

Dereflection Any Image with Diffusion Priors and Diversified Data

Jichen Hu^{1*}, Chen Yang^{1*}, Zanwei Zhou¹, Jiemin Fang^{2†}, Xiaokang Yang¹, Qi Tian², Wei Shen^{1✉†}
¹ MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University
² Huawei Inc.



Figure 1. Our model demonstrates strong and general reflection removal capabilities in diverse real-world scenarios. **Upper:** Original images with reflections. **Bottom:** Results generated by our model. The demonstrated scenarios encompass multiple reflection types, including glass, plastic materials, water surfaces, and digital displays, *etc.*

Abstract

Reflection removal of a single image remains a highly challenging task due to the complex entanglement between target scenes and unwanted reflections. Despite significant progress, existing methods are hindered by the scarcity of high-quality, diverse data and insufficient restoration priors, resulting in limited generalization across various real-world scenarios. In this paper, we propose *Dereflection Any Image*¹, a comprehensive solution with an efficient data preparation pipeline and a generalizable model for robust reflection removal. First, we introduce a dataset named *Diverse Reflection Removal (DRR)* created by randomly rotating reflective mediums in target scenes, enabling variation of reflection angles and intensities, and setting a new benchmark in scale, quality, and diversity. Second, we propose a diffusion-based framework with one-step diffusion for deterministic outputs and fast inference. To ensure stable learning, we design a three-stage progressive training strategy, including reflection-invariant finetuning to encourage consistent outputs across varying reflection patterns that characterize our dataset. Extensive experiments

show that our method achieves SOTA performance on both common benchmarks and challenging in-the-wild images, showing superior generalization across diverse real-world scenes. Project page: <https://abuuu122.github.io/DAI.github.io/>.

1. Introduction

Capturing images through glass or other reflective mediums often introduces unwanted reflections, which degrade the visibility of the underlying target scene, resulting in a mixed image with two layers superposed. These reflections significantly degrade both visual aesthetics and usages in downstream tasks [20, 34]. Developing strong reflection-removal (**dereflection** for short) methods is essential for practical applications.

Traditional methods [18, 28, 32] for reflection removal rely on empirical assumptions like reflections being blurred or exhibiting ghosting effects, which often fail in real-world scenarios. Learning-based methods [2, 7–10, 30, 53, 54] attempt to learn dereflection capabilities through paired data. However, such paired data is difficult to obtain in real-world settings, *e.g.*, capturing museum exhibits with and without reflections typically requires physically removing the glass from display cases, which is often impractical. Consequently, existing real-world datasets [16, 33, 52, 54] are

* Equal contribution.

† Project lead.

✉ Corresponding author.

¹“Dereflection” may not be an accurate word, which is used to represent removing reflection for simplicity. This name is to respect [13, 48].

limited in scale and diversity, failing to capture a wide range of reflection patterns encountered in practice. Moreover, synthetic datasets [9, 12, 42] often suffer from significant domain gaps, as their reflection patterns and image characteristics do not fully align with real-world scenarios.

To address these limitations, we propose an efficient data collection pipeline and introduce Diverse Reflection Removal (DRR), a 4K dataset with diverse reflection patterns. Specifically, for each view of the capturing device, reflective mediums are randomly rotated with the target scene. Then we capture mixed videos and decompose them into frames to form data pairs. This approach allows us to flexibly vary reflection angles, intensities, and scene diversity, significantly enhancing the realism and variety of reflection patterns compared to existing datasets [16, 33, 52, 54], as shown in Table 1 and Fig. 2. We also construct synthetic pairs to strengthen the data sufficiency, which are filtered with the CLIP score [25] to guarantee the data realism. Besides extending the diversity of training data, we propose to leverage the generalization capabilities of diffusion models, which excel in image-to-image translation tasks even tuned with finite data [31, 44]. Our model integrates ControlNet [50] to use the mixed image as a conditioning signal. To ensure deterministic outputs for reliable transmission layer recovery, we build our generative prior-based model on one-step diffusion [46, 49], enabling both deterministic results and fast inference. To ensure stable learning of translation from complex mixed images to dereflected ones, we design a three-stage progressive training strategy. Initially, we employ foundation training with image pairs to get basic performance. We further exploit the unique characteristics of our dataset through a reflection-invariant finetuning strategy. Since our dataset contains identical transmission scenes with varying reflection patterns, we train the model to produce consistent outputs despite these variations, enhancing generalization by focusing on invariant properties of transmission scenes rather than variable reflection characteristics. Finally, a cross-latent decoder is trained to mitigate blurriness and preserve details. The new dataset, code and model checkpoints will be released.

Our key contributions are summarized as follows:

- We introduce an efficient pipeline for data collection, and present DRR, a high-quality dataset featuring diverse reflections with varying angles, fostering future advancements in the field.
- We design a novel diffusion-based framework with a progressive training strategy ensuring both stable optimization and strong generalization across diverse reflection types, as shown in Fig. 1.
- Extensive experiments demonstrate that our method achieves SOTA performance not only on benchmark datasets but also on challenging in-the-wild images captured by mobile devices, exhibiting superior generaliza-

tion across diverse real-world reflection scenarios.

2. Related Work

2.1. Reflection Removal

Reflection removal is an ill-posed problem that aims to separate the reflection layer from a mixed image and recover the underlying transmission layer. It is well-established that additional information can significantly simplify this task. For instance, multi-view [5, 17, 24] information can guide robust reflection removal, but capturing image sequences is often redundant for users. Flash-based methods [15, 37] utilize pairs of images—one with flash and one without—since the transmission layer typically appears brighter under flash illumination. Similarly, polarization-based techniques [14] exploit images captured at different polarization angles. However, these approaches rely on specialized equipment, which limits their general applicability. As a result, we focus on the most challenging yet widely applicable task: single-image reflection removal (SIRR).

Traditional methods for reflection removal rely on hand-crafted priors, such as the assumption that reflections are blurred [18, 32] or exhibit ghosting effects [28]. However, these assumptions frequently break down in complex real-world scenarios, resulting in suboptimal performance. Universal image restoration techniques [1, 11, 40] could serve as a potential solution, however, they cannot outperform methods designed specifically for single image reflection removal. In recent years, learning-based methods [2, 4, 7–10, 30, 36, 41, 47, 53] have become the mainstream approach, enabling models to remove reflections by training on large datasets. While some real datasets [16, 33, 52, 54] have been collected, their scale and quality remain insufficient. To address this, some methods employ empirical formulations [9], physics-based rendering [12], and tailored models [42] to generate synthetic data. Nevertheless, these approaches often struggle to bridge the gap between synthetic and real-world data. In this work, we introduce a comprehensive data collection pipeline to acquire a large volume of real-world data DRR and a data purification strategy to enhance the quality of synthetic data.

