

Improving Image Restoration through Removing Degradations in Textual Representations

Jingbo Lin¹, Zhilu Zhang¹, Yuxiang Wei¹, Dongwei Ren¹, Dongsheng Jiang², Wangmeng Zuo^{1,*}

¹Harbin Institute of Technology ²Huawei Cloud Computing Co., Ltd.

jblincls1996@gmail.com, cszlzhang@outlook.com, yuxiang.wei.cs@gmail.com,
rendongweihiit@gmail.com, dongsheng-jiang@outlook.com, cswmzuo@gmail.com

Abstract

In this paper, we introduce a new perspective for improving image restoration by removing degradation in the textual representations of a given degraded image. Intuitively, restoration is much easier on text modality than image one. For example, it can be easily conducted by removing degradation-related words while keeping the content-aware words. Hence, we combine the advantages of images in detail description and ones of text in degradation removal to perform restoration. To address the cross-modal assistance, we propose to map the degraded images into textual representations for removing the degradations, and then convert the restored textual representations into a guidance image for assisting image restoration. In particular, We ingeniously embed an image-to-text mapper and text restoration module into CLIP-equipped text-to-image models to generate the guidance. Then, we adopt a simple coarse-to-fine approach to dynamically inject multi-scale information from guidance to image restoration networks. Extensive experiments are conducted on various image restoration tasks, including deblurring, dehazing, deraining, and denoising, and all-in-one image restoration. The results showcase that our method outperforms state-of-the-art ones across all these tasks. The codes and models are available at <https://github.com/mrluin/TextualDegRemoval>.

1. Introduction

Image restoration aims to reconstruct a high-quality clean image from its degraded observations. Most existing methods design deep networks for specific restoration tasks, including image denoising [57, 114, 118, 120, 122], deblurring [14, 15, 19, 95, 113], deraining [17, 95, 100, 113], dehazing [9, 20, 37, 90, 92], etc. Recent works [16, 51, 71, 73, 117, 127] expect to explore a unified model for multiple

degradations. However, the severely ill-posed nature of the task makes it non-trivial to separate degradations and desired image content. Especially for unified models, the potential conflicts in dealing with various degradations bring more uncertainty.

From a broader perspective, the purpose of the restoration is to enhance the clarity of scenes for human perception and recognition, while performing it on image modality is just one option. Moreover, degradations in image modality are tightly coupled with desirable content, making their removal challenging. When recording the scenes in some other modalities, this problem can be alleviated. Let us take the text modality as an example. Assume the textual description of a clean scene is represented by ‘a scene of *’. When this scene is subjected to rainfall, the description can be converted into ‘a rainy scene of *’. Thus, deraining can be readily achieved by simply removing the rain-related text ‘rainy’. Furthermore, it is also convenient to remove multiple degradations in textual space with one unified model. Note that text modality may only describe rough aspects and ignore some details. We can further introduce image modality to combine their complementary potential in degradation removal and image restoration.

Motivated by the above, we propose to perform restoration first in a modality where degradations and content are loosely coupled, and then utilize the corrected content to guide image restoration. In particular, we adopt the commonly used textual modality. Note that the text corresponding to the scene is not always available and the cross-modal assistance is difficult to achieve directly. The gap between text and image should be bridged. First, degraded images should be mapped into textual space for removing degradations. Second, the restored text should be mapped into a guidance image to assist restoration of the degraded image. As shown in Fig. 1, for the former, we have tried to convert degraded images to textual captions by image-to-text (I2T) models (e.g., BLIP [52, 53]), and find the captions can explicitly express both degradation types (e.g., blur, rain, haze, and noise) and content information. For the latter, we can

*Correspondence author.

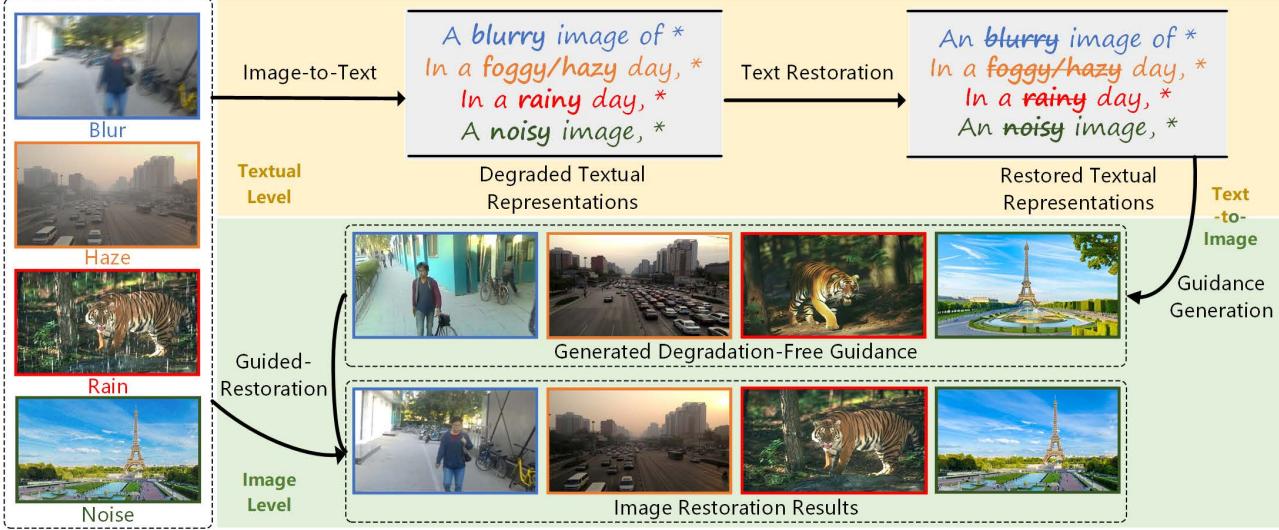


Figure 1. **Overview of the proposed method.** We propose to improve image restoration by performing restoration on the textual level, in which content and degradation information are loosely coupled. We first encode image concepts into textual space and then remove degradation-related information. To achieve cross-modal assistance, we employ pre-trained T2I models to generate clean guidance for the image restoration process.

utilize the **text-to-image (T2I) models** (e.g., Imagen [84]), which have demonstrated impressive image generation capability from the given texts.

Although such a scheme is conceptually feasible, converting to explicit text may lose a lot of information from images, and additional degraded-clean text pairs need to be collected to train the text restoration module. Fortunately, Contrastive Language-Image Pre-training (CLIP) [76] implicitly aligns concepts of images and text. And we can leverage the CLIP-equipped T2I models (e.g., Stable Diffusion [83]) to build an end-to-end framework for generating clear guidance from degraded images gracefully. Specifically, CLIP has been demonstrated as an effective I2T converter [24]. For converting degraded images into degraded textual representations, we adopt the image encoder of CLIP, and add an I2T mapper after the encoder. For restoring degraded representations, we further append a textual restoration module after I2T mapper. Then, I2T mapper and textual restoration module can be sequentially trained by employing paired data from different image restoration datasets, without additional text data.

Due to the ease of removing degradations on textual level, we train the I2T mapper and textual restoration module with multiple degradations at once. When training is done, it can serve multiple tasks, generating content-related and degradation-free guidance images from degraded images. Finally, we adopt a simple coarse-to-fine approach to dynamically inject multi-scale information from guidance to classic image restoration networks. Extensive experiments are conducted on multiple tasks, including all-in-one restoration, image deburring, dehazing, deraining, and denoising. The results show our method achieves better uni-

versally than corresponding state-of-the-art methods. Especially for all-in-one restoration, 0.5 dB PSNR improvement is obtained in comparison with PromptIR [73].

The main contributions can be summarized as follows:

- We introduce a new perspective for image restoration, *i.e.*, performing restoration first in textual space where degradations and content are loosely coupled, and then utilizing the restored content to guide image restoration.
- To address the cross-modal assistance, we propose to embed an image-to-text mapper and textual restoration module into CLIP-equipped text-to-image models to generate clear guidance from degraded images.
- Extensive experiments on multiple tasks demonstrate that our method improves the performance of state-of-the-art image restoration networks.

2. Related Work

2.1. Image Restoration

Image Restoration for Specific Tasks. Image restoration is a fundamental computer vision problem, which aims to reconstruct high-quality images from corresponding degraded inputs. In recent years, deep learning has achieved great progress in image restoration. Starting from some simple convolutional neural networks (CNNs) [25, 118], the introduction of channel-attention [124], spatial-attention [36, 112], non-local operation [59, 125], skip-connection architectures [55, 112] and multi-stage scheme [95, 113] enables image restoration performance continuously improve. With the emergence of vision transformers [27, 62], the capability of capturing long-range dependencies in the image allows transformer-based

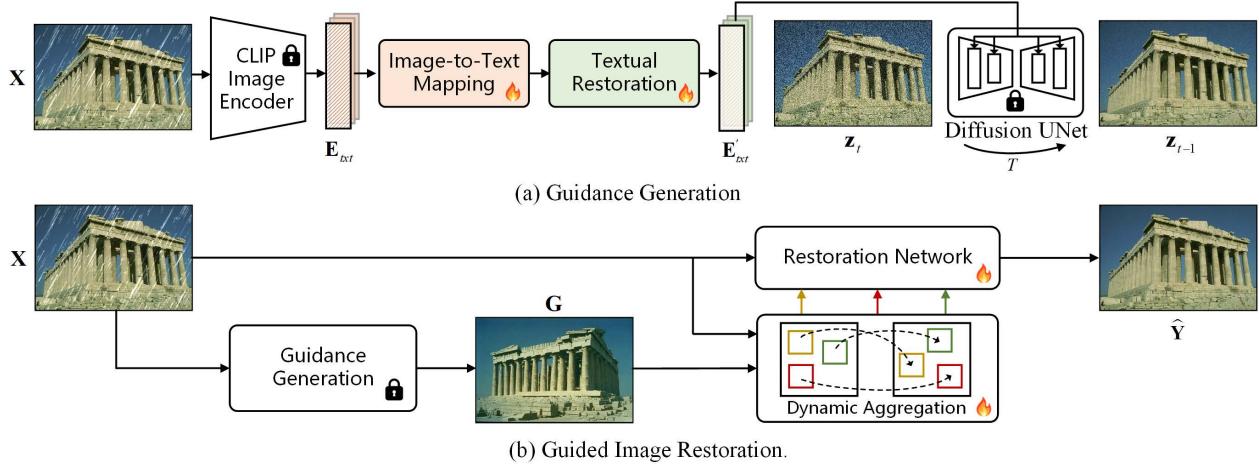


Figure 2. **Illustration of the proposed pipeline.** (a) We sequentially train image-to-text mapper \mathcal{M}_{i2t} and textual restoration module \mathcal{M}_{clean} to convert image concepts into textual representations and remove textual degradation information, respectively. (b) The guidance image is used to assist the image restoration process.

methods to achieve better performance, gradually replacing previous CNN-based methods. To balance computational cost, window-based attention [18, 57] and transposed attention [114] are also introduced into image restoration tasks. Albeit image restoration performance has benefited from various advanced architecture designs, most of the works only focus on one specific degradation, due to the difficulty of learning to remove multiple degradations.

All-in-One Image Restoration. Promoting by the development of unified model, some works focus on the solution to the all-in-one image restoration [16, 51, 71, 73, 117, 127]. The first all-in-one image restoration work, *i.e.*, IPT [13], employs ViT-based backbone with multi-heads and multi-tails for inputs with different degradations. However, the method is limited in handling specific synthesized degradations and cannot directly generalize to unknown tasks. AirNet [51] develops the unified model for denoising-deraining-dehazing, which utilizes contrastive learning to capture degradation representations of different tasks and adaptively injects the degradation priors into backbone restoration framework for aiding in learning better results. ADMS [71] exploits FAIG [101] in all-in-one image restoration, by learning specific filters and degradation classifiers, achieving better performance on deraining-denoising-deblurring and deraining-desnowing-dehazing tasks. IDR [117] proposes a two-stage training strategy, which first learns separate task-oriented hubs for each degradation. Then, in the second stage, it reformulates the learned hubs into a single ingredient-oriented hub by learnable PCA and adopts the reformulated hub as prior to adaptively aid in restoring corrupted images. PromptIR [73] learns to encode and adopt degradation information as prompts for dehazing-deraining-denoising task. Although existing works can unify a set of image restoration tasks into one unified model, the performance is still limited

caused by the large gap among different degradations.

Prior-Based Image Restoration. Prior-based restoration aims at improving performance by introducing external priors, including structures, images, pre-trained models, *etc*. For instance, some works [10, 43, 63] adopt information from the high-resolution reference image to help improve the performance of super-resolution. Recent works [58, 94, 103] utilize pre-trained generative priors to restore more realistic and natural results. TextIR [7] develops a text-driven image restoration framework by incorporating information on textual representation.

In this work, we provide a new perspective, *i.e.*, utilizing restoration in textual space to assist image restoration. We also combine the advantages of both text-based prior (by learning degradations in textual level) and image-based prior (by giving clean guidance) for better performance.

2.2. Text-to-Image Generation

Text-to-image generation models [23, 32, 69, 77, 78, 85, 98] have attracted intensive attention in recent years due to its ability to generate high-quality and diverse images based on given text descriptions. A variety of techniques, including generative adversarial networks (GAN) [35], autoregressive models, and diffusion models have been investigated. Initial studies [49, 85, 98] mainly rely on GAN-based architectures, and train a conditional model from given paired image-caption datasets to generate samples. Some efforts focus on autoregressive models [11, 22, 32, 77, 107] have also shown exciting results. These models, such as CogView [22] and Muse [11] first learn a discrete codebook through training an autoencoder, and then adopt an autoregressive transformer to predict the tokens sequentially. With the development of diffusion models [39, 89], text conditioned image synthesis has shown remarkable improvement. By training with huge corpora, large diffusion

models, such as DALLE-2 [78], Imagen [84], Stable Diffusion [82], and DALLE-3 [8] have demonstrated excellent semantic understanding, and can generate diverse and photo-realistic images.

3. Proposed Method

In this section, we present our proposed method for image restoration. As illustrated in Fig. 2, we suggest first conducting **restoration in the textual modality**, in which degradation is **loosely coupled with the content and can be easily removed**. Then we in turn utilize the restored results as guidance to improve image restoration. Specifically, in Sec. 3.2, we propose an **image-to-text mapper** \mathcal{M}_{i2t} and a **textual restoration module** \mathcal{M}_{clean} to extract the **degradation-free textual representations** \mathbf{E}'_{txt} from **degraded images** $\mathbf{X} \in \mathbb{R}^{H \times W \times 3}$. Then with a **pre-trained T2I model**, we can **generate a guidance image** $\mathbf{G} \in \mathbb{R}^{H \times W \times 3}$ that is **related to the content of \mathbf{X} but free of degradation**. In Sec. 3.3, to **leverage clean content information from \mathbf{G}** for **enhancing image restoration performance**, we introduce a **coarse-to-fine approach** that **dynamically incorporates multi-scale information from guidance \mathbf{G}** into restoration frameworks. Besides, we first give a preliminary knowledge of the T2I model in Sec. 3.1.

3.1. Preliminary

In this work, we employ **Stable Diffusion** [83] as our **text-to-image model**. Stable Diffusion pretrains an autoencoder $(\mathcal{E}(\cdot), \mathcal{D}(\cdot))$ to map the input image \mathbf{X} into a lower dimensional latent space by $\mathbf{z} = \mathcal{E}(\mathbf{X})$. The decoder $\mathcal{D}(\cdot)$ learns to map the latent code back to the image as $\mathcal{D}(\mathcal{E}(\mathbf{X})) \approx \mathbf{X}$. Then, the conditional diffusion model $\epsilon_\theta(\cdot)$ is trained on the latent space to generate latent codes based on text condition \mathbf{p} . To train the diffusion model, simple mean-squared loss is adopted as,

$$L_{LDM} = \mathbb{E}_{\mathbf{z} \sim \mathcal{E}(\mathbf{X}), \mathbf{p}, \epsilon \sim \mathcal{N}(0, 1), t} \left[\|\epsilon - \epsilon_\theta(\mathbf{z}_t, t, \tau_\theta^t(\mathbf{p}))\|_2^2 \right], \quad (1)$$

where ϵ denotes the unscaled noise, t is the timestep, \mathbf{z}_t is the latent noised to time t , and $\tau_\theta^t(\cdot)$ represents the pre-trained CLIP [76] text encoder. During inference, a random Gaussian noise \mathbf{z}_T is iteratively denoised to \mathbf{z}_0 , and the final image is obtained through the decoder $\mathbf{X}' = \mathcal{D}(\mathbf{z}_0)$.

3.2. Degradation-Free Guidance Generation

Image captioning models (e.g., **BLIP2** [53]) can generate text descriptions of degraded images, which can be taken into the text restoration module and text-to-image models to reconstruct clean guidance images. However, **such descriptions may lose too many details, resulting in severely inconsistent content between the guidance and the input images**. To address these problems, we propose an **image-to-text mapper** \mathcal{M}_{i2t} to encode the degraded image \mathbf{X} as im-

plicit textual representations \mathbf{E}_{txt} rather than explicit texts, which can be used to **reconstruct the content more faithfully**. Then, we utilize a **textual restoration module** \mathcal{M}_{clean} to **remove the degradation from \mathbf{E}_{txt} and obtain degradation-free textual representations** \mathbf{E}'_{txt} .

Image-to-Text Mapping. Following [96], we adopt the **textual word embedding space of CLIP** [76] as the **target space** in **image-to-text mapping**. Benefiting from the aligned vision-language feature space of CLIP, we first adopt **pre-trained CLIP image encoder** $\tau_\theta^i(\cdot)$ to encode input \mathbf{X} into CLIP image embedding \mathbf{E}_{img} . Then, we propose an **image-to-text mapper** \mathcal{M}_{i2t} to project CLIP image embedding into textual word embedding \mathbf{E}_{txt} ,

$$\mathbf{E}_{txt} = \mathcal{M}_{i2t}(\tau_\theta^i(\mathbf{X})), \quad (2)$$

where $\mathbf{E}_{txt} \in \mathbb{R}^{N \times D}$ and D is the dimension of textual word embedding. N is the number of learned words. To make the obtained word embedding describe more details of the input image, we set N to a large number (e.g., $N=20$).

