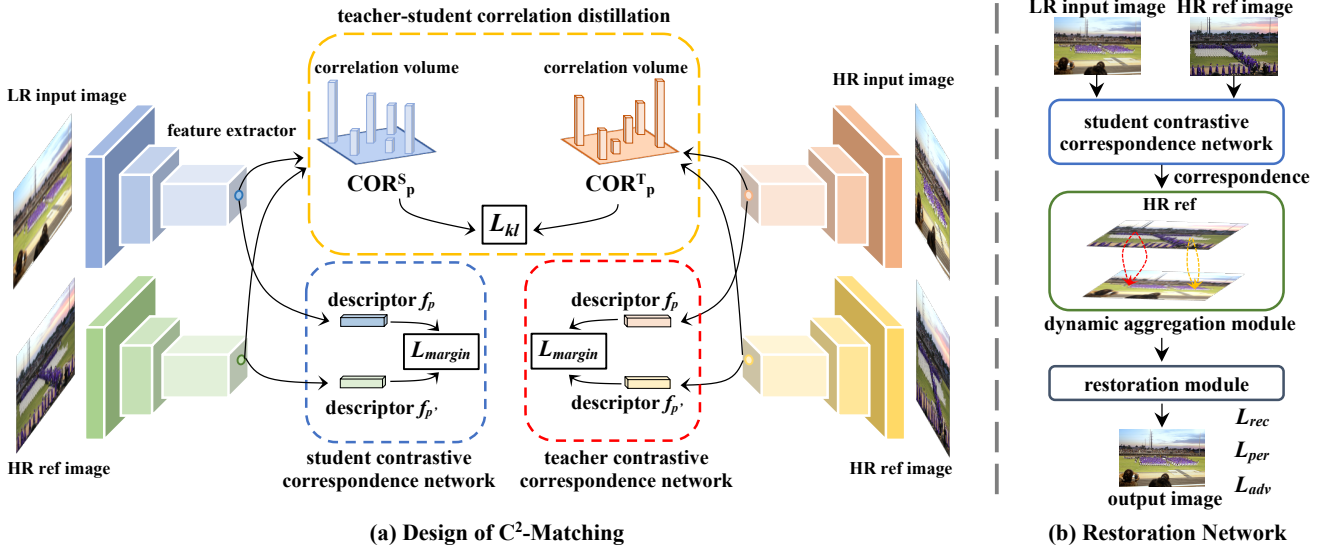


Figure 2. (a) **The overview of our proposed  $C^2$ -Matching.** The contrastive correspondence network is designed for transformation-robust correspondence matching. The student contrastive correspondence network takes both the LR input image and HR reference image (the transformed version of the HR input image serves as the HR reference image during training) as input. The descriptors before and after transformations are pushed closer while distances of the irrelevant descriptors are maximized by  $L_{margin}$ . To enable the student LR-HR contrastive correspondence network to perform correspondence matching better on highly textured regions, we embed a teacher-student correlation distillation process to distill the knowledge of the easier HR-HR teacher matching network to the student model by  $L_{kl}$ . (b) **The overall pipeline of restoration network.** The correspondences are first computed by the trained student contrastive correspondence network, after which the correspondences are used for subsequent dynamic aggregation module and restoration module.

texture regions. The correspondences obtained from the contrastive correspondence network are used for the subsequent restoration network (Fig. 2(b)). In the restoration network, to better aggregate the information of reference images, a dynamic aggregation module is proposed to handle the underlying misalignment problem. Then the aggregated features are used for the later restoration task.

### 3.1. Contrastive Correspondence Network

To transfer textures from reference images, correspondences should first be computed between input images and reference images, *i.e.* specify the location of similar regions in reference images for every region in input images.

Existing methods [39, 34, 32] computed correspondences according to the content and appearance similarities between degraded input images and reference images. For example, Zhang *et al.* [39] used VGG features [27] for the correspondence matching while [34, 32] used features trained end-to-end together with the downstream tasks. The drawback of this scheme is that it cannot handle the matching well if there are scale and rotation transformations between input images and reference images. An inaccurate correspondence would lead to an imperfect texture transfer for restoration. In this paper, we propose a learnable contrastive correspondence network to extract features that are robust to scale and rotation transformations.

In contrastive correspondence network, we deal with the correspondence matching between LR input images and HR reference images. Since the resolutions of the input image and the reference image are different, we adopt two networks with the same architecture but non-shared weights for feature extractions. The architecture of feature extractors will be explained in supplementary files.

For training, we synthesize HR reference images by applying homography transformations to original HR input images. By doing so, for every position  $p$  in the LR input image  $I$ , we can compute its ground-truth correspondence point  $p'$  in the transformed image  $I'$  according to the homography transformation matrix. We regard point  $p$  and its corresponding point  $p'$  as a positive pair. During optimization, we push the distances between feature representations of positive pairs closer, while maximizing the distances between other irrelevant but confusing negative pairs (defined as Eq. (3)). Similar to [8], we use the triplet margin ranking loss as follows:

$$L_{margin} = \frac{1}{N} \sum_{p \in I} \max(0, m + \text{Pos}(p) - \text{Neg}(p)), \quad (1)$$

where  $N$  is the total number of points in image  $I$  and  $m$  is the margin value.

The positive distance  $\text{Pos}(p)$  between the descriptor  $f_p$  of position  $p$  and its corresponding descriptor  $f_{p'}$  is defined

as follows:

$$\text{Pos}(p) = \|f_p - f_{p'}\|_2^2. \quad (2)$$

As for the calculation of negative distance  $\text{Neg}(p)$ , to avoid easy negative samples dominating the loss, we only select the hardest sample. The negative distance is defined as follows:

$$\text{Neg}(p) = \min\left(\min_{k \in I', \|k-p'\|_\infty > T} \|f_p - f_k\|_2^2, \min_{k \in I, \|k-p\|_\infty > T} \|f_{p'} - f_k\|_2^2\right), \quad (3)$$

where  $T$  is a threshold to filter out neighboring points of the ground-truth correspondence point.

Thanks to the triplet margin ranking loss and the transformed versions of input image pairs, the contrastive correspondence network can compute correspondences more robust to scale and rotation transformations. Once trained, the original LR input image and the HR reference image are fed into the contrastive correspondence network to compute correspondences.

### 3.2. Teacher-Student Correlation Distillation

Since a lot of information is lost in LR input images, correspondence matching is difficult, especially for highly textured regions. Matching between two HR images has a better performance than LR-HR matching. To mitigate the gap, we employ the idea of knowledge distillation [10] into our framework. Traditional knowledge distillation tasks [10, 18] deal with model compression issues, while we aim to transfer the matching ability of HR-HR matching to LR-HR matching. Thus, different from the traditional knowledge distillation models that have the same inputs but with different model capacities, in our tasks, the teacher HR-HR matching model and the student LR-HR matching have the exact same designs of architecture, but with different inputs.