2.2. Diffusion Model

Diffusion models [6, 26, 29] have emerged as a powerful framework for image generation and restoration, leveraging an iterative denoising process that learns the underlying data distribution. Their ability to produce high-quality, diverse outputs has garnered significant attention in recent years. Notably, diffusion models have been successfully applied to a wide range of image restoration tasks [11], such as deblurring [43], super-resolution [35, 38] and inpainting [23, 45], demonstrating exceptional performance in recovering fine-grained details. These characteristics make

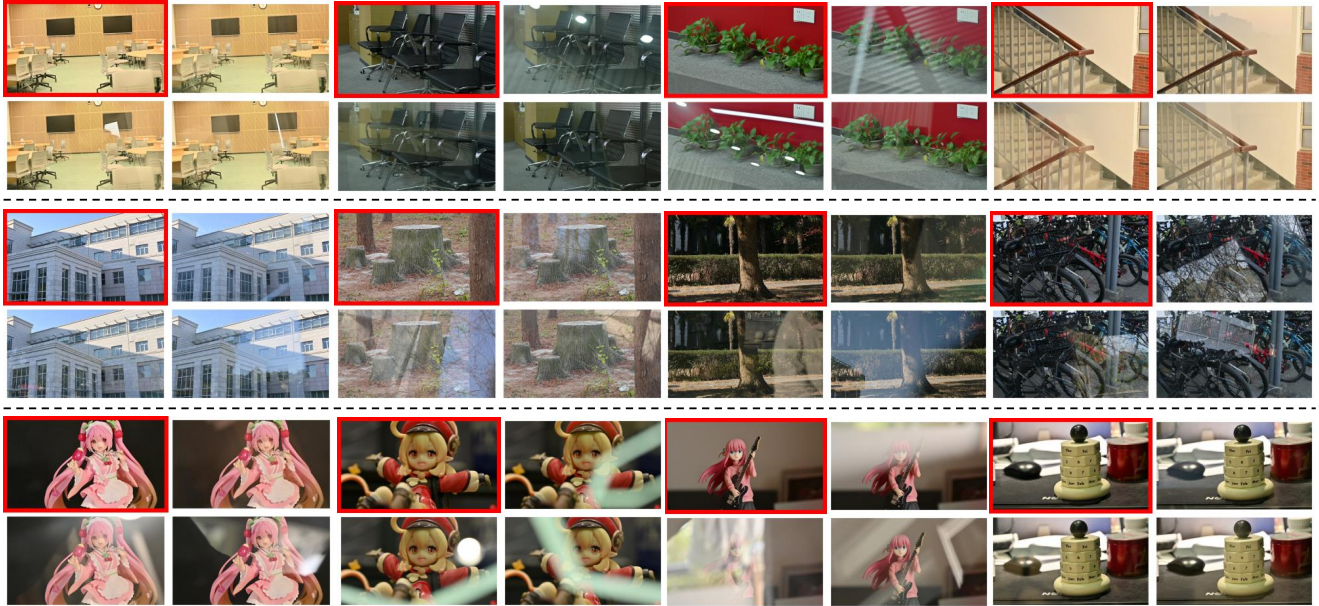


Figure 2. Our dataset contains a diverse collection of scenes, each accompanied by multiple reflection images. As illustrated in the figure, the ground truth transmission layer is highlighted in red boxes, while the remaining images represent various mixed images. The dataset demonstrates remarkable diversity, encompassing **indoor**, **outdoor**, and **object-centric scenes**. All image pairs maintain high resolution with rich textual details. (Best viewed on screen.)

Table 1. Summary of existing real reflection datasets. Compared to these datasets, our proposed DRR dataset demonstrates significant advantages in three key aspects: (1) greater diversity in reflection (varying glass angles, reflected contents, reflection intensity) and scenes (various environment conditions), (2) a substantially larger collection of image pairs, and (3) superior image quality with higher resolution.

| | SIR ² [33] | Real[52] | Nature[16] | RRW[54] | DRR (Ours) |
|-----------|-----------------------|------------|------------|-----------|------------------|
| Year | 2017 | 2018 | 2020 | 2023 | 2025 |
| Videos | × | × | × | ✓ | ✓ |
| Angles | × | × | × | × | ✓ |
| Usage | Test | Train/Test | Train/Test | Train | Train/Test |
| Pairs | 500 | 89/20 | 200/20 | 14952 | 23303/400 |
| Scenes | 126 | 89/20 | 68 | 150 | 217/40 |
| Avg. res. | 540*400 | 1152*930 | 598*398 | 2580*1460 | 3840*2160 |

diffusion models particularly well-suited for the challenging task of reflection removal, where disentangling complex visual structures and maintaining image fidelity are critical. Recent work, such as L-Differ [7], utilizes ControlNet [50] to inject information from the mixed image and iteratively denoise it. However, this approach requires multiple steps to recover the transmission layer and text prompt to guide the process. To address this limitation, we propose an alternative solution based on a recently proposed one-step denoising strategy [46, 49]. It enables stable and deterministic results while significantly accelerating the process, making it more practical for real-world applications.

3. Method

In this section, we first present our data collection pipeline, which captures diverse real-world data and generates high-quality synthetic data. Next, we introduce our diffusion-based framework specifically designed for reflection removal. Finally, we propose a reflection-invariant finetuning strategy to fully leverage the diversity of our dataset and enhance robustness for generalizing to real-world reflection scenarios.

3.1. Data Collection Pipeline

Formally, we define the target scene as transmission layer \mathbf{T} , which is superposed with undesired reflection layer \mathbf{R} , resulting in the mixed image \mathbf{M} . Given that the core objective is to translate mixed image \mathbf{M} to transmission layer \mathbf{T} , we design a comprehensive data collection pipeline to acquire aligned image pairs of \mathbf{M} and \mathbf{T} , including our diverse real-world data, and synthetic data serving as a supplementary, demonstrated in Fig. 3.

Real Data Our dataset is captured using a Nikon Z50 camera and three mobile phones mounted on a fixed tripod, with a portable glass slab positioned in front of the lens. By rotating the glass at different angles within a sequence, we generate reflection images with varying reflection and intensity. The corresponding ground-truth transmission images are acquired by removing the glass entirely. To ensure

diversity and robustness, we carefully vary key parameters, including scenes (indoor, outdoor, object-centric), lighting conditions (skylight, sunlight, incandescent), glass thickness (3 mm and 8 mm), camera-to-glass distance, viewing angles, exposure values, and aperture settings.

The dataset consists of 257 unique scenes, each captured with two glass thicknesses to ensure reflection diversity. It is partitioned into a training set (217 scenes, 23,303 image pairs) and a testing set (40 scenes, 400 image pairs). All training images are captured in 4K resolution (3840×2160 pixels) to provide high visual fidelity for model training. The testing set is further divided into two subsets: DRR-S, containing standard reflections captured with the same camera used in training, and DRR-C, containing challenging reflections captured using three mobile phones. This division enables a comprehensive evaluation of model performance across varying reflection complexities and real-world scenarios.

To address spatial shifts caused by glass refraction, we employ a robust post-processing pipeline. Inspired by [33], we use Scale-Invariant Feature Transform (SIFT) [22] for feature point detection and Random Sample Consensus (RANSAC) for precise alignment between reflection images and their ground-truth transmission pairs.

Compared to existing datasets, our collection offers superior image quality, substantial quantity, and exceptional diversity, derived from varying glass angles, reflected contents, reflection intensity and various environment conditions. The existing RRW dataset focuses on object movement and occlusions in front of the glass, lacking the variation in angular perspectives and diverse reflection content from the real scenes that characterize our dataset. A detailed comparison is provided in Table 1.

Synthetic Data Supplementary In addition to datasets collected from the real world, synthetic data can serve as a significant supplementary resource to enrich the diversity of training data. In our work, we adopt the following formulation from DSRNet [9] to generate synthetic data:

$$\mathbf{M} = \gamma_1 \mathbf{T} + \gamma_2 \mathbf{R} - \gamma_1 \gamma_2 \mathbf{T} \circ \mathbf{R}, \quad (1)$$

where \mathbf{T} , \mathbf{R} , \mathbf{I} represent the transmission, reflection and mixed layers, respectively. To enhance the diversity of the synthesized data, we randomly sample transmission and reflection layers from the COCO [19] and PASCAL VOC [3] datasets and assign random values to the coefficients $\gamma_1 \in [0.8, 1.0]$ and $\gamma_2 \in [0.4, 1.0]$ during synthesis. An intuitive observation is that the synthetic data exhibit a wide range of quality: while some images closely resemble real-world reflections, others appear less realistic. To address this, we leverage CLIP similarity [25] with the text prompt “image with glass reflection” to evaluate the synthetic data.