Textual Restoration. Although \mathcal{M}_{i2t} can project images into the textual word embedding space, the projected textual word embedding will include corresponding degradation information if input \mathbf{X} is degraded. Condition on the degraded textual word embedding, the synthesized guidance inevitably reflects the corresponding degradation pattern. Therefore, we further propose a **textual restoration module** \mathcal{M}_{clean} to **remove the degradation information from the textual word embedding** \mathbf{E}_{txt} ,

$$\mathbf{E}'_{txt} = \mathcal{M}_{clean}(\mathbf{E}_{txt}), \quad (3)$$

where \mathbf{E}'_{txt} denote the **restored textual representations**. When feeding \mathbf{E}'_{txt} into **Stable Diffusion**, we can obtain a **degradation-free image** \mathbf{G} , which has **similar content as \mathbf{X} but is free of degradation**, as shown in Fig. 1.

Model Optimization. During training, the **image-to-text mapper** \mathcal{M}_{i2t} and **textual restoration module** \mathcal{M}_{clean} are optimized sequentially. In the **first training stage**, we collect both **clean and degraded images** from **different restoration tasks as training data**. We adopt Eq. (1) as a loss function, which **constrains the diffusion model** to reconstruct clean and degraded images conditioning **on their own projected embedding** \mathbf{E}_{txt} . In the **second** training stage, we use **pairs of degraded-clean images** from different image restoration datasets as training data. Based on pre-trained \mathcal{M}_{i2t} , we **further deploy** \mathcal{M}_{clean} to remove degradation-related concepts from degraded embedding \mathbf{E}_{txt} , obtaining \mathbf{E}'_{txt} . We still adopt Eq. (1), which **constrains the diffusion model** to reconstruct **clean images** conditioning **on restored embedding** \mathbf{E}'_{txt} . Note that, due to the ease of removing degradations in the textual space, it is feasible to train \mathcal{M}_{i2t} and \mathcal{M}_{clean} with multiple degradations simultaneously. When training is done, it can serve multiple image restoration tasks, **generating content-related and**

LQ 跟 HQ 不需要 pair, 利用圖像轉文字, 各自生成各自的文字再丟進 Diff 再用 loss 約束生成圖像的一致, train 的是 i2t mapper

不動 i2t · paired data · Diff loss 訓練 clean module

Table 1. **All-in-one image restoration results.** Following PromptIR [73], we train and evaluate the proposed method in all-in-one image restoration task, our method outperforms PromptIR across all the benchmark datasets.

Method	Dehazing	Derain	Denoise on BSD68			Average
	on SOTS	on Rain100L	$\sigma = 15$	$\sigma = 25$	$\sigma = 50$	
BRDNet [91]	23.23/0.895	27.42/0.895	32.26/0.898	29.74/0.836	26.34/0.836	27.80/0.843
LPNet [34]	20.84/0.828	24.88/0.784	26.47/0.778	24.77/0.748	21.26/0.552	23.64/0.738
FDGAN [33]	24.71/0.924	29.89/0.933	30.25/0.910	28.81/0.868	26.43/0.776	28.02/0.883
MPRNet [113]	25.28/0.954	33.57/0.954	33.54/0.927	30.89/0.880	27.56/0.779	30.17/0.899
DL [28]	26.92/0.391	32.62/0.931	33.05/0.914	30.41/0.861	26.90/0.740	29.98/0.875
AirNet [51]	27.94/0.962	34.90/0.967	33.92/0.933	31.26/0.888	28.00/0.797	31.20/0.910
PromptIR [73]	30.58/0.974	36.37/0.972	33.98/0.933	31.31/0.888	28.06/0.799	32.06/0.913
Ours	31.63/0.980	37.58/0.979	34.01/0.933	31.39/0.890	28.18/0.802	32.56/0.916

degradation-free guidance images \mathbf{G} from degraded images with one unified model (as shown in Fig. 1). More training details can be seen in the *Suppl.*

3.3. Guided Restoration

Although guidance image \mathbf{G} is clean, it is **inevitable that there are content differences between \mathbf{G} and degraded image \mathbf{X} .** Simply fusing the features of \mathbf{G} and \mathbf{X} is not enough to improve image restoration. Instead, we suggest a **dynamic aggregation module** to extract and exploit helpful information from \mathbf{G} .

Dynamic Aggregation. First, we perform feature matching [63] between the guidance \mathbf{G} and given degraded image \mathbf{X} . Specifically, we use a shared image encoder to extract multi-scale features for \mathbf{X} and \mathbf{G} , named \mathbf{F}_x and \mathbf{F}_g , respectively. According to the similarity score between \mathbf{F}_x and \mathbf{F}_g , we search for the most useful feature from \mathbf{F}_g for each small patch in \mathbf{F}_x in a coarse-to-fine manner. After searching, we get a set of guidance features $\hat{\mathbf{F}}_g$, whose content is aligned spatially with the input content. Second, we integrate the matched feature $\hat{\mathbf{F}}_g$ into the image restoration backbone. Both CNN-based and transformer-based restoration networks can be adopted. Without bells and whistles, we modulate the original input features as,

$$\mathbf{F}_x = \mathbf{F}_x + \alpha \cdot \mathcal{B}([\mathbf{F}_x, \hat{\mathbf{F}}_g]), \quad (4)$$

where $[\cdot, \cdot]$ represents concatenation operation, \mathcal{B} represents one CNN-based block or transformer-based block, α is the hyper-parameter to trade-off feature integration. Moreover, we perform it over multiple layers or levels. More details can be seen in the *Suppl.*

Loss Function. To train guided restoration framework, we simply adopt ℓ_1 loss as the reconstruction loss between predicted results $\hat{\mathbf{Y}}$ and target \mathbf{Y} , i.e.,

$$\ell_1 = \|\hat{\mathbf{Y}} - \mathbf{Y}\|_1. \quad (5)$$

4. Experiments

We evaluate the proposed method on five image restoration tasks, including **(1) all-in-one image restoration** [73], **(2)**

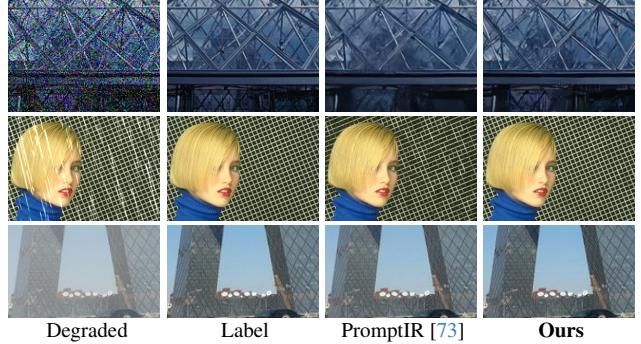


Figure 3. **All-in-one image restoration results.** Top: image dehazing, mid: image deraining, bottom: image dehazing.

image deblurring [15, 114] (*i.e.*, motion image deblurring and defocus image deblurring), **(3) image dehazing** [20], **(4) image deraining** [17], and **(5) image denoising** [114] (*i.e.*, grayscale and color image denoising on Gaussian noise, and real-world image denoising).

Implementation Details. We adopt publicly released **stable-diffusion-v2.1-base** as our T2I generative model. N is set to 20 in image-to-text mapper. Moreover, we employ the state-of-the-art methods as our image restoration backbone in different image restoration tasks: **(1)** we adopt PromptIR [73] in all-in-one image restoration, **(2)** we adopt NAFNet [15] and Restormer [114] in single-image motion deblurring, and Restormer [114] in defocus deblurring, **(3)** we adopt SFNet [20] in image dehazing, **(4)** we adopt DRSformer [17] in image deraining, **(5)** we adopt Restormer [114] in image denoising. The training settings of our guided restoration network are consistent with those of the corresponding backbone methods. Details of datasets, network architecture, and training hyperparameters are provided in *Suppl.*

4.1. All-in-One Restoration Results

Following [73], we adopt BSD400 [6] and WED [64] for color image denoising, Rain100L [28] for image deraining, and SOTS [50] for image dehazing. We train and evaluate the proposed method on these three tasks in all-in-one image restoration setting. The large difference among different degradations makes unifying multiple image restoration tasks into one unified network difficult, and thus the performance of existing all-in-one image restoration methods is limited. Compared with performing all-in-one restoration in the image level, it is much easier for us to conduct all-in-one restoration in the textual level. Once trained, our method can generate clean textual representation for different image restoration tasks, and thus we can synthesize guidance images for various types of degradation. Comparison results in Table 1 show that our method outperforms PromptIR [73] across all the benchmark datasets, achieving +1.05 dB PSNR on dehazing task, +1.21 dB PSNR on de-

Table 2. **Motion image deblurring results.** We train models with GoPro training data. We evaluate our method on GoPro, HIDE, RealBlur benchmark datasets. PSNR and SSIM scores are calculated on RGB-channels.

Method	GoPro [68]		HIDE [86]		RealBlur-R [81]		RealBlur-J [81]	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
DBGAN [121]	31.10	0.942	28.94	0.915	33.78	0.909	24.93	0.745
MT-RNN [70]	31.15	0.945	29.15	0.918	35.79	0.951	28.44	0.862
DMPHN [116]	31.20	0.940	29.09	0.924	35.70	0.948	28.42	0.860
SPAIR [74]	32.06	0.953	30.29	0.931	-	-	28.81	0.875
MIMO-Unet+ [19]	32.45	0.957	29.99	0.930	35.54	0.947	27.63	0.837
IPT [13]	32.52	-	-	-	-	-	-	-
MPRNet [113]	32.66	0.959	30.96	0.939	35.99	0.952	28.70	0.873
HINet [14]	32.71	0.959	30.32	0.932	-	-	-	-
Uformer [95]	32.97	0.967	-	-	-	-	-	-
Restormer [114]	32.92	0.961	31.22	0.942	36.19	0.957	28.96	0.879
Ours-Restormer	33.11	0.962	31.26	0.943	36.47	0.959	29.17	0.875
NAFNet [15]	33.69	0.966	31.32	0.943	33.62	0.944	26.33	0.856
Ours-NAFNet	33.97	0.968	31.57	0.946	33.87	0.950	26.76	0.861

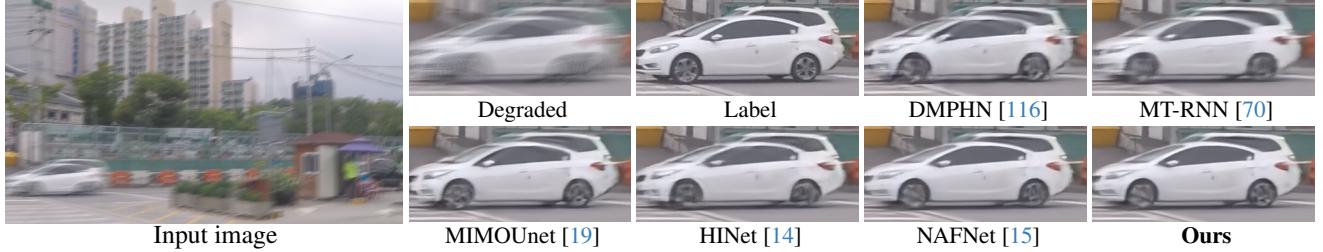


Figure 4. **Motion image deblurring results.** Compared with others, our method can predict clearer results with clearer boundaries.

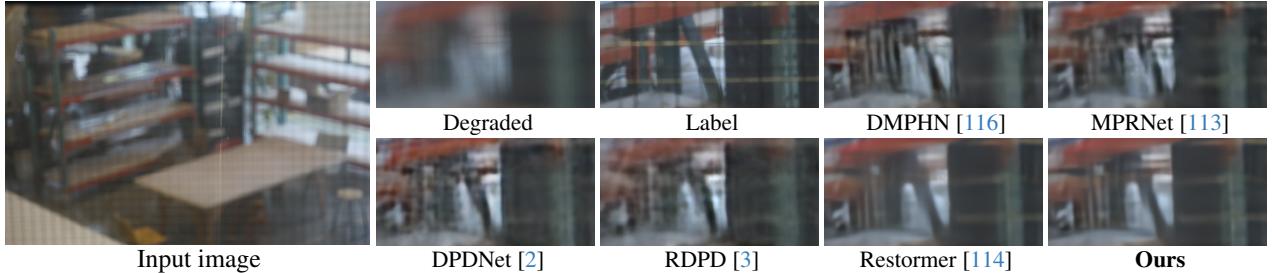


Figure 5. **Defocus image deblurring results.** Compared with others, our method can predict clearer results with clearer boundaries.

raining task, and +0.5 dB PSNR on average improvement. As shown in Fig. 3, the prediction of the proposed method is finer and more detailed compared with PromptIR.

4.2. Image Deblurring Results

Motion Deblurring. Following [15, 114], we train the proposed method on GoPro training data and evaluate our method on GoPro [68], HIDE [86], and real-world datasets (RealBlur-R [81] and RealBlur-J [81]). As shown in Table 2, **benefiting from synthesized high-quality guidance**, Ours-Restormer and Ours-NAFNet outperform the state-of-the-art methods on all four benchmark datasets. Our method achieves +0.28 dB PSNR improvement on GoPro testing data. The visual results in Fig. 4 show that our method can

Table 3. **Defocus image deblurring results.** We train and evaluate methods on DPDD dataset [2]. S denotes single-image defocus deblurring model. D denotes dual-pixel defocus deblurring. PSNR and SSIM scores are calculated on RGB channels.

Method	Indoor Scenes		Outdoor Scenes		Combined	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
EBDB _S [44]	25.77	0.772	21.25	0.599	23.45	0.683
DMENets _S [46]	25.50	0.788	21.43	0.644	23.41	0.714
JNB _S [87]	26.73	0.828	21.10	0.608	23.84	0.715
DPDNet _S [2]	26.54	0.816	22.25	0.682	24.34	0.747
KPAC _S [88]	27.97	0.852	22.62	0.701	25.22	0.774
IFAN _S [47]	28.11	0.861	22.76	0.720	25.37	0.789
Restormer _S [114]	28.87	0.882	23.24	0.743	25.98	0.811
Ours_S	29.11	0.889	23.35	0.748	26.15	0.817
DPDNet _D [2]	27.48	0.849	22.90	0.726	25.13	0.786
RDPD _D [3]	28.10	0.843	22.82	0.704	25.39	0.772
Uformer _D [95]	28.23	0.860	23.10	0.728	25.65	0.795
IFAN _D [47]	28.66	0.868	23.46	0.743	25.99	0.804
Restormer _D [114]	29.48	0.895	23.97	0.773	26.66	0.833
Ours_D	29.62	0.899	24.16	0.775	26.82	0.835

restore images with **sharper boundaries and details**, demonstrating its effectiveness.

Defocus Deblurring. Following [114], we train and evaluate the proposed method on DPDD [2] dataset for single-image defocus deblurring and dual-pixel defocus deblurring. From Table 3, compared with Restormer [114], our method achieves +0.17 dB PSNR and +0.16 dB PSNR gains on single-image defocus deblurring and dual-pixel defocus deblurring, respectively. As illustrated in Fig. 5, the prediction of our method shows a **clearer structure**.

4.3. Image Dehazing Results

For image dehazing task, we conduct experiments on the synthetic benchmark RESIDE [50] dataset. We train the

Table 4. **Image dehazing results.** We separately train and evaluate our method indoor scene and outdoor scene. PSNR and SSIM scores are calculated on RGB-channels.

Method	SOTS-Indoor [50]		SOTS-Outdoor [50]	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑
DehazeNet [9]	19.82	0.821	24.75	0.927
AOD-Net [48]	20.51	0.861	24.14	0.920
GridDehazeNet [61]	32.16	0.984	30.86	0.982
MSBDN [26]	33.67	0.985	33.48	0.982
FFA-Net [75]	36.39	0.989	33.57	0.984
ACER-Net [97]	37.17	0.990	-	-
DeHamer [37]	36.63	0.988	35.18	0.986
MAXIM-2S [92]	38.11	0.991	34.19	0.985
PMNet [105]	38.41	0.990	34.74	0.985
DehazeFormer-L [90]	40.05	0.996	-	-
SFNet [20]	41.24	0.996	40.05	0.996
Ours	41.48	0.996	40.29	0.996

Table 6. **Grayscale image denoising on Gaussian noise.**

Upper-bracket: models are trained on a range of noise levels.
Lower-bracket: models are trained on the fixed noise level.

Method	Set12 [118]			BSD68 [65]			Urban100 [40]		
	$\sigma=15$	$\sigma=25$	$\sigma=50$	$\sigma=15$	$\sigma=25$	$\sigma=50$	$\sigma=15$	$\sigma=25$	$\sigma=50$
DnCNN [118]	32.67	30.35	27.18	31.62	29.16	26.23	32.28	29.80	26.35
FFDNet [120]	32.75	30.43	27.32	31.63	29.19	26.29	32.40	29.90	26.50
IRCNN [119]	32.76	30.37	27.12	31.63	29.15	26.19	32.46	29.80	26.22
DRUNet [122]	33.25	30.94	27.90	31.91	29.48	26.59	33.44	31.11	27.96
Restormer [114]	33.35	31.04	28.01	31.95	29.51	26.62	33.67	31.39	28.33
Ours	33.35	31.30	28.13	31.98	29.58	26.77	33.62	31.47	28.46
FOCNet [41]	33.07	30.73	27.68	31.83	29.38	26.50	33.15	30.64	27.40
MWCNN [60]	33.15	30.79	27.74	31.86	29.41	26.53	33.17	30.66	27.42
NLRN [59]	33.16	30.80	27.64	31.88	29.41	26.47	33.45	30.94	27.49
RNAN [125]	-	-	27.70	-	-	26.48	-	-	27.65
DeamNet [79]	33.19	30.81	27.74	31.91	29.44	26.54	33.37	30.85	27.53
DAGL [67]	33.28	30.93	27.81	31.93	29.46	26.51	33.79	31.39	27.97
SwinIR [57]	33.36	31.01	27.91	31.97	29.50	26.58	33.70	31.30	27.98
Restormer [114]	33.42	31.08	28.00	31.96	29.52	26.62	33.79	31.46	28.29
Ours	33.47	31.15	28.12	31.92	29.67	26.78	33.78	31.58	28.38



Figure 6. **Color image denoising results on Gaussian noise.** Compared with others, the proposed method can obtain finer results.