The distances between the descriptors of HR input images and reference images provide additional supervision for knowledge distillation. Thus, we propose to push closer the correlation volume (a matrix represents the distances between descriptors of input images and reference images) of teacher model to that of student model. For an input image, we have  $N$  descriptors, and its reference image has  $M$  descriptors. By computing correlations between descriptors of input images and reference images, we can obtain an  $N \times M$  matrix to represent the correlation volume, and view it as a probability distribution by applying a softmax function with temperature  $\tau$  over it. To summarize, the correlation of the descriptor of input image at position  $p$  and the descriptor of reference image at position  $q$  is computed as follows:

$$\text{cor}_{pq} = \frac{e^{\frac{f_p \cdot f_q}{\|f_p\| \cdot \|f_q\|} / \tau}}{\sum_{k \in I'} e^{\frac{f_p \cdot f_k}{\|f_p\| \cdot \|f_k\|} / \tau}}. \quad (4)$$

By computing the correlations  $\text{cor}_{pq}$  for every pair of descriptor  $p$  and  $q$ , we can obtain the correlation volume. We denote  $\text{COR}^T$  and  $\text{COR}^S$  as the teacher correlation volume and student correlation volume, respectively. For every descriptor  $p$  of input image, the divergence of teacher model's correlation and student model's correlation can be measured by Kullback Leibler divergence as follows:

$$\begin{aligned} \text{Div}_p &= \text{KL}(\text{COR}_p^T \parallel \text{COR}_p^S) \\ &= \sum_{k \in I'} \text{cor}_{pk}^T \log\left(\frac{\text{cor}_{pk}^T}{\text{cor}_{pk}^S}\right). \end{aligned} \quad (5)$$

The correlation volume contains the knowledge of relationship between descriptors. By minimizing the divergence between two correlation volumes, the matching ability of teacher model can be transferred to the student model. This objective is defined as follows:

$$L_{kl} = \frac{1}{N} \sum_{p \in I} \text{Div}_p. \quad (6)$$

With the teacher-student correlation distillation, the total loss used for training the contrastive correspondence network is:

$$L = L_{\text{margin}} + \alpha_{kl} \cdot L_{kl}, \quad (7)$$

where  $\alpha_{kl}$  is the weight for the KL-divergence loss.

### 3.3. Dynamic Aggregation Module

After obtaining correspondences, we fuse textures from reference images by a dynamic aggregation module. For every position  $p$  in input images, we compute its correspondence point  $p'$  in reference images as follows:

$$p' = \underset{q \in I'}{\text{argmax}} \text{cor}_{pq}. \quad (8)$$

In order to transfer the texture of reference image, we need to aggregate the information around the position  $p'$ . We denote  $p_0$  as the spatial difference between position  $p$  and  $p'$ , i.e.  $p_0 = p' - p$ . Then the aggregated reference feature  $y$  at position  $p$  is computed by fusing original reference feature  $x$  with a modified DCN as follows:

$$y(p) = \sum_{k=1}^K w_k \cdot x(p + p_0 + p_k + \Delta p_k) \cdot \Delta m_k, \quad (9)$$

where  $p_k \in \{(-1, 1), (-1, 0), \dots, (1, 1)\}$ ,  $w_k$  denotes the convolution kernel weight,  $\Delta p_k$  and  $\Delta m_k$  denote the learnable offset and modulation scalar, respectively.

Compared to the reference feature aggregation operation used in [39] that cropped patches with a fixed size around corresponding points, our dynamic aggregation module dynamically utilizes the information around the precomputed corresponding points with learnable offset  $\Delta p_k$ .



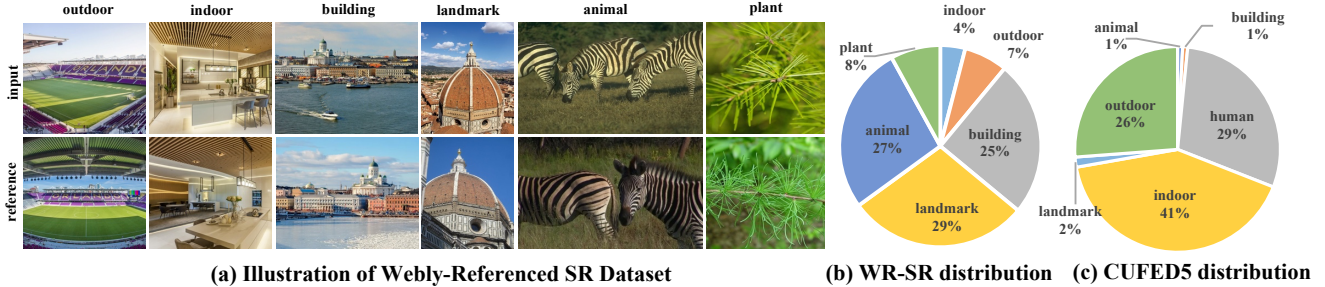


Figure 3. (a) **Illustration of Webly-Referenced SR dataset.** The contents of WR-SR dataset include outdoor scenes, indoor scenes, building, famous landmarks, animals and plants. The first line shows the HR input images and the second line is their reference images searched by Google Image. (b) **WR-SR Dataset Distribution.** (c) **CUFED5 Dataset Distribution.** The WR-SR dataset contains a more diverse category of images. It has more animal, building and landmark images than CUFED5 dataset.

### 3.4. Restoration Module

The restoration module takes the aggregated reference features and input features as input. These features are first concatenated and then fed into residual blocks to generate the desired output SR image. We employ the commonly used reconstruction loss  $L_{rec}$ , perceptual loss  $L_{per}$  and adversarial loss  $L_{adv}$  for the restoration network. The reconstruction loss we adopted is the  $\ell_1$ -norm. The perceptual loss is calculated on relu5-1 VGG features.

### 3.5. Implementation Details

The overall network is trained in two-stage: **1)** training of  $C^2$ -Matching, *i.e.* contrastive correspondence network accompanied with teacher-student correlation distillation. **2)** training of restoration network.

**Training of  $C^2$ -Matching.** We synthesize the image pairs by applying synthetic homography transformation to input images. Homography transformation matrix is obtained by `cv2.getPerspectiveTransform`. The margin value  $m$  in Eq. (1) is set as 1.0, the threshold value  $T$  in Eq. (3) is set as 4.0, the temperature  $\tau$  in Eq. (4) is set as 0.15, and the weight  $\alpha_{kl}$  for KL-divergence loss is 15. The learning rate is set as  $10^{-3}$ .

**Training of Restoration Network.** In this stage, correspondences obtained from the student contrastive correspondence network are used for the calculation of  $p_0$  specified in Eq. (9). The weights for  $L_{rec}$ ,  $L_{per}$  and  $L_{adv}$  are 1.0,  $10^{-4}$  and  $10^{-6}$ , respectively. The learning rate for the training of restoration network is set as  $10^{-4}$ . During training, the input sizes for LR images and HR reference images are  $40 \times 40$  and  $160 \times 160$ , respectively.