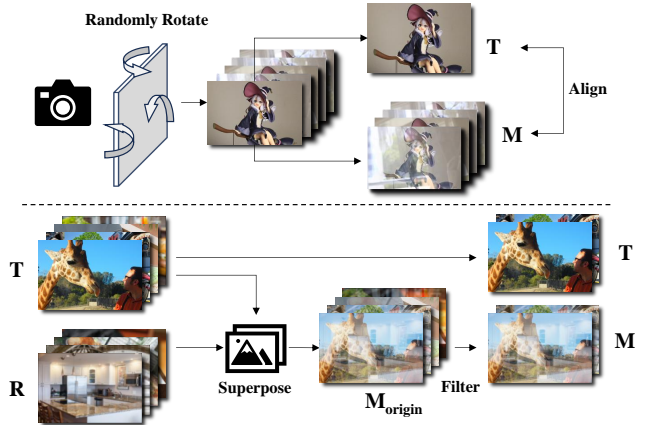


Figure 3. Data collection pipeline of real (above) and synthetic (below) data. Real data is captured by recording videos while rotating a glass panel at various angles, then processed to align mixed images with their ground truth transmission layers. Synthetic data is generated by randomly chosen coefficients and filtered to produce high-quality image pairs.

Based on this metric, we filtered synthetic data from initially 69,443 pairs to 20,833 high quality pairs. Similar to our real-world dataset, our synthetic data is also capable of generating multiple reflection for a single scene, thereby enhancing the diversity and practicality of the training data.

3.2. Framework

Our proposed dataset offers a comprehensive and diverse training resource, however, it is impractical to encompass all possible reflection types. Fortunately, diffusion models provide an effective solution. By leveraging the powerful generative priors of diffusion models, we can maximize the utility of limited data and address more challenging reflection scenarios. Our diffusion-based model is composed of three key parts, a U-net for one-step denoising, a ControlNet to input the mixed image, and a cross-latent decoder to preserve details, Shown in Fig. 4.

Traditional diffusion models generate results through iterative denoising from Gaussian noise. But we use a one-step denoising strategy as an alternative, the target latent is predicted with only one step by the U-net, resulting in deterministic results and fast inference. ControlNet [50] is initially designed to impose structural constraints on generated images. In our work, we adapt it as a mechanism to inject information from mixed images. To address the inherent limitations of one-step diffusion, such as over-smoothed results and potential shifts in details during the denoising process, we add a cross-latent decoder \mathcal{D} to preserve the high-frequency details of the input mixed image. Inspired by [27, 37], the multi-scale latent features extracted by the encoder are directly connected to the decoder with zero convolution, creating a shortcut that preserve the initial infor-

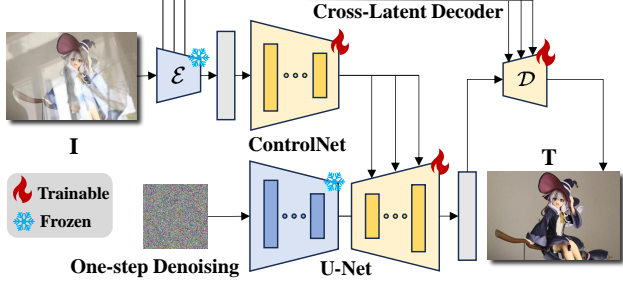


Figure 4. Our proposed framework. It consists of a U-net with one-step denoising strategy, a ControlNet to input the mixed image processed by the encoder \mathcal{E} , and a cross-latent decoder \mathcal{D} to mitigate blurriness and preserve details.

mation of the mixed image.

3.3. Progressive Training

To effectively train our reflection removal model, we adopt a progressive training strategy that decomposes the learning process into three distinct yet interconnected phases, as illustrated in Fig. 5. This hierarchical approach enables stable optimization and ensures that each component of our model is properly trained to handle the complex task of reflection removal.

Foundation Training The initial stage focuses on establishing the fundamental reflection removal capability by jointly training the ControlNet and the upsampling blocks of the U-Net. We employ the following one-step diffusion loss function:

$$\mathcal{L}_{diff} = \mathbb{E}_{z_t, \mathbf{M}} \left[\|z_t - \mu_{\theta, \phi}^{z_t}(\mathbf{M})\|_2^2 \right], \quad (2)$$

where θ and ϕ refer to the parameters of U-Net and ControlNet, z_t is the latent representation after adding noise over t steps with $t \in (0, T)$. (in Stable Diffusion, T is set to 1,000). For notational simplicity, we leave the full formulation in the supplementary material.

Reflection-Invariant Fine-tuning Building upon the foundation, we introduce a novel reflection-invariant fine-tuning strategy that leverages the unique characteristics of our dataset, where each scene contains multiple mixed images with varying reflections. Our key insight is that the model should produce consistent results regardless of the specific reflection patterns present in the input. This is achieved by incorporating a consistent loss:

$$\mathcal{L}_{con} = \mathbb{E}_{\mathbf{M}_1, \mathbf{M}_2} \left[\|\mu_{\theta, \phi}^{z_t}(\mathbf{M}_1) - \mu_{\theta, \phi}^{z_t}(\mathbf{M}_2)\|_2^2 \right], \quad (3)$$

$$\mathcal{L} = \mathcal{L}_{diff}(\mathbf{M}_1) + \mathcal{L}_{diff}(\mathbf{M}_2) + \mathcal{L}_{con}(\mathbf{M}_1, \mathbf{M}_2). \quad (4)$$

The reflection-invariant finetuning enhances generalization by focusing on invariant properties of transmission

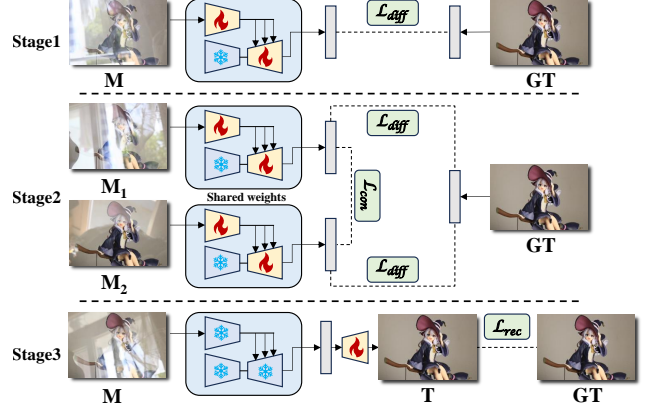


Figure 5. The three stages of progressive training. First, we train the ControlNet and the upsampling blocks of the U-Net using the basic one-step diffusion loss. Second, we finetune these components by incorporating the consistent loss. Finally, we train the cross-latent decoder using the image reconstruction loss.

scenes rather than variable reflection characteristics, making the model more robust and generalize to complex reflection scenarios.

Cross-Latent Decoder In the final stage, we freeze the previously trained components and focus on training the cross-latent decoder to mitigate blurriness and preserve high-frequency details. This is accomplished through a comprehensive image reconstruction loss:

$$\mathcal{L}_{rec} = \mathcal{L}_1 + \mathcal{L}_{SSIM} + \mathcal{L}_{LPIPS}. \quad (5)$$

4. Experiments

4.1. Implementation Details

Our model is trained on a combination of datasets, including our newly introduced real dataset DRR, our synthetic dataset, and the previously established Real and Nature datasets. To enhance the robustness and generalization capability of our model, we implement a comprehensive suite of data augmentation techniques. These techniques include random cropping to a resolution of 768, random flipping, and random adjustments of contrast, hue, brightness, and saturation applied uniformly to all training pairs.

