

ICM-SR: Image-Conditioned Manifold Regularization for Image Super-Resolution

Junoh Kang^{*1} Donghun Ryu^{*2} Bohyung Han^{1,2}
 Computer Vision Laboratory, ¹ECE & ²IPAI, Seoul National University
 {junoh.kang, dhryou, bhhan}@snu.ac.kr

Abstract

Real world image super-resolution (Real-ISR) often leverages the powerful generative priors of text-to-image diffusion models by regularizing the output to lie on their learned manifold. However, existing methods often overlook the importance of the regularizing manifold, typically defaulting to a text-conditioned manifold. This approach suffers from two key limitations. Conceptually, it is misaligned with the Real-ISR task, which is to generate high quality (HQ) images directly tied to the low quality (LQ) images. Practically, the teacher model often reconstructs images with color distortions and blurred edges, indicating a flawed generative prior for this task. To correct these flaws and ensure conceptual alignment, a more suitable manifold must incorporate information from the images. While the most straightforward approach is to condition directly on the raw input images, their high information densities make the regularization process numerically unstable. To resolve this, we propose image-conditioned manifold regularization (ICM), a method that regularizes the output towards a manifold conditioned on the sparse yet essential structural information: a combination of colormap and Canny edges. ICM provides a task-aligned and stable regularization signal, thereby avoiding the instability of dense-conditioning and enhancing the final super-resolution quality. Our experiments confirm that the proposed regularization significantly enhances super-resolution performance, particularly in perceptual quality, demonstrating its effectiveness for real-world applications. We will release the source code of our work for reproducibility.

1. Introduction

Image Super-Resolution (ISR), which aims to restore a high-quality (HQ) image from its low-quality (LQ) counterpart, is a classical problem in computer vision. While recent advances in deep learning have significantly improved the

^{*}indicates equal contribution.

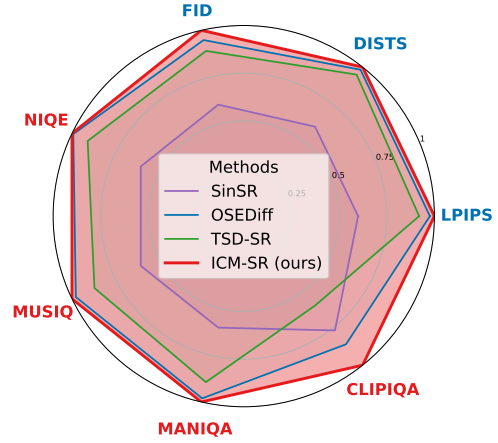


Figure 1. Performance comparison on DRealSR benchmark [34]. The red and blue metrics are no-reference and reference perceptual metrics, respectively. ICM-SR (ours) stands out for perceptual metrics, highlighting its strong performance in practical scenarios.

ISR performance [5, 8, 14, 18, 19], they often fail to generalize to the diverse and unknown degradations encountered in real-world scenarios. To address this limitation, real-world image super-resolution (Real-ISR) [31, 44] aims to achieve practical super-resolution by applying significantly more diverse and complex degradation pipelines. Since training models solely with a pixel-wise reconstruction loss inevitably leads to blurry and oversmoothed results, training frameworks from generative models such as Generative Adversarial Networks (GANs) [7] and diffusion models [10, 25, 26] are adopted for Real-ISR. Both GAN-based methods [15, 30, 31] and diffusion-based methods [20, 43] enable the generation of more realistic and sharp images with superior perceptual quality compared to the models trained only with reconstruction loss.

The emergence of generative foundation models has opened new avenues for Real-ISR, leveraging the powerful generative priors of pretrained text-to-image diffusion models [25]. One approach [36, 42] is to adapt pretrained diffu-

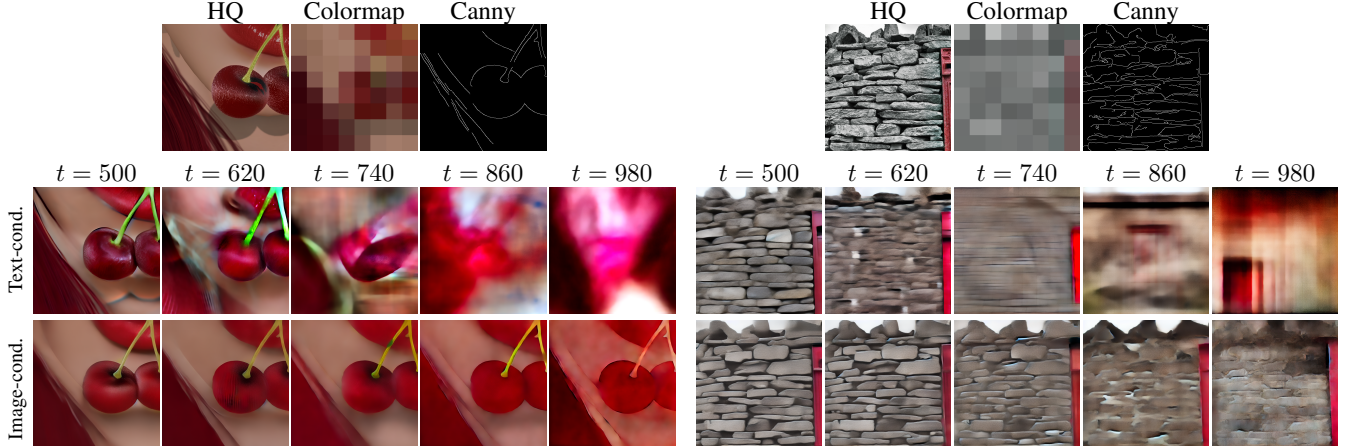


Figure 2. Visualization of denoised latents from the teacher diffusion model. We add noise corresponding to timestep t to the ground-truth latent and visualize the model’s single-step denoised prediction. **(Text-cond.)** The standard text-conditioned prior struggles to reconstruct the image from noisy latents, especially at large t . It produces oversaturated colors (left) and fails to recover edges (right). **(Image-cond.)** In contrast, our proposed image-conditioned prior, guided by a colormap and Canny edges, provides a much more accurate prediction. It consistently reconstructs latents with faithful color and sharp structural details, demonstrating a more stable and task-aligned generative prior.

sion models to the Real-ISR task by training LoRA [11] or ControlNet [46]. This approach preserves the powerful generative priors inherent in the text-to-image models, thereby achieving superior generalization capabilities. Despite their impressive performance, standard diffusion models for ISR often require high computational cost due to their iterative sampling process. For efficient inference, researches [6, 16, 35, 41] have explored distilling generative priors into one-step super-resolution models. They adopt distillation techniques for diffusion models such as distribution matching [22, 40] and consistency trajectory matching [13, 27]. Among these efficient models, OSEDiff [35] regularizes the super-resolution outputs towards the natural image prior embedded within the pretrained diffusion models, facilitated by Variational Score Distillation (VSD) [33].

However, many existing one-step Real-ISR methods focus on applying and enhancing distillation techniques developed for general-purpose image generation. This leads them to overlook a more fundamental aspect: choosing a target manifold that aligns with the characteristics of the Real-ISR task. Specifically, these methods regularize output towards a manifold conditioned on text prompts. This approach creates a conceptual mismatch, as the Real-ISR task requires generating an HQ image that is faithful to the LQ input, not just plausible based on a text description. Moreover, as visualized in Figure 2 (Text cond.), the text-conditioned teacher models often reconstruct images with saturated color and blunt boundaries. This indicates that the generative prior is not only misaligned but also practically flawed for the precise task of image restoration. An intuitive solution to resolve this mismatch is conditioning

the target manifold on the LQ images. However, we prove that conditioning on the information-dense signal causes VSD [33] to become numerically unstable and degenerate towards SDS [24], thereby harming the distillation performance.

To address the dilemma between conceptual alignment and distillation stability, we propose Image-Conditioned Manifold (ICM) regularization. Our method conditions the target manifold on core structural information, which we compose from a low-resolution colormap and Canny edges. This combination is specifically designed to resolve the aforementioned practical failure of the text-conditioned prior; the colormap provides global guidance to prevent color shifts, while Canny edges enforce sharp structural details. We implement this conditioning using a pretrained T2I-Adapter [23], and this clearly mines more stable and accurate prior from diffusion models as shown in Figure 2 (bottom row).

ICM regularization offers two key advantages. Conceptually, ICM provides a regularization manifold that is fundamentally better aligned with the objectives of Real-ISR. Practically, the structural conditioning improves score estimation accuracy, especially at large diffusion timesteps. Consequently, this synergy of conceptual alignment and practical stability allows ICM regularization to yield superior one-step diffusion models for Real-ISR.