## 4. Webly-Referenced SR Dataset

In Ref-SR tasks, the performance relies on similarities between input images and reference images. Thus, the quality of reference images is vital. Currently, Ref-SR methods are trained and evaluated on the CUFED5 dataset [39],

where each input image is accompanied by five references of different levels of similarity to the input image. A pair of input and reference image in CUFED5 dataset is selected from a same event album. Constructing image pairs from albums ensures a high similarity between the input and reference image. However, in realistic settings, it is not always the case that we can find the reference images from off-the-shelf albums.

In real-world applications, given an LR image, users may find possible reference images through web search engines like Google Image. Motivated by this, we propose a more reasonable dataset named Webly Referenced SR (WR-SR) Dataset to evaluate Ref-SR methods. The WR-SR dataset is much closer to the practical usage scenarios, and it is set up as follows:

**Data Collection.** We select about 150 images from BSD500 dataset [21] and Flickr website. These images are used as query images to search for their visually similar images through Google Image. For each query image, the top 50 similar images are saved as reference image pools for the subsequent Data Cleaning procedure.

**Data Cleaning.** Images downloaded from Google Image are of different levels of quality and similarity. Therefore, we manually select the most suitable reference image for each query image. Besides, since some reference images are significantly larger than input images, we rescale the reference images to a comparable scale as HR input images. We also abandon the images with no proper reference images found.

**Data Organization.** A total number of 80 image pairs are collected for WR-SR dataset. Fig. 3 (a) illustrates our WR-SR dataset. The contents of the input images in our dataset include outdoor scenes, indoor scenes, building images, famous landmarks, animals and plants. We analyze the distributions of our WR-SR dataset and CUFED5 dataset. As shown in Fig. 3 (b), compared to CUFED5 dataset (Fig. 3 (c)), we have a more diverse category, and we include more animal, landmark, building and plant images.

Table 1. **Quantitative Comparisons.** PSNR / SSIM are used for evaluation. We group methods by SISR and Ref-SR. We mark the best results in **bold**. The models trained with GAN loss are marked in gray. The suffix ‘-rec’ means only reconstruction loss is used for training.

	Method	CUFED5	Sun80	Urban100	Manga109	WR-SR
SISR	SRCNN [6]	25.33 / .745	28.26 / .781	24.41 / .738	27.12 / .850	27.27 / .767
	EDSR [17]	25.93 / .777	28.52 / .792	25.51 / .783	28.93 / .891	28.07 / .793
	RCAN [36]	26.06 / .769	29.86 / .810	25.42 / .768	29.38 / .895	28.25 / .799
	SRGAN [15]	24.40 / .702	26.76 / .725	24.07 / .729	25.12 / .802	26.21 / .728
	ENet [23]	24.24 / .695	26.24 / .702	23.63 / .711	25.25 / .802	25.47 / .699
	ESRGAN [30]	21.90 / .633	24.18 / .651	20.91 / .620	23.53 / .797	26.07 / .726
	RankSRGAN [35]	22.31 / .635	25.60 / .667	21.47 / .624	25.04 / .803	26.15 / .719
Ref-SR	CrossNet [40]	25.48 / .764	28.52 / .793	25.11 / .764	23.36 / .741	-
	SRNTT	25.61 / .764	27.59 / .756	25.09 / .774	27.54 / .862	26.53 / .745
	SRNTT-rec [39]	26.24 / .784	28.54 / .793	25.50 / .783	28.95 / .885	27.59 / .780
	TTSR	25.53 / .765	28.59 / .774	24.62 / .747	28.70 / .886	26.83 / .762
	TTSR-rec [34]	27.09 / .804	30.02 / .814	25.87 / .784	30.09 / .907	27.97 / .792
	SSEN	25.35 / .742	-	-	-	-
	SSEN-rec [26]	26.78 / .791	-	-	-	-
	E2ENT <sup>2</sup>	24.01 / .705	28.13 / .765	-	-	-
	E2ENT <sup>2</sup> -rec [32]	24.24 / .724	28.50 / .789	-	-	-
	CIMR	26.16 / .781	29.67 / .806	25.24 / .778	-	-
CIMR-rec [33]	26.35 / .789	30.07 / .813	25.77 / .792	-	-	
Ours	$C^2$ -Matching	27.16 / .805	29.75 / .799	25.52 / .764	29.73 / .893	27.80 / .780
	$C^2$ -Matching-rec	<b>28.24 / .841</b>	<b>30.18 / .817</b>	<b>26.03 / .785</b>	<b>30.47 / .911</b>	<b>28.32 / .801</b>

To summarize, our WR-SR dataset has two advantages over CUFED5 dataset: 1) The pairs of input images and reference images are collected in a more realistic way. 2) Our contents are more diverse than CUFED5 dataset.

## 5. Experiments

### 5.1. Datasets and Evaluation Metrics

**Training Dataset.** We train our models on the training set of CUFED5 dataset [39], which has 11,871 image pairs and each image pair has an input image and a reference image.

**Testing Datasets.** The performance are evaluated on the testing set of CUFED5 dataset, SUN80 dataset [28], Urban100 dataset [11], Manga109 dataset [22] and our WR-SR dataset. The CUFED5 dataset has 126 input images and each has 5 reference images with different similarity levels. The SUN80 dataset has 80 images with 20 reference images for each input image. WR-SR dataset has been introduced in Section 4. As for the SISR datasets, we adopt the same evaluation setting as [39, 34]. The Urban100 dataset contains 100 building images, and the LR versions of images serve as reference images. The Manga109 dataset has 109 manga images and the reference images are randomly selected from the dataset.

**Evaluation Metrics.** PSNR and SSIM on Y channel of YCrCb space are adopted as evaluation metrics. Input LR images for evaluation are constructed by bicubic downsampling  $4\times$  from HR images.

### 5.2. Results Comparisons

**Quantitative Comparison.** We compare the proposed  $C^2$ -Matching with SISR methods and Ref-SR methods. For SISR methods, we include SRCNN [6], EDSR [17], RCAN [36], SRGAN [15], ENet [23], ESRGAN [30] and RankSRGAN [35]. For Ref-SR methods, CrossNet [40], SRNTT [39], SSEN [26], TTSR [34], E2ENT<sup>2</sup> [32] and CIMR [33] are included.

Table 1 shows the quantitative comparison results. We mark the methods trained with GAN loss in gray. On the standard CUFED5 benchmark, our proposed method outperforms state of the arts by a large margin. We also achieve the best results on the Sun80, Urban100, Manga109 and WR-SR dataset, which demonstrates the great generalizability of  $C^2$ -Matching. Notably, CIMR [33] is a multiple reference-based SR method, which transfers the HR textures from a collection of reference images. Our proposed method performs better than CIMR, which further verifies the superiority of our method.