We employ a progressive training strategy to divide the whole process into three distinct stages. All three stages are executed on 3 NVIDIA GeForce RTX 3090 GPUs with a batch size of 1. The pretrained weights of Stable Diffusion V2.1 are used to initialize the U-net and ControlNet. We employ the AdamW [21] optimizer with a fixed learning rate of 3×10^{-4} except during the reflection-invariant finetuning phase, where a learning rate of 1×10^{-4} is used. To optimize memory usage during the reflection-invariant finetuning in the second phase, we alternate the training of



Figure 6. Qualitative comparison of our method with other approaches fine-tuned under our data settings. The benchmark datasets, listed from top to bottom, are *Nature*, *Real*, *SIR²*, *DRR-S*, and *DRR-C*. The red rectangles highlight key regions for comparison.

the ControlNet and the upsampling blocks of the U-Net every 100 steps. The durations for the three parts are 2 days, 1 day, and 1 day, respectively. Our method achieves fast inference speed, processing a 768-resolution image in approximately 1 second on a single NVIDIA GeForce RTX 3090 GPU.

We evaluate the performance of our model on three established benchmarks: *Nature* [16], *Real* [52], and *SIR²* [33]. Additionally, we introduce a new benchmark derived from our DRR dataset to provide a more comprehensive evaluation. Our DRR dataset includes one standard

set and another challenging set, each containing 200 image pairs sampled from 20 different scenes. This design further validates the robustness and generalization capabilities of our approach across diverse real-world conditions.

4.2. Comparison to State-of-the-arts

We compare our method with several recently proposed approaches, including RobustSIRR [30], DSRNet [9], L-Differ [7], RRW [54], and DSIT [10]. All evaluations are conducted using the same testing scripts and the same testing data to ensure fairness. We finetune the above meth-

Table 2. Quantitative comparisons on the existing three reflection benchmarks, and our new dataset DRR, consisting a standard set (DRR-S) and a challenging set (DRR-C). The scores after finetuning on our data setting are labeled as †. The best results are in **bold**, and the second-best results are underlined.

| Methods | Venue | <i>Nature</i> (20) | | <i>Real</i> (20) | | <i>SIR</i> ² (500) | | <i>DRR-S</i> (200) | | <i>DRR-C</i> (200) | |
|------------------|--------------|--------------------|--------------|------------------|--------------|-------------------------------|--------------|--------------------|--------------|--------------------|--------------|
| | | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| RobustSIRR [30] | CVPR 2023 | 20.94 | 0.770 | 22.71 | 0.787 | 22.61 | 0.872 | 19.68 | 0.756 | 20.24 | 0.692 |
| DSRNet [9] | ICCV 2023 | 24.86 | 0.823 | 23.31 | 0.791 | 25.65 | 0.919 | 22.33 | 0.846 | 21.93 | 0.820 |
| RRW [54] | CVPR 2024 | 25.79 | 0.833 | 21.51 | 0.767 | 25.31 | 0.907 | 22.39 | 0.857 | 21.84 | 0.820 |
| L-Differ [7] | ECCV 2024 | 23.95 | 0.831 | 23.77 | 0.821 | 25.18 | 0.911 | - | - | - | - |
| DSIT [10] | NeurIPS 2024 | 26.25 | 0.833 | 24.54 | 0.814 | <u>26.34</u> | 0.922 | 23.48 | 0.869 | 22.46 | 0.817 |
| RobustSIRR† [30] | CVPR 2023 | 23.74 | 0.818 | 23.00 | 0.793 | 24.14 | 0.897 | 23.25 | 0.793 | 20.94 | 0.709 |
| DSRNet† [9] | ICCV 2023 | 24.99 | 0.827 | 23.77 | 0.809 | 24.65 | 0.911 | 24.08 | 0.872 | 22.00 | 0.826 |
| RRW† [54] | CVPR 2024 | 26.37 | <u>0.842</u> | 23.82 | 0.802 | 25.08 | 0.908 | 23.00 | 0.872 | 22.19 | 0.826 |
| DSIT† [10] | NeurIPS 2024 | <u>26.44</u> | 0.836 | <u>24.94</u> | <u>0.823</u> | 26.18 | <u>0.927</u> | <u>26.50</u> | <u>0.890</u> | <u>22.77</u> | <u>0.832</u> |
| Ours | - | 26.81 | 0.843 | 25.21 | 0.841 | 27.19 | 0.930 | 27.25 | 0.902 | 23.77 | 0.843 |

Table 3. Quantitative results of ablation study, CLD refers to Cross-Latent Decoder, DRR is our dataset Diverse Reflection Removal, RIF refers to Reflection-Invariant Finetuning.

| CLD | DRR | RIF | <i>Nature</i> (20) | | <i>Real</i> (20) | | <i>SIR</i> ² (500) | | <i>DRR-S</i> (200) | | <i>DRR-C</i> (200) | |
|-----|-----|-----|--------------------|--------------|------------------|--------------|-------------------------------|--------------|--------------------|--------------|--------------------|--------------|
| | | | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| × | × | × | 25.54 | 0.782 | 23.67 | 0.762 | 24.96 | 0.861 | 22.32 | 0.769 | 20.50 | 0.666 |
| ✓ | × | × | 26.03 | 0.829 | 24.39 | 0.830 | 26.00 | 0.919 | 23.67 | 0.878 | 22.17 | 0.827 |
| ✓ | ✓ | × | 26.46 | 0.836 | 24.78 | 0.836 | 26.75 | 0.926 | 26.47 | 0.898 | 23.39 | 0.837 |
| ✓ | ✓ | ✓ | 26.81 | 0.843 | 25.21 | 0.841 | 27.19 | 0.930 | 27.25 | 0.902 | 23.77 | 0.843 |

ods using our data setting, with the exception of L-Differ, for which pre-trained weights and implementation codes are not publicly available. Qualitative results in Fig. 6 show that our diffusion-based framework show SOTA performance and strong generalization across diverse reflection types. The quantitative results on the three established benchmarks are demonstrated in Table 2. Our method achieves the highest average PSNR and SSIM scores, surpassing all recent approaches both before and after finetuning on our dataset. This demonstrates the overall superiority of our approach.

In addition, we evaluate the performance on our two testing benchmarks. The quantitative results are also presented in Table 2. We observe that existing methods perform poorly on these benchmarks, suggesting potential overfitting to older datasets. In contrast, our approach exhibits superior generalization capabilities, achieving strong performance on both the standard and challenging sets. Although the performance of existing methods improves after finetuning on our dataset, our method consistently maintains the highest scores, highlighting its effectiveness and robustness in handling diverse and challenging reflection scenarios.