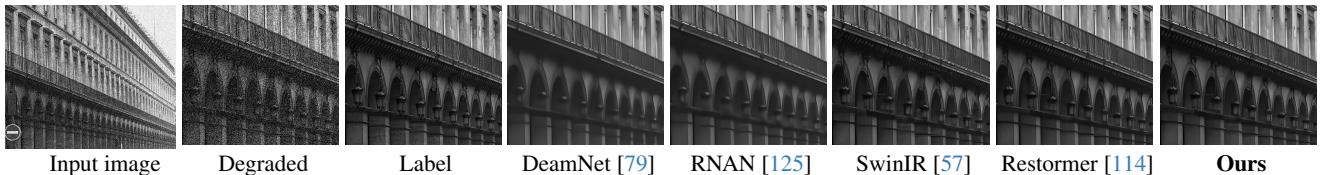


Figure 7. **Grayscale image denoising results on Gaussian noise.** Our method can restore detailed results especially in texture regions.

proposed method on indoor scene and outdoor scene data, and evaluate the performance on SOTS-indoor and SOTS-outdoor testing dataset [50] separately. Table 4 reports the quantitative comparison results. One can see that our

Table 5. **Image deraining results.** We separately train and evaluate our method on Rain200H, Rain200L, DID-Data, and DDN-Data. PSNR and SSIM scores are calculated on Y channel in YCbCr color space.

Method	Rain200L [104]		Rain200H [104]		DID-Data [115]		DDN-Data [30]	
	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑	PSNR↑	SSIM↑
DDN [29]	34.68	0.967	26.05	0.805	30.97	0.911	30.00	0.904
RESCAN [54]	36.09	0.967	26.75	0.835	33.38	0.941	31.94	0.935
PReNet [80]	37.80	0.981	29.04	0.899	33.17	0.948	32.60	0.946
MSPFN [42]	38.53	0.983	29.36	0.903	33.72	0.955	32.99	0.933
RCDNet [93]	39.17	0.989	30.24	0.904	34.08	0.953	33.04	0.947
MPRNet [113]	39.47	0.982	30.67	0.911	33.99	0.959	33.10	0.935
DualGCN [31]	40.73	0.989	31.15	0.912	34.37	0.962	33.01	0.949
SPDNet [106]	40.50	0.988	31.28	0.920	34.57	0.956	33.15	0.946
Uformer [95]	40.20	0.986	30.80	0.910	35.02	0.962	33.95	0.955
Restormer [114]	40.99	0.989	32.00	0.932	35.29	0.964	34.20	0.957
IDT [100]	40.74	0.988	32.10	0.934	34.89	0.962	33.84	0.955
DRSformer [17]	41.21	0.989	32.16	0.933	35.24	0.962	34.23	0.955
Ours	41.59	0.990	31.97	0.931	35.46	0.964	34.57	0.958

Table 7. **Color image denoising on Gaussian noise.** Upper-bracket: models are trained on a range of noise levels. Lower-bracket: models are trained on the fixed noise level. PSNR is calculated on RGB channels.

Method	CBSD68 [66]			Kodak24			McMaster [123]			Urban100 [40]		
	$\sigma=15$	$\sigma=25$	$\sigma=50$	$\sigma=15$	$\sigma=25$	$\sigma=50$	$\sigma=15$	$\sigma=25$	$\sigma=50$	$\sigma=15$	$\sigma=25$	$\sigma=50$
IRCNN [119]	33.86	31.16	27.86	34.69	32.18	28.93	34.58	32.18	28.91	33.78	31.20	27.70
FFDNet [120]	33.87	31.21	27.96	34.63	32.13	28.98	34.66	32.35	29.18	33.83	31.40	28.05
DnCNN [118]	33.90	31.24	27.95	34.60	32.14	28.95	33.45	31.52	28.62	32.98	30.81	27.59
DSNet [72]	33.91	31.28	28.05	34.63	32.16	29.05	34.67	32.40	29.28	-	-	-
DRUNet [122]	34.30	31.69	28.51	35.31	32.89	29.86	35.40	33.14	30.08	34.81	32.60	29.61
Restormer [114]	34.39	31.78	28.59	35.44	33.02	30.00	35.55	33.31	30.29	35.06	32.91	30.02
Ours	34.37	31.87	28.68	35.52	33.13	30.15	35.62	33.38	30.40	35.03	32.97	30.19
RPCNN [99]	-	31.24	28.06	-	32.34	29.25	-	32.33	29.33	-	31.81	28.62
BRDNet [91]	34.10	31.43	28.16	34.88	32.41	29.22	35.08	32.75	29.52	34.42	31.99	28.56
RNAN [125]	-	-	28.27	-	-	29.58	-	-	29.72	-	-	29.08
RDN [126]	-	-	28.31	-	-	29.66	-	-	-	-	-	29.38
IPT [13]	-	-	28.39	-	-	29.64	-	-	29.98	-	-	29.71
SwinIR [57]	34.42	31.78	28.56	35.34	32.89	29.79	35.61	33.20	30.22	35.13	32.90	29.82
Restormer [114]	34.40	31.79	28.60	35.47	33.04	30.01	35.61	33.34	30.30	35.13	32.96	30.02
Ours	34.48	31.97	28.83	35.58	33.21	30.23	35.75	33.56	30.46	35.11	33.13	30.27

Table 8. **Real-world image denoising results.** We train and evaluate our method on SIDD [1] datasets. * denotes methods using additional training data. PSNR and SSIM scores are calculated on RGB channels.

Dataset	Method	DnCNN [118]	BM3D [21]	CBDNet* [38]	RIDNet* [5]	AINDNet* [45]	VDN [108]	SADNet* [12]	DANet+* [109]	CycleISP* [110]	MIRNet [111]	DeamNet* [79]	MPRNet [113]	DAGL [67]	Uformer [95]	Restormer [114]	Ours
SIDD [1]	PSNR \uparrow	23.66	25.65	30.78	38.71	39.08	39.28	39.46	39.47	39.52	39.72	39.47	39.71	38.94	39.77	40.02	40.09
	SSIM \uparrow	0.583	0.685	0.801	0.951	0.954	0.956	0.957	0.957	0.957	0.959	0.957	0.958	0.953	0.959	0.960	0.960

Table 9. Effect of condition information.

Method	baseline	$N=5$	$N=10$	$N=20$	$N=30$	$N=40$
PSNR \uparrow	30.16	31.13	31.36	31.57	31.51	31.60
SSIM \uparrow	0.932	0.941	0.945	0.947	0.947	0.948

Table 10. Effect of integration strategy.

Method	baseline	Enc.	Dec.	Enc. & Dec.
PSNR \uparrow	30.16	31.37	30.31	31.57
SSIM \uparrow	0.932	0.946	0.934	0.947

Table 11. Effect of generated guidance.

Method	baseline	Degr.	Ours
PSNR \uparrow	30.16	30.13	31.57
SSIM \uparrow	0.932	0.931	0.947

4.4. Image Deraining Results

Following [17], we train and evaluate the proposed method separately on Rain200L [104], Rain200H [104], DID-Data [115], and DDN-Data [30] datasets. From Table 5, our method achieves +0.31 dB PSNR average improvement on Rain200L, DID-Data, and DDN-Data against DRS-former [17]. On the Rain200H dataset, we only have comparable performance, as its severe rain streak may interfere with the dynamic aggregation process between the guidance information and the degraded images. Visual comparisons of deraining results will be provided in the *Suppl.*

4.5. Image Denoising Results

Image Denoising on Gaussian Noise. Following [114], we train the blind denoising model (σ with range of [0, 50]) and non-blind denoising models ($\sigma=15, 25, 50$) for grayscale and color image denoising on Gaussian noise. We train the proposed method on DFBW dataset with synthetic noise degradation. For grayscale image denoising, we evaluate our method on Set12 [118], BSD68 [66], and Urban100 [40]. For color image denoising, we evaluate our method on CBSD68 [66], Kodak24, McMaster [123], and Urban100 [40]. As shown in Table 6 and Table 7, our method achieves comparable performance with baseline on low noise-level ($\sigma=15$), while **our method outperforms baseline when noise becomes heavier**, e.g., achieving +0.25 dB PSNR gain on Urban100 [40] with Gaussian color noise $\sigma=50$). Qualitative comparisons in Fig. 6 and Fig. 7 show that our method can restore finer texture and sharper boundaries while other methods tend to generate smoother results.

Real-World Image Denoising. Furthermore, We train and evaluate the proposed method on real-world denoising dataset SIDD [1]. From Table 8, although Restormer [114] has already achieved superior performance on real-world denoising task, our method can also bring +0.07 dB PSNR improvement, demonstrating its effectiveness. Visual comparison results can be found in *Suppl.*

4.6. Ablation Studies

In ablation studies, we mainly conduct discussion from three perspectives: (1) the effect of the **textural representation size**, (2) the effect of **guidance integration strategy** and (3) the effect of **generated guidance**. During experiments, we adopt NAFNet with 32 channels as our **baseline**. We train each method on GoPro datasets for 30k iterations and evaluate the performance on GoPro testing data.

Effect of Textural Representation Size. In our method, we use the number of words N to control the representation ability of \mathbf{E}_{txt} , which is used as condition information to generate guidance. We conduct the experiments on five settings, i.e., $N = 5, 10, 20, 30, 40$. As shown in Table 9, the restoration performance **generally improves from $N=5$ to $N=40$** , as more informative textual descriptions can help T2I model generate more faithful guidance. Since the performance **only improves slightly when N is larger than 20**, we finally **set N to 20** to make a **trade-off** between performance and computational cost.

Effect of Integrating Strategy. To evaluate the effect of different guidance integrating strategies, we conduct experiments by **integrating guidance information into different modules of image restoration networks**, i.e., **Enc.** (encoder only), **Dec.** (decoder only), and **Enc. & Dec.** (both encoder and decoder). As shown in Table 10, the experimental results show that **integrating guidance information into both encoder and decoder achieves better performance**.

Effect of Generated Guidance. To assess the effect of generated guidance, we conduct an experiment by **replacing it with degraded input (named ‘Degr.’)**. From Table 11, when replacing the generated guidance with degraded input, it **cannot obtain better results compared with baseline**, which shows dynamically aggregating the degraded image itself cannot help the image restoration. In contrast, the performance can improve effectively with our guidance.

5. Conclusion

In this paper, we provide a new perspective for image restoration. Considering that **content and degradation are tightly coupled in image representation**, while **decoupled in textual representation**. Rather than direct restoration on image level, we suggest conducting **restoration on textual level and in turn utilizing the restored textual content to assist image restoration**. To this end, we propose a plug-

and-play approach that first encodes degraded images into degraded textual representations and then removes textual degradation information to obtain restored textual representations. Conditioned on the representations, we employ a pre-trained T2I model to synthesize clean guidance images for improving image restoration. We evaluate our method on various image restoration tasks. The experimental results demonstrate it outperforms the state-of-the-art ones.

References

- [1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S. Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, pages 1692–1700, 2018. [8](#), [15](#)
- [2] Abdullah Abuolaim and Michael S. Brown. Defocus deblurring using dual-pixel data. In *ECCV*, pages 111–126, 2020. [6](#), [14](#), [31](#), [32](#)
- [3] Abdullah Abuolaim, Mauricio Delbracio, Damien Kelly, Michael S. Brown, and Peyman Milanfar. Learning to reduce defocus blur by realistically modeling dual-pixel data. In *ICCV*, pages 2269–2278, 2021. [6](#), [31](#), [32](#)
- [4] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *CVPRW*, 2017. [15](#)
- [5] Saeed Anwar and Nick Barnes. Real image denoising with feature attention. In *ICCV*, pages 3155–3164, 2019. [8](#)
- [6] Pablo Arbelaez, Michael Maire, Charless C. Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(5):898–916, 2011. [5](#), [14](#), [15](#)
- [7] Yunpeng Bai, Cairong Wang, Shuzhao Xie, Chao Dong, Chun Yuan, and Zhi Wang. Textir: A simple framework for text-based editable image restoration. *arXiv preprint arXiv:2302.14736*, 2023. [3](#)
- [8] James Betker, Gabriel Goh, Li Jing, † TimBrooks, Jianfeng Wang, Linjie Li, † LongOuyang, † JuntangZhuang, † JoyceLee, † YufeiGuo, † WesamManassra, † PrafullaDhariwal, † CaseyChu, † YunxinJiao, and Aditya Ramesh. Improving image generation with better captions, 2023. <https://cdn.openai.com/papers/dall-e-3.pdf>
- [9] Bolun Cai, Xiangmin Xu, Kui Jia, Chunmei Qing, and Dacheng Tao. Dehazenet: An end-to-end system for single image haze removal. *IEEE Trans. Image Process.*, 25(11):5187–5198, 2016. [1](#), [7](#)
- [10] Jiezhang Cao, Jingyun Liang, Kai Zhang, Yawei Li, Yulun Zhang, Wenguan Wang, and Luc Van Gool. Reference-based image super-resolution with deformable attention transformer. In *ECCV*, pages 325–342, 2022. [3](#)
- [11] Huiwen Chang, Han Zhang, Jarred Barber, AJ Maschinot, Jose Lezama, Lu Jiang, Ming-Hsuan Yang, Kevin Murphy, William T Freeman, Michael Rubinstein, et al. Muse: Text-to-image generation via masked generative transformers. *arXiv preprint arXiv:2301.00704*, 2023. [3](#)
- [12] Meng Chang, Qi Li, Huajun Feng, and Zhihai Xu. Spatial-adaptive network for single image denoising. In *ECCV*, pages 171–187, 2020. [8](#)
- [13] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *CVPR*, pages 12299–12310, 2021. [3](#), [6](#), [7](#)
- [14] Liangyu Chen, Xin Lu, Jie Zhang, Xiaojie Chu, and Chengpeng Chen. Hinet: Half instance normalization network for image restoration. In *CVPRW*, pages 182–192, 2021. [1](#), [6](#), [29](#), [30](#)
- [15] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *ECCV*, pages 17–33, 2022. [1](#), [5](#), [6](#), [14](#), [29](#), [30](#)
- [16] Wei-Ting Chen, Zhi-Kai Huang, Cheng-Che Tsai, Hao-Hsiang Yang, Jian-Jiun Ding, and Sy-Yen Kuo. Learning multiple adverse weather removal via two-stage knowledge learning and multi-contrastive regularization: Toward a unified model. In *CVPR*, pages 17632–17641, 2022. [1](#), [3](#)
- [17] Xiang Chen, Hao Li, Mingqiang Li, and Jinshan Pan. Learning A sparse transformer network for effective image deraining. In *CVPR*, pages 5896–5905, 2023. [1](#), [5](#), [7](#), [8](#), [14](#), [15](#), [34](#), [35](#), [36](#)
- [18] Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer. In *CVPR*, pages 22367–22377, 2023. [3](#)
- [19] Sung-Jin Cho, Seo-Won Ji, Jun-Pyo Hong, Seung-Won Jung, and Sung-Jea Ko. Rethinking coarse-to-fine approach in single image deblurring. In *ICCV*, pages 4621–4630, 2021. [1](#), [6](#), [29](#), [30](#)
- [20] Yuning Cui, Yi Tao, Zhenshan Bing, Wenqi Ren, Xinwei Gao, Xiaochun Cao, Kai Huang, and Alois Knoll. Selective frequency network for image restoration. In *ICLR*, 2023. [1](#), [5](#), [7](#), [14](#), [15](#), [33](#)
- [21] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen O. Egiazarian. Image denoising by sparse 3-d transform-domain collaborative filtering. *IEEE Trans. Image Process.*, 16(8):2080–2095, 2007. [8](#)
- [22] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *NeurIPS*, 34:19822–19835, 2021. [3](#)
- [23] Ming Ding, Wendi Zheng, Wenyi Hong, and Jie Tang. Cogview2: Faster and better text-to-image generation via hierarchical transformers. *arXiv preprint arXiv:2204.14217*, 2022. [3](#)
- [24] Yuxuan Ding, Chunna Tian, Haoxuan Ding, and Lingqiao Liu. The clip model is secretly an image-to-prompt converter. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. [2](#)
- [25] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoxu Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, pages 184–199, 2014. [2](#)
- [26] Hang Dong, Jinshan Pan, Lei Xiang, Zhe Hu, Xinyi Zhang, Fei Wang, and Ming-Hsuan Yang. Multi-scale boosted dehazing network with dense feature fusion. In *CVPR*, pages 2154–2164, 2020. [7](#)
- [27] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner,

- Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [28] Qingnan Fan, Dongdong Chen, Lu Yuan, Gang Hua, Nenghai Yu, and Baoquan Chen. A general decoupled learning framework for parameterized image operators. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(1):33–47, 2021. 5
- [29] Xueyang Fu, Jiaxin Huang, Delu Zeng, Yue Huang, Xinghao Ding, and John W. Paisley. Removing rain from single images via a deep detail network. In *CVPR*, pages 1715–1723, 2017. 7
- [30] Xueyang Fu, Jiaxin Huang, Delu Zeng, Yue Huang, Xinghao Ding, and John W. Paisley. Removing rain from single images via a deep detail network. In *CVPR*, pages 1715–1723, 2017. 7, 8, 15
- [31] Xueyang Fu, Qi Qi, Zheng-Jun Zha, Yurui Zhu, and Xinghao Ding. Rain streak removal via dual graph convolutional network. In *AAAI*, pages 1352–1360, 2021. 7
- [32] Oran Gafni, Adam Polyak, Oron Ashual, Shelly Sheynin, Devi Parikh, and Yaniv Taigman. Make-a-scene: Scene-based text-to-image generation with human priors. In *ECCV*, pages 89–106, 2022. 3
- [33] Hongyun Gao, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Dynamic scene deblurring with parameter selective sharing and nested skip connections. In *CVPR*, pages 3848–3856, 2019. 5
- [34] Hongyun Gao, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Dynamic scene deblurring with parameter selective sharing and nested skip connections. In *CVPR*, pages 3848–3856, 2019. 5
- [35] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 27, 2014. 3
- [36] Shuhang Gu, Yawei Li, Luc Van Gool, and Radu Timofte. Self-guided network for fast image denoising. In *ICCV*, pages 2511–2520, 2019. 2
- [37] Chunle Guo, Qixin Yan, Saeed Anwar, Runmin Cong, Wenqi Ren, and Chongyi Li. Image dehazing transformer with transmission-aware 3d position embedding. In *CVPR*, pages 5802–5810, 2022. 1, 7, 33
- [38] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *CVPR*, pages 1712–1722, 2019. 8
- [39] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. 3
- [40] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *CVPR*, pages 5197–5206, 2015. 7, 8, 37
- [41] Xixi Jia, Sanyang Liu, Xiangchu Feng, and Lei Zhang. Focnet: A fractional optimal control network for image denoising. In *CVPR*, pages 6054–6063, 2019. 7
- [42] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Baojin Huang, Yimin Luo, Jiayi Ma, and Junjun Jiang. Multi-scale progressive fusion network for single image deraining. In *CVPR*, pages 8343–8352, 2020. 7
- [43] Yuming Jiang, Kelvin C. K. Chan, Xintao Wang, Chen Change Loy, and Ziwei Liu. Robust reference-based super-resolution via c2-matching. In *CVPR*, pages 2103–2112, 2021. 3
- [44] Ali Karaali and Cláudio R. Jung. Edge-based defocus blur estimation with adaptive scale selection. *IEEE Trans. Image Process.*, 27(3):1126–1137, 2018. 6
- [45] Yoosik Kim, Jae Woong Soh, Gu Yong Park, and Nam Ik Cho. Transfer learning from synthetic to real-noise denoising with adaptive instance normalization. In *CVPR*, pages 3479–3489, 2020. 8
- [46] Junyoung Lee, Sungkil Lee, Sunghyun Cho, and Seungyong Lee. Deep defocus map estimation using domain adaptation. In *CVPR*, pages 12222–12230, 2019. 6
- [47] Junyoung Lee, Hyeongseok Son, Jaesung Rim, Sunghyun Cho, and Seungyong Lee. Iterative filter adaptive network for single image defocus deblurring. In *CVPR*, pages 2034–2042, 2021. 6
- [48] Boyi Li, Xiulian Peng, Zhangyang Wang, Jizheng Xu, and Dan Feng. Aod-net: All-in-one dehazing network. In *ICCV*, pages 4780–4788, 2017. 7
- [49] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip Torr. Controllable text-to-image generation. *NeurIPS*, 32, 2019. 3
- [50] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. Benchmarking single-image dehazing and beyond. *IEEE Trans. Image Process.*, 28(1):492–505, 2019. 5, 6, 7, 14, 27, 33
- [51] Boyun Li, Xiao Liu, Peng Hu, Zhongqin Wu, Jiancheng Lv, and Xi Peng. All-in-one image restoration for unknown corruption. In *CVPR*, pages 17431–17441, 2022. 1, 3, 5
- [52] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, pages 12888–12900, 2022. 1
- [53] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, pages 19730–19742, 2023. 1, 4, 16
- [54] Xia Li, Jianlong Wu, Zhouchen Lin, Hong Liu, and Hongbin Zha. Recurrent squeeze-and-excitation context aggregation net for single image deraining. In *ECCV*, pages 262–277, 2018. 7, 34, 35, 36
- [55] Xia Li, Jianlong Wu, Zhouchen Lin, Hong Liu, and Hongbin Zha. Recurrent squeeze-and-excitation context aggregation net for single image deraining. In *ECCV*, pages 262–277, 2018. 2
- [56] Yawei Li, Kai Zhang, Jingyun Liang, Jiezhang Cao, Ce Liu, Rui Gong, Yulun Zhang, Hao Tang, Yun Liu, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. LSDIR: A large scale dataset for image restoration. In *CVPRW*, pages 1775–1787, 2023. 15
- [57] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *Int. Conf. Comput. Vis. Worksh.*, pages 1833–1844, 2021. 1, 3, 7, 37

- [58] Xinqi Lin, Jingwen He, Ziyan Chen, Zhaoyang Lyu, Ben Fei, Bo Dai, Wanli Ouyang, Yu Qiao, and Chao Dong. Diffbir: Towards blind image restoration with generative diffusion prior. *arXiv preprint arXiv:2308.15070*, 2023. 3
- [59] Ding Liu, Bihan Wen, Yuchen Fan, Chen Change Loy, and Thomas S. Huang. Non-local recurrent network for image restoration. In *NeurIPS*, pages 1680–1689, 2018. 2, 7
- [60] Pengju Liu, Hongzhi Zhang, Kai Zhang, Liang Lin, and Wangmeng Zuo. Multi-level wavelet-cnn for image restoration. In *CVPRW*, pages 773–782, 2018. 7
- [61] Xiaohong Liu, Yongrui Ma, Zhihao Shi, and Jun Chen. Griddehazenet: Attention-based multi-scale network for image dehazing. In *ICCV*, pages 7313–7322. IEEE, 2019. 7
- [62] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 2
- [63] Liying Lu, Wenbo Li, Xin Tao, Jiangbo Lu, and Jiaya Jia. MASA-SR: matching acceleration and spatial adaptation for reference-based image super-resolution. In *CVPR*, pages 6368–6377, 2021. 3, 5, 15
- [64] Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei Zhang. Waterloo exploration database: New challenges for image quality assessment models. *IEEE Trans. Image Process.*, 26(2):1004–1016, 2017. 5, 14, 15
- [65] David R. Martin, Charless C. Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, pages 416–425, 2001. 7
- [66] David R. Martin, Charless C. Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, pages 416–425, 2001. 7, 8, 28
- [67] Chong Mou, Jian Zhang, and Zhuoyuan Wu. Dynamic attentive graph learning for image restoration. In *ICCV*, pages 4308–4317, 2021. 7, 8
- [68] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, pages 257–265, 2017. 6, 14, 29, 30
- [69] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3
- [70] Dongwon Park, Dong Un Kang, Jisoo Kim, and Se Young Chun. Multi-temporal recurrent neural networks for progressive non-uniform single image deblurring with incremental temporal training. In *ECCV*, pages 327–343, 2020. 6, 29, 30
- [71] Dongwon Park, Byung Hyun Lee, and Se Young Chun. All-in-one image restoration for unknown degradations using adaptive discriminative filters for specific degradations. In *CVPR*, pages 5815–5824, 2023. 1, 3
- [72] Yali Peng, Lu Zhang, Shigang Liu, Xiaojun Wu, Yu Zhang, and Xili Wang. Dilated residual networks with symmetric skip connection for image denoising. *Neurocomputing*, 345:67–76, 2019. 7
- [73] Vaishnav Potlapalli, Syed Waqas Zamir, Salman Khan, and Fahad Shahbaz Khan. Promptir: Prompting for all-in-one blind image restoration. In *NeurIPS*, 2023. 1, 2, 3, 5, 14, 25, 26, 27, 28
- [74] Kuldeep Purohit, Maitreya Suin, A. N. Rajagopalan, and Vishnu Naresh Boddeti. Spatially-adaptive image restoration using distortion-guided networks. In *ICCV*, pages 2289–2299, 2021. 6
- [75] Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia. Ffa-net: Feature fusion attention network for single image dehazing. In *AAAI*, pages 11908–11915, 2020. 7
- [76] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. 2, 4
- [77] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3
- [78] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 3, 4
- [79] Chao Ren, Xiaohai He, Chuncheng Wang, and Zhibo Zhao. Adaptive consistency prior based deep network for image denoising. In *CVPR*, pages 8596–8606, 2021. 7, 8
- [80] Dongwei Ren, Wangmeng Zuo, Qinghua Hu, Pengfei Zhu, and Deyu Meng. Progressive image deraining networks: A better and simpler baseline. In *CVPR*, pages 3937–3946, 2019. 7, 34, 35, 36
- [81] Jaesung Rim, Haeyun Lee, Juchol Won, and Sunghyun Cho. Real-world blur dataset for learning and benchmarking deblurring algorithms. In *ECCV*, pages 184–201, 2020. 6, 14
- [82] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 4
- [83] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10674–10685, 2022. 2, 4
- [84] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2, 4

- [85] Axel Sauer, Tero Karras, Samuli Laine, Andreas Geiger, and Timo Aila. Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. *arXiv preprint arXiv:2301.09515*, 2023. 3
- [86] Ziyi Shen, Wenguan Wang, Xiankai Lu, Jianbing Shen, Haibin Ling, Tingfa Xu, and Ling Shao. Human-aware motion deblurring. In *ICCV*, pages 5571–5580, 2019. 6, 14
- [87] Jianping Shi, Li Xu, and Jiaya Jia. Just noticeable defocus blur detection and estimation. In *CVPR*, pages 657–665, 2015. 6
- [88] Hyeongseok Son, Junyoung Lee, Sunghyun Cho, and Seungyong Lee. Single image defocus deblurring using kernel-sharing parallel atrous convolutions. In *ICCV*, pages 2622–2630, 2021. 6
- [89] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 3
- [90] Yuda Song, Zhuqing He, Hui Qian, and Xin Du. Vision transformers for single image dehazing. *IEEE Trans. Image Process.*, 32:1927–1941, 2023. 1, 7, 33
- [91] Chunwei Tian, Yong Xu, and Wangmeng Zuo. Image denoising using deep CNN with batch renormalization. *Neural Networks*, 121:461–473, 2020. 5, 7
- [92] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan C. Bovik, and Yinxiao Li. MAXIM: multi-axis MLP for image processing. In *CVPR*, pages 5759–5770, 2022. 1, 7
- [93] Hong Wang, Qi Xie, Qian Zhao, and Deyu Meng. A model-driven deep neural network for single image rain removal. In *CVPR*, pages 3100–3109, 2020. 7
- [94] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin C. K. Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. *arXiv preprint arXiv:2305.07015*, 2023. 3
- [95] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *CVPR*, pages 17662–17672, 2022. 1, 2, 6, 7, 8, 34, 35, 36
- [96] Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. ELITE: encoding visual concepts into textual embeddings for customized text-to-image generation. In *ICCV*, 2023. 4
- [97] Haiyan Wu, Yanyun Qu, Shaohui Lin, Jian Zhou, Ruizhi Qiao, Zhizhong Zhang, Yuan Xie, and Lizhuang Ma. Contrastive learning for compact single image dehazing. In *CVPR*, pages 10551–10560, 2021. 7
- [98] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2256–2265, 2021. 3
- [99] Zhihao Xia and Ayan Chakrabarti. Identifying recurring patterns with deep neural networks for natural image denoising. In *WACV*, pages 2415–2423, 2020. 7
- [100] Jie Xiao, Xueyang Fu, Aiping Liu, Feng Wu, and Zheng-Jun Zha. Image de-raining transformer. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(11):12978–12995, 2023. 1, 7
- [101] Liangbin Xie, Xintao Wang, Chao Dong, Zhongang Qi, and Ying Shan. Finding discriminative filters for specific degradations in blind super-resolution. In *NeurIPS*, pages 51–61, 2021. 3
- [102] Qinhong Yang, Dongdong Chen, Zhentao Tan, Qiankun Liu, Qi Chu, Jianmin Bao, Lu Yuan, Gang Hua, and Nenghai Yu. HQ-50K: A large-scale, high-quality dataset for image restoration. *arXiv preprint arXiv:2306.05390*, 2023. 15
- [103] Tao Yang, Peiran Ren, Xuansong Xie, and Lei Zhang. Pixel-aware stable diffusion for realistic image super-resolution and personalized stylization. *arXiv preprint arXiv:2308.14469*, 2023. 3
- [104] Wenhan Yang, Robby T. Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan. Deep joint rain detection and removal from a single image. In *CVPR*, pages 1685–1694, 2017. 7, 8, 14, 15, 25, 26
- [105] Tian Ye, Yunchen Zhang, Mingchao Jiang, Liang Chen, Yun Liu, Sixiang Chen, and Erkang Chen. Perceiving and modeling density for image dehazing. In *ECCV*, pages 130–145. Springer, 2022. 7
- [106] Qiaosi Yi, Juncheng Li, Qinyan Dai, Faming Fang, Guixu Zhang, and Tieyong Zeng. Structure-preserving deraining with residue channel prior guidance. In *ICCV*, pages 4218–4227, 2021. 7
- [107] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 3
- [108] Zongsheng Yue, Hongwei Yong, Qian Zhao, Deyu Meng, and Lei Zhang. Variational denoising network: Toward blind noise modeling and removal. In *NeurIPS*, pages 1688–1699, 2019. 8
- [109] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. Dual adversarial network: Toward real-world noise removal and noise generation. In *ECCV*, pages 41–58, 2020. 8
- [110] Syed Waqas Zamir, Aditya Arora, Salman H. Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Cycleisp: Real image restoration via improved data synthesis. In *CVPR*, pages 2693–2702, 2020. 8
- [111] Syed Waqas Zamir, Aditya Arora, Salman H. Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In *ECCV*, pages 492–511, 2020. 8
- [112] Syed Waqas Zamir, Aditya Arora, Salman H. Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In *ECCV*, pages 492–511, 2020. 2
- [113] Syed Waqas Zamir, Aditya Arora, Salman H. Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *CVPR*, pages 14821–14831, 2021. 1, 2, 5, 6, 7, 8, 31, 32, 34, 35, 36
- [114] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang.

- Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, pages 5718–5729, 2022. [1](#), [3](#), [5](#), [6](#), [7](#), [8](#), [14](#), [15](#), [31](#), [32](#), [37](#)
- [115] He Zhang and Vishal M. Patel. Density-aware single image de-raining using a multi-stream dense network. In *CVPR*, pages 695–704, 2018. [7](#), [8](#), [15](#), [34](#), [35](#), [36](#)
- [116] Hongguang Zhang, Yuchao Dai, Hongdong Li, and Piotr Koniusz. Deep stacked hierarchical multi-patch network for image deblurring. In *CVPR*, pages 5978–5986, 2019. [6](#), [29](#), [30](#), [31](#), [32](#)
- [117] Jinghao Zhang, Jie Huang, Mingde Yao, Zizheng Yang, Hu Yu, Man Zhou, and Feng Zhao. Ingredient-oriented multi-degradation learning for image restoration. In *CVPR*, pages 5825–5835, 2023. [1](#), [3](#)
- [118] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep CNN for image denoising. *IEEE Trans. Image Process.*, 26(7):3142–3155, 2017. [1](#), [2](#), [7](#), [8](#)
- [119] Kai Zhang, Wangmeng Zuo, Shuhang Gu, and Lei Zhang. Learning deep CNN denoiser prior for image restoration. In *CVPR*, pages 2808–2817, 2017. [7](#)
- [120] Kai Zhang, Wangmeng Zuo, and Lei Zhang. Ffdnet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE Trans. Image Process.*, 27(9):4608–4622, 2018. [1](#), [7](#)
- [121] Kaihao Zhang, Wenhan Luo, Yiran Zhong, Lin Ma, Björn Stenger, Wei Liu, and Hongdong Li. Deblurring by realistic blurring. In *CVPR*, pages 2734–2743, 2020. [6](#)
- [122] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(10):6360–6376, 2022. [1](#), [7](#)
- [123] Lei Zhang, Xiaolin Wu, Antoni Buades, and Xin Li. Color demosaicking by local directional interpolation and non-local adaptive thresholding. *J. Electronic Imaging*, 20(2):023016, 2011. [7](#), [8](#)
- [124] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *ECCV*, pages 294–310, 2018. [2](#)
- [125] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. In *ICLR*, 2019. [2](#), [7](#), [37](#)
- [126] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image restoration. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(7):2480–2495, 2021. [7](#)
- [127] Yurui Zhu, Tianyu Wang, Xueyang Fu, Xuanyu Yang, Xin Guo, Jifeng Dai, Yu Qiao, and Xiaowei Hu. Learning weather-general and weather-specific features for image restoration under multiple adverse weather conditions. In *CVPR*, pages 21747–21758, 2023. [1](#), [3](#)

Improving Image Restoration through Removing Degradations in Textual Representations

Supplementary Material

The content of the supplementary material involves:

- Experimental details in Sec. A.
- Effect of textual restoration in Sec. B.
- Explicit textual representation v.s. implicit textual representation in Sec. C.
- Guidance visualization in Sec. D.
- More visual comparisons in Sec. E.

A. Experimental Details

In this section, we list experimental details for different image restoration task (*i.e.*, all-in-one image restoration [73], image deblurring [15, 114], image dehazing [20], image deraining [17], and image denoising [114]), the proposed degradation-free guidance generation process (*i.e.*, image-to-text mapping and textual restoration), and guided-restoration.

All-in-One Image Restoration. We adopt **PromptIR** [73] as our **backbone** in **all-in-one image restoration**. Following [73], network has 8 stages (the first 7 stages as main network, the last stage as refinement), the number of blocks for each stages is [4, 6, 6, 8, 6, 6, 4, 4], network width is 48, the number of heads for each stages is [1, 2, 4, 8, 4, 2, 1, 1]. In perspective of training data, we adopt concatenation of 400 images from BSD [6] and 4,744 images from WED [64] dataset as denoising training data, 200 images from Rain100L [104] for deraining task, 72,135 images from SOTS for dehazing task. Considering dataset size gap among different tasks, we properly enlarge deraining data and denoising data as [73]. To train, we adopt AdamW optimizer with CosineAnnealing learning rate scheduler, the initial learning rate of the main restoration network and dynamic aggregation is set to 2e-4 and 1e-4, respectively. We train all-in-one image restoration on 4 Tesla-V100 GPUs with training patch size 128, batch size 48. Performance reported on Table 1 is referred to [73]. PSNR and SSIM scores are calculated on RGB channels, except which of deraining task are calculated on Y-channel in YCbCr color space.