**Qualitative Evaluation.** Fig. 4 shows the qualitative comparison with state of the arts. We compare our method with ESRGAN [30], RankSRGAN [35], SRNTT [39], and TTSR [34]. The results of our method have the best visual quality containing many realistic details and are closer to their respective HR ground-truths. As shown in the top left example,  $C^2$ -Matching successfully recovers the exact word “EVERY” while other methods fail. Besides, as shown in

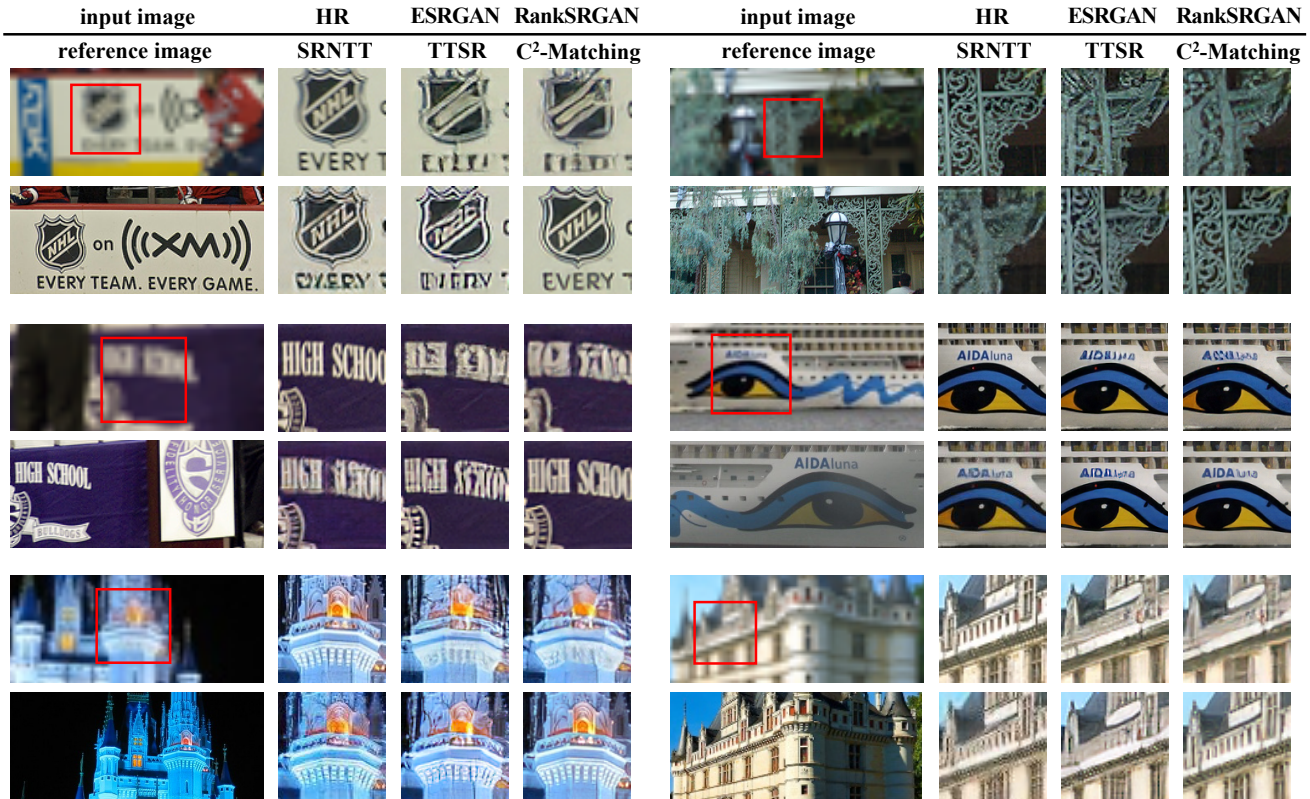


Figure 4. **Qualitative Comparisons.** We compare our results with ESRGAN [30], RankSRGAN [35], SRNTT [39], and TTSR [34]. All these methods are trained with GAN loss. Our results have better visual quality with more texture details.

the bottom left example, our approach can deal with reasonable color and illumination changes to a certain extent without deliberate augmentation.

### 5.3. Ablation Study

We perform ablation studies to assess the effectiveness of each module. To evaluate the effectiveness of our proposed modules on texture regions, we select a subset of CUFED5 dataset that contains images with complicated textures; we name it “texture subset”. On top of the baseline model, we progressively add the dynamic aggregation module, contrastive correspondence network and teacher-student correlation distillation to show their effectiveness. The ablation study results are shown in Fig. 5.

**Dynamic Aggregation Module.** We first analyze the effectiveness of the Dynamic Aggregation (Dyn-Agg) module because it deals with the reference texture transfer problem. Only with a better texture transfer module would the improvements of correspondence matching module be reflected, *i.e.* the baseline accompanied with Dyn-Agg module provides a stronger backup. The Dyn-Agg module dynamically fuses the HR information from reference images. Compared to the previous scheme that cropped patches of a fixed size from HR reference features, the Dyn-Agg module

has a more flexible patch size. With Dyn-Agg module, we observe an increment in PSNR by 0.08dB in the full dataset.

**Contrastive Correspondence Network.** We further replace the fixed VGG feature matching module with the contrastive correspondence (Contras) network. With the learnable contrastive correspondence network, the PSNR value increases by about 0.2dB. This result demonstrates the contrastive correspondence network computes more accurate correspondences and further boosts the performance of restoration. Fig. 5 shows one example of visual comparisons. With the contrastive correspondence network, the output SR images have more realistic textures.

**Teacher-Student Correlation Distillation.** With the teacher-student correlation (TS Corr) distillation, the performance further increases by 0.09dB on the whole dataset. For the texture subset, the performance increases by 0.16dB. The TS Corr module aims to push closer the correspondence of LR-HR student Contras network and that of HR-HR teacher Contras network. Since HR-HR teacher matching model is more capable of matching texture regions, the TS Corr module mainly boosts the performance on texture regions. As indicated in Fig. 5, with TS Corr module, the textures of output SR images are enriched.