Recognizing that all compared methods exhibit improved performance after finetuning on our dataset, we present a comprehensive comparison of their results before and after this adaptation in Fig. 7. It indicates that our dataset significantly enhances the robustness of reflection removal models. To illustrate this enhancement, we specif-

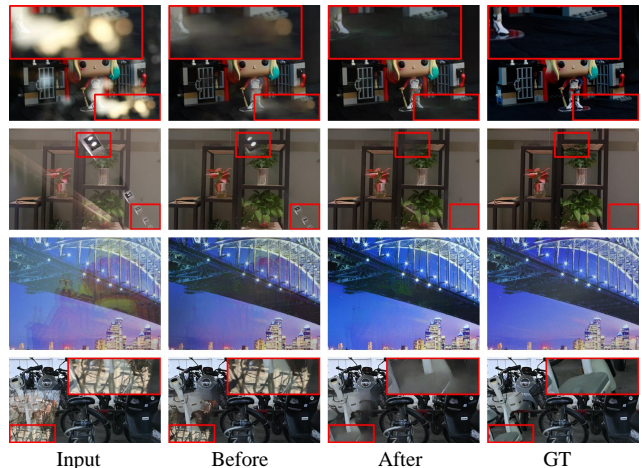


Figure 7. Visual comparison of DSIT performance before and after fine-tuning on our dataset. The results demonstrate significant improvements in reflection removal quality after adaptation to our data setting.

ically select the second-best performing method DSIT as a representative example.

4.3. Ablation Study

To validate the effectiveness of each component in our proposed method, we conduct a comprehensive ablation study.

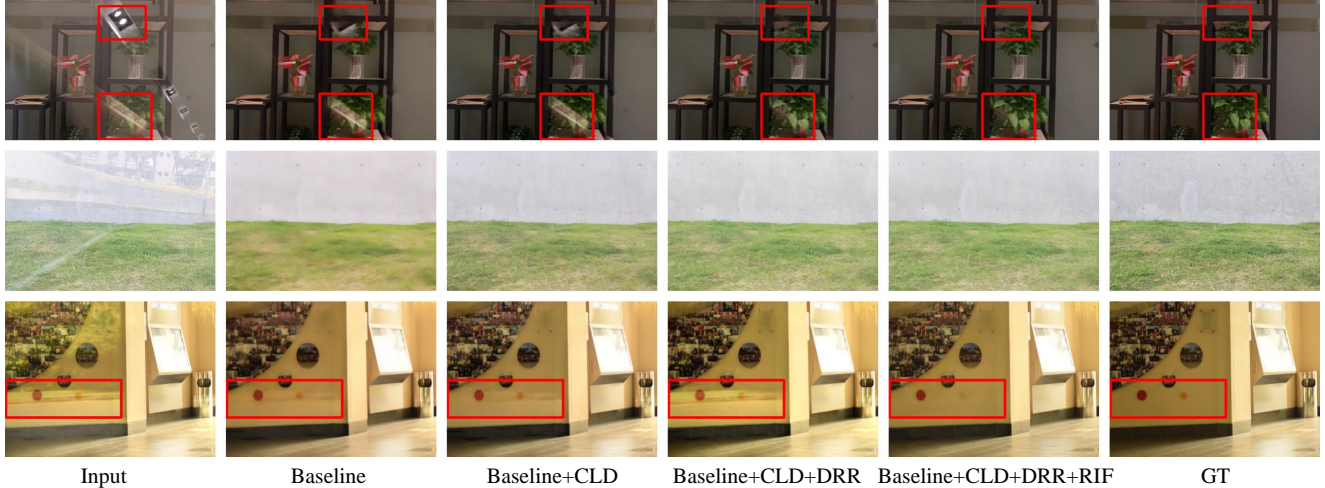


Figure 8. Qualitative results of the ablation study. CLD is Cross-Latent Decoder, DRR is our dataset Diverse Reflection Removal, RIF refers to Reflection-Invariant Finetuning. Baseline is the basic setting without CLD, DRR and CF.

We evaluate the impact of 1) the cross-latent decoder module, 2) our new dataset DRR, and 3) the reflection-invariant finetuning strategy by adding the component gradually from a baseline. The results are summarized in Table 3 and visual quality is demonstrated in Fig. 8.

Cross-Latent Decoder The baseline employs a vanilla decoder from Stable Diffusion without connecting it with encoder, thus failing to restore high-frequency details. To address the over-smooth issue of one-step diffusion process, we consider cross-latent decoder as an essential component. After incorporating this module, more original information is directly injected during the decoding process, effectively mitigating blurring effects and preserving details. This enhancement results in more faithful reconstruction quality, with improved PSNR and SSIM scores.

Influence of Dataset When training without our newly proposed dataset DRR, the model relies only on the old real datasets (Nature and Real) and synthetic data. Our experiments reveal that incorporating DRR enhances the performance across all benchmarks, including both established benchmarks and our proposed ones. This demonstrates that DRR introduces diverse and challenging scenarios, which are crucial for improving the generalization and robustness of reflection removal models.

Reflection-Invariant Finetuning Initially, we train the model on the entire dataset by treating all image pairs equally, without any distinction. However, experiments confirm that it is not a ideal solvement, thus leading to sub-optimal performance. After incorporating our reflection-invariant finetuning strategy, our model become more robust



Figure 9. Visualization of the generalization capability across diverse real-world reflection patterns. Our proposed method exhibits robust generalization to diverse reflection types unseen during training, including reflections from water surfaces, glossy plastics, digital displays, and even stylized scenes such as anime.

in handling diverse reflection scenarios, ultimately achieving more stable and reliable results.

4.4. Dereflection Any Image

To further demonstrate the versatility and generalization capability of our proposed method, we conduct additional experiments on challenging real-world images exhibiting various complex reflection patterns, such as reflections on water surfaces, glossy plastics, digital screens, and even stylized anime scenes. Due to the absence of ground-truth transmission layers in these scenarios, our evaluation is qualitative. Visualization results in Fig. 9 clearly illustrate the robustness and effectiveness of our approach in handling diverse, uncontrolled reflection conditions.

5. Conclusion

We present a comprehensive framework for single-image reflection removal, comprising a high-quality dataset (DRR) featuring diverse reflection scenarios and an effi-

cient capture pipeline, a robust diffusion-based model tailored for reflection removal, and a novel reflection-invariant finetuning strategy to enhance generalization. Our dataset significantly improves training quality by addressing the limitations of existing datasets. The proposed model, combined with reflection-invariant finetuning, achieves state-of-the-art performance on both common and our newly introduced benchmarks. Extensive experiments and ablation studies validate the effectiveness of each component, demonstrating superior generalization capabilities and robustness across diverse real-world scenarios. This work not only advances reflection removal technology but also provides a valuable resource and a strong baseline for future research in the field.

Limitation When the reflection and target scene are superposed with similar intensities, it is hard to distinguish which one is the target to be preserved. To address this limitation, incorporating semantic or user guidance offers a viable solution and a promising direction for future research.

