

OmniSSR: Zero-shot Omnidirectional Image Super-Resolution using Stable Diffusion Model

Runyi Li[†], Xuhan Sheng[†], Weiqi Li, and Jian Zhang[✉]

School of Electronic and Computer Engineering, Peking University

Abstract. Omnidirectional images (ODIs) are commonly used in real-world visual tasks, and high-resolution ODIs help improve the performance of related visual tasks. Most existing super-resolution methods for ODIs use end-to-end learning strategies, resulting in inferior realness of generated images and a lack of effective out-of-domain generalization capabilities in training methods. Image generation methods represented by diffusion model provide strong priors for visual tasks and have been proven to be effectively applied to image restoration tasks. Leveraging the image priors of the **Stable Diffusion (SD)** model, we achieve **omnidirectional image super-resolution** with both fidelity and realness, dubbed as **OmniSSR**. Firstly, we transform the equirectangular projection (ERP) images into tangent projection (TP) images, whose distribution approximates the planar image domain. Then, we use SD to iteratively sample initial high-resolution results. At each denoising iteration, we further correct and update the initial results using the proposed Octadecaplex Tangent Information Interaction (OTII) and Gradient Decomposition (GD) technique to ensure better consistency. Finally, the TP images are transformed back to obtain the final high-resolution results. Our method is zero-shot, requiring no training or fine-tuning. Experiments of our method on two benchmark datasets demonstrate the effectiveness of our proposed method.

Keywords: Omnidirectional Imaging · Super-Resolution · Latent Diffusion Model

1 Introduction

Omnidirectional images (ODIs) capture the entire scene in all directions, exceeding the narrow field of view (FOV) offered by planar images. Super-Resolution (SR) techniques enhance the visual quality of ODIs by increasing their resolution, thereby revealing finer details and enabling more accurate scene analysis and interpretation. This becomes particularly crucial in applications like virtual

[†] means equal contribution.

[✉] means corresponding author.

reality and surveillance, where high-resolution ODIs are essential for precise perception and decision-making.

Current research in omnidirectional image super-resolution (ODISR) explores various methodologies to enhance the resolution of ODIs [15, 38]. SphereSR [60] addresses non-uniformity in different projections by learning upsampling processes and ensuring information consistency using LIIF [5]. OSRT [61] designs a distortion-aware Transformer to modulate equirectangular projection (ERP) distortions continuously and self-adaptively. Without a cumbersome process, OSRT outperforms previous methods remarkably. However, existing ODISR methods face the following challenges: (1) The majority are end-to-end models that can only produce a deterministic output, always better data fidelity but worse visual perception quality [18]. It’s promising to develop a generation-based model, but requiring high data demands, yet high-resolution ODIs are high cost to collect [56, 57]. (2) Most methods directly perform SR on ERP format ODIs, while users usually watch ODIs in a narrow FOV using tangent projection (TP). So another promising direction is to use off-the-shelf planar models on TP images. Recent times have witnessed the introduction and widespread application of diffusion models [24, 45], especially Stable Diffusion (SD) [40], which have provided a robust backbone for visual tasks [22, 25, 58, 62], including SR [32, 42, 49, 53, 54, 63]. However, if TP images are trivially one-by-one processed using diffusion-based SR models, they will exhibit discrepancies in the overlapping region when re-projected onto the ERP image. As a result, the global continuity is compromised.

Leveraging the strong image prior provided by SD, we propose the first diffusion-based zero-shot method for ODISR, named OmniSSR. Specifically, we propose Octadecaplex Tangent Information Interaction (OTII). OTII entails iterative conversion of intermediate SR results between ERP and TP representations, bridging the domain gap between ODIs and planar images. Building upon OTII, we further employ an approximate analytical solution of gradient descent, namely as Gradient Decomposition, to guide high-fidelity, high-quality omnidirectional image super-resolution. By capitalizing on SD’s effective image prior, our approach strikes a balance between *fidelity* and *realness*, ensuring that the restored ODIs exhibit both fidelity to the input data and realistic visual details. This method shows potential for advancing the current state of ODISR, providing improved resolution and visual quality across various applications. Fig. 1 showcases results fully demonstrating the superiority and performance of our proposed methods.