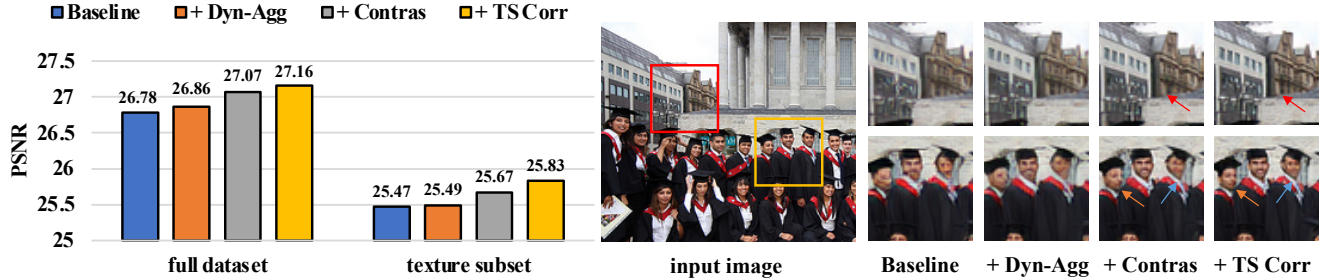


Figure 5. **Ablation Study.** We evaluate the effectiveness of each component on the full CUFED5 dataset and the texture region subset. Since the contrastive correspondence (Contrás) network and teacher-student correlation distillation (TS Corr) focus on improving texture details, we also add visual comparisons with the component added.

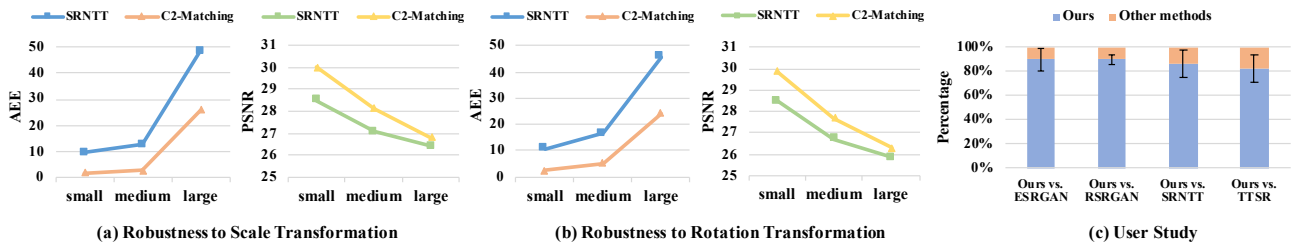


Figure 6. **Further Analysis.** (a) Robustness to scale transformation. (b) Robustness to rotation transformations. The proposed  $C^2$ -Matching is more robust to scale and rotation transformation compared to SRNTT. (c) User study. Compared to other state of the arts, over 80% users prefer our results.

## 5.4. Further Analysis

**Robustness to Scale and Rotation Transformations.** We perform further analysis on the robustness of our  $C^2$ -Matching to scale and rotation transformations. We build a transformation-controlled dataset based on CUFED5 dataset. The scaled and rotated versions of input images serve as reference images. We adopt two metrics to measure the robustness: Average End-to-point Error (AEE) for matching accuracy and PSNR for restoration performance.

Fig. 6 shows the robustness to the scale and transformations in AEE and PSNR. We separately analyze the impact of scale and rotation. We classify the degrees of scale and rotations into three groups: small, medium and large. The AEE rises as the degrees of transformations increases, which indicates a larger degree of transformation makes it harder to perform correspondence matching. Based on AEE, our proposed  $C^2$ -Matching computes more accurate correspondences than SRNTT under scale and rotation transformations. With the features that are more robust to scale and rotation transformations, according to the PSNR, the restoration performance of our proposed  $C^2$ -Matching is also more robust than that of SRNTT [39]. It should be noted that large transformations are not included during training but our proposed  $C^2$ -Matching still exhibits superior performance compared to SRNTT.

**User Study.** We perform a user study to further demonstrate the superiority of our method qualitatively. A total number of 20 users are asked to compare the visual quality

of our method and state of the arts on the CUFED5 dataset, including ESRGAN [30], RankSRGAN [35], SRNTT [39] and TTSR [34]. We present images in pairs, of which one is the result of our method, and ask users to choose the one offering better visual quality. As shown in Fig. 6, over 80% of the users felt that the result of our method is superior compared to that of state of the arts.

## 6. Conclusion

In this paper, we propose a novel  $C^2$ -Matching for robust reference-based super-resolution, which consists of contrastive correspondence network, teacher-student correlation distillation and dynamic aggregation module. The motivation of contrastive correspondence network is to perform scale and rotation robust matching between input images and reference images. The teacher-student correlation distillation is proposed to distill the teacher HR-HR matching to guide the student LR-HR matching to improve the visual quality of texture regions. After obtaining the correspondences, we fuse the information of reference images through a dynamic aggregation module. We achieve over 1dB improvement in PSNR over state of the arts. To facilitate a more realistic evaluation of Ref-SR tasks, we also contribute a new benchmark named WR-SR dataset, which is collected in a more realistic way.

**Acknowledgement.** This research was conducted in collaboration with SenseTime. This work is supported by NTU NAP and A\*STAR through the Industry Alignment Fund - Industry Collaboration Projects Grant.



## References

- [1] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. Patchmatch: A randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24, 2009. 2
- [2] Kelvin C.K. Chan, Xintao Wang, Ke Yu, Chao Dong, and Chen Change Loy. Understanding deformable alignment in video super-resolution. In *AAAI*, 2021. 2
- [3] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Int. Conf. Comput. Vis.*, pages 764–773, 2017. 2
- [4] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 11065–11074, 2019. 1, 2
- [5] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superpoint: Self-supervised interest point detection and description. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, pages 224–236, 2018. 2
- [6] 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. 1, 2, 6
- [7] Shivam Duggal, Shenlong Wang, Wei-Chiu Ma, Rui Hu, and Raquel Urtasun. Deeppruner: Learning efficient stereo matching via differentiable patchmatch. In *Int. Conf. Comput. Vis.*, pages 4384–4393, 2019. 2
- [8] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler. D2-net: A trainable cnn for joint description and detection of local features. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8092–8101, 2019. 2, 3
- [9] Qinquan Gao, Yan Zhao, Gen Li, and Tong Tong. Image super-resolution using knowledge distillation. In *Asian Conference on Computer Vision*, pages 527–541. Springer, 2018. 2
- [10] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 4
- [11] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5197–5206, 2015. 6
- [12] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Eur. Conf. Comput. Vis.*, pages 694–711. Springer, 2016. 2, 11
- [13] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1646–1654, 2016. 1
- [14] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1637–1645, 2016. 1, 2
- [15] 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, 11
- [16] Wonkyung Lee, Junghyup Lee, Dohyung Kim, and Bumsub Ham. Learning with privileged information for efficient image super-resolution. In *European Conference on Computer Vision*, pages 465–482. Springer, 2020. 2
- [17] 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. 1, 2, 6
- [18] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured knowledge distillation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2604–2613, 2019. 4
- [19] Or Lotan and Michal Irani. Needle-match: Reliable patch matching under high uncertainty. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 439–448, 2016. 2
- [20] David G Lowe. Object recognition from local scale-invariant features. In *Int. Conf. Comput. Vis.*, volume 2, pages 1150–1157. IEEE, 1999. 2
- [21] D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *Int. Conf. Comput. Vis.*, volume 2, pages 416–423, July 2001. 5
- [22] Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications*, 76(20):21811–21838, 2017. 6
- [23] Mehdi SM Sajjadi, Bernhard Scholkopf, and Michael Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. In *Int. Conf. Comput. Vis.*, pages 4491–4500, 2017. 2, 6
- [24] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature matching with graph neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 4938–4947, 2020. 2
- [25] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1874–1883, 2016. 1, 2
- [26] Gyumin Shim, Jinsun Park, and In So Kweon. Robust reference-based super-resolution with similarity-aware deformable convolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, June 2020. 1, 2, 6
- [27] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *Int. Conf. Learn. Represent.*, 2015. 2, 3, 11