References

- [1] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *European conference on computer vision*, pages 17–33. Springer, 2022. 2
- [2] Zheng Dong, Ke Xu, Yin Yang, Hujun Bao, Weiwei Xu, and Rynson WH Lau. Location-aware single image reflection removal. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5017–5026, 2021. 1, 2
- [3] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111:98–136, 2015. 4
- [4] Qingnan Fan, Jiaolong Yang, Gang Hua, Baoquan Chen, and David Wipf. A generic deep architecture for single image reflection removal and image smoothing. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3238–3247, 2017. 2
- [5] Xiaojie Guo, Xiaochun Cao, and Yi Ma. Robust separation of reflection from multiple images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2187–2194, 2014. 2
- [6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [7] Yuchen Hong, Haofeng Zhong, Shuchen Weng, Jinxiu Liang, and Boxin Shi. L-differ: Single image reflection removal with language-based diffusion model. In *European Conference on Computer Vision*, pages 58–76. Springer, 2024. 1, 2, 3, 6, 7
- [8] Qiming Hu and Xiaojie Guo. Trash or treasure? an interactive dual-stream strategy for single image reflection separation. *Advances in Neural Information Processing Systems*, 34:24683–24694, 2021.
- [9] Qiming Hu and Xiaojie Guo. Single image reflection separation via component synergy. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13138–13147, 2023. 2, 4, 6, 7
- [10] Qiming Hu, Hainuo Wang, and Xiaojie Guo. Single image reflection separation via dual-stream interactive transformers. *Advances in Neural Information Processing Systems*, 37:55228–55248, 2024. 1, 2, 6, 7
- [11] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. *Advances in Neural Information Processing Systems*, 35:23593–23606, 2022. 2
- [12] Soomin Kim, Yuchi Huo, and Sung-Eui Yoon. Single image reflection removal with physically-based training images. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5164–5173, 2020. 2
- [13] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4015–4026, 2023. 1
- [14] Chenyang Lei, Xuhua Huang, Mengdi Zhang, Qiong Yan, Wenxiu Sun, and Qifeng Chen. Polarized reflection removal with perfect alignment in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1750–1758, 2020. 2
- [15] Chenyang Lei, Xudong Jiang, and Qifeng Chen. Robust reflection removal with flash-only cues in the wild. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(12):15530–15545, 2023. 2
- [16] Chao Li, Yixiao Yang, Kun He, Stephen Lin, and John E Hopcroft. Single image reflection removal through cascaded refinement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3565–3574, 2020. 1, 2, 3, 6
- [17] Tingtian Li, Yuk-Hee Chan, and Daniel PK Lun. Improved multiple-image-based reflection removal algorithm using deep neural networks. *IEEE Transactions on Image Processing*, 30:68–79, 2020. 2
- [18] Yu Li and Michael S Brown. Single image layer separation using relative smoothness. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2752–2759, 2014. 1, 2
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer vision—ECCV 2014: 13th European conference, zurich, Switzerland, September 6–12, 2014, proceedings, part v 13*, pages 740–755. Springer, 2014. 4
- [20] Yunfei Liu, Xingjun Ma, James Bailey, and Feng Lu. Reflection backdoor: A natural backdoor attack on deep neural networks. In *Computer vision—ECCV 2020: 16th European conference, Glasgow, UK, August 23–28, 2020, proceedings, part X 16*, pages 182–199. Springer, 2020. 1
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 5
- [22] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004. 4

- [23] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11461–11471, 2022. 2
- [24] Simon Niklaus, Xuaner Cecilia Zhang, Jonathan T Barron, Neal Wadhwa, Rahul Garg, Feng Liu, and Tianfan Xue. Learned dual-view reflection removal. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3713–3722, 2021. 2
- [25] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021. 2, 4
- [26] 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. 2
- [27] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 4
- [28] YiChang Shih, Dilip Krishnan, Fredo Durand, and William T Freeman. Reflection removal using ghosting cues. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3193–3201, 2015. 1, 2
- [29] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2
- [30] Zhenbo Song, Zhenyuan Zhang, Kaihao Zhang, Wenhan Luo, Zhaoxin Fan, Wenqi Ren, and Jianfeng Lu. Robust single image reflection removal against adversarial attacks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24688–24698, 2023. 1, 2, 6, 7
- [31] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1921–1930, 2023. 2
- [32] Renjie Wan, Boxin Shi, Tan Ah Hwee, and Alex C Kot. Depth of field guided reflection removal. In *2016 IEEE International Conference on Image Processing (ICIP)*, pages 21–25. IEEE, 2016. 1, 2
- [33] Renjie Wan, Boxin Shi, Ling-Yu Duan, Ah-Hwee Tan, and Alex C Kot. Benchmarking single-image reflection removal algorithms. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3922–3930, 2017. 1, 2, 3, 4, 6
- [34] Renjie Wan, Boxin Shi, Haoliang Li, Ling-Yu Duan, and Alex C Kot. Face image reflection removal. *International Journal of Computer Vision*, 129:385–399, 2021. 1
- [35] Jianyi Wang, Zongsheng Yue, Shangchen Zhou, Kelvin CK Chan, and Chen Change Loy. Exploiting diffusion prior for real-world image super-resolution. *International Journal of Computer Vision*, 132(12):5929–5949, 2024. 2
- [36] Tao Wang, Wanglong Lu, Kaihao Zhang, Wenhan Luo, Tae-Kyun Kim, Tong Lu, Hongdong Li, and Ming-Hsuan Yang. Promptrr: Diffusion models as prompt generators for single image reflection removal. *arXiv preprint arXiv:2402.02374*, 2024. 2
- [37] Tianfu Wang, Mingyang Xie, Haoming Cai, Sachin Shah, and Christopher A Metzler. Flash-split: 2d reflection removal with flash cues and latent diffusion separation. *arXiv preprint arXiv:2501.00637*, 2024. 2, 4
- [38] Yufei Wang, Wenhan Yang, Xinyuan Chen, Yaohui Wang, Lanqing Guo, Lap-Pui Chau, Ziwei Liu, Yu Qiao, Alex C Kot, and Bihan Wen. Sinsr: diffusion-based image super-resolution in a single step. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 25796–25805, 2024. 2
- [39] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 12
- [40] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. Uformer: A general u-shaped transformer for image restoration. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 17683–17693, 2022. 2
- [41] Kaixuan Wei, Jiaolong Yang, Ying Fu, David Wipf, and Hua Huang. Single image reflection removal exploiting misaligned training data and network enhancements. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8178–8187, 2019. 2
- [42] Qiang Wen, Yinjie Tan, Jing Qin, Wenxi Liu, Guoqiang Han, and Shengfeng He. Single image reflection removal beyond linearity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3771–3779, 2019. 2
- [43] Jay Whang, Mauricio Delbracio, Hossein Talebi, Chitwan Saharia, Alexandros G Dimakis, and Peyman Milanfar. Deblurring via stochastic refinement. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16293–16303, 2022. 2
- [44] Bin Xia, Yulun Zhang, Shiyin Wang, Yitong Wang, Xinglong Wu, Yapeng Tian, Wenming Yang, and Luc Van Gool. Diffir: Efficient diffusion model for image restoration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13095–13105, 2023. 2
- [45] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22428–22437, 2023. 2
- [46] Guangkai Xu, Yongtao Ge, Mingyu Liu, Chengxiang Fan, Kangyang Xie, Zhiyue Zhao, Hao Chen, and Chunhua Shen. What matters when repurposing diffusion models for general dense perception tasks? *arXiv preprint arXiv:2403.06090*, 2024. 2, 3, 12

- [47] Jie Yang, Dong Gong, Lingqiao Liu, and Qinfeng Shi. Seeing deeply and bidirectionally: A deep learning approach for single image reflection removal. In *Proceedings of the european conference on computer vision (ECCV)*, pages 654–669, 2018. [2](#)
- [48] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10371–10381, 2024. [1](#)
- [49] Chongjie Ye, Lingteng Qiu, Xiaodong Gu, Qi Zuo, Yushuang Wu, Zilong Dong, Liefeng Bo, Yuliang Xiu, and Xiaoguang Han. Stablenormal: Reducing diffusion variance for stable and sharp normal. *ACM Transactions on Graphics (TOG)*, 43(6):1–18, 2024. [2](#), [3](#), [12](#)
- [50] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 3836–3847, 2023. [2](#), [3](#), [4](#), [12](#)
- [51] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. [12](#)
- [52] Xuaner Zhang, Ren Ng, and Qifeng Chen. Single image reflection separation with perceptual losses. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4786–4794, 2018. [1](#), [2](#), [3](#), [6](#)
- [53] Haofeng Zhong, Yuchen Hong, Shuchen Weng, Jinxiu Liang, and Boxin Shi. Language-guided image reflection separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24913–24922, 2024. [1](#), [2](#)
- [54] Yurui Zhu, Xueyang Fu, Peng-Tao Jiang, Hao Zhang, Qibin Sun, Jinwei Chen, Zheng-Jun Zha, and Bo Li. Revisiting single image reflection removal in the wild. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. [1](#), [2](#), [3](#), [6](#), [7](#)