Our main contributions are summarized as follows:

- We propose OmniSSR, the first zero-shot ODISR method, using an off-the-shelf diffusion-based model, requiring no training or fine-tuning, leveraging existing image generation model priors to solve ODISR task.
- To bridge the domain gap between ODIs and planar images, we introduce Octadecaplex Tangent Information Interaction by repeatedly transforming ODIs between ERP format and TP format, enabling ODISR task with pre-trained diffusion models on planar images.

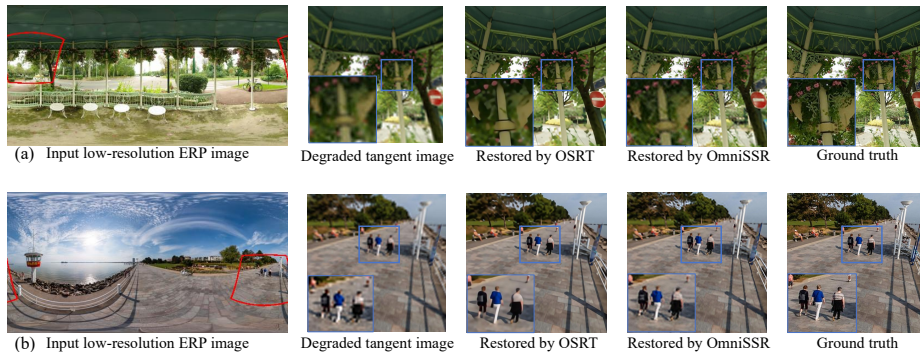


Fig. 1: We address omnidirectional image super-resolution in a *zero-shot* manner via OmniSSR. Presented above are select outcomes that sketch the essence of OmniSSR compared with current state-of-the-art approach OSRT [61]. Part (a) and (b) illustrate that OmniSSR upholds fidelity and visual realism at the same time, providing vivid and realistic details, while OSRT outputs over-smoothed and distorted results. Zoom in for more details.

- By iteratively updating images using the developed Gradient Decomposition technique, we introduce consistency constraints into the sampling process of the latent diffusion model, ensuring a trade-off between fidelity and realism in the reconstructed results.
- Extensive experiments are conducted on the benchmark datasets, demonstrating the superior performance of our method over existing state-of-the-art approaches, which validate the effectiveness of OmniSSR.

2 Related Work

2.1 Single Image Super-Resolution (SISR)

Image super-resolution methods based on deep learning have undergone significant development over an extended period. Currently, they can be broadly classified into two categories of solutions. The first category involves end-to-end network training methods, which utilize image pairs consisting of low-resolution degraded images and high-resolution ground truth images for network training [6–8, 29, 33, 34, 64, 67]. The network architectures employed in this category include CNN [17], Transformers [48], etc. The second category employs image generation models as priors, such as GAN [21], diffusion models [24, 45], etc., where low-resolution images are used as conditions to generate high-resolution images. We will mainly introduce the methods using generative prior.

Single Image SR using GAN prior In SR works utilizing GAN priors [4, 13, 35, 39, 59], including real-world scenarios [8, 51, 52, 66], pre-trained GAN networks are employed to transform image features into latent space, where the corresponding latent code for the high-resolution image is searched, ultimately yielding the reconstructed high-resolution result.

Single Image SR using diffusion prior The diffusion model provides a powerful image prior, and the diffusion sampling process can generate highly realistic images. This strong prior distribution can be applied to various image restoration tasks, including super-resolution [9, 10, 20, 42, 44, 53]. Image-domain diffusion models directly provide prior distributions of image-domain data. DDNM [53] based on the mathematical method of Range-Nullspace Decomposition, iteratively refines content on the zero space, combining image prior content in the value domain to achieve image restoration. DDRM [26] uses SVD decomposition to obtain restoration results, which is similar to DDNM. DPS [9] transforms the image super-resolution problem into an optimization problem with consistency constraints, using gradient descent algorithms to guide the generation of image-domain diffusion models. GDP [20] further uses such gradient to update the degradation operator to tackle blind image inverse problems. Other methods including MCG [10], DDS [9] and unified control of diffusion generation [20, 44] use same strategy for image restoration, especially image super-resolution.