Image Deblurring. We adopt **NAFNet** [15] as our backbone in **single-image motion deblurring**, **Restormer** [114] as backbone in **defocus deblurring**. In single-image motion deblurring task, we follow [15], main restoration network has 9 stages, and the number of blocks for each stage is [1, 1, 1, 28, 1, 1, 1, 1, 1], network width is 64. We

adopt GoPro [68] as our training data and directly evaluated the trained model on GoPro validation set, HIDE [86] testing set, and Realblur [81] dataset. GoPro dataset has 4,214 blur-sharp paris of data (2,103 for training and 1,111 for validation), testset of HIDE dataset consist of 2,025 images, and two subsets of Realblur both have 980 images. To train single-image motion deblurring, we adopt AdamW optimizer with CosineAnnealing learning rate scheduler, the initial learning rate of the main restoration network and dynamic aggregation module is 1e-4 and 5e-5, respectively. We train single-image motion deblurring on 8 Tesla-V100 GPUs with training patch size of 256, batch size of 16. In defocus deblurring task, we follow [114], main restoration network has 8 stages (the last stage as refinement stage), the number of blocks of each stage is [4, 6, 6, 8, 6, 6, 4, 4], network width is 48, the number of heads of each stage is [1, 2, 4, 8, 4, 2, 1, 1]. Note that input channels is 6 in dual-pixel defocus deblurring, the output channels is 3 in both single-image defocus deblurring and dual-pixel defocus deblurring. We adopt DPDD [2] dataset as our training data. DPDD dataset contains 500 indoor & outdoor scenes captured by DSLR camera. Each scene includes three defocus input images and a corresponding all-in-focus ground-truth image. Three input images are labeled as left, right and center views. The left and right defocused sub-aperture views are acquired with a wide camera aperture setting, and the corresponding all-in-focus ground-truth image captured with a narrow aperture. Following Restormer [114], we use sub-aperture data to train dual-pixel defocus deblurring, and we use center input image to train single-image defocus deblurring. To perform evaluation, we separately evaluate trained model in indoor & outdoor scene testing data and the average performance is calculated by weighted combination. To train defocus deblurring, we maintain progressive learning in official implementation, we adopt AdamW optimizer with CosineAnnealing learning rate scheduler, the initial learning rate of main restoration network and dynamic aggregation module is set to 3e-4 and 1e-4, respectively. Performance reported in Table 2 and Table 3 is referred to [15] and [114]. PSNR and SSIM scores are calculated on RGB channels.

Image Dehazing. We adopt **SFNet** [20] as our **backbone** in **image dehazing**. Following [20], network width is set to 32, the number of resblocks is 16. We train image dehazing on RESIDE [50] dataset, we train and evaluate method separately on indoor scene and outdoor scene data. To train image dehazing, we use Adam optimizer with CosineAn-

nealing learning rate scheduler, the initial learning rate of main network and dynamic aggregation is set to 1e-4 and 5e-5, respectively. Performance in Table 4 is referred to [20]. PSNR and SSIM scores are calculated on RGB channels.

Image Deraining. We adopt DRSformer [17] as our backbone in image deraining. Following [17], we adopt 7 stages for main restoration network, the number of blocks for each stage is [4, 6, 6, 8, 6, 6, 4], network width is set to 48, the number of heads for each stage is [1, 2, 4, 8, 4, 2, 1]. For training, we separately train and evaluate the proposed method on four datasets with synthetic rainstreak degradation, including Rain200L [104], Rain200H [104], DID-Data [115], and DDN-Data [30]. Rain200H and Rain200L dataset contain 1,800 pairs of rainy-clean images for training and 200 pairs images for testing. In DID-Data and DDN-Data, the synthetic rainstreak has different directions and different levels. DID-Data contains 12,000 pairs of images for training and 1,200 pairs of images for testing. DDN-Data contains 12,600 pairs of images for training and 1,400 images for testing. We employ MEFC [17] module for Rain200H, DID-Data, and DDN-Data. During training, we adopt AdamW optimizer with CosineAnnealing learning rate scheduler, patch size and batch size is set to 128 and 16 with 4 Tesla-V100 GPUs, the initial learning rate of main network and dynamic aggregation is set to 1e-4 and 5e-5, respectively. Performance in Table 5 is referred to [17]. PSNR and SSIM scores are calculated on Y channel in YCbCr color space.

Image Denoising. We adopt Restormer [114] as our backbone in image denoising. Following [114], we employ the bias-free network with 8 stages, the number of blocks for each stage is [4, 6, 6, 8, 6, 6, 4, 4], network width is set to 48, the number of heads of each stage is [1, 2, 4, 8, 4, 2, 1, 1]. We adopt concatenation data of Div2k [4] (800 images for training), Flickr2k (2,650 images for training), BSD [6] (400 images for training), and WED [64] (4,744 images for training) to train Gaussian grayscale denoising and Gaussian color denoising. We adopt 320 high-resolution images in SIDD [1] dataset for real-world denoising. During training, we also maintain the progressive learning strategy, we adopt AdamW optimizer with CosinAnnealing learning rate scheduler, the initial learning rate of main restoration network and dynamic aggregation is set to 3e-4 and 1e-4, respectively. Performance in Table 6, Table 7, and Table 8 is referred to [114]. PSNR and SSIM scores are calculated on RGB channels.

Image-to-Text Mapping. To enable our image-to-text mapping network can project both clean images and degraded images into textual space, we use the collection of high-quality data, degraded data from different image restoration tasks as our training data. High-quality data includes LSDIR [56] dataset and HQ-50K [102] dataset, LS-

DIR dataset contains 84,991 high-quality images for training, HQ-50K dataset contains 50,000 high-quality images for training. Degraded data includes GoPro, RESIDE, Rain200H, Rain200L, DID-Data, DDN-Data, and DFBW data with synthetic Gaussian noise. During training, we crop high-quality high-resolution data (LSDIR and HQ-50K) into 512×512 as input, for others we centerly crop images along shorter side and resize them to 512×512 as input. The mapping network is implemented as four-layer MLP network, and we adopt $N=20$ words to control representation capability of textual word embedding. To encode image concepts into textual space, feature from the last layer of CLIP image encoder is selected as input to image-to-text mapping network. The learning rate is set to 1e-6 and batch size is set to 4.

Textual Restoration. To train textual restoration network, we use concatenation of training dataset used in different image restoration tasks as our training data. We adopt the same strategy to preprocess pairs of degraded-clean data to 512×512 patches as input. The same with image-to-text mapping network, the textual restoration network is also implemented by four-layer MLP network. The learning rate is set to 1e-6 and training batch size is set to 4. During guidance generation, we use 200 steps of DDIM scheduler with scale of 5.

Guided-Restoration. Following [63], the dynamic aggregation includes two steps: feature matching and feature aggregation. In feature matching, we adopt a shared n -stages encoder to extract multi-scale feature from degraded input and clean guidance, n depends on total downsampling ratio of the main restoration network, each stage is with 4 residual blocks, the width of encoder is the same to the width of main network. We then adopt a coarse-to-fine manner to match useful information for each patch of degraded input, e.g., we first match in coarse block level then match in fine patch level. In coarse matching, feature block size is set to 8, dilation ratio is set to [1, 2, 3]. In fine matching, patch size is set to 3. For feature aggregation, we employ a more general way. We simply use concatenation & residual/self-attention blocks with adaptive scaling factor α to fuse guidance information to main restoration network, i.e., Eq. (4).

B. Effect of Textual Restoration

In this section, we demonstrate the effectiveness of our textual restoration. We discard textual restoration and directly use the output of image-to-text mapping network as conditional input for diffusion model, and the visual results of the synthetic guidance images are shown in Fig. A, denoted as w/o. textual restoration.

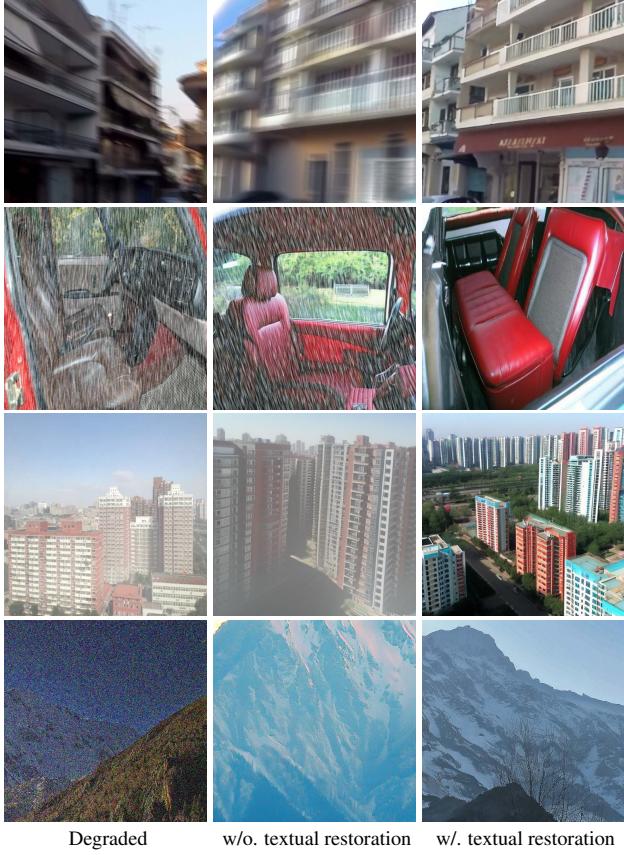


Figure A. Visual comparison of w/o. textual restoration and w/. textual restoration.

C. Explicit Textual Representation v.s. Implicit Textual Representation

In this section, we compare synthetic guidance images conditioned on explicit text representation and implicit textual representation, 1) **explicit text representation**: we first convert degraded images into image caption by BLIPv2 [53], then we manually discarding degradation-related text in image caption, finally we use the processed image caption as text prompt input to StableDiffusion to get synthetic guidance images. Denoted as **Explicit**. 2) **implicit textual representation**: our method, which is denoted as **Ours**. As shown in Fig. B, Fig. C, Fig. D, and Fig. E, we illustrate visual comparison for image deblurring, image deraining, image dehazing, and image denoising tasks. We can found though explicit text representation can describe content of degraded image properly, the synthetic results cannot maintain style, details and texture of original content. And in image denoising task, explicitly converting degraded noise image into image caption usually leads to wrong captions and thus cannot provide useful guidance image for restoration.

D. Guidance Visualization

We illustrate synthesized guidance images for each image restoration: image deblurring shows in Fig. F, image deraining shows in Fig. G, image dehazing shows in Fig. H, image denoising shows in Fig. I.

E. More Visual Comparisons

We provide visual comparison for different image restoration tasks:

- All-in-one image restoration: image deraining results show in Fig. J and Fig. K, image dehazing results show in Fig. L, image denoising results show in Fig. M.
- Image deblurring results: single-image motion deblurring results show in Fig. N and Fig. O, defocus deblurring results show in Fig. P and Fig. Q.
- Image dehazing results: Fig. R.
- Image deraining results: Fig. S, Fig. T, and Fig. U.
- Image denoising results: Fig. V.



Figure B. Visual comparison of synthetic guidance by explicit and implicit textual representation on image deblurring task.



Figure C. Visual comparison of synthetic guidance by explicit and implicit textual representation on image deraining task.



Figure D. Visual comparison of synthetic guidance by explicit and implicit textual representation on image dehazing task.



Figure E. Visual comparison of synthetic guidance by explicit and implicit textual representation on image denoising task.

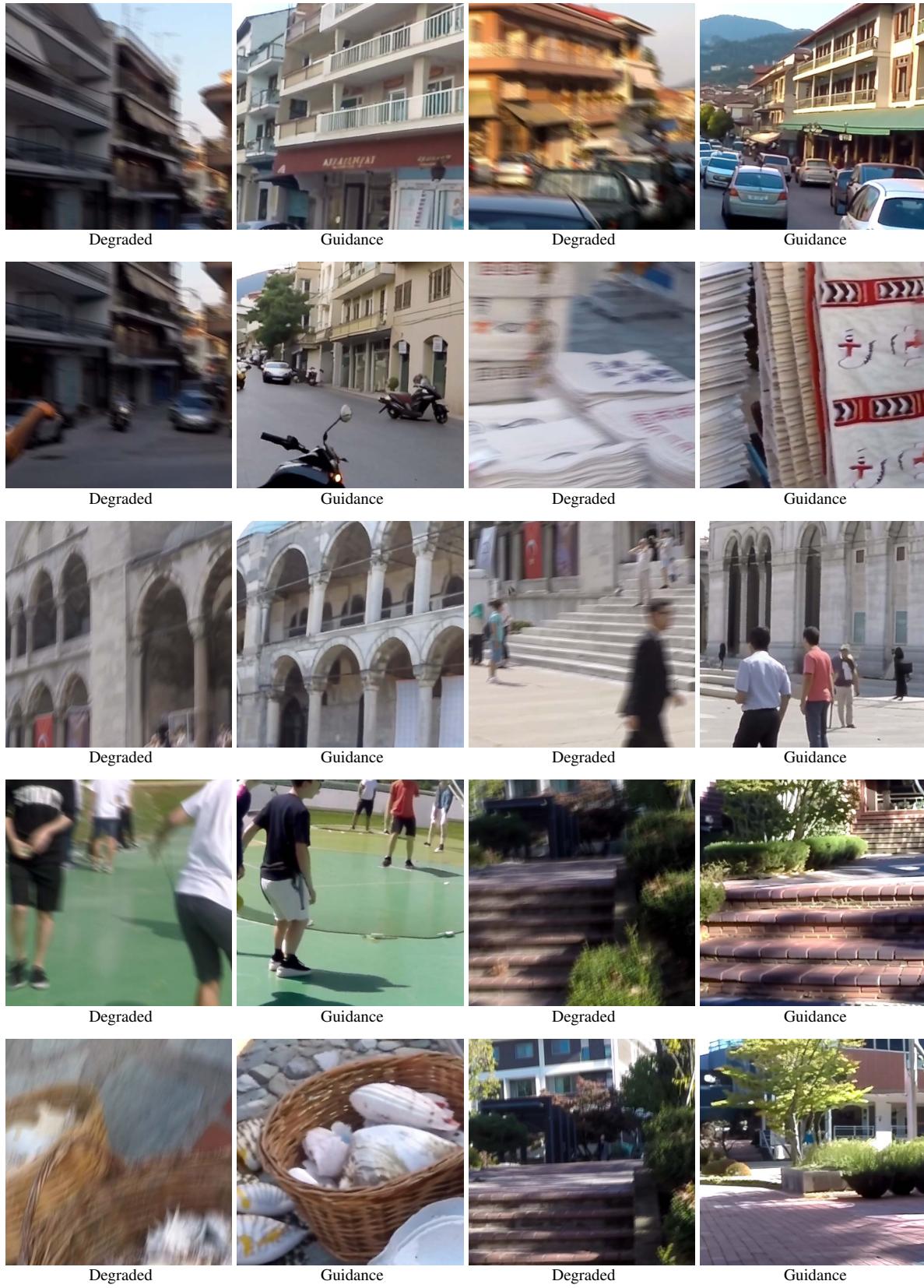


Figure F. Illustration of guidance images for image deblurring task.

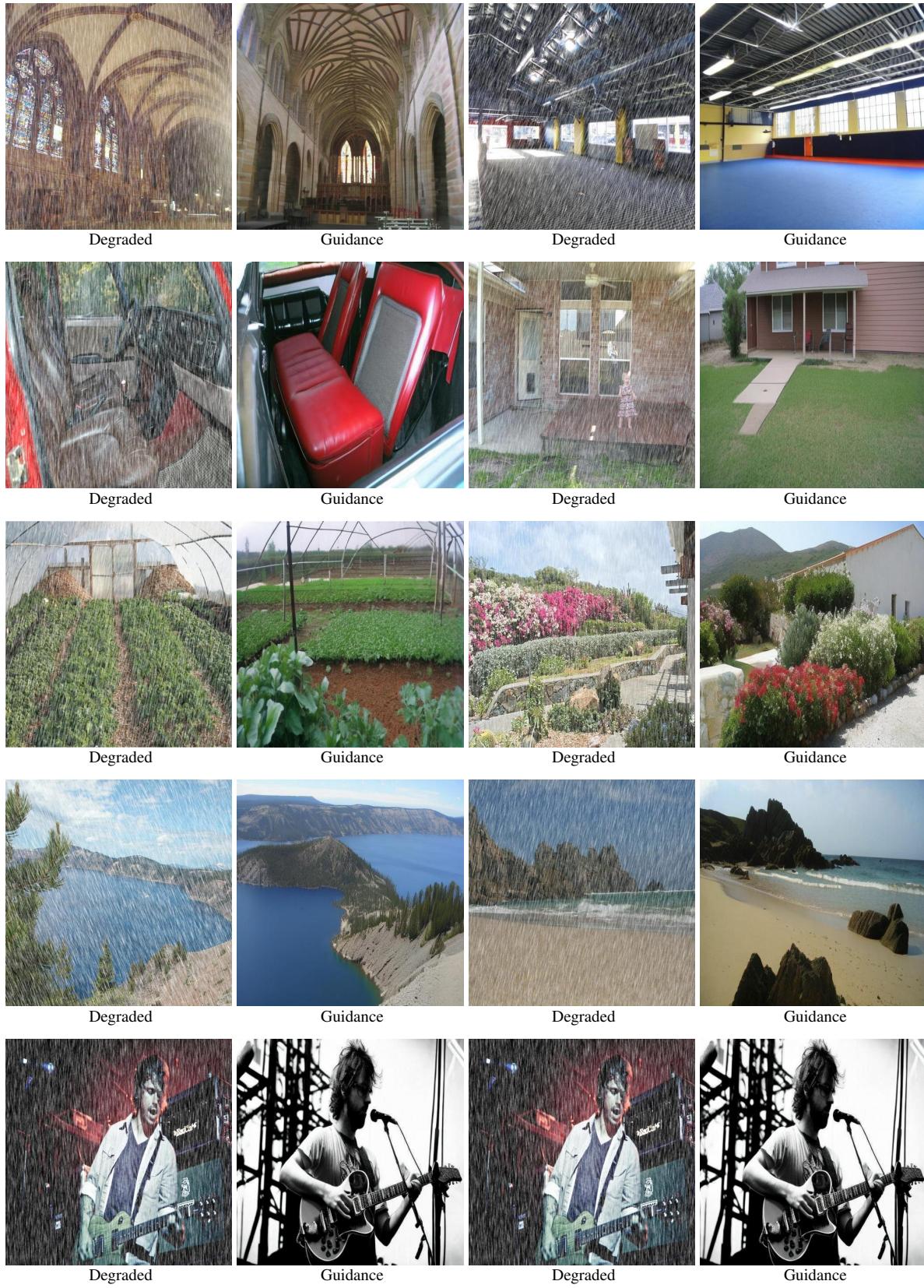


Figure G. Illustration of guidance images for image deraining task.