Supplementary materials

A. Details of Loss Functions

In the forward diffusion process, Gaussian noise ε is incrementally applied to the latent representation z of the training image, resulting in a noisy latent representation defined as: $z_t = \sqrt{\alpha_t}z + \sqrt{1 - \alpha_t}\varepsilon$. In the backward denoising process, the traditional diffusion models learn to predict the noise at randomly sampled step t , and denoise with multi-step for inference. The corresponding multi-step loss is formulated as:

$$\mathcal{L}_{multi-step} = \mathbb{E}_{\varepsilon \sim \mathcal{N}(0,1), c, t} \left[\|\varepsilon - \mu_{\theta}^{\varepsilon}(z_t, c, t)\|_2^2 \right], \quad (6)$$

where z_t is the latent representation after adding noise over t steps with $t \in (0, T)$. In Stable Diffusion, T is set to 1,000. The embedding condition c is derived from the text prompt “remove glass reflection”. μ is the U-net with parameters θ , predicting the noise ε as the training target.

The noise introduced during the multi-step diffusion process creates challenges in producing stable results, making traditional diffusion models unsuitable for our reflection removal task, which requires deterministic outputs. In contrast, our approach employs one-step denoising [46, 49], which significantly accelerates inference speed while producing deterministic results. The denoising loss function for our method is formulated as:

$$\mathcal{L}_{one-step} = \mathbb{E}_{z_t, c, t} \left[\|z_t - \mu_{\theta}^{z_t}(z_T, c, t)\|_2^2 \right], \quad (7)$$

During the training process, the noise is added to the max step T , the target is predicting latents at t step with $t \in (0, T)$. By setting $t = 0$ during inference, we obtain the final deterministic result.

We adapt ControlNet [50] as a mechanism to inject information from mixed images. By jointly training the ControlNet with the upsampling blocks of the U-Net architecture, we enable the model to effectively perform reflection removal. Our loss function is formulated as:

$$\mathcal{L}_{diff} = \mathbb{E}_{z_t, c, \mathbf{M}, t} \left[\|z_t - \mu_{\theta}^{z_t}(z_T, c, t, f_{\phi}(\mathcal{E}(\mathbf{M})))\|_2^2 \right], \quad (8)$$

where μ and f denote U-Net and ControlNet with parameters θ and ϕ , respectively. The mixed image \mathbf{M} is first processed by the encoder \mathcal{E} , and the resulting encoded features are subsequently fed as input into the ControlNet. For simplicity, we express the formulation as follows:

$$\mathcal{L}_{diff} = \mathbb{E}_{z_t, \mathbf{M}} \left[\|z_t - \mu_{\theta, \phi}^{z_t}(\mathbf{M})\|_2^2 \right], \quad (9)$$

We introduce a novel reflection-invariant fine-tuning strategy that encourages consistent outputs across varying reflection patterns. Specifically, we first randomly select two different mixed images, \mathbf{M}_1 and \mathbf{M}_2 from our real

dataset DRR and synthetic data. then we use two models with shared weights to process them respectively. The outputs of these models are used to calculate the basic diffusion loss show in Eq. 9. Additionally, we design a novel consistent loss to constrain the outputs of the two models to be similar in the latent space, formulated as:

$$\mathcal{L}_{con} = \mathbb{E}_{\mathbf{M}_1, \mathbf{M}_2} \left[\|\mu_{\theta, \phi}^{z_t}(\mathbf{M}_1) - \mu_{\theta, \phi}^{z_t}(\mathbf{M}_2)\|_2^2 \right], \quad (10)$$

We combine the consistent loss with the two basic diffusion losses using equal weights, formulated as:

$$\mathcal{L} = \mathcal{L}_{diff}(\mathbf{M}_1) + \mathcal{L}_{diff}(\mathbf{M}_2) + \mathcal{L}_{con}(\mathbf{M}_1, \mathbf{M}_2), \quad (11)$$

At the final stage of our progressive training, the cross-latent decoder is trained using a combination of L1 loss, SSIM [39] loss, and LPIPS [51] loss, with equal weighting applied to each component:

$$\mathcal{L}_{rec} = \mathcal{L}_1 + \mathcal{L}_{SSIM} + \mathcal{L}_{LPIPS}. \quad (12)$$

B. More Qualitative Comparison

We provide additional qualitative comparisons between our method and other approaches, as illustrated in Fig. 10 and Fig. 11. Notably, the other approaches are not finetuned on our data setting, which highlights the superior performance of our comprehensive method.

C. Additional examples of DRR

To show the diversity of our dataset, we randomly select several mixed images of one scene, demonstrated in Fig. 12. The ground truth transmission layer is the first image while the remaining images represent various mixed images.

D. Downstream application

Our model can be applied to various downstream tasks, including semantic segmentation, object detection, depth estimation, and normal estimation. The results in Fig. 13 demonstrate that our model can effectively remove reflections and restore the original scene, which is beneficial for subsequent tasks.