Our key contributions are summarized as follows:

- We propose Image-Conditioned Manifold (ICM) regularization, highlighting the overlooked importance of the target manifold for regularization. Our key idea is to condition the manifold on image information, which conceptually

ally aligns the generative prior to the Real-ISR task.

- We prove that conditioning on information-dense signals causes the VSD loss to degenerate towards SDS, leading to numerical instability. Our method, ICM regularization, resolves this by conditioning only on core structural information. This achieves the desired alignment of the prior while ensuring a stable distillation process.
- We show the effectiveness of ICM regularization through extensive experiments on Real-ISR benchmarks, confirming ICM-SR’s superior performance, particularly in perceptual quality, validating the practicality of the proposed approach for real-world applications.

2. Related Work

2.1. Real-ISR for Perceptual Enhancement

Early deep learning-based super-resolution methods focused on pixel-wise accuracy, often resulting in blurry outputs for real-world images. To address this, the field of Real-World Image Super-Resolution (Real-ISR) emerged, prioritizing perceptual quality [31, 44]. A key advancement was the adoption of Generative Adversarial Networks (GANs), pioneered by SRGAN [15]. Subsequent methods like Real-ESRGAN [31] and BSRGAN [44] further improved photo-realism by training on more complex and realistic degradation models.

2.2. Multi-step Diffusion-based Real-ISR

Recently, diffusion models have set a new standard for perceptual quality in Real-ISR. These methods leverage powerful priors from pre-trained text-to-image models, such as Stable Diffusion [25], by conditioning the generation process on the low-resolution (LQ) input. Representative works include StableSR [29], DiffBIR [20], SeeSR [36], and SUPIR [42], which employ various fine-tuning or adapter-based strategies for conditioning. While demonstrating impressive perceptual quality, a major drawback of these multi-step diffusion methods is their high inference cost, requiring numerous sampling steps and resulting in slow processing speeds.

2.3. Efficient Diffusion Models for Real-ISR

The considerable inference time of multi-step diffusion models necessitates acceleration techniques. Distillation has emerged as a primary strategy to reduce sampling steps, often targeting one-step generation for diffusion-based ISR. Noteworthy efficient methods include SinSR [32] using consistency preserving distillation, AddSR [37] employing adversarial diffusion distillation, and OSSEDiff [35] which proposes a Variational Score Distillation (VSD) loss in the latent space. More recently, TSD-SR [6] introduced Target Score Distillation (TSD) and a Distribution-Aware Sampling Method (DASM) to improve score estimation stabil-

ity. Although these methods significantly reduce inference steps, challenges persist in achieving perfect score distillation across all timesteps and preventing artifacts, highlighting the need for further improvements in efficient diffusion-based Real-ISR.

3. Preliminary

Real world image super-resolution Real-ISR is a task to reconstruct high quality (HQ) image \mathbf{x}_H from low quality (LQ) image \mathbf{x}_L . Given a dataset $\mathcal{D} = \{(\mathbf{x}_L, \mathbf{x}_H)\}_{i=1}^N$, a super-resolution model G_θ is trained to minimize

$$\mathbb{E}_{(\mathbf{x}_L, \mathbf{x}_H) \sim \mathcal{D}} [\mathcal{L}_{\text{Rec}}(G_\theta(\mathbf{x}_L), \mathbf{x}_H) + \mathcal{L}_{\text{Reg}}(G_\theta(\mathbf{x}_L))], \quad (1)$$

where \mathcal{L}_{Rec} is the reconstruction loss such as MSE or LPIPS [47], and \mathcal{L}_{Reg} is the regularization loss forcing $G_\theta(\mathbf{x}_L)$ to lie on the desired real HQ image manifold.

Regularizing towards a generative prior To ensure the output to be photorealistic, the regularization loss \mathcal{L}_{Reg} regularizes the model’s output distribution towards the powerful generative prior learned by large-scale text-to-image diffusion models. Following prior distillation works [33, 40], OSSEDiff [35] implements this by defining the desired manifold as the distribution of high-quality images conditioned on text prompts \mathbf{c}_t . The regularization loss is formalized using the Variational Score Distillation (VSD) loss:

$$\begin{aligned} \mathcal{L}_{\text{Reg}}(G_\theta(\mathbf{x}_L)) &= \mathcal{L}_{\text{VSD}}(G_\theta(\mathbf{x}_L), \mathbf{c}_t) \\ &= \int_0^T w(t) D_{\text{KL}}(q_t^\theta(\hat{\mathbf{x}}_H | \mathbf{c}_t) || p_t^{\text{real}}(\mathbf{x}_H | \mathbf{c}_t)) dt, \end{aligned} \quad (2)$$

where \mathbf{c}_t is text prompt describing \mathbf{x}_H , q_t^θ is a perturbed distribution of $\hat{\mathbf{x}}_H = G_\theta(\mathbf{x}_L)$ and p_t^{real} is a noisy distribution learned by pretrained diffusion models. As the computation of \mathcal{L}_{VSD} is intractable, the model is optimized by the gradient of \mathcal{L}_{VSD} with respect to θ . The gradient is expressed with two score functions estimated via diffusion models:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{\text{VSD}} &= \nabla_\theta \int_0^T w(t) D_{\text{KL}}(q_t^\theta(\hat{\mathbf{x}}_H | \mathbf{c}_t) || p_t^{\text{real}}(\mathbf{x}_H | \mathbf{c}_t)) dt \\ &= \nabla_\theta \int_0^T w(t) \mathbb{E}_{\hat{\mathbf{z}}_t \sim q_t^\theta} [\log q_t^\theta(\hat{\mathbf{z}}_t | \mathbf{c}_t) - \log p_t^{\text{real}}(\hat{\mathbf{z}}_t | \mathbf{c}_t)] dt \\ &= \int_0^T w(t) \mathbb{E}_{\hat{\mathbf{z}}_t} \left[(\nabla \log q_t^\theta(\hat{\mathbf{z}}_t | \mathbf{c}_t) - \nabla \log p_t^{\text{real}}(\hat{\mathbf{z}}_t | \mathbf{c}_t)) \frac{\partial \hat{\mathbf{z}}_t}{\partial \theta} \right] dt \\ &\approx \int_0^T w'(t) \mathbb{E}_{\hat{\mathbf{z}}_t} [\epsilon_\phi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t) - \epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t)] \frac{\partial \hat{\mathbf{z}}_t}{\partial \theta} dt, \end{aligned} \quad (3)$$

where $\hat{\mathbf{z}}_t$ denotes perturbed latents (encoded image), ϵ_ψ is a diffusion model trained to learn distribution of $G_\theta(\mathbf{x}_L)$, and ϵ_ϕ is a pretrained diffusion model that estimates the score function of real-world, high-fidelity image distribution.