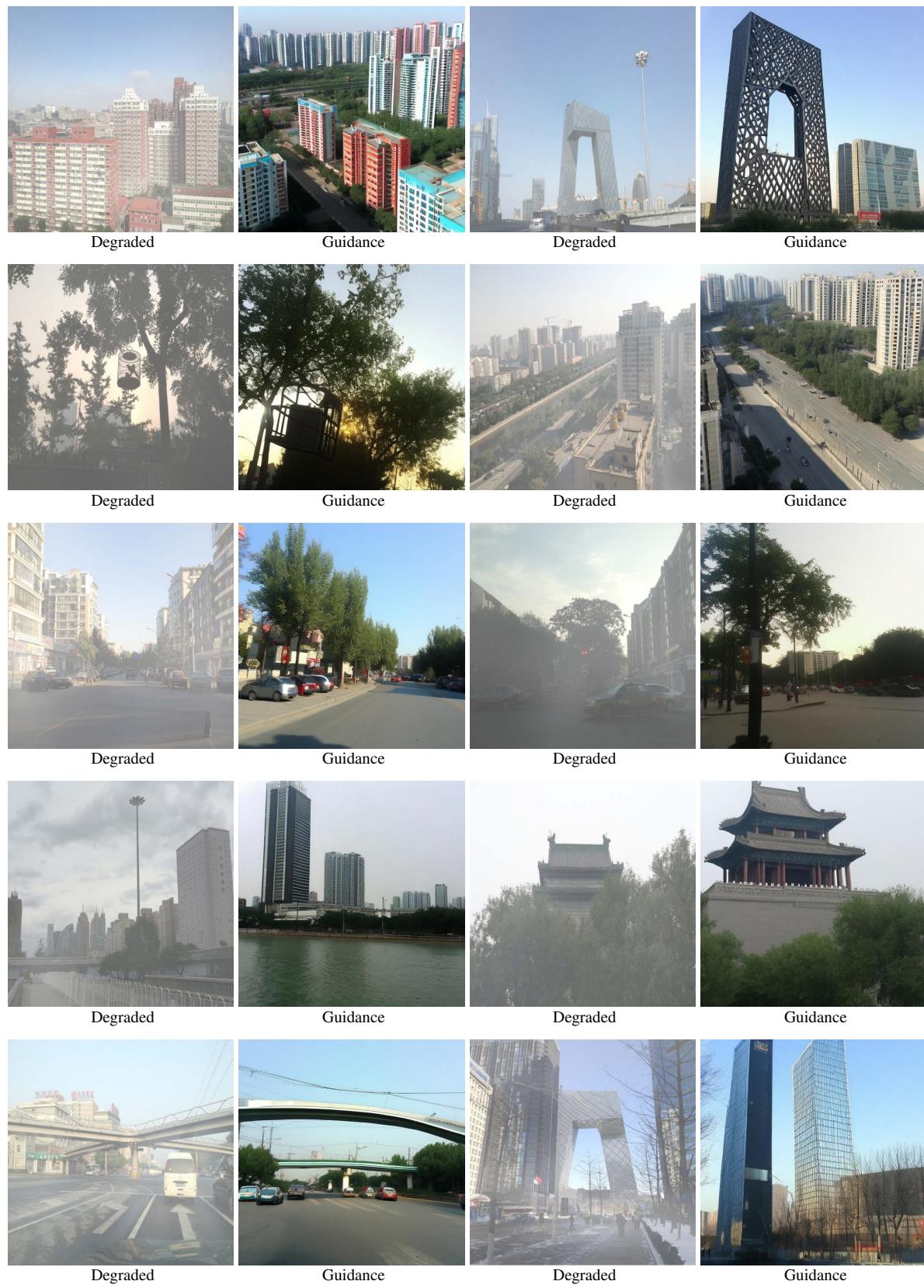


Figure H. Illustration of guidance images for image dehazing task.

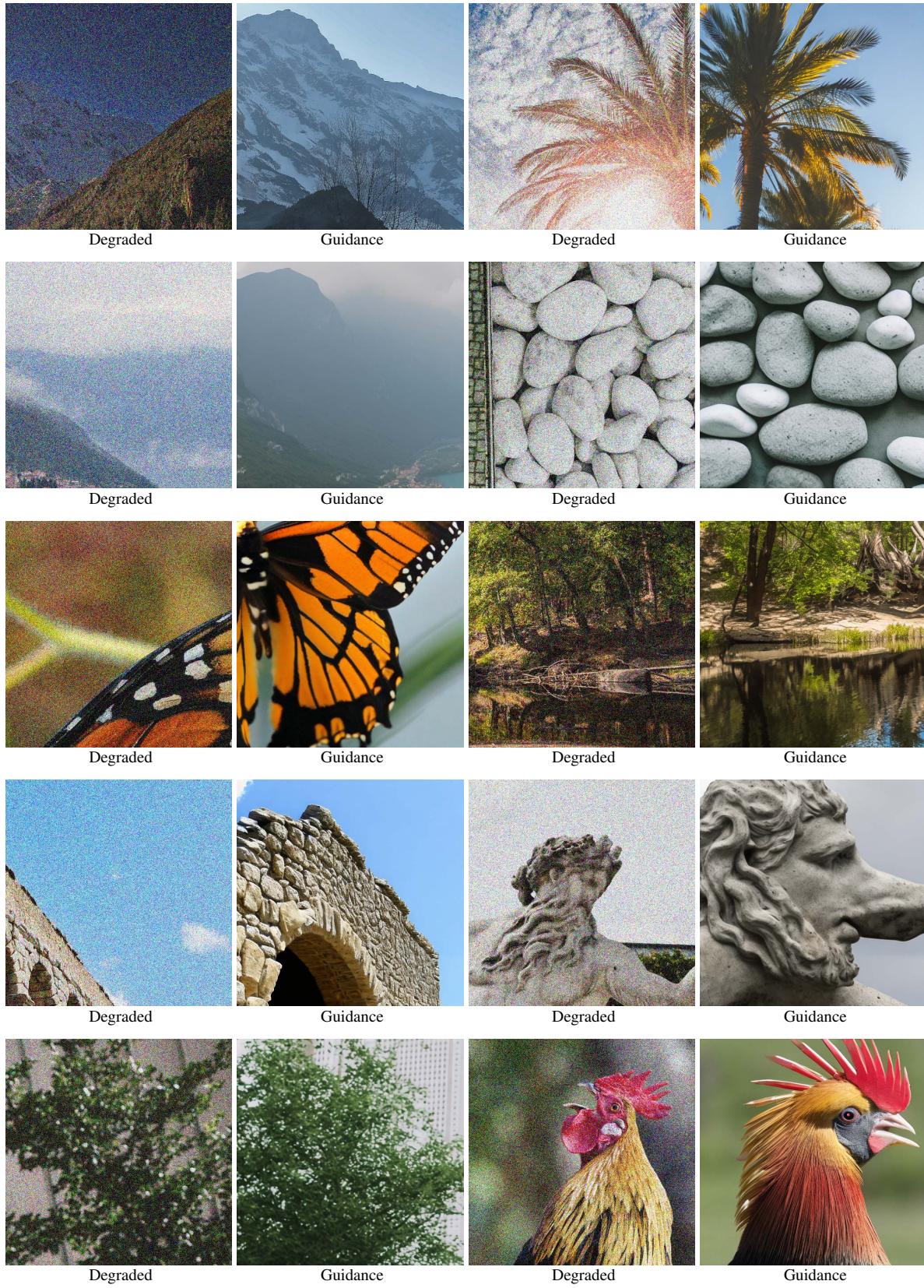


Figure I. Illustration of guidance images for image denoising task.

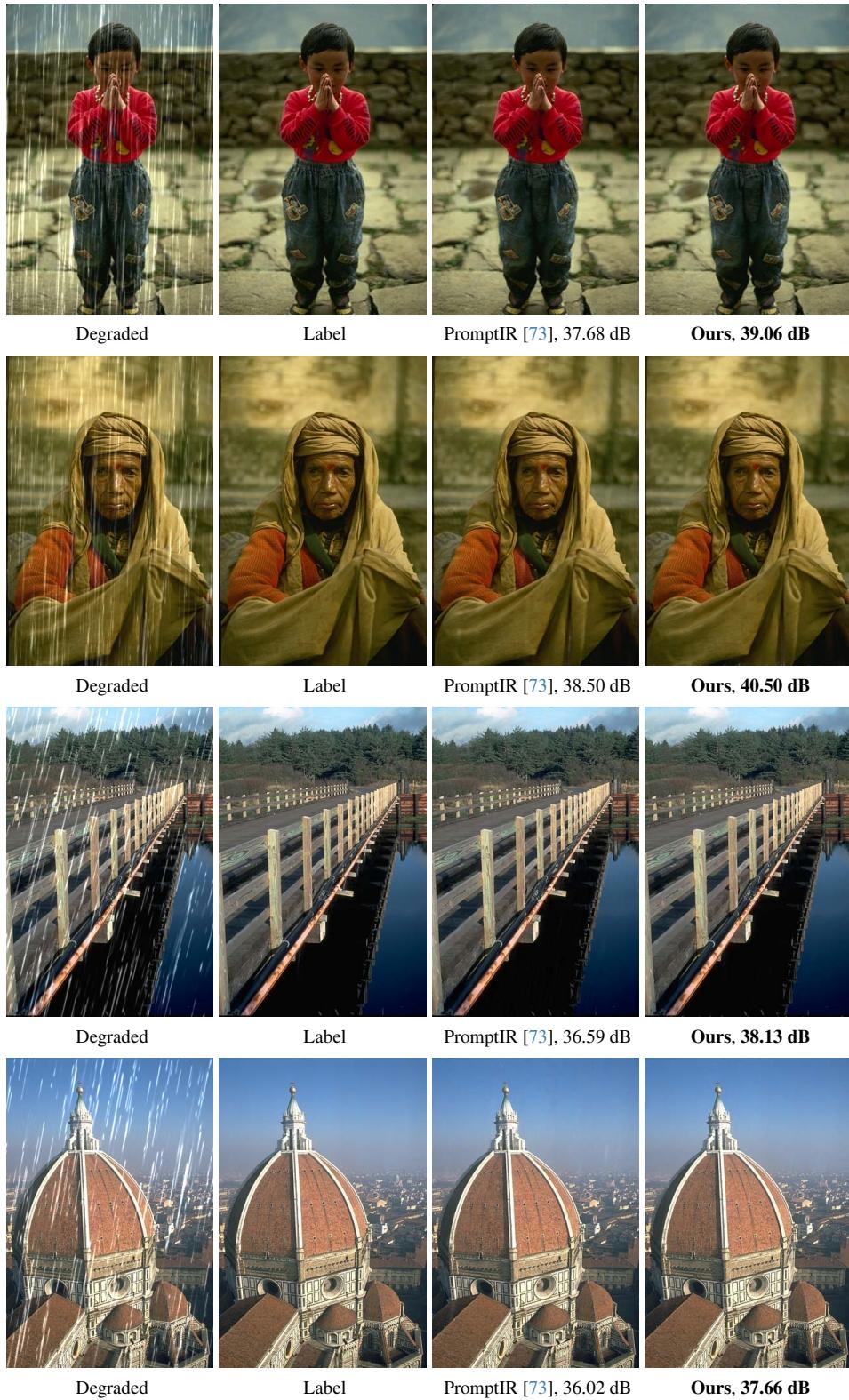


Figure J. Image Deraining on Rain100L [104].

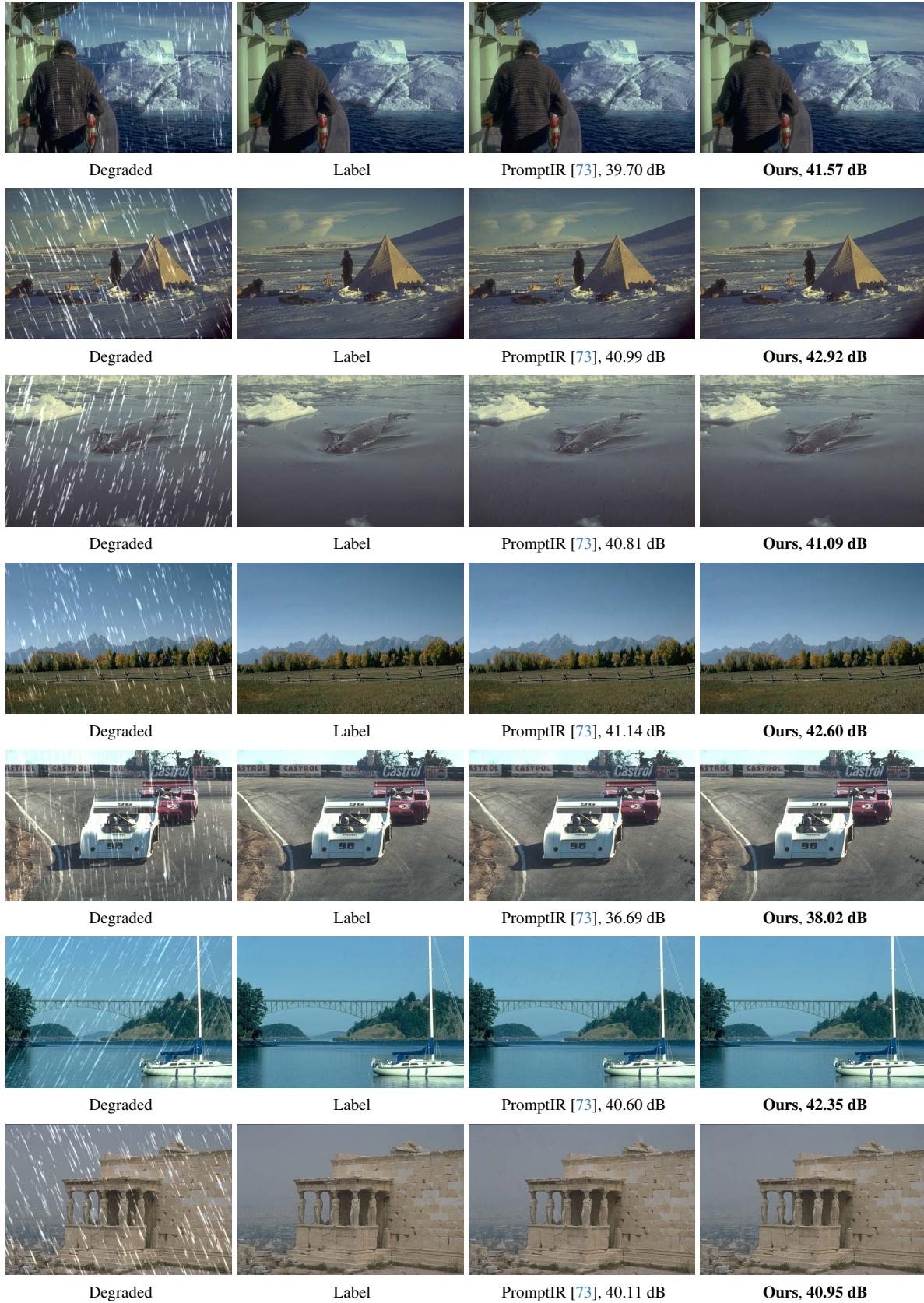


Figure K. Image Deraining on Rain100L [104].

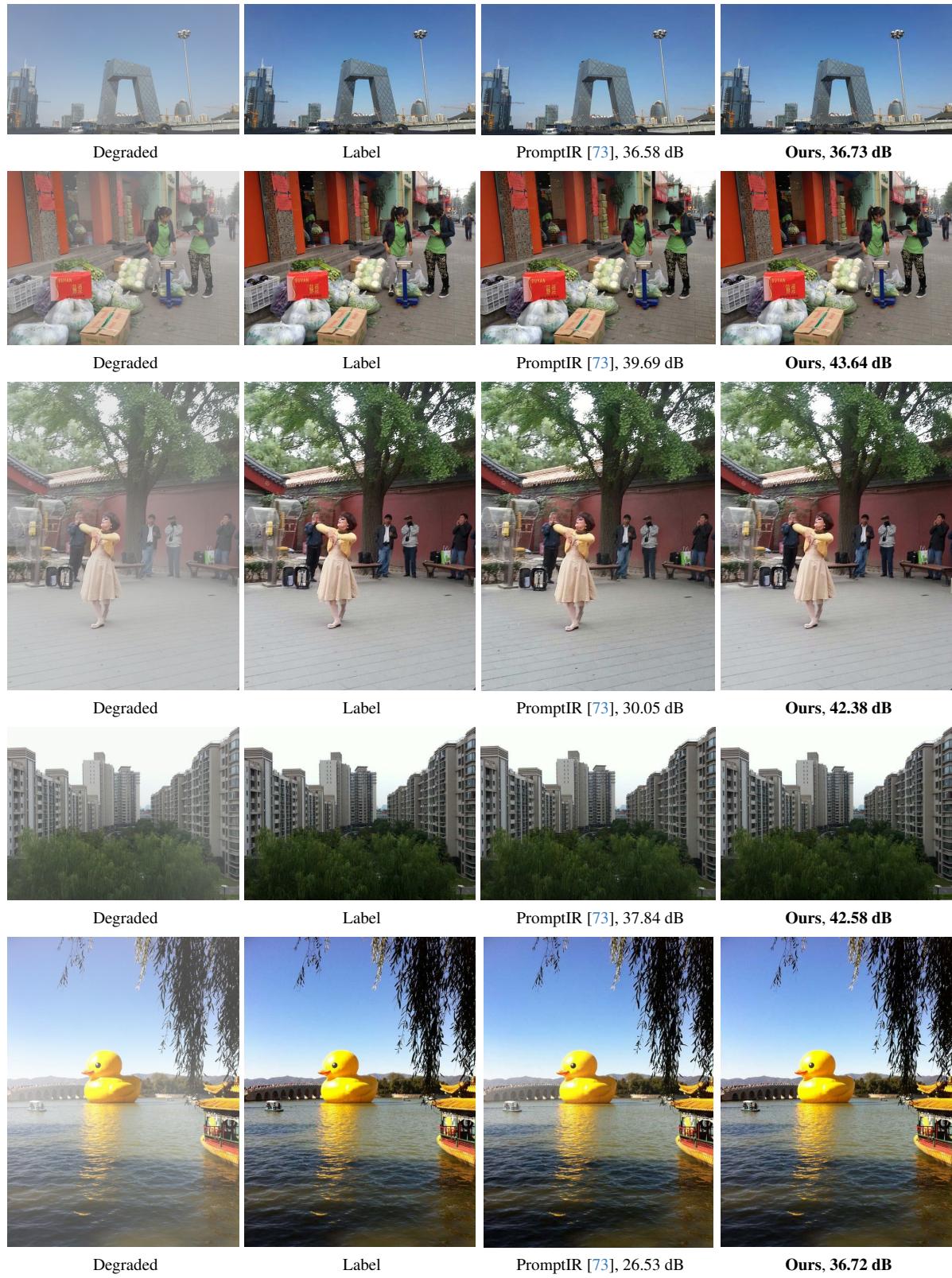


Figure L. Image Dehazing on SOTS-outdoor [50].