- [28] Libin Sun and James Hays. Super-resolution from internet-scale scene matching. In *IEEE Int. Conf. Comput. Photo.*, pages 1–12. IEEE, 2012. [6](#)
- [29] Xintao Wang, Kelvin C.K. Chan, Ke Yu, Chao Dong, and Chen Change Loy. Edvr: Video restoration with enhanced deformable convolutional networks. In *IEEE Conf. Comput. Vis. Pattern Recog. Worksh.*, June 2019. [2](#)
- [30] 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.*, 2018. [2](#), [6](#), [7](#), [8](#), [12](#), [13](#), [14](#)
- [31] Olivia Wiles, Sebastien Ehrhardt, and Andrew Zisserman. D2d: Learning to find good correspondences for image matching and manipulation. *arXiv preprint arXiv:2007.08480*, 2020. [2](#)
- [32] Yanchun Xie, Jimin Xiao, Mingjie Sun, Chao Yao, and Kaizhu Huang. Feature representation matters: End-to-end learning for reference-based image super-resolution. In *Eur. Conf. Comput. Vis.*, pages 230–245. Springer, 2020. [1](#), [2](#), [3](#), [6](#)
- [33] Xu Yan, Weibing Zhao, Kun Yuan, Ruimao Zhang, Zhen Li, and Shuguang Cui. Towards content-independent multi-reference super-resolution: Adaptive pattern matching and feature aggregation. In *Eur. Conf. Comput. Vis.*, volume 2, 2020. [2](#), [6](#)
- [34] Fuzhi Yang, Huan Yang, Jianlong Fu, Hongtao Lu, and Bainig Guo. Learning texture transformer network for image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5791–5800, 2020. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [12](#), [13](#), [14](#)
- [35] Wenlong Zhang, Yihao Liu, Chao Dong, and Yu Qiao. Ranksrgan: Generative adversarial networks with ranker for image super-resolution. In *Int. Conf. Comput. Vis.*, pages 3096–3105, 2019. [2](#), [6](#), [7](#), [8](#), [12](#), [13](#), [14](#)
- [36] 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. [2](#), [6](#), [12](#)
- [37] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image super-resolution. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2472–2481, 2018. [2](#)
- [38] Yulun Zhang, Zhifei Zhang, Stephen DiVerdi, Zhaowen Wang, Jose Echevarria, and Yun Fu. Texture hallucination for large-factor painting super-resolution. In *Eur. Conf. Comput. Vis.*, 2020. [2](#)
- [39] Zhifei Zhang, Zhaowen Wang, Zhe Lin, and Hairong Qi. Image super-resolution by neural texture transfer. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 7982–7991, 2019. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [11](#), [12](#), [13](#), [14](#)
- [40] Haitian Zheng, Mengqi Ji, Haoqian Wang, Yebin Liu, and Lu Fang. Crossnet: An end-to-end reference-based super resolution network using cross-scale warping. In *Eur. Conf. Comput. Vis.*, pages 88–104, 2018. [1](#), [2](#), [6](#), [12](#)
- [41] Shangchen Zhou, Jiawei Zhang, Wangmeng Zuo, and Chen Change Loy. Cross-scale internal graph neural network for image super-resolution. In *Adv. Neural Inform. Process. Syst.*, 2020. [2](#)
- [42] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9308–9316, 2019. [2](#)

## Supplementary

In this supplementary file, we will explain the network structures (*i.e.* Contrastive Correspondence Network and Restoration Network) and training details in Section A. Then we will introduce training losses we used in Section B. In Section C, the model size will be analyzed. In Section D, we will provide more visual comparisons with state-of-the-art methods. Finally, we will show more visual comparisons of ablation study in Section E.

### A. Network Structures and Training Details

#### A.1. Contrastive Correspondence Network

**Network Structure.** Table 2 shows the detailed feature extractor structure of contrastive correspondence network. Since the resolutions of input image and reference image are different, we adopt two feature extractors for LR input image and HR reference image, respectively.

Table 2. **The feature extractor structure of contrastive correspondence network.** The kernel size of convolution layers is  $3 \times 3$  and the MaxPool operation is with kernel size of  $2 \times 2$ .

#	Layer name(s)
0	Conv (3, 64), ReLU
1	Conv (64, 64), ReLU
2	MaxPool ( $2 \times 2$ )
3	Conv (64, 128), ReLU
4	Conv (128, 128), ReLU
5	MaxPool ( $2 \times 2$ )
6	Conv (128, 256)

**Training Details.** To enable teacher-student correlation distillation, a teacher contrastive correspondence network should be first trained. The hyperparameters for the training of teacher model are set as follows: the margin value  $m$  is 1.0, the threshold value  $T$  is 4.0, the batch size is set as 8, and the learning rate is  $10^{-3}$ . We use the pretrained weights of VGG-16 to initialize the feature extractor. Then the student contrastive correspondence network is trained with the teacher network fixed. The margin value  $m$ , threshold value  $T$ , batch size and learning rate are the same as the teacher network. The temperature value  $\tau$  is 0.15, and the weight  $\alpha_{kl}$  for KL-divergence loss is 15.

#### A.2. Restoration Network

**Network Structure.** The restoration network consists of dynamic aggregation module and restoration module. For each image, three reference features (*i.e.* pretrained VGG relu3\_1, relu2\_1, relu1\_1 feature [27]) are aggregated by dynamic aggregation module, and the aggregated reference features are denoted as Aggregated Reference Feature1, Aggregated Reference Feature2 and Aggregated Reference

Feature3, respectively. The structure of restoration module is illustrated in Table. 3.

Table 3. **The structure of restoration module.** The kernel size of convolution layers is  $3 \times 3$ . PixelShuffle layers are  $2 \times$ . RB denotes residual block. Aggregated Reference Feature denotes the reference feature aggregated by the dynamic aggregation module.