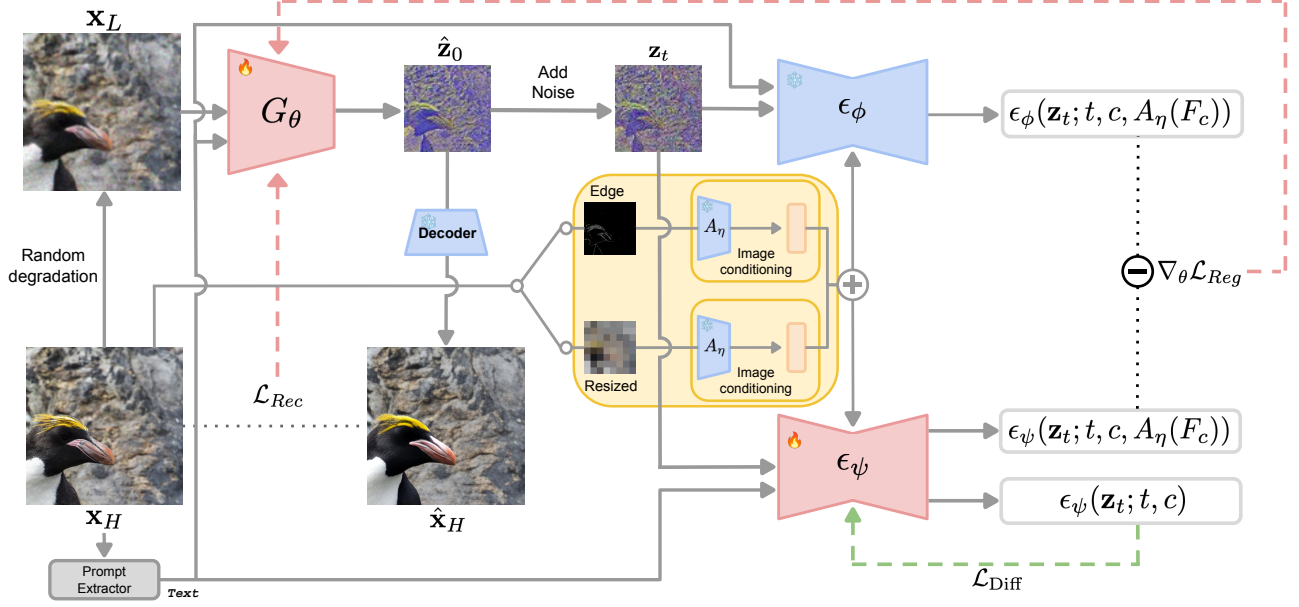


Figure 3. Training framework of ICM-SR. Our framework trains a one-step super-resolution generator using two main losses. A reconstruction loss \mathcal{L}_{Rec} ensures fidelity to the ground-truth \mathbf{x}_H . For realism, a VSD-based regularization loss \mathcal{L}_{Reg} is applied, which involves a frozen pre-trained diffusion model ϵ_ϕ and a trainable auxiliary model ϵ_ψ . The key innovation of our method is to condition the target manifold on structural information \mathbf{F}_c (e.g., edges, resized image) from the HQ image \mathbf{x}_H . These conditions are encoded by T2I-Adapter A_η and then injected into both ϵ_ϕ and ϵ_ψ to guide the generator G_θ towards producing outputs that are not only realistic but also structurally aligned with the target image.

4. Method

4.1. Image-conditioned manifold for Real-ISR

The dilemma of conditioning for Real-ISR A core principle for successful regularization is the proper alignment of prior knowledge with the task’s objectives. However, existing regularization methods for Real-ISR have not addressed the choice of the manifold, often defaulting to a text-conditioned manifold inherited from text-to-image models. This default choice creates a fundamental mismatch with the objective of Real-ISR: generating an HQ image that is a direct enhancement of the given LQ image, not an arbitrary one that merely matches a text description.

This mismatch naturally suggests that a suitable manifold must be conditioned on image information. While the most straightforward solution is to condition directly on the raw LQ image, our analysis reveals that this approach leads to a critical problem. Conditioning on such an information-dense signal over-constrains the target manifold, causing the VSD to become numerically unstable and degenerate towards SDS [24], thereby harming distillation performance. We formalize this instability in the following lemma:

Lemma 1. *Let \mathbf{c} be a strong condition such that the latent variable \mathbf{z}_0 is deterministic, i.e., $\mathbf{z}_0|\mathbf{c} = \mu(\mathbf{c})$. If the perturbed distribution $q_t(\mathbf{z}_t|\mathbf{c})$ is generated by $\mathbf{z}_t = a_t\mathbf{z}_0 + b_t\epsilon$ with $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then the auxiliary denoiser ϵ_ψ collapses*

to the sampled noise ϵ :

$$\epsilon_\psi(\mathbf{z}_t; t, \mathbf{c}) = \epsilon. \quad (4)$$

Consequently, the gradient of the VSD loss $\nabla_\theta \mathcal{L}_{\text{VSD}}$ degenerates to the gradient of the SDS loss $\nabla_\theta \mathcal{L}_{\text{SDS}}$. In other words, the score prediction difference becomes:

$$\epsilon_\phi(\mathbf{z}_t; t, \mathbf{c}) - \epsilon_\psi(\mathbf{z}_t; t, \mathbf{c}) = \epsilon_\phi(\mathbf{z}_t; t, \mathbf{c}) - \epsilon, \quad (5)$$

where ϵ_ϕ is the denoiser of real high-quality images.

Proof. $\mathbf{z}_0|\mathbf{c} \sim \delta(\mathbf{z}_0 - \mu(\mathbf{c}))$ implies $\mathbf{z}_t|\mathbf{c} \sim \mathcal{N}(\mu(\mathbf{c}), b_t^2\mathbf{I})$. Hence, $\nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t|\mathbf{c}) = -\frac{\mathbf{z}_t - a_t\mu(\mathbf{c})}{b_t^2}$ holds, which leads to

$$\begin{aligned} \epsilon_\psi(\mathbf{z}_t; t, \mathbf{c}) &= -b_t \nabla_{\mathbf{z}_t} \log q_t(\mathbf{z}_t|\mathbf{c}) \\ &= (-b_t) \times \left(-\frac{\mathbf{z}_t - a_t\mu(\mathbf{c})}{b_t^2} \right) \\ &= \frac{a_t\mathbf{z}_0 + b_t\epsilon - a_t\mu(\mathbf{c})}{b_t} = \epsilon. \end{aligned}$$

□

An effective manifold for Real-ISR regularization must satisfy two criteria: (1) it must be conditioned on image information to be conceptually aligned, and (2) it must not be overly restrictive to ensure distillation stability.

The design of core structural information \mathbf{F}_c To satisfy both criteria simultaneously, our guiding principle is to extract and condition on the core identity of the target image, while discarding fine-grained, high-density information. The design of this core structural information, \mathbf{F}_c , is motivated by the observed practical failures of the text-conditioned prior. As visualized in Figure 2, the teacher model guided only by text often reconstructs outputs with oversaturated colors and blurred edges. This indicates that the text-conditioned prior is not only conceptually misaligned but also practically flawed for the image restoration. To address the failures while adhering to the two criteria, we design the core structural information as a combination of a low-resolution colormap and Canny edges. The 8x8 colormap guides the teacher model to extract priors aligned with the target’s color distribution, while the Canny edges enforce the elicitation of priors rich in structural details. This combination of sparse yet essential information provides the necessary conceptual alignment from image conditioning (criterion 1) without the distillation instability caused by the high information density (criterion 2).

Image-conditioned manifold regularization We formalize the concept of conditioning on the core structural information with our proposed Image-Conditioned Manifold (ICM) regularization loss. This loss regularizes the super-resolution outputs to lie on the manifold conditioned on both the structural information and the text prompt, $p_t^{\text{real}}(\mathbf{x}_H|\mathbf{F}_c, \mathbf{c})$:

$$\mathcal{L}_{\text{ICM}} = \int_0^T w(t) D_{\text{KL}}(q_t^\theta(\hat{\mathbf{x}}_H|\mathbf{F}_c, \mathbf{c}_t) || p_t^{\text{real}}(\mathbf{x}_H|\mathbf{F}_c, \mathbf{c}_t)) dt. \quad (6)$$

The gradient of this loss with respect to θ can be derived similarly to Equation (3), yielding:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{\text{ICM}} \\ \approx \int_0^T w'(t) \mathbb{E}_{\hat{\mathbf{z}}_t \sim q_t^\theta} [\epsilon_\phi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, \mathbf{F}_c) - \epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, \mathbf{F}_c)] dt. \end{aligned} \quad (7)$$

4.2. Training framework of ICM-SR

The overall training framework of ICM-SR is designed to regularize the super-resolution output towards the image-conditioned manifold. As depicted in Figure 3 and Algorithm 1, the process involves training G_θ with two main losses and training an auxiliary diffusion model, ϵ_ψ , which learns the distribution of current super-resolution output, $G_\theta(\mathbf{x}_L)$. The key components of our framework are detailed below.

4.2.1. Training super-resolution model G_θ

The super-resolution model G_θ is trained with a combined loss function, which consists of a reconstruction loss and a regularization loss.

Algorithm 1 ICM-SR

```

1: Require:  $G_\theta, \epsilon_\phi, \epsilon_\psi, A_\eta, \text{Dec}, \text{Scheduler}, F$ 
2: while train do
3:    $(\mathbf{x}_L, \mathbf{x}_H, \mathbf{c}_t) \sim \mathcal{D}$ 
4:    $\hat{\mathbf{z}}_0 \leftarrow G_\theta(\mathbf{x}_L), \mathbf{F}_c \leftarrow F(\mathbf{x}_H)$ 
5:   /* Compute reconstruction loss
6:    $\hat{\mathbf{x}}_H \leftarrow \text{Dec}(\hat{\mathbf{z}}_0)$ 
7:    $\mathcal{L}_{\text{Rec}} \leftarrow \mathcal{L}_2(\hat{\mathbf{x}}_H, \mathbf{x}_H) + \mathcal{L}_{\text{LPIPS}}(\hat{\mathbf{x}}_H, \mathbf{x}_H)$ 
8:   /* Compute regularization gradient
9:    $t \sim \mathcal{U}(20, 980), \epsilon \sim \mathcal{N}(0, \mathbf{I})$ 
10:   $a_t, b_t \leftarrow \text{Scheduler}(t)$ 
11:   $\hat{\mathbf{z}}_t \leftarrow a_t \hat{\mathbf{z}}_0 + b_t \epsilon$ 
12:   $\epsilon_{\text{fake}} \leftarrow \text{stopgrad}(\epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, A_\eta(\mathbf{F}_c)))$ 
13:   $\epsilon_{\text{real}} \leftarrow \text{stopgrad}(\text{cfg}(\epsilon_\phi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, A_\eta(\mathbf{F}_c))))$ 
14:   $\nabla_\theta \mathcal{L}_{\text{Reg}} \leftarrow w(t)(\epsilon_{\text{fake}} - \epsilon_{\text{real}})$ 
15:  /* Compute auxiliary diffusion loss
16:   $\hat{\mathbf{z}}_t \leftarrow \text{stopgrad}(\hat{\mathbf{z}}_t)$ 
17:   $\mathcal{L}_{\text{Aux}} \leftarrow \|\epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t) - \epsilon\|_2^2$ 
18:  /* Update generator  $G_\theta$  and
    auxiliary diffusion model  $\epsilon_\psi$ 
19:  Update  $\theta$  with  $\mathcal{L}_{\text{Rec}}$  and  $\nabla_\theta \mathcal{L}_{\text{Reg}}$ 
20:  Update  $\psi$  with  $\mathcal{L}_{\text{Aux}}$ 
21: end while
```