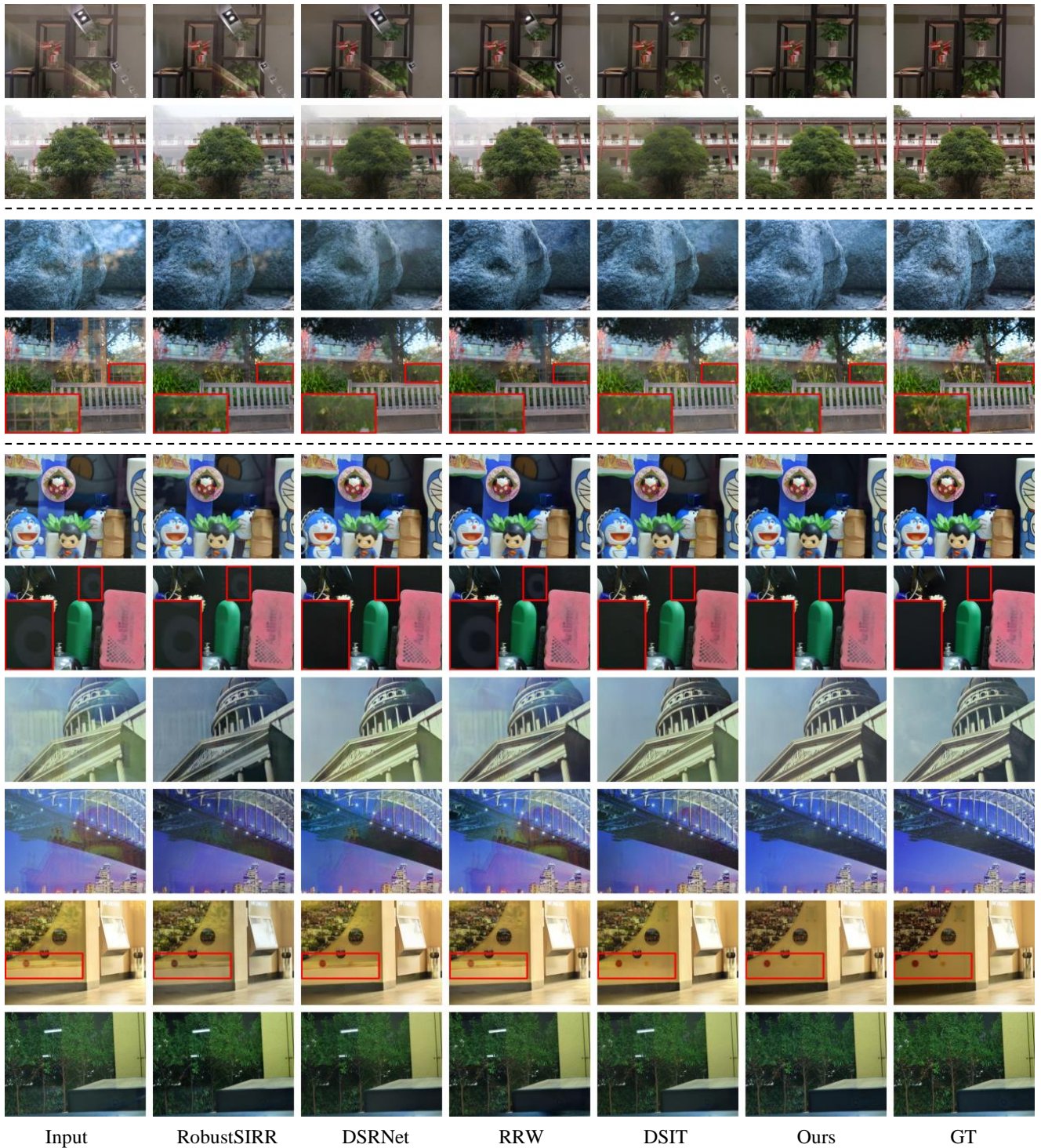


Figure 10. Qualitative comparison of our method with other approaches (not finetuned on our data setting). The benchmark datasets, listed from top to bottom, are *Nature*, *Real*, SIR^2 .

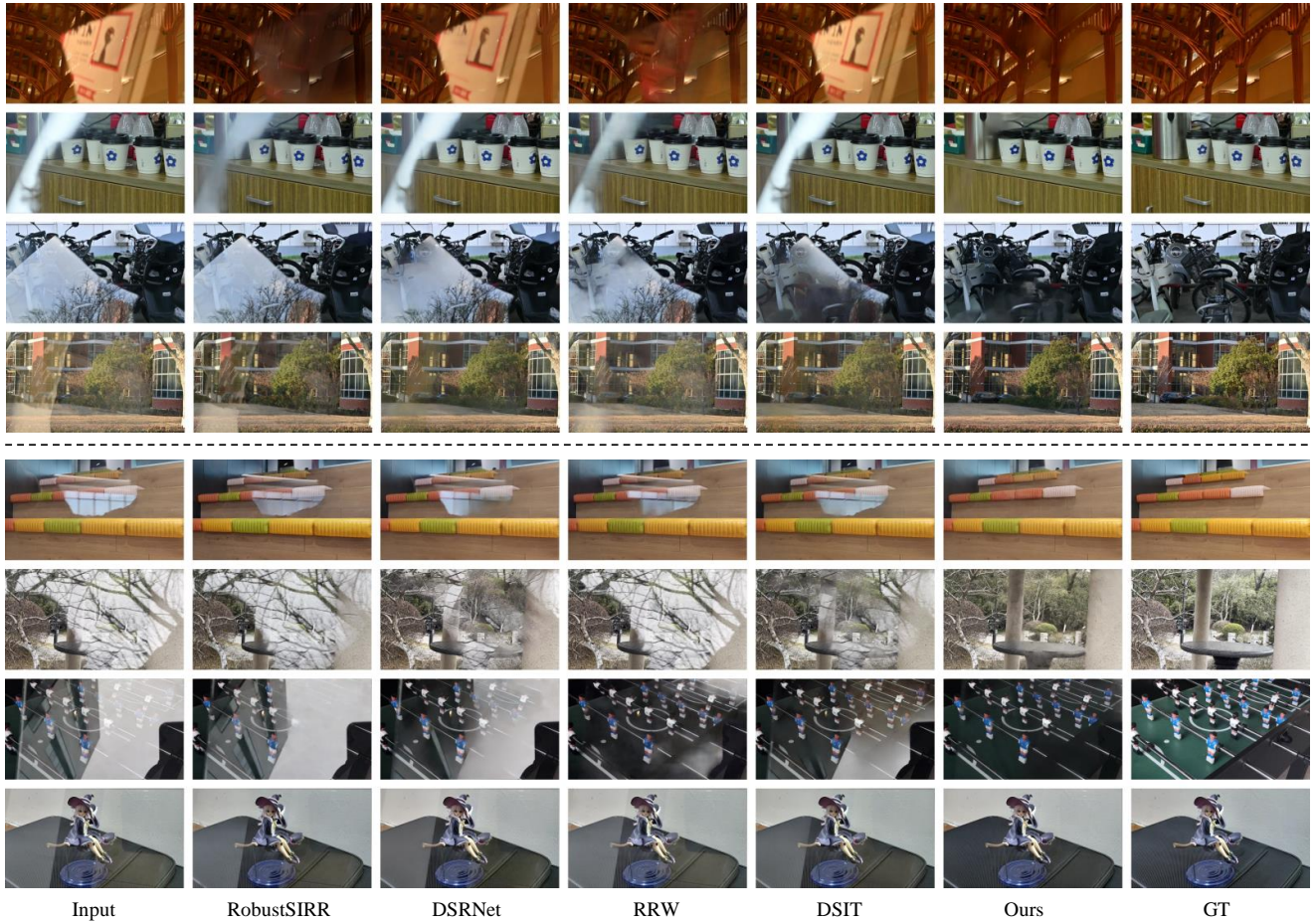


Figure 11. Qualitative comparison of our method with other approaches (not finetuned on our data setting). The benchmark datasets, listed from top to bottom, are *DRR-S* and *DRR-C*.



Figure 12. Additional visual demonstration of DRR. Each scene has diverse reflection patterns, varying glass angles and reflection intensity. The ground truth transmission layer is the first image while the remaining images represent various mixed images.

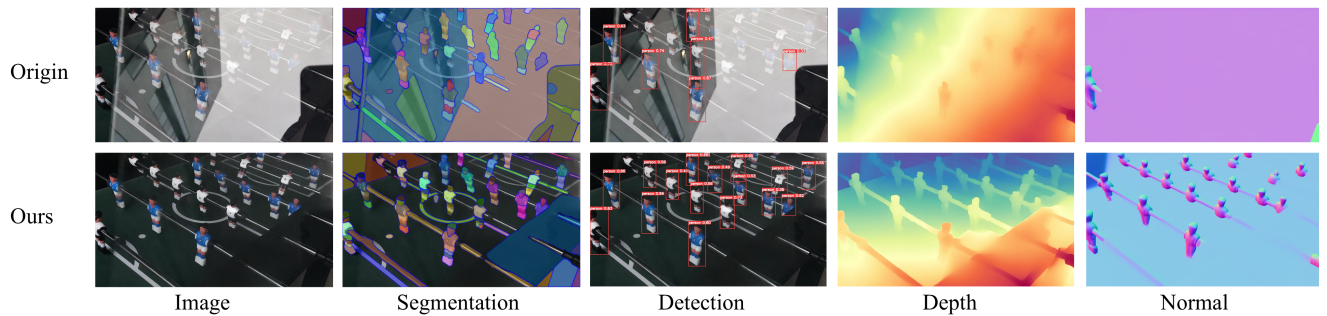


Figure 13. Various downstream tasks, including semantic segmentation, object detection, depth estimation, and normal estimation.