#	Layer name(s)
0	Conv(3, 64), LeakyReLU
1	RB [Conv(64, 64), ReLU, Conv(64, 64)] $\times$ 16
2	Concat [#1, Aggregated Reference Feature1]
3	Conv(320, 64), LeakyReLU
4	RB [Conv(64, 64), ReLU, Conv(64, 64)] $\times$ 16
5	ElementwiseAdd(#1, #4)
6	Conv(64, 256), PixelShuffle, LeakyReLU
7	Concat [#6, Aggregated Reference Feature2]
8	Conv(192, 64), LeakyReLU
9	RB [Conv(64, 64), ReLU, Conv(64, 64)] $\times$ 16
10	ElementwiseAdd(#6, #9)
11	Conv(64, 256), PixelShuffle, LeakyReLU
12	Concat [#11, Aggregated Reference Feature3]
13	Conv(128, 64), LeakyReLU
14	RB [Conv(64, 64), ReLU, Conv(64, 64)] $\times$ 16
15	ElementwiseAdd(#11, #14)
16	Conv(64, 32), LeakyReLU
17	Conv(32, 3)

**Training Details.** The learning rate is set as  $10^{-4}$ . For the training of the network with adversarial loss and perceptual loss, we adopt the same setting as [39] (*i.e.* the network is trained with only reconstruction loss for the first 10K iterations).

### B. Loss Functions

**Reconstruction Loss.** The  $\ell_1$ -norm is adopted to keep the spatial structure of the LR images. It is defined as follows:

$$L_{rec} = \|I^{HR} - I^{SR}\|_1. \quad (10)$$

**Perceptual Loss.** The perceptual loss [12] is employed to improve the visual quality. It is defined as follows:

$$L_{per} = \frac{1}{V} \sum_{i=1}^C \|\phi_i(I^{HR}) - \phi_i(I^{SR})\|_F, \quad (11)$$

where  $V$  and  $C$  denotes the volume and channel number of feature maps.  $\phi$  denotes the relu5\_1 features of VGG19 model [27].  $\|\cdot\|_F$  denotes the Frobenius norm.

**Adversarial Loss.** The adversarial loss [15] is defined as follows:

$$L_{adv} = -D(I^{SR}). \quad (12)$$

The loss for training discriminator  $D$  is defined as follows:

$$L_D = D(I^{SR}) - D(I^{HR}) + \lambda(\|\nabla_{\hat{I}}D(\hat{I})\|_2 - 1)^2. \quad (13)$$

where  $\hat{I}$  is the random convex combination of  $I^{SR}$  and  $I^{HR}$ .

### C. Comparison of Model Size

The comparison of model size (*i.e.* the number of trainable parameters) is shown in Table 4. Our proposed  $C^2$ -Matching has a total number of 8.9M parameters and achieves a PSNR of 28.24dB. For a fair comparison in terms of model size, we build a light version of  $C^2$ -Matching, which has fewer trainable parameters. The  $C^2$ -Matching-*light* is built by setting the number of residual blocks of layer #9 and layer #14 to 8 and 4, respectively, and removing the Aggregated Reference Feature1. The  $C^2$ -Matching-*light* has a total number of 4.8M parameters. The light version has fewer parameters than TTSR [34] but significantly better performance.

Table 4. **Model sizes of different methods.** PSNR / SSIM are adopted as the evaluation metrics.

Method	Params	PSNR/SSIM
RCAN [36]	16M	26.06 / .769
RankSRGAN [35]	1.5M	22.31 / .635
CrossNet [40]	33.6M	25.48 / .764
SRNTT [39]	4.2M	26.24 / .784
TTSR [34]	6.4M	27.09 / .804
$C^2$ -Matching- <i>light</i>	4.8M	28.14 / .839
$C^2$ -Matching	8.9M	28.24 / .841

### D. More Visual Comparisons with State-of-the-art Methods

In Fig. 7 and Fig. 8, more visual comparisons with ES-RGAN [30], RankSRGAN [35], SRNTT [39] and TTSR [34] are provided. The images restored by our proposed  $C^2$ -Matching have better visual quality.

### E. More Visual Comparisons of Ablation Study

In this paper, the proposed  $C^2$ -Matching consists of three major components: Dynamic Aggregation Module (Dyn-Agg), Contrastive Correspondence Network (Contras) and Teacher-Student Correlation Distillation (TS Corr). On top of the baseline model, we progressively add the Dyn-Agg module, Contras network and TS Corr distillation. In Fig. 9, we show more visual comparisons with these proposed modules progressively added.