Reconstruction loss Following OSediff [35], the reconstruction loss measures the pixel-wise and perceptual difference between the generated and ground-truth images. It is defined as the sum of ℓ_2 and LPIPS losses:

$$\mathcal{L}_{\text{Rec}}(\mathbf{x}_H, \hat{\mathbf{x}}_H) = \|\hat{\mathbf{x}}_H - \mathbf{x}_H\|_2^2 + \mathcal{L}_{\text{LPIPS}}(\hat{\mathbf{x}}_H, \mathbf{x}_H), \quad (8)$$

for $\hat{\mathbf{x}}_H = \text{Dec}(G_\theta(\mathbf{x}_L))$.

ICM regularization loss The ICM regularization encourages the output to lie on the proposed manifold, $\mathbf{x}_H|\mathbf{F}_c, \mathbf{c}_t$. As derived in Equation (7), the gradient used for the optimization is:

$$\begin{aligned} \nabla_\theta \mathcal{L}_{\text{Reg}}(\hat{\mathbf{z}}_t, \mathbf{x}_L, t, \mathbf{c}_t, \mathbf{F}_c) \\ = w(t) [\epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, A_\eta(\mathbf{F}_c)) - \epsilon_\phi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, A_\eta(\mathbf{F}_c))], \end{aligned} \quad (9)$$

where $\hat{\mathbf{z}}_t = a_t \hat{\mathbf{z}}_0 + b_t \epsilon$ for a_t and b_t are predefined noise scheduling, ϵ_ϕ is the pretrained teacher model, ϵ_ψ is the auxiliary student model, and the core structural information \mathbf{F}_c is injected using the T2I-Adapter A_η .

4.2.2. Training auxiliary diffusion model ϵ_ψ

The auxiliary diffusion model with trainable LoRA [11] learns the distribution of the generator’s current output. Therefore, its objective is solely based on this output without conditioning on the external map:

$$\mathcal{L}_{\text{Aux}}(\hat{\mathbf{z}}_t, t, \mathbf{c}_t) = \|\epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t) - \epsilon\|_2^2. \quad (10)$$

Table 1. Quantitative comparison with state-of-the-art one-step super-resolution methods on both synthetic (DIV2K-Val) and real-world benchmarks (DrealSR, RealSR). ‘†’ indicates models trained by us for fair comparison. The best and second best results are highlighted in **red** and **blue**, respectively.

Datasets	Methods	Perceptual w/ ref.			Perceptual w/o ref.						Fidelity w/ ref.	
		LPIPS↓	DISTS↓	FID↓	NIQE↓	MUSIQ↑	MANIQ↑	CLIPQ↑	TOPIQ↑	LIQE↑	PSNR↑	SSIM↑
DIV2K-Val	SinSR	0.3220	0.2044	37.29	5.7861	63.28	0.5411	0.6537	0.5796	3.5015	24.29	0.6012
	AddSR	0.3779	0.2174	30.58	5.2324	58.91	0.5260	0.5337	0.5218	3.3877	23.10	0.5928
	OSDiff†	0.2847	0.1905	26.15	4.4918	67.73	0.6081	0.6394	0.6014	4.1704	23.40	0.6160
	TSD-SR†	0.2759	0.1894	25.45	4.6859	65.06	0.5935	0.6306	0.6252	3.6920	24.68	0.6257
	ICM-SR	0.2799	0.1861	24.72	4.4411	68.00	0.6169	0.6440	0.6138	4.2094	23.77	0.6173
DrealSR	SinSR	0.3537	0.2495	177.37	6.7022	57.14	0.5021	0.6590	0.5369	3.2267	28.12	0.7533
	AddSR	0.3079	0.2207	154.24	7.4512	53.47	0.4770	0.5451	0.4840	2.8429	27.92	0.7880
	OSDiff†	0.2974	0.2183	130.50	6.4852	65.46	0.5975	0.6840	0.5961	4.0183	25.57	0.7705
	TSD-SR†	0.2869	0.2134	127.29	6.5985	62.13	0.5741	0.6614	0.6028	3.5815	28.25	0.7885
	ICM-SR	0.2871	0.2142	125.30	6.4163	65.96	0.6051	0.6929	0.6088	4.0977	26.85	0.7763
RealSR	SinSR	0.3050	0.2325	135.51	6.0516	60.82	0.5423	0.6212	0.5321	3.1935	26.15	0.7385
	AddSR	0.3141	0.2198	135.58	6.1957	62.65	0.5620	0.5141	0.5532	3.3821	24.22	0.7032
	OSDiff†	0.2637	0.2023	112.26	5.5930	68.16	0.6246	0.6271	0.5998	4.0785	24.22	0.7276
	TSD-SR†	0.2699	0.2049	116.17	5.6907	66.05	0.6057	0.6106	0.6210	3.6515	25.92	0.7366
	ICM-SR	0.2611	0.2009	108.74	5.5842	68.59	0.6288	0.6360	0.6094	4.1336	24.99	0.7309

5. Experiment

5.1. Experimental Settings

Training details For training dataset, we utilize a combined dataset comprising DIV2K [1] and LSDIR [17]. Following common practice in real-world super-resolution, low-quality (LQ) images are synthesized from their high-quality (HQ) counterparts using the degradation pipeline of Real-ESRGAN [31]. During the training process, we randomly crop 512×512 patches from the HR images. The model is optimized using the AdamW optimizer [21] for a total of 100K iterations with a batch size of 4. The initial learning rate is set to 5×10^{-5} and is adjusted using a cosine learning rate scheduler. For the LoRA [11] components integrated into our model, we set the rank to 4. We utilize pretrained Stable Diffusion v1-4¹ and T2I-Adapter² trained by TencentARC for ICM-SR. The conditioning function, F_c , generates guidance either through Canny edge detection or by downsampling an HQ image. For the downsampling condition, the image is first resized to 8×8 pixels and subsequently upsampled to 512×512 pixels via nearest-neighbor interpolation to align with the T2I-Adapter’s input requirements. All trainings were conducted using four NVIDIA RTX A6000 GPUs.

Evaluation details We evaluate the performance of our method on three benchmark datasets: DIV2K [1] validation dataset, RealSR [2], and DRealSR [34] provided by Sta-

bleSR [29]. For a thorough evaluation, we utilize a diverse set of both full-reference and no-reference image quality metrics. For full-reference evaluation, we use PSNR and SSIM for fidelity, LPIPS [47] and DISTS [4] for perceptual similarity, and FID [9] for distribution similarity. For no-reference evaluation, we utilize NIQE [45], MUSIQ [12], MANIQA [38], and CLIPQA [28], TOPIQ [3], and LIQE [48].

5.2. Quantitative comparison

Tab. 1 presents a quantitative comparison of the proposed method, ICM-SR, with the state-of-the-art one-step super-resolution models on both synthetic and real-world benchmarks. For a fair comparison, we reproduce OSDiff [35] and TSD-SR [6] using the same backbone (Stable Diffusion v1-4) and training dataset.

As evident in Tab. 1, ICM-SR consistently excels across all benchmarks, particularly for both full-reference perceptual metrics (*e.g.*, LPIPS, DISTS, FID) and no-reference quality metrics (*e.g.*, MUSIQ, MANIQA). ICM-SR’s strong performance in both metric categories highlights its superior ability to harmonize fidelity with perceptual realism. In contrast, while TSD-SR also improves upon OSDiff in full-reference metrics, it struggles with no-reference metrics. This suggests that TSD-SR’s approach, while enhancing fidelity, may fail to fully leverage the generative prior’s capacity for producing aesthetically pleasing images.

The strong performance of ICM-SR extends to the challenging real-world datasets, such as DrealSR and RealSR, where low-quality (LQ) images are sourced from practical scenarios. This robust performance on practical, unknown degradations confirms the effectiveness of image-

¹<https://huggingface.co/CompVis/stable-diffusion-v1-4>

²https://huggingface.co/TencentARC/t2iadapter_color_sd14v1

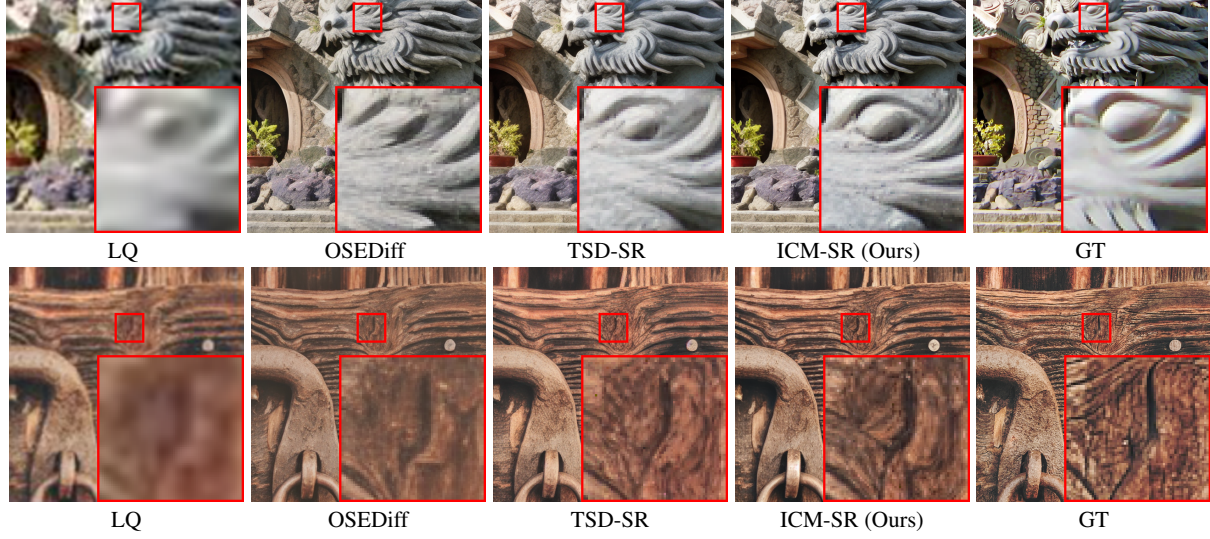


Figure 4. Qualitative results of our method compared to OSERDiff and TSD-SR on the Div2k validation dataset. Our method demonstrates superior performance in recovering fine details. Zoom in for better visualization.

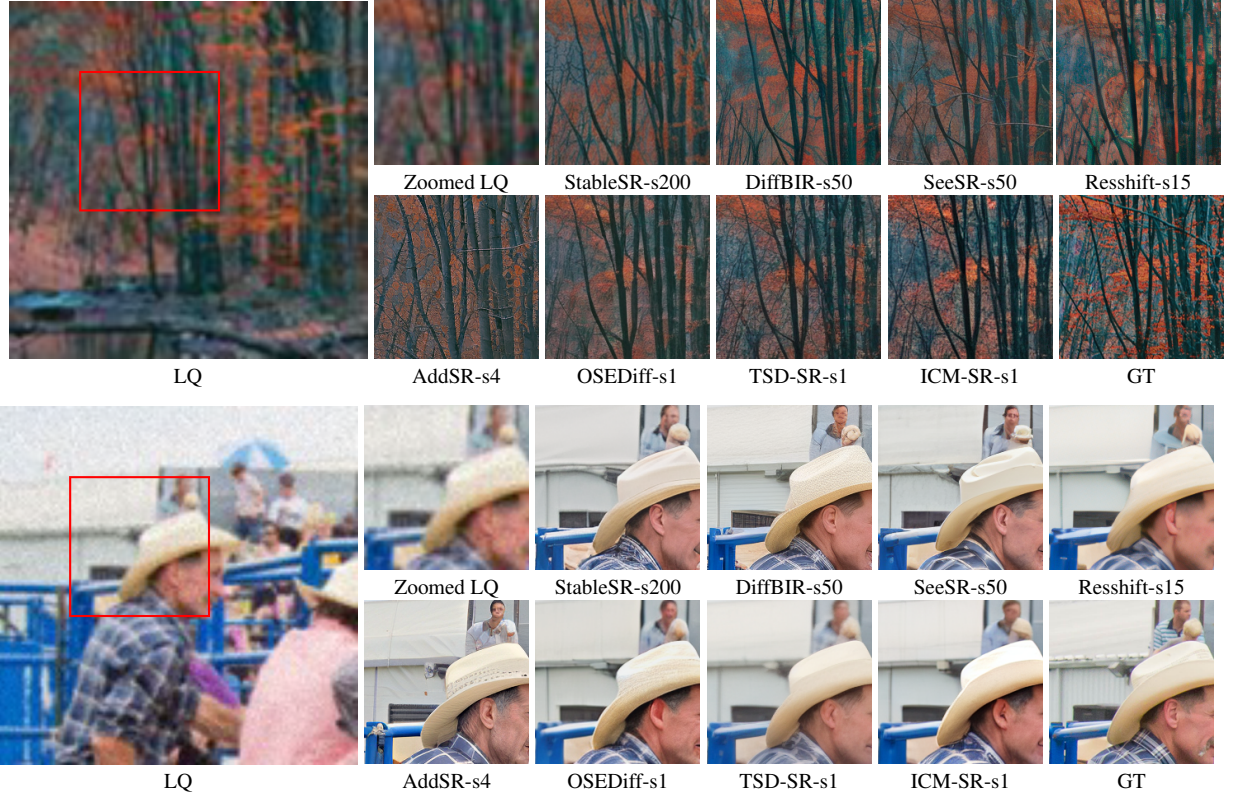


Figure 5. Qualitative comparison of our method with various multi-step and one-step diffusion-based methods. ‘s’ denotes the number of network inferences in the method. Zoom in for better visualization.

conditioned manifold regularization. By providing a more stable and task-aligned regularization signal, ICM-SR successfully generates faithful and visually appealing super-resolution results. For a comparison to multi-step methods, please refer to the supplementary material.

5.3. Qualitative comparison

Figure 4 provides a qualitative comparison of ICM-SR, OSERDiff, and TSD-SR, highlighting ICM-SR’s ability to produce highly expressive textures that are faithful to the

Table 2. Ablation studies of ICM-SR on the choice of structural image condition \mathbf{F}_c . The best and second best results of each metric are highlighted in **red** and **blue**, respectively.

Datasets	Methods	Perceptual w/ ref.			NIQE↓	MUSIQ↑	Perceptual w/o ref.				Fidelity w/ ref.	
		LPIPS↓	DISTS↓	FID↓			MANIQ↑	CLIPQ↑	TOPIQ↑	LIQE↑	PSNR↑	SSIM↑
DIV2K-Val	OSDiff†	0.2847	0.1905	26.15	4.4918	67.73	0.6081	0.6394	0.6014	4.1704	23.40	0.6160
	ICM-SR-lq	0.3042	0.1977	28.32	4.7224	67.47	0.6135	0.6446	0.5794	4.1874	22.72	0.6044
	ICM-SR-color	0.2849	0.1877	25.49	4.4361	68.64	0.6192	0.6520	0.6235	4.2722	23.43	0.6158
	ICM-SR-canny	0.2791	0.1856	24.66	4.4630	67.84	0.6143	0.6457	0.6089	4.1967	23.66	0.6180
	ICM-SR-multi	0.2799	0.1861	24.72	4.4411	68.00	0.6169	0.6440	0.6138	4.2094	23.77	0.6173

ground truth. This qualitative superiority aligns with our strong performance in perceptual quality measures, as detailed in Tab. 1. In the first super-resolution example in Figure 4, ICM-SR excels at recovering fine details from the LQ input, clearly reconstructing the contour of the eyes while OSDiff and TSD-SR fail. Notably, ICM-SR’s reconstruction inherits naturalness from pretrained diffusion models, enabling it to surpass the quality of the ground-truth image, which containing some noise. In the second example in Figure 4, ICM-SR demonstrates its superior ability to restore complex textures. While OSDiff produces a smoothed-out, blurred texture that loses the intricate details of the wood grain, ICM-SR accurately recovers the deep grooves and unique patterns present in the ground truth. This results in an output that is far more faithful to the original material and perceptually more realistic. Additionally, Figure 5 provides a broader qualitative comparison against various state-of-the-art diffusion-based SR methods, spanning both multi-step and one-step models. Across all examples, ICM-SR yields results that are simultaneously the sharpest and most consistent with the ground truth. For example, our method successfully recovers the clear edges of the leaves and branches in the first example, and the fine details of the ear in the second example. In contrast, other methods tend to either lose fine details to blurring or generate textures that deviate from the original. For more qualitative results, please see supplementary material.

5.4. Ablation study

Tab. 2 presents a quantitative comparison of the proposed method, ICM-SR, with various structural image information. As discussed in Sec. 4.1, directly conditioning the manifold with LQ images shows poor performance. Interestingly, we observe that conditioning on the colormap, which preserves broad color distributions, results in higher scores for no-reference perceptual metrics. Conversely, conditioning on canny edges, which explicitly defines structural outlines, leads to better performance on full-reference metrics that reward structural fidelity. Notably, a combination of both the colormap and canny edges proves the most effective, outperforming OSDiff across all metrics and datasets.

5.5. Analysis on image conditioned score

Figure 2 shows the impact of image-conditioning on score estimation across various diffusion timesteps. At large timestep ($t = 980$), denoised output without image-conditioning struggle to retain image structure and produce oversaturated image with totally different content. This indicates that the extracted generative prior is inappropriate for Real-ISR regularization. In contrast, predictions conditioned on the structural image information (colormap and Canny edges) effectively capture overall color tone and image structure, even at large timestep. This indicates that image-conditioned manifold is task-aligned and it enhances score estimation accuracy.

6. Conclusion and future work

We introduced Image-Conditioned Manifold (ICM) regularization which addressed the fundamental problem of conceptual mismatch in one-step Real-ISR methods that rely on text-conditioned generative priors. We identified that text-conditioning is often misaligned with the image fidelity goal of Real-ISR and naïve integration of dense image conditions may lead to numerical instability in the VSD loss. The proposed ICM regularization resolves this dilemma by conditioning the target manifold on core structural information, colormaps and Canny edges. This allows the teacher model to provide a stable and accurate prior with faithful colors and rich structural details.

Our comprehensive experiments demonstrated that ICM-SR achieved outstanding performance, particularly in perceptual quality metrics, highlighting its ability to generate visually appealing images that are faithful to LQ images. While maintaining the efficiency of single-step inference, ICM-SR produced results with vivid textures and sharp details, often rivaling or surpassing the perceptual quality of multi-step diffusion methods.

While ICM-SR demonstrates strong performance, its reliance on large pretrained models results in a significant parameter count, and there remains opportunity for further refinement in restoring very fine details. Future work will explore model compression techniques and advanced conditioning strategies to further enhance the practical utility and performance in diverse real-world applications.

References

- [1] Eirikur Agustsson and Radu Timofte. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *CVPRW*, 2017. 6, 11, 13
- [2] Jianrui Cai, Hui Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. In *ICCV*, 2019. 6, 11
- [3] Chaofeng Chen, Jiadi Mo, Jingwen Hou, Haoning Wu, Liang Liao, Wenxiu Sun, Qiong Yan, and Weisi Lin. Topiq: A top-down approach from semantics to distortions for image quality assessment. *IEEE Transactions on Image Processing*, 2024. 6
- [4] Keyan Ding, Kede Ma, Shiqi Wang, and Eero P Simoncelli. Image quality assessment: Unifying structure and texture similarity. *IEEE transactions on pattern analysis and machine intelligence*, 44(5):2567–2581, 2020. 6
- [5] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, 2014. 1
- [6] Linwei Dong, Qingnan Fan, Yihong Guo, Zhonghao Wang, Qi Zhang, Jinwei Chen, Yawei Luo, and Changqing Zou. Tsd-sr: One-step diffusion with target score distillation for real-world image super-resolution. In *CVPR*, 2025. 2, 3, 6, 13
- [7] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *NeurIPS*, 27, 2014. 1
- [8] Hang Guo, Jinmin Li, Tao Dai, Zhihao Ouyang, Xudong Ren, and Shu-Tao Xia. Mambair: A simple baseline for image restoration with state-space model. In *ECCV*. Springer, 2024. 1
- [9] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS*, 2017. 6
- [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1
- [11] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 2022. 2, 5, 6
- [12] Junjie Ke, Qifei Wang, Yilin Wang, Peyman Milanfar, and Feng Yang. Musiq: Multi-scale image quality transformer. In *ICCV*, 2021. 6
- [13] Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency trajectory models: Learning probability flow ode trajectory of diffusion. In *ICLR*, 2024. 2
- [14] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *CVPR*, 2016. 1
- [15] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 1, 3
- [16] Jianze Li, Jiezhong Cao, Yong Guo, Wenbo Li, and Yulun Zhang. One diffusion step to real-world super-resolution via flow trajectory distillation. *arXiv preprint arXiv:2502.01993*, 2025. 2
- [17] Yawei Li, Kai Zhang, Jingyun Liang, Jiezhong Cao, Ce Liu, Rui Gong, Yulun Zhang, Hao Tang, Yun Liu, Denis Demandolx, et al. Lsdir: A large scale dataset for image restoration. In *CVPR*, 2023. 6
- [18] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration using swin transformer. In *ICCV*, 2021. 1
- [19] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *CVPRW*, 2017. 1
- [20] Xinqi Lin, Jingwen He, Ziyang Chen, Zhaoyang Lyu, Bo Dai, Fanghua Yu, Yu Qiao, Wanli Ouyang, and Chao Dong. Diffbir: Toward blind image restoration with generative diffusion prior. In *ECCV*. Springer, 2024. 1, 3, 11, 13
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 6
- [22] Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff-instruct: A universal approach for transferring knowledge from pre-trained diffusion models, 2023. 2
- [23] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 2
- [24] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 2, 4
- [25] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 3
- [26] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021. 1
- [27] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *ICML*, 2023. 2
- [28] Jianyi Wang, Kelvin CK Chan, and Chen Change Loy. Exploring clip for assessing the look and feel of images. In *AAAI*, 2023. 6
- [29] 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. 3, 6, 11, 13
- [30] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. Esrgan: Enhanced super-resolution generative adversarial networks. In *ECCVW*, 2018. 1
- [31] Xintao Wang, Liangbin Xie, Chao Dong, and Ying Shan. Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In *ICCV*, 2021. 1, 3, 6

- [32] 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 *CVPR*, 2024. 3
- [33] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. In *NeurIPS*, 2023. 2, 3
- [34] Pengxu Wei, Ziwei Xie, Hannan Lu, Zongyuan Zhan, Qixiang Ye, Wangmeng Zuo, and Liang Lin. Component divide-and-conquer for real-world image super-resolution. In *ECCV*. Springer, 2020. 1, 6, 11
- [35] Rongyuan Wu, Lingchen Sun, Zhiyuan Ma, and Lei Zhang. One-step effective diffusion network for real-world image super-resolution. *NeurIPS*, 37:92529–92553, 2024. 2, 3, 5, 6, 12, 13
- [36] Rongyuan Wu, Tao Yang, Lingchen Sun, Zhengqiang Zhang, Shuai Li, and Lei Zhang. Seesr: Towards semantics-aware real-world image super-resolution. In *CVPR*, 2024. 1, 3, 11
- [37] Rui Xie, Chen Zhao, Kai Zhang, Zhenyu Zhang, Jun Zhou, Jian Yang, and Ying Tai. Addsr: Accelerating diffusion-based blind super-resolution with adversarial diffusion distillation. *arXiv preprint arXiv:2404.01717*, 2024. 3, 11, 13
- [38] Sidi Yang, Tianhe Wu, Shuwei Shi, Shanshan Lao, Yuan Gong, Mingdeng Cao, Jiahao Wang, and Yujiu Yang. Maniqa: Multi-dimension attention network for no-reference image quality assessment. In *CVPR*, 2022. 6
- [39] Tao Yang, Rongyuan Wu, Peiran Ren, Xuansong Xie, and Lei Zhang. Pixel-aware stable diffusion for realistic image super-resolution and personalized stylization. In *European Conference on Computer Vision*, pages 74–91. Springer, 2024. 11
- [40] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *CVPR*, pages 6613–6623, 2024. 2, 3
- [41] Weiye You, Mingyang Zhang, Leheng Zhang, Kexuan Shi, Xingyu Zhou, and Shuhang Gu. Consistency trajectory matching for one-step generative super-resolution. *arXiv preprint arXiv:2503.20349*, 2025. 2
- [42] Fanghua Yu, Jinjin Gu, Zheyuan Li, Jinfan Hu, Xiangtao Kong, Xintao Wang, Jingwen He, Yu Qiao, and Chao Dong. Scaling up to excellence: Practicing model scaling for photo-realistic image restoration in the wild. In *CVPR*, 2024. 1, 3, 11
- [43] Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. *Advances in Neural Information Processing Systems*, 36:13294–13307, 2023. 1, 11, 13
- [44] Kai Zhang, Jingyun Liang, Luc Van Gool, and Radu Timofte. Designing a practical degradation model for deep blind image super-resolution. In *ICCV*, 2021. 1, 3
- [45] Lin Zhang, Lei Zhang, and Alan C Bovik. A feature-enriched completely blind image quality evaluator. *IEEE Transactions on Image Processing*, 24(8):2579–2591, 2015. 6
- [46] 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
- [47] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 3, 6
- [48] Weixia Zhang, Guangtao Zhai, Ying Wei, Xiaokang Yang, and Kede Ma. Blind image quality assessment via vision-language correspondence: A multitask learning perspective. In *CVPR*, 2023. 6

Appendix

A. Quantitative comparison to multi-step super-resolution models

Compared to multi-step methods [20, 29, 36, 37, 39, 42, 43], ICM-SR outperforms in perceptual quality metrics with reference and delivers comparable results in no-reference perceptual metrics across all evaluated datasets [1, 2, 34]. Given that ICM-SR operates as a one-step model, unlike the multi-step approaches, it achieves impressive perceptual quality with significantly lower computational overhead.

Table 3. Quantitative comparison with state-of-the-art methods on both synthetic and real-world benchmarks. ‘s’ denotes the number of network inferences in the method. The best and second best results of each metric are highlighted in **red** and **blue**, respectively.

Datasets	Methods	Perceptual w/ ref.			Perceptual w/o ref.				Fidelity w/ ref.	
		LPIPS↓	DISTS↓	FID↓	NIQE↓	MUSIQ↑	MANIQ↑	CLIPQ↑	PSNR↑	SSIM↑
DIV2K-Val	StableSR-s200	0.3114	0.2048	24.44	4.7581	65.92	0.6190	0.6771	23.26	0.5726
	DiffBIR-s50	0.3669	0.2209	32.70	4.9903	69.87	0.6461	0.7299	23.14	0.5441
	SeeSR-s50	0.3194	0.1966	25.82	4.7927	68.42	0.6219	0.6867	23.73	0.6056
	SUPIR-s50	0.3919	0.2312	31.40	5.6767	63.86	0.5903	0.7146	22.13	0.5279
	PASD-s20	0.3779	0.2305	39.12	4.8587	67.36	0.6121	0.6327	24.00	0.6041
	ResShift-s15	0.3077	0.2136	30.79	6.9152	58.89	0.5283	0.5717	24.59	0.6232
	AddSR-s4	0.3816	0.2340	34.71	5.8441	69.18	0.6324	0.7532	22.38	0.5557
	ICM-SR-s1	0.2799	0.1861	24.72	4.4411	68.00	0.6169	0.6440	23.77	0.6173
DrealSR	StableSR-s200	0.3284	0.2269	148.95	6.5239	58.51	0.5586	0.6357	28.03	0.7536
	DiffBIR-s50	0.4669	0.2882	180.52	6.3293	66.15	0.6230	0.7068	25.91	0.6245
	SeeSR-s50	0.3142	0.2299	146.85	6.4825	64.74	0.6005	0.6895	28.14	0.7711
	SUPIR-s50	0.4243	0.2795	169.48	7.3918	58.79	0.5471	0.6749	25.09	0.6460
	PASD-s20	0.3579	0.2524	171.03	6.7661	63.23	0.5919	0.6242	27.79	0.7495
	ResShift-s15	0.3870	0.2632	160.16	8.6344	51.23	0.4644	0.5399	27.05	0.7404
	AddSR-s4	0.3709	0.2662	169.34	7.9004	65.23	0.6014	0.7153	26.66	0.7406
	ICM-SR-s1	0.2871	0.2142	125.30	6.4163	65.96	0.6051	0.6929	26.85	0.7763
RealSR	StableSR-s200	0.3002	0.2139	128.49	5.8809	65.88	0.6249	0.6234	24.65	0.7080
	DiffBIR-s50	0.3650	0.2399	130.67	5.8335	69.28	0.6511	0.7051	24.83	0.6501
	SeeSR-s50	0.3004	0.2218	125.09	5.3938	69.69	0.6453	0.6674	25.20	0.7215
	SUPIR-s50	0.3541	0.2488	130.38	6.1099	62.09	0.5780	0.6707	23.65	0.6620
	PASD-s20	0.3144	0.2304	134.18	5.7616	68.33	0.6323	0.5783	25.68	0.7273
	ResShift-s15	0.3279	0.2475	128.03	8.0708	56.88	0.5100	0.5362	25.66	0.7360
	AddSR-s4	0.3820	0.2688	153.35	6.4357	71.87	0.6767	0.7306	22.53	0.6452
	ICM-SR-s1	0.2611	0.2009	108.74	5.5842	68.59	0.6511	0.6360	24.99	0.7309

B. Ablation studies

B.1. Adapter scale

Table 4. Ablation study on the scale of the T2I Adapter. We investigate the effect of varying the adapter scale values ($\{0.0, 0.5, 1.0\}$) on the super-resolution performance evaluated on the DIV2K validation dataset.

Methods	Scale	Perceptual w/ ref.			Perceptual w/o ref.						Fidelity w/ ref.	
		LPIPS↓	DISTS↓	FID↓	NIQE↓	MUSIQ↑	MANIQ↑	CLIPQI↑	TOPIQ↑	LIQE↑	PSNR↑	SSIM↑
OSDiff	0	0.2847	0.1905	26.15	4.4918	67.73	0.6081	0.6394	0.6014	4.1704	23.40	0.6160
ICM-SR	0.5	0.2804	0.1862	25.09	4.5090	67.92	0.6178	0.6478	0.6110	4.2117	23.68	0.6193
	1.0 (Ours)	0.2799	0.1861	24.72	4.4411	68.00	0.6169	0.6440	0.6138	4.2094	23.77	0.6173

To verify the impact of the structural guidance strength, we conduct an ablation study on the scale of the T2I-Adapter. We evaluate the super-resolution performance on the DIV2K validation dataset by varying the adapter scale values within $\{0.0, 0.5, 1.0\}$. Note that setting the scale to 0.0 is equivalent to the baseline method, OSDiff [35], which relies solely on text conditioning. As presented in Tab. 4, introducing the core structural information (\mathbf{F}_c) via the T2I-Adapter significantly improves performance across all metrics compared to the baseline. However, we observe that the performance benefits saturate above a certain threshold (e.g., 0.5) and, we adopt 1.0 as the default setting for our main experiments.

B.2. Updating strategy for the auxiliary diffusion model, ϵ_ψ

Table 5. Ablation study on the training strategy for ϵ_ψ . The results are evaluated on the DIV2K validation dataset.

Methods	T2I-Adapter	Perceptual w/ ref.			Perceptual w/o ref.						Fidelity w/ ref.	
		LPIPS↓	DISTS↓	FID↓	NIQE↓	MUSIQ↑	MANIQ↑	CLIPQ↑	TOPIQ↑	LIQE↑	PSNR↑	SSIM↑
ICM-SR	✓	0.2785	0.1852	24.18	4.5000	67.15	0.6126	0.6307	0.5991	4.1165	23.96	0.6202
	✗ (Ours)	0.2799	0.1861	24.72	4.4411	68.00	0.6169	0.6440	0.6138	4.2094	23.77	0.6173

We conduct an ablation study on the training strategy for the auxiliary diffusion model (ϵ_ψ), specifically investigating the effect of incorporating core structural information. Formally, the objective of ϵ_ψ with T2I-Adapter is

$$\|\epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t, A_\eta(\mathbf{F}_c)) - \epsilon\|_2^2,$$

which learns the conditional distribution of $G_\theta(\mathbf{x}_L)|\mathbf{F}_c$. On the other hand, the objective of ϵ_ψ trained without T2I-Adapter (our choice) is

$$\|\epsilon_\psi(\hat{\mathbf{z}}_t; t, \mathbf{c}_t) - \epsilon\|_2^2,$$

which learns the distribution of generator’s output conditioned only on the text, *i.e.*, $G_\theta(\mathbf{x}_L)s$.