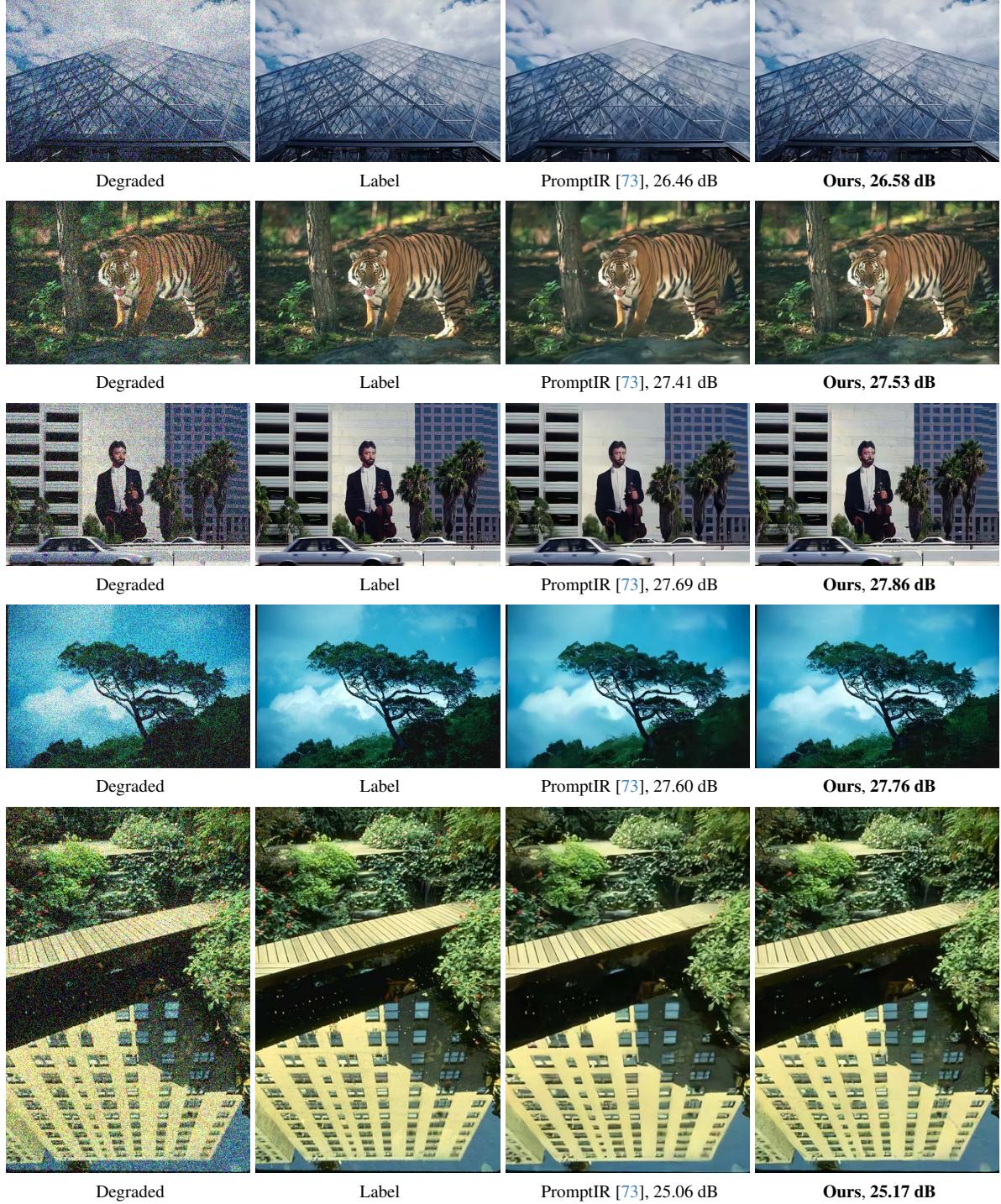


Figure M. Image Denoising on CBSD68 [66].

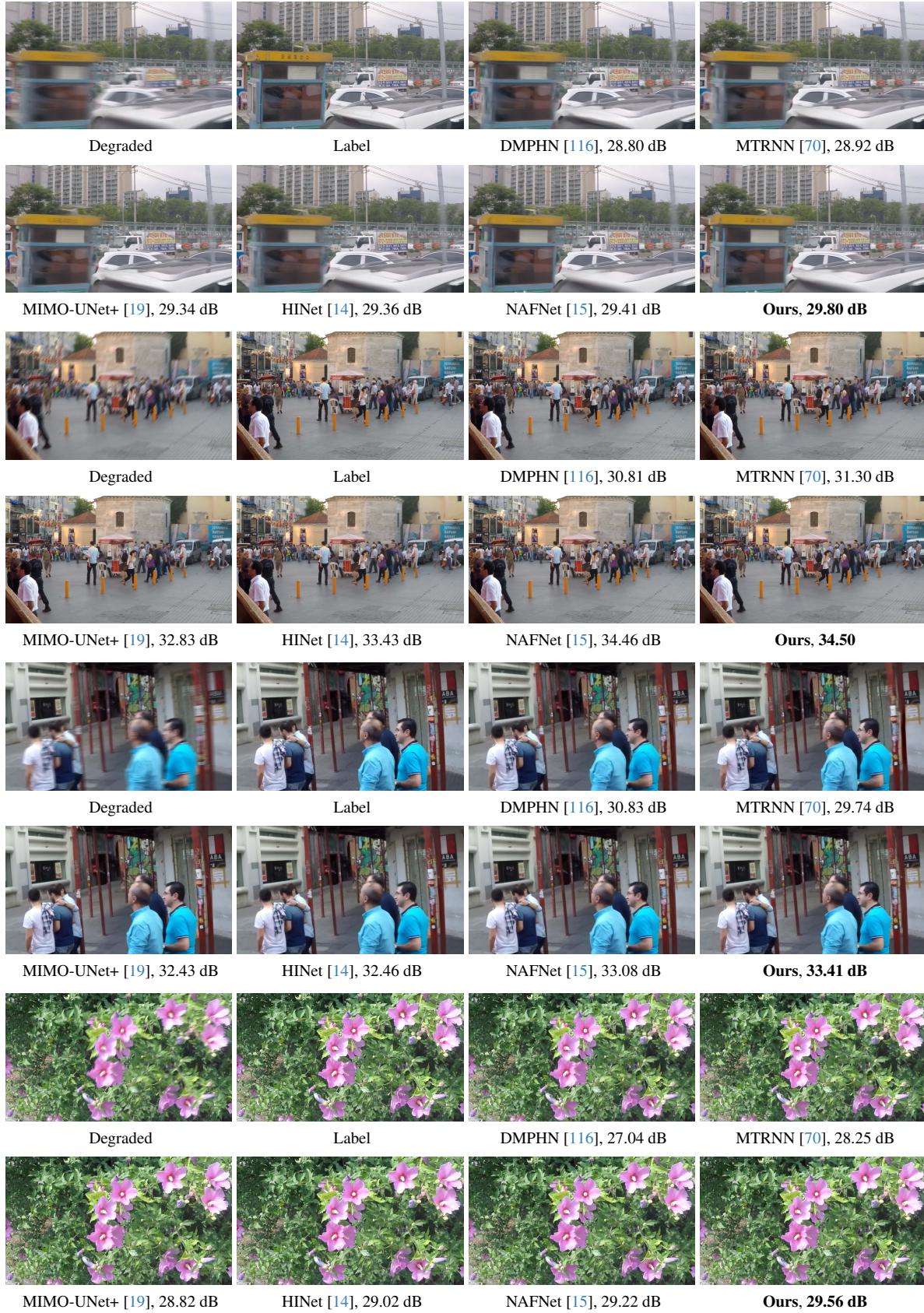


Figure N. Single-image motion deblurring on GoPro [68].

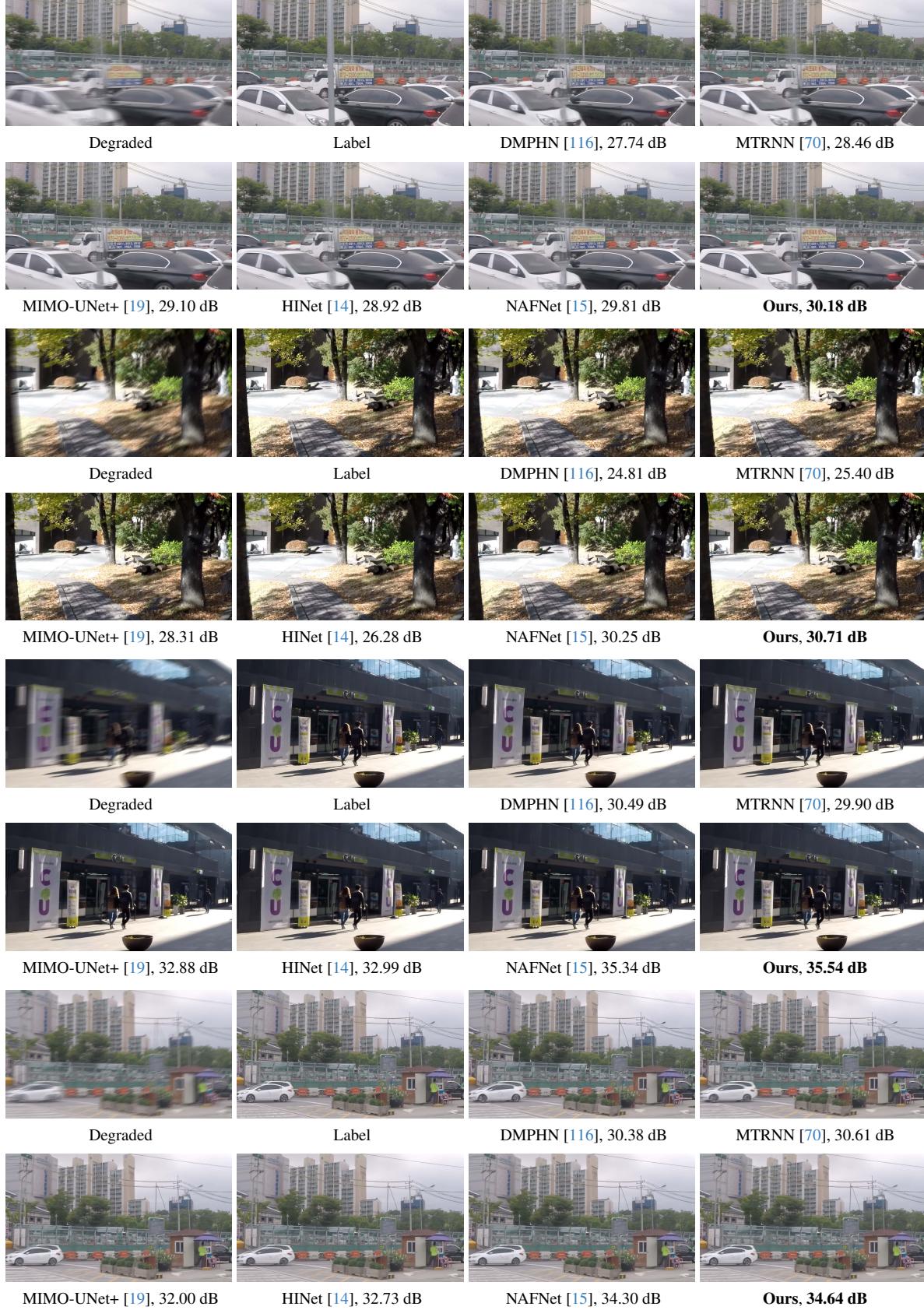


Figure O. Single-image motion deblurring on GoPro [68].



Figure P. Defocus deblurring on DPDD [2].

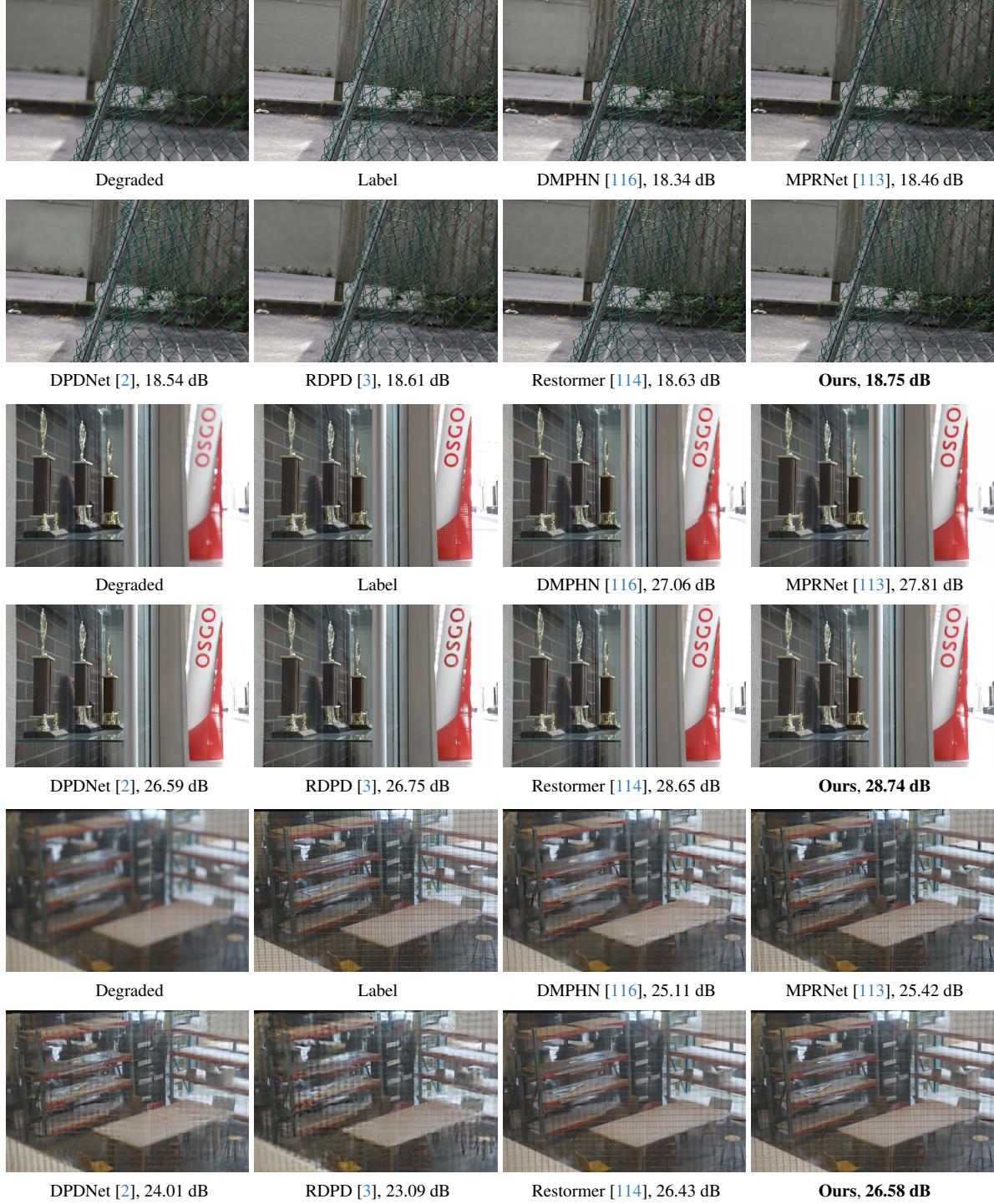


Figure Q. Defocus deblurring on DPDD [2].