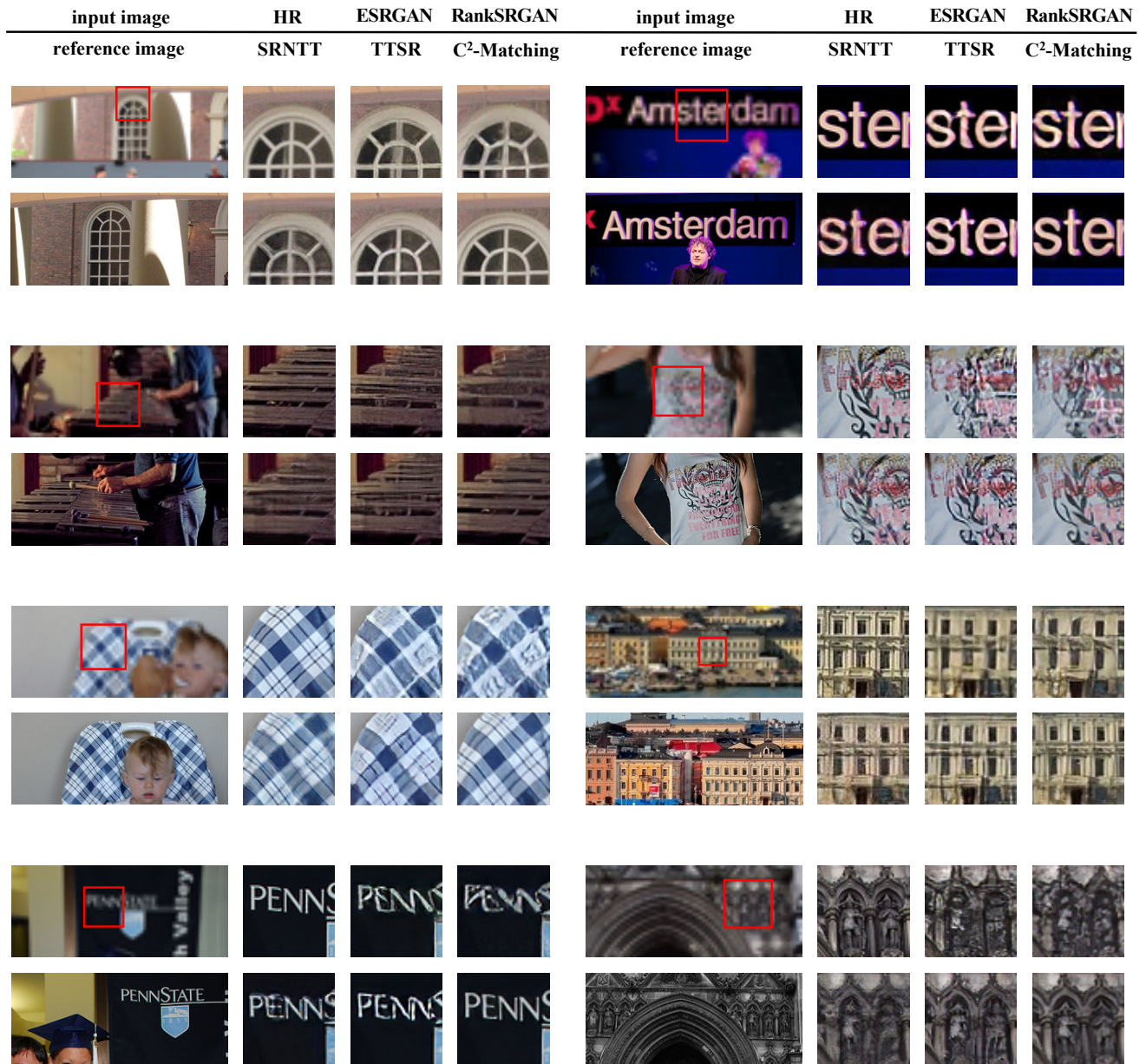


Figure 7. **More visual comparisons.** We compare our results with ESRGAN [30], RankSRGAN [35], SRNTT [39], and TTSR [34]. All these methods are trained with GAN loss. Our results have better visual quality with more texture details.

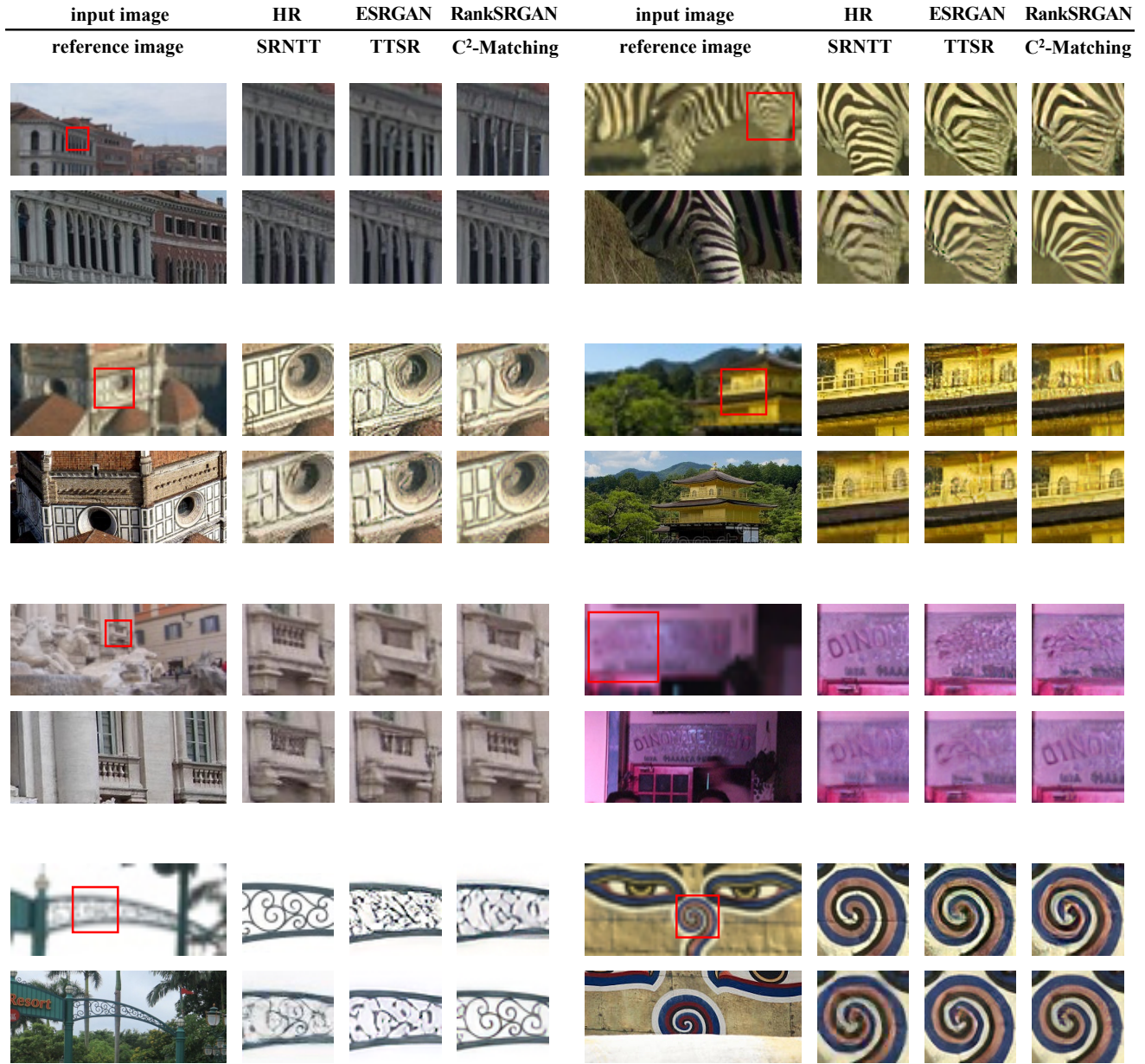
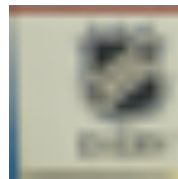


Figure 8. **More visual comparisons.** We compare our results with ESRGAN [30], RankSRGAN [35], SRNTT [39], and TTSR [34]. All these methods are trained with GAN loss. Our results have better visual quality with more texture details.

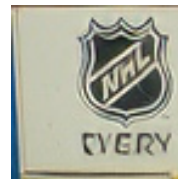




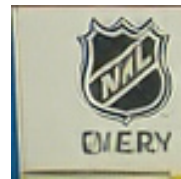
input image



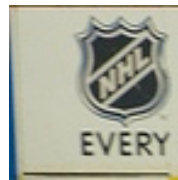
LR



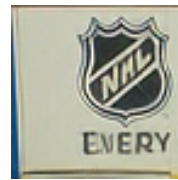
Baseline



+ Dyn-Agg



HR



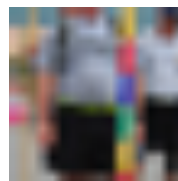
+ Contrast



+ TS Corr



input image



LR



Baseline



+ Dyn-Agg



HR



+ Contrast



+ TS Corr

Figure 9. **More visual comparisons of ablation study.** On top of the baseline model, Dynamic Aggregation Module (Dyn-Agg), Contrastive Correspondence Network (Contrast) and Teacher-Student Correlation Distillation (TS Corr) are progressively added.