As shown in Tab. 5, the strategy incorporating the T2I-Adapter into the auxiliary model (✓) yields slightly higher fidelity. However, our chosen strategy without the adapter (✗) demonstrates substantial gains in no-reference perceptual metrics with only marginal compromise in reference-based metrics. In addition to this favorable performance trade-off, we consider the strategy without the adapter to be conceptually more appropriate. While ICM utilizes the T2I-Adapter to extract high-fidelity signal from pretrained diffusion models, the auxiliary model ϵ_ψ is tasked with learning the distribution of the generated output itself. Providing the ground-truth structural condition (\mathbf{F}_c) to ϵ_ψ makes the learning task trivial, potentially causing the model to over-rely on the condition instead of accurately capturing the distribution of the generated images. Therefore, considering both the significant perceptual gains and the conceptual alignment, we adopt the strategy without the T2I-Adapter for ϵ_ψ .

C. Qualitative comparisons

In this section, we provide more qualitative examples in Figure 7 on DIV2K [1] validation dataset, compared with other methods [6, 20, 29, 35–37, 43]. For various LQ images, our method, ICM-SR, produces visually appealing super-resolution results.

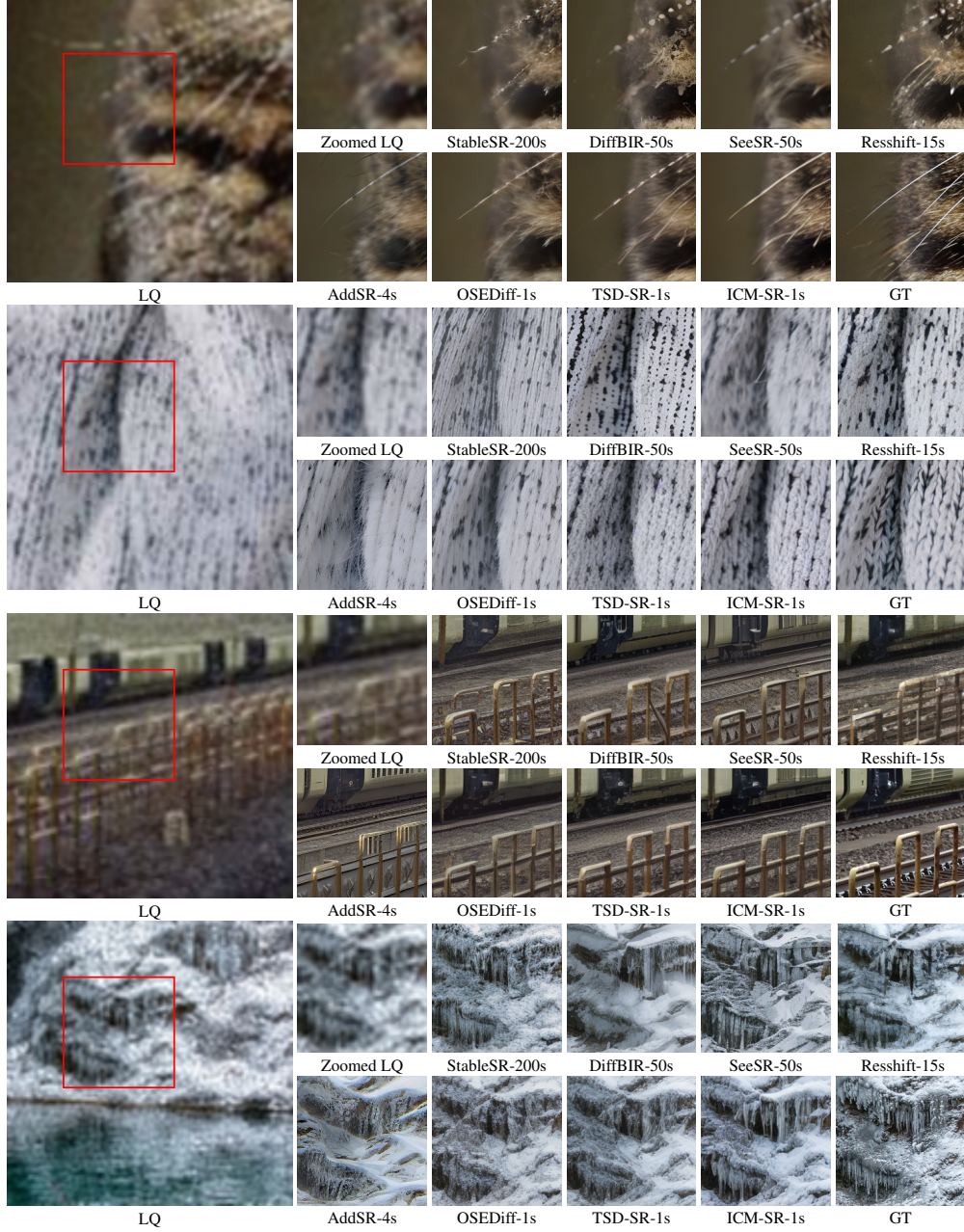


Figure 6. Qualitative results. Zoom in for better visualization.

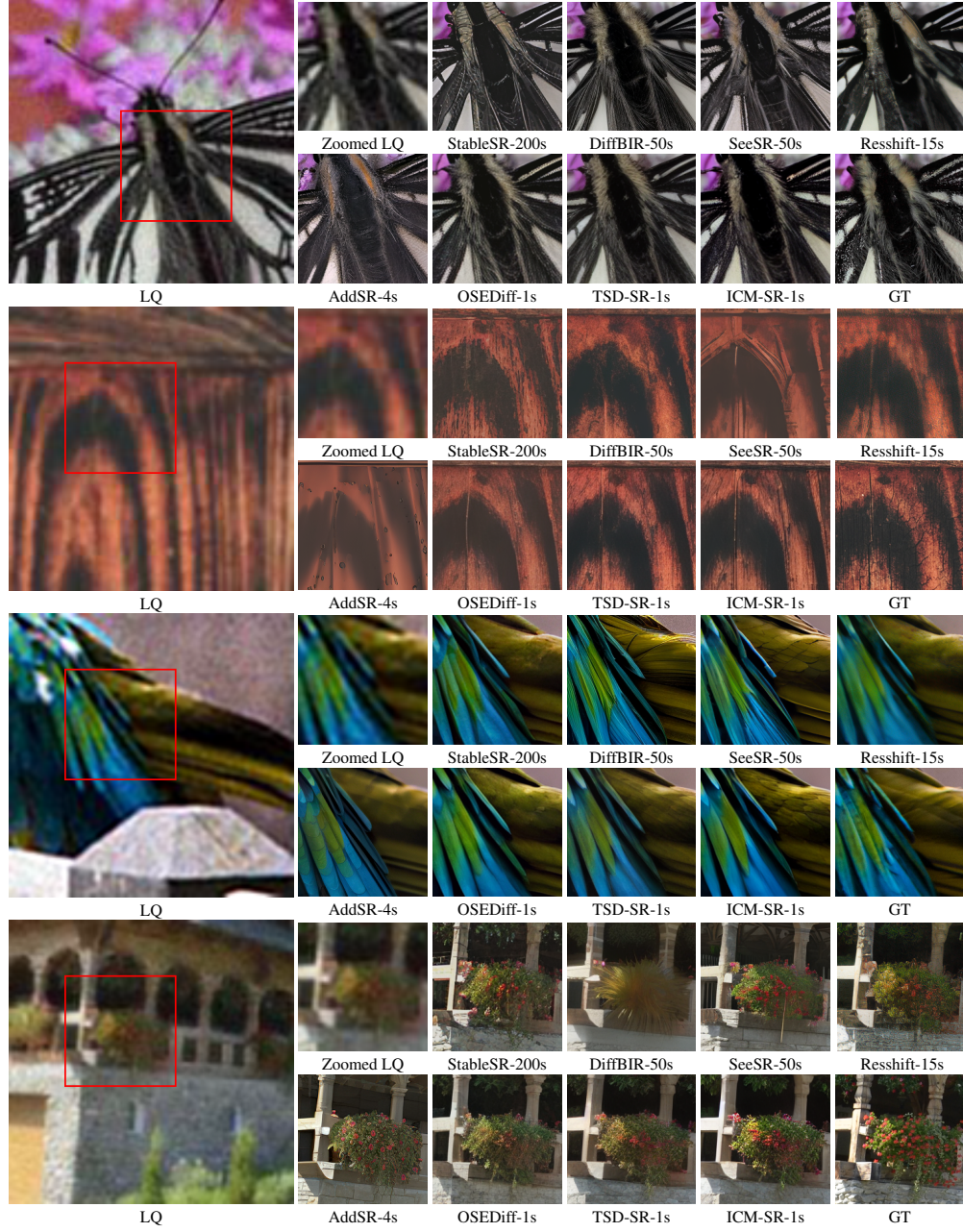


Figure 7. Qualitative results. Zoom in for better visualization.