The latent space diffusion model encodes data from various modalities into a latent space, samples its distribution, and then decodes it into the target domain. Image super-resolution works based on latent space domain include PSLD [41], P2L [11] and TextReg [28]. PSLD transfers the gradient-guided method of DPS [9] to the latent space diffusion model, while P2L furthermore considers prompt design, iteratively optimizing the prompt embedding of SD to improve the quality and visual effects of image reconstruction. TextReg applies the textual description of the preconception of the solution of the image inverse problem during the reverse sampling phase, of which the description is dynamically reinforced through null-text optimization for adaptive negation.

2.2 Omnidirectional Image Super-Resolution

Omnidirectional image super-resolution (ODISR) aims to enhance the resolution of omnidirectional or 360-degree images, which are commonly captured by cameras with a wide field of view. This field has garnered increasing attention due to its applications in virtual reality, omnidirectional video streaming, and surveillance. Several approaches have been proposed to address the unique challenges of ODISR [1–3, 37, 46]. For instance, Kämäräinen et al. [19] propose a deep learning-based approach for omnidirectional super-resolution, leveraging convolutional neural networks to effectively upscale low-resolution omnidirectional images while preserving spatial details. Similarly, Smolic et al. [38] introduce a novel omnidirectional super-resolution algorithm utilizing generative adversarial networks (GANs) to enhance the visual quality of omnidirectional images by hallucinating high-frequency details.

For evaluation purposes, researchers commonly utilize datasets such as the ODI-SR dataset from LAU-Net [14], and SUN 360 Panorama dataset [55]. These datasets enable the quantitative assessment of ODISR algorithms across various scenarios and facilitate fair comparisons between different methods.

3 Method

In this section, we first briefly introduce the preliminary background of our method (Sec. 3.1), and give an overall view of our proposed OmniSSR (Sec. 3.2). Then, we discuss the designs of Octadecaplex Tangent Information Interaction, which transform ODIs between ERP and TP formats with pre-upsampling strategy (Sec 3.3), and the Gradient Decomposition correction (Sec. 3.4).

3.1 Preliminaries

ERP↔TP Transformation The essence of projection transformations between ERP and TP lie in determining the positions of target image pixels within the source image and computing their corresponding pixel values using interpolation algorithms, as digital images are always stored discretely [30]. Therefore, the ERP→TP transformation involves locating the TP image pixels on the ERP imaging plane, and vice versa. Gnomonic projection [12] provides the correspondence between ERP image pixels and TP image pixels.

For a pixel $P_e(x_e, y_e)$ within the ERP image, we first find its corresponding pixel $P_s(\theta, \phi)$ on the unit sphere using Eq. 1:

$$\theta = 2\pi x_e/W, \quad \phi = \pi y_e/H, \quad (1)$$

where H and W are the height and width of the ERP image. The Cartesian coordinates of the ERP image and the angular coordinates on the unit sphere exhibit a straightforward one-to-one linear relationship, suggesting a conceptual equivalence between these two projection formats.

Given the spherical coordinates of the tangent plane center (θ_c, ϕ_c) , The transformation from $P_s(\theta, \phi)$ to $P_t(x_t, y_t)$, i.e. ERP→TP, is defined as:

$$\begin{aligned} x_t &= (\cos(\phi) \sin(\theta - \theta_c)) / \zeta, \\ y_t &= (\cos(\phi_c) \sin(\phi) - \sin(\phi_c) \cos(\phi) \cos(\theta - \theta_c)) / \zeta, \end{aligned} \quad (2)$$

where $\zeta = \sin(\phi_c) \sin(\phi) + \cos(\phi_c) \cos(\phi) \cos(\theta - \theta_c)$.

The corresponding inverse transformation, i.e. TP→ERP, is:

$$\begin{aligned} \theta &= \theta_c + \arctan((x_t \sin(c)) / (\rho \cos(\phi_1) \cos(c) - y_t \sin(\phi_c) \sin(c))), \\ \phi &= \arcsin(\cos(c) \sin(\phi_c) + y_t \sin(c) \cos(\phi_c) / \rho), \end{aligned} \quad (3)$$

where $\rho = \sqrt{x_t^2 + y_t^2}$ and $c = \arctan(\rho)$.

With Eq. 2 and Eq. 3, we can build one-to-one forward and inverse mapping functions between pixels on the ERP image and pixels on the TP images. An illustration of the ERP→TP transformation is shown in Fig. 2(a).

Iterative Denoising for Super-Resolution Utilizing the rich image priors provided by SD, we can super-resolve planar images. During initialization, the images are passed through the encoder \mathcal{E} of SD to obtain latent codes, which are

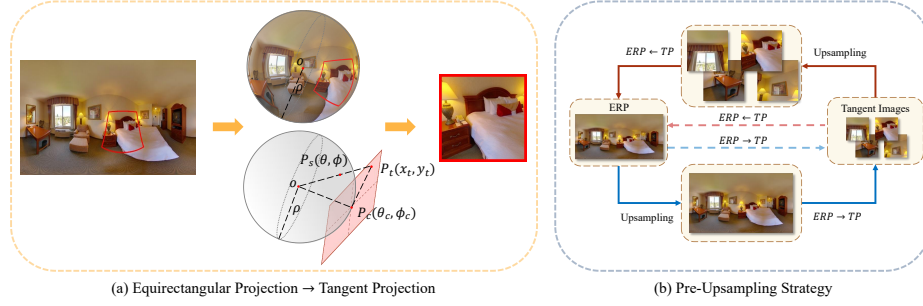


Fig. 2: Details about gnomonic transformations. (a) conversion from ERP to TP. (b) pre-upsampling proposed in Octadecaplex Tangent Information Interaction (Sec. 3.3) mitigating loss during transformation.

then added to pure noise to generate initial noise \mathbf{z}_T . Following the approach proposed by StableSR [49], we pass the images through a time-aware adapter \mathcal{T} . This adapter network structure is similar to the down-sampling part in denoising UNet, taking the image and the time step t of diffusion sampling as inputs to obtain the latent code feature for step t . This feature, along with the latent code \mathbf{z}_t for each step and the time step t , is then passed through denoising UNet to calculate the denoised result $\mathbf{z}_0|_t$ and the latent code \mathbf{z}_{t-1} for the next sampling step. By iterating this process T times, we can obtain the final super-resolution result via decoder \mathcal{D} of SD, yielding high-resolution images.

3.2 Overview

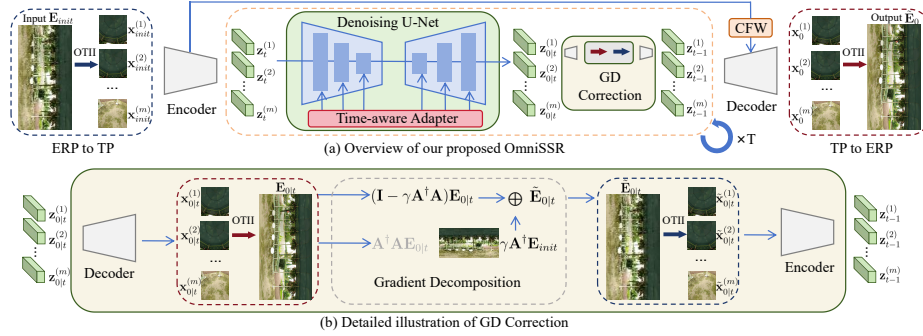


Fig. 3: Overview of our proposed OmniSSR. Input low-resolution omnidirectional image \mathbf{E}_{init} in ERP format is first projected onto Tangent Projection (TP) images $\mathbf{x}_{init}^{(1)}, \mathbf{x}_{init}^{(2)}, \dots, \mathbf{x}_{init}^{(m)}$, then iteratively refined via Stable Diffusion (SD) with a time-aware adapter and controllable feature wrapping (CFW) module. In each step of diffusion sampling, we adopt the Gradient Decomposition (GD) correction technique to introduce consistency constraints for the restored intermediate results. After T steps of sampling, we obtain the final result $\hat{\mathbf{E}}_0$ with high resolution and better visual quality.

Algorithm 1: OmniSSR Pipeline

Input: $\mathbf{E}_{init}, \mathcal{F}, \mathcal{F}^{-1}, \mathbf{A}, \mathbf{A}^\dagger, \mathcal{E}, \mathcal{D}, T$

Output: SR result $\tilde{\mathbf{E}}_0$

- 1 $\{\mathbf{x}_{init}^{(1)}, \mathbf{x}_{init}^{(2)}, \dots, \mathbf{x}_{init}^{(m)}\} = \mathcal{F}(\mathbf{E}_{init})$
- 2 **for** $i = 1$ **to** m **do**
- 3 $\mathbf{z}_{init}^{(i)} = \mathcal{E}(\mathbf{x}_{init}^{(i)})$
- 4 $\boldsymbol{\epsilon}^{(i)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5 $\mathbf{z}_T^{(i)} = \sqrt{\alpha_T} \mathbf{z}_{init}^{(i)} + \sqrt{1 - \alpha_T} \boldsymbol{\epsilon}^{(i)}$
- 6 **end**
- 7 Get $\{\mathbf{z}_0^{(1)}, \mathbf{z}_0^{(2)}, \dots, \mathbf{z}_0^{(m)}\}$ from Algo. 2
- 8 **for** $i = 1$ **to** m **do**
- 9 $\mathbf{x}_0^{(i)} = \mathcal{D}(\mathbf{z}_0^{(i)})$
- 10 **end**
- 11 $\mathbf{E}_0 = \mathcal{F}^{-1}(\{\mathbf{x}_0^{(1)}, \mathbf{x}_0^{(2)}, \dots, \mathbf{x}_0^{(m)}\})$
- 12 $\tilde{\mathbf{E}}_0 = \mathbf{E}_0 + \gamma_p \mathbf{A}^\dagger (\mathbf{E}_{init} - \mathbf{A} \mathbf{E}_0)$
- 13 **return** $\tilde{\mathbf{E}}_0$

Algorithm 2: Iterative Denoising with GD Correction

Input: $\mathbf{E}_{init}, \mathcal{F}, \mathcal{F}^{-1}, \mathbf{A}, \mathbf{A}^\dagger, \mathcal{E}, \mathcal{D}, \mathcal{T}, T$

Output: Latent code $\{\mathbf{z}_0^{(1)}, \mathbf{z}_0^{(2)}, \dots, \mathbf{z}_0^{(m)}\}$

- 1 **for** $t = T$ **to** 1 **do**
- 2 **for** $i = 1$ **to** m **do**
- 3 $\boldsymbol{\epsilon}_t = \boldsymbol{\epsilon}_\theta(\mathbf{z}_t^{(i)}, \mathcal{T}(\mathbf{z}_{init}^{(i)}, t), t)$
- 4 $\mathbf{z}_{0|t}^{(i)} = \frac{1}{\sqrt{\alpha_t}} (\mathbf{z}_t^{(i)} - \boldsymbol{\epsilon}_t \sqrt{1 - \alpha_t})$
- 5 $\mathbf{x}_{0|t}^{(i)} = \mathcal{D}(\mathbf{z}_{0|t}^{(i)})$
- 6 **end**
- 7 $\mathbf{E}_{0|t} = \mathcal{F}^{-1}(\{\mathbf{x}_{0|t}^{(1)}, \mathbf{x}_{0|t}^{(2)}, \dots, \mathbf{x}_{0|t}^{(m)}\})$
- 8 $\tilde{\mathbf{E}}_{0|t} = \mathbf{E}_{0|t} + \gamma_e \mathbf{A}^\dagger (\mathbf{E}_{init} - \mathbf{A} \mathbf{E}_{0|t})$
- 9 $\{\tilde{\mathbf{x}}_{0|t}^{(1)}, \tilde{\mathbf{x}}_{0|t}^{(2)}, \dots, \tilde{\mathbf{x}}_{0|t}^{(m)}\} = \mathcal{F}(\tilde{\mathbf{E}}_{0|t})$
- 10 **for** $i = 1$ **to** m **do**
- 11 $\tilde{\mathbf{z}}_{0|t}^{(i)} = (1 - \gamma_l) \mathbf{z}_{0|t}^{(i)} + \gamma_l \mathcal{E}(\tilde{\mathbf{x}}_{0|t}^{(i)})$
- 12 $\mathbf{z}_{t-1}^{(i)} \sim p(\mathbf{z}_{t-1}^{(i)} | \mathbf{z}_t^{(i)}, \tilde{\mathbf{z}}_{0|t}^{(i)})$
- 13 **end**
- 14 **end**
- 15 **return** $\{\mathbf{z}_0^{(1)}, \mathbf{z}_0^{(2)}, \dots, \mathbf{z}_0^{(m)}\}$

Our approach can be divided into three parts. The first part is pre-processing, where we initially up-sample the low-resolution ERP images \mathbf{E}_{init} , then project them onto tangent planes to obtain a series of TP images. These TP images are transformed to the latent space by the SD encoder, iteratively processed through denoising UNet and time-aware adapter network, and then decoded to obtain high-resolution TP images. During each denoising step, these TP images are transformed back via inverse transformation to ERP images, employing the Gradient Decomposition correction to ensure consistency constraints in diffusion sampling. After T iterations, the final super-resolution result is obtained. A formulaic description for OmniSSR pipeline is shown in Algo. 1. Fig. 3 shows the overview of our proposed pipeline.

3.3 Octadecaplex Tangent Information Interaction (OTII)

Motivation To apply SD for ODISR, a straightforward way is to perform the ERP \rightarrow TP transformation on the input ERP image. Then, each obtained TP image is fed into the SD-based model for SR. Finally, the TP \rightarrow ERP transformation yields the ultimate super-resolved ERP image. OmniFusion [30] employs a similar approach for depth estimation. However, this simplistic strategy fractures the inherent global coherence of ODIs, leading to pixel-level discontinuities in the fused ERP images. Moreover, interpolation algorithms cause significant information loss in the original projection transformations, resulting in more blurred images. If applied multiple times, this exacerbates the information loss even further. To mitigate this, a trivial solution is to increase the pixel count (resolution) of the intermediate projection imaging plane. However, excessively high resolutions in TP images can introduce unnecessary computational overhead during the denoising stage and potentially compromise the denoising performance.

Information Interaction and Pre-upsampling Based on the observations and analysis presented above, we propose OTII by alternating the intermediate results between ERP and TP formats at each denoising step, where a single ERP image is represented by 18 TP images. From Sec. 3.1, we can achieve the ERP→TP transformation (denoted as $\mathcal{F}(\cdot)$) and the TP→ERP transformation (denoted as $\mathcal{F}^{-1}(\cdot)$). Through the ERP→TP transformation, we can convert distorted ERP images into TP images with content distributions that approximate those of planar images. This enables the use of the original SD super-resolution method for planar images. Conversely, through the TP→ERP transformation, we can fuse information between different TP images holistically, while providing ERP-format input for the subsequent GD Correction in Sec. 3.4. To handle information loss during projection transformation, we further propose to pre-upsample the source image before projection transformations, as shown in Fig. 2(b). Our experiments in Sec. 4.4 demonstrate that this pre-upsampling strategy can significantly mitigate the information loss caused by projection transformations.

3.4 Gradient Decomposition (GD) Correction for Fidelity

SD-based methods, as introduced in Sec. 3.1, can perform SR on sliced TP images. However, relying solely on the SR results from SD may lack consistency and fail to accurately preserve the original semantic information and details of the low-resolution image.¹ To enhance the consistency of the SR results from SD, we opt to use convex optimization methods to iteratively refine them. Modeling the SR task as an image inverse problem, the following equation is formulated:

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}, \quad \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (4)$$

where \mathbf{x} represents the original image, \mathbf{y} denotes the degraded result, \mathbf{A} is the degradation operator (e.g., bicubic downsampling for super-resolution), and \mathbf{n} is random noise. The objective we aim to solve can be expressed as the following convex optimization problem:

$$\operatorname{argmin}_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda\mathcal{R}(\mathbf{x}), \quad (5)$$

where the first term is the data-fidelity term, ensuring the consistency of image reconstruction, and the second term is the regulation term, ensuring the sparsity of the reconstruction result, thus making the image more realistic. The regularization term can be the 1-norm, Total Variation, etc. The aforementioned convex optimization problem can be solved using a series of algorithms, such as gradient descent, ADMM, etc. Considering the trade-off between time and performance, we turn to find a solution based on gradient descent, and provide an approximate analytical solution composed of a *fidelity* term and a *realness* term, named

¹ This claim will be further illustrated in subsequent experiments.

"Gradient Decomposition (GD)":

$$\begin{aligned}\tilde{\mathbf{E}}_{0|t} &= \mathbf{E}_{0|t} + \alpha \nabla_{\mathbf{E}_{0|t}} \|\mathbf{E}_{init} - \mathbf{A}\mathbf{E}_{0|t}\|_F = \mathbf{E}_{0|t} + \alpha \times 2(\mathbf{A}^\dagger \mathbf{E}_{init} - \mathbf{A}^\dagger \mathbf{A} \mathbf{E}_{0|t}) \\ &= \mathbf{E}_{0|t} + \gamma \mathbf{A}^\dagger (\mathbf{E}_{init} - \mathbf{A} \mathbf{E}_{0|t}) = \gamma \mathbf{A}^\dagger \mathbf{E}_{init} + (\mathbf{I} - \gamma \mathbf{A}^\dagger \mathbf{A}) \mathbf{E}_{0|t}\end{aligned}\tag{6}$$

where \mathbf{A}^\dagger denotes pseudo-inverse of degradation operator \mathbf{A} , \mathbf{E}_{init} denotes initial low-resolution ERP input, $\mathbf{E}_{0|t}$ denotes restored result by SD, $\tilde{\mathbf{E}}_{0|t}$ denotes corrected result by GD, α denotes the learning rate of gradient descent, and γ denotes the simplified hyper-parameter which is further tuned using grid search. The final setting of γ on different stages is shown in Sec. 4.1, and the ablation studies of parameter choice are in Sec. 4.4.

This technique could be seen as a step of gradient descent optimization, and the optimized result could be decomposed of (1) $\gamma \mathbf{A}^\dagger \mathbf{E}_{init}$, which ensures the consistency of the generated result, and (2) $(\mathbf{I} - \gamma \mathbf{A}^\dagger \mathbf{A}) \mathbf{E}_{0|t}$, which serves as the iteratively updated generated result to improve its realness; γ is a hyper-parameter balancing the data fidelity and visual quality. For a better diversity and generality of the SR process, we expand this solution to latent space, and obtain the denoising result from both denoising UNet and corrected TP images (Algo. 2 line 11). A more detailed understanding of the iterative denoising process and application of GD correction could be referred to Algo. 2.

4 Experiments

4.1 Implementation Details

Datasets and Pretrained Models We choose the test set of ODI-SR dataset from LAU-Net [14] and SUN 360 Panorama dataset [55], comprising 97 and 100 omnidirectional images respectively, for experimental evaluation. The ground truth images are of size 1024×2048 pixels. In SR methods such as GDP [20] and PSLD [41] for planar images, we partitioned the images into several 256×256 patches and performed super-resolution on each patch individually.

For pre-trained models, we adopt from StableSR [49], which provided a SR network for planar images based on SD. This network architecture includes a time-aware adapter, a controllable feature wrapping (CFW) module, and the original SD structure from HuggingFace. All of them are kept untrained in our proposed OmniSSR.

Settings We set diffusion sampling steps to 200, which is the same as StableSR. The steps for other diffusion-based methods are set the same as their default settings (e.g. 1000 steps for PSLD). The degradation for low-resolution ERP images is bicubic down-sampling, and the implementation of its pseudo-inverse can be referred from code of DDRM [26]². For choices of hyper-parameter γ in GD correction, we set $\gamma_p = 1.0$, $\gamma_e = 1.0$, $\gamma_l = 0.5$. Our code is developed via PyTorch on NVIDIA 3090Ti GPU.³

² <https://github.com/bahjat-kawar/ddrm>

³ Code will be made available.

Table 1: SR results under bicubic downsampling on ODI-SR and SUN 360 Panorama datasets. For tasks not implemented in those papers, we mark N/A in corresponding results. Best results are shown in **Red**, and second best results are shown in **Blue**.

Method	Scale	ODI-SR				SUN 360 Panorama			
		WS-PSNR \uparrow	WS-SSIM \uparrow	FID \downarrow	LPIPS \downarrow	WS-PSNR \uparrow	WS-SSIM \uparrow	FID \downarrow	LPIPS \downarrow
Bicubic	$\times 2$	28.14	0.8343	24.00	0.2164	28.67	0.8537	29.25	0.1933
DDRM [26]		27.90	0.8317	12.28	0.1661	29.55	0.8670	13.10	0.1426
DPS [9]		20.99	0.6194	148.30	0.5249	21.44	0.6598	148.83	0.5175
GDP [20]		27.89	0.8157	26.56	0.2724	28.60	0.8376	28.02	0.2445
PSLD [41]		N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
DiffIR [54]		23.77	0.6583	57.23	0.4687	23.54	0.6775	58.06	0.4658
StableSR [49]		22.70	0.6458	44.87	0.3039	23.30	0.6907	43.49	0.2858
OmniSSR		28.57	0.8540	13.01	0.1575	29.69	0.8781	12.99	0.1459
Bicubic	$\times 4$	25.43	0.7059	50.84	0.3755	25.49	0.7229	55.99	0.3656
DDRM [26]		25.43	0.7367	32.69	0.3206	25.83	0.7443	32.93	0.3304
DPS [9]		24.75	0.6594	120.74	0.4911	21.09	0.6119	175.2143	0.5541
GDP [20]		23.16	0.6692	77.43	0.4260	23.75	0.6569	90.23	0.4240
PSLD [41]		21.72	0.5498	107.99	0.5329	21.75	0.5828	141.49	0.5461
DiffIR [54]		24.01	0.6770	54.14	0.4367	23.90	0.7014	50.37	0.4235
StableSR [49]		23.33	0.6577	49.95	0.3135	23.99	0.6998	46.03	0.3023
OmniSSR		25.77	0.7279	30.97	0.2977	26.01	0.7481	34.58	0.2963

4.2 Comparison of OmniSSR with diffusion-based methods

To evaluate the performance of proposed OmniSSR, we compare our method with recent state-of-the-art zero-shot methods for single image SR task: DPS [9], DDRM [26], GDP [20] which are based on the image-domain diffusion model, and PSLD [41], which is based on latent diffusion model. We also choose supervised diffusion-based super-resolution approaches including StableSR [49] and DiffIR [54] for comparison. We conduct experiments on $\times 2$ and $\times 4$ SR with ERP bicubic downsampling, on ODI-SR test-set and SUN test-set. We choose WS-PSNR [47], WS-SSIM [68], FID [23], and LPIPS [65] as the main metrics.

Quantitative results are presented in Tab. 1. With proposed OTII and GD correction, OmniSSR out-performs previous methods in terms of both *Fidelity* (from WS-PSNR and WS-SSIM) and *Realness* (from FID, LPIPS), which shows superior performance to existing diffusion-based methods for ODISR tasks on different scales.

Qualitative results are shown in Fig. 4 and Fig. 5, which illustrates the visualization of SR results on SUN test set and ODI-SR test set with $\times 2$ and $\times 4$ scales, by different methods. The visual results indicate that our OmniSSR exhibits superior capability for detail recovery compared to other methods, particularly evident in textual elements (e.g., the text "flapping" in upper part of Fig. 4), complex objects (e.g., the black desk with a screen in lower part of Fig. 4, patterns above the white door in lower part of Fig. 5), and small-scale objects (e.g., the person and clock behind the desk in upper part of Fig. 5). OmniSSR demonstrates the ability to recover highly detailed and realistic visual effects from TP images.



Fig. 4: Visualized comparison of $\times 2$ and $\times 4$ SR results on SUN 360 testset. 001 and 009 is the id number in testset filenames. We also calculate the PSNR and SSIM to HR ground truth of each SR result and downsampled image.

4.3 Comparison with end-to-end supervised methods

The experiments of comparison in Sec. 4.2 are mainly focused on zero-shot image super-resolution methods, and supervised single image super-resolution methods, where the approaches are not trained or fine-tuned on omnidirectional images. In this part, we will compare OmniSSR to supervised end-to-end methods with end-to-end training on ODI datasets, including SwinIR and OSRT. Besides the main metrics in Sec. 4.2, we also use NIQE [36] and DISTS [16] to evaluate the visual perception of SR outputs. Results are presented in Tab. 2, which shows that although our OmniSSR exhibits inferior fidelity metrics compared to end-to-end supervised methods trained directly on ODI datasets, it demonstrates notable improvements in the visual quality and authenticity of super-resolved images. Notably, end-to-end methods often produce smoothed reconstructions with distortions, whereas our approach preserves finer details and adheres more closely to the realistic distribution. Considering that our method has never been trained or tuned on ODI datasets, nor having omnidirectional images prior, this result is acceptable.