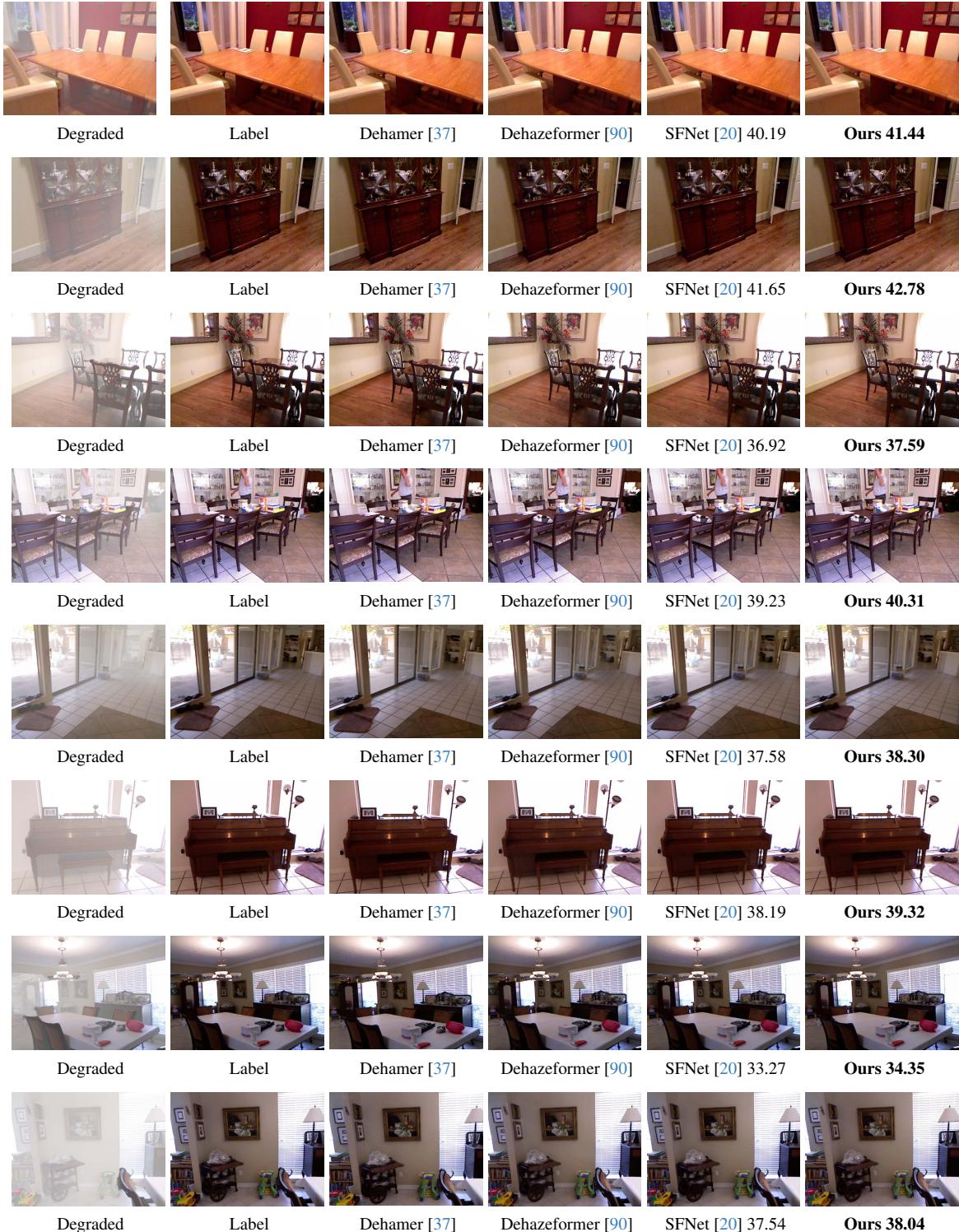


Figure R. Image dehazing results on SOTS [50].

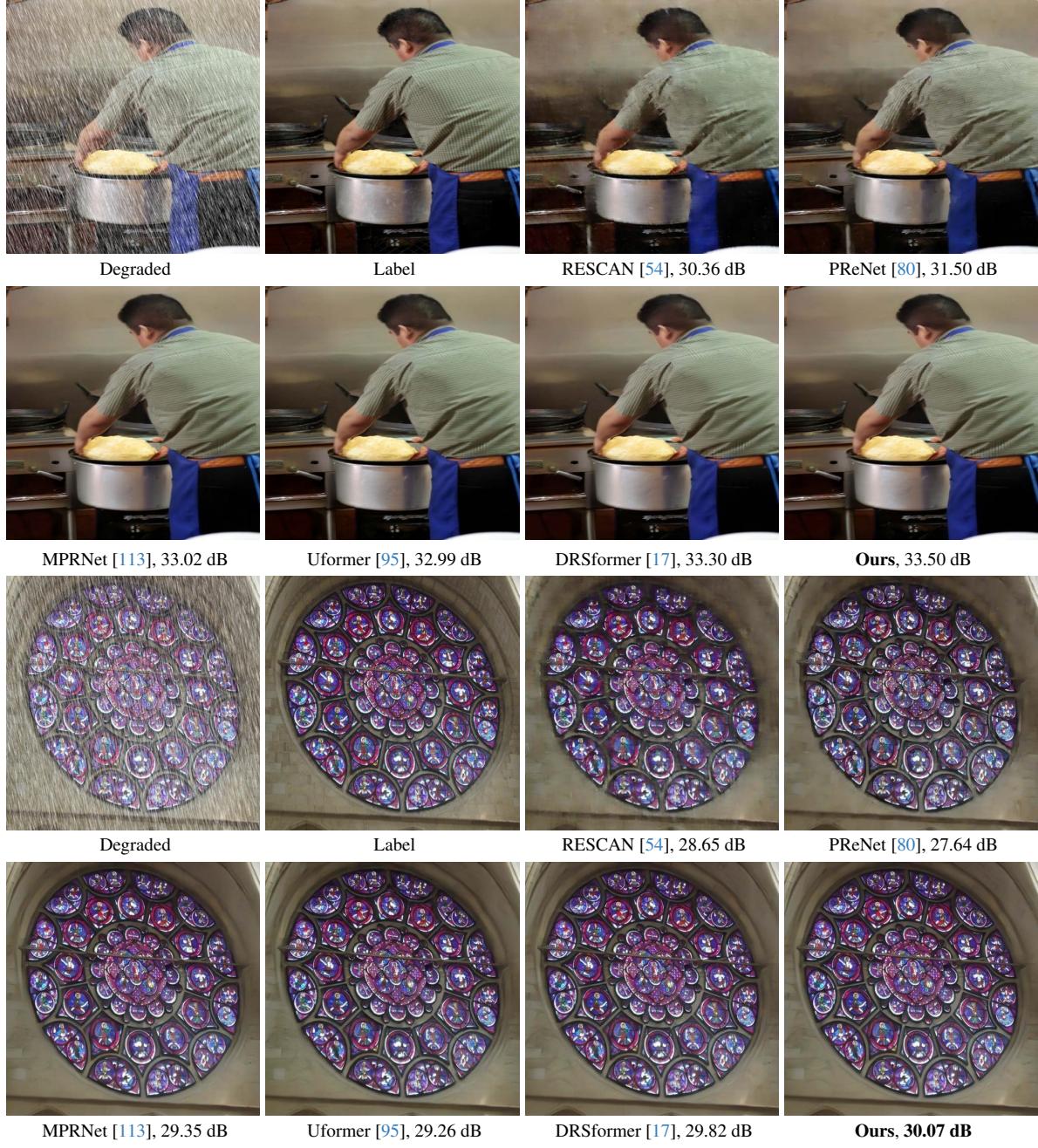


Figure S. Image deraining results on DID-Data [115]

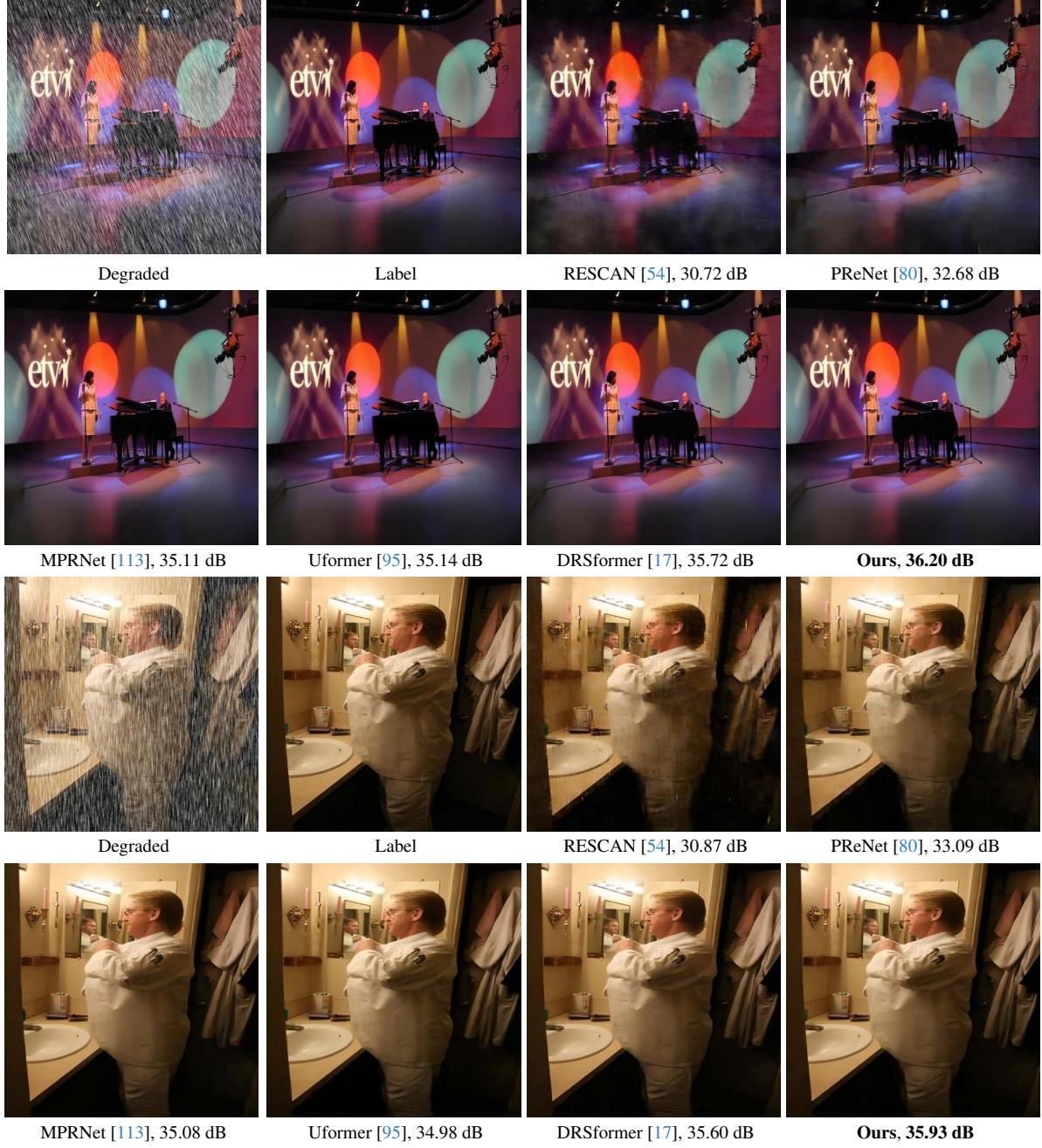


Figure T. Image deraining results on DID-Data [115]

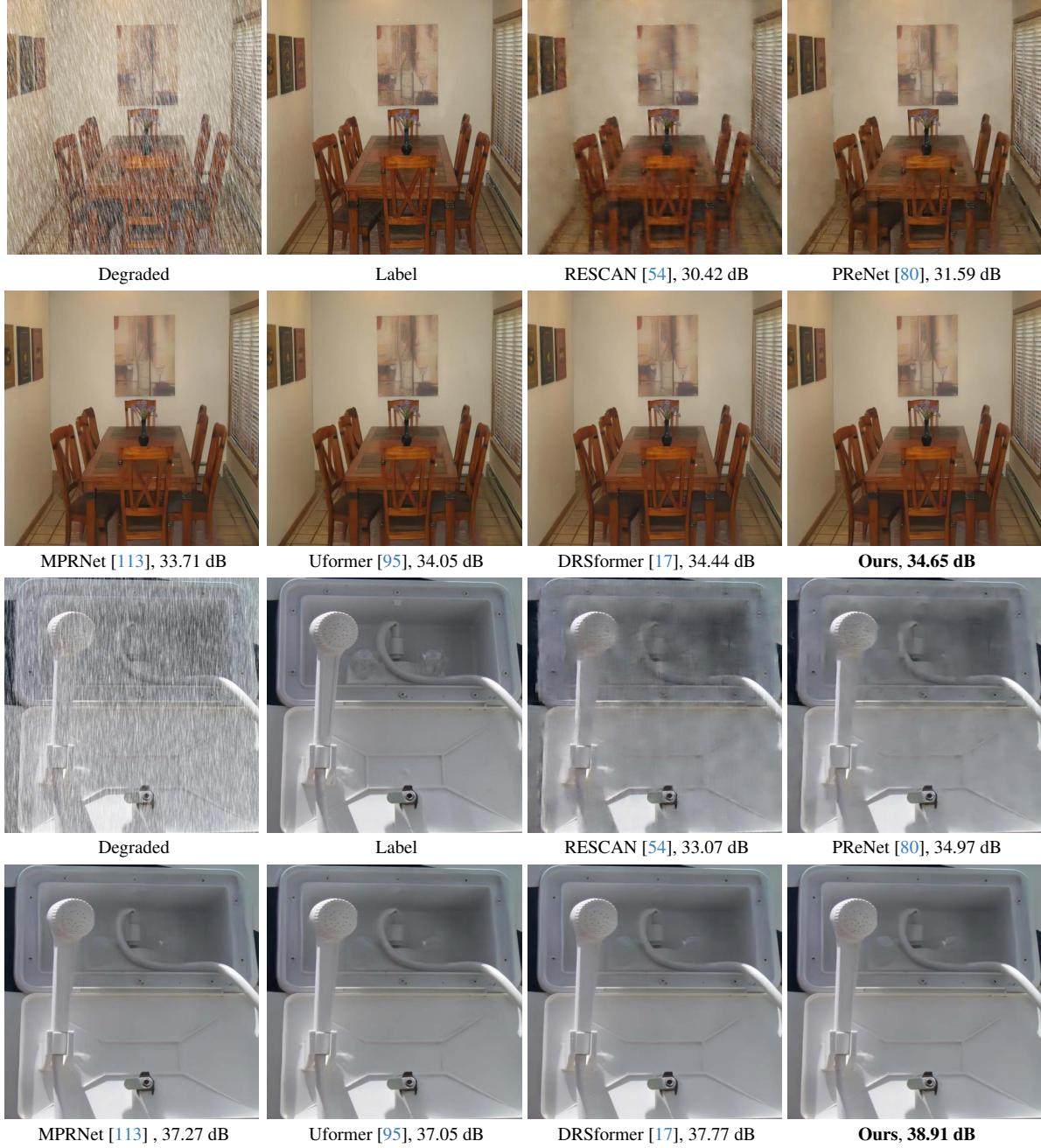


Figure U. Image deraining results on DID-Data [115]

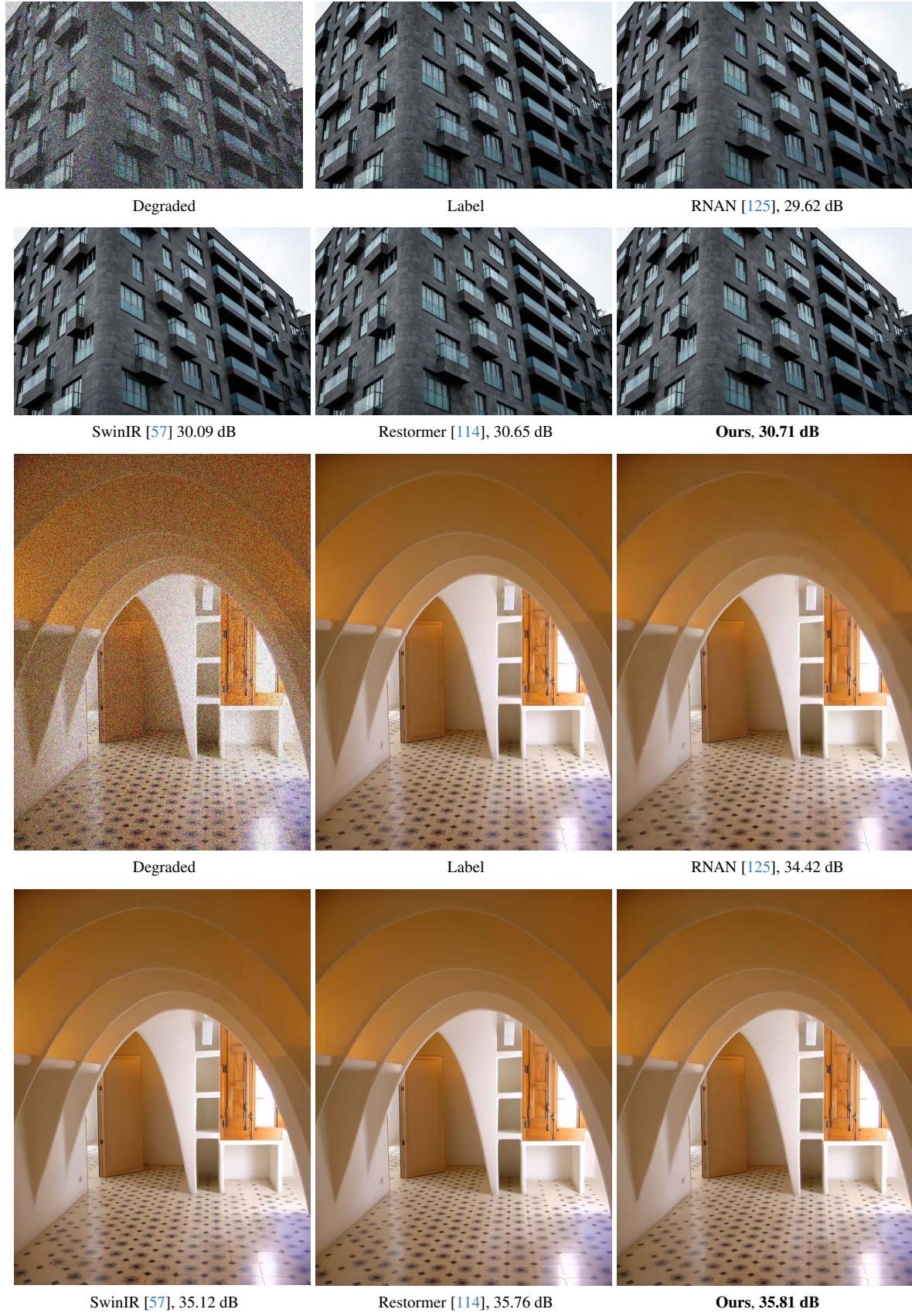


Figure V. Gaussian color denoising results on Urban100 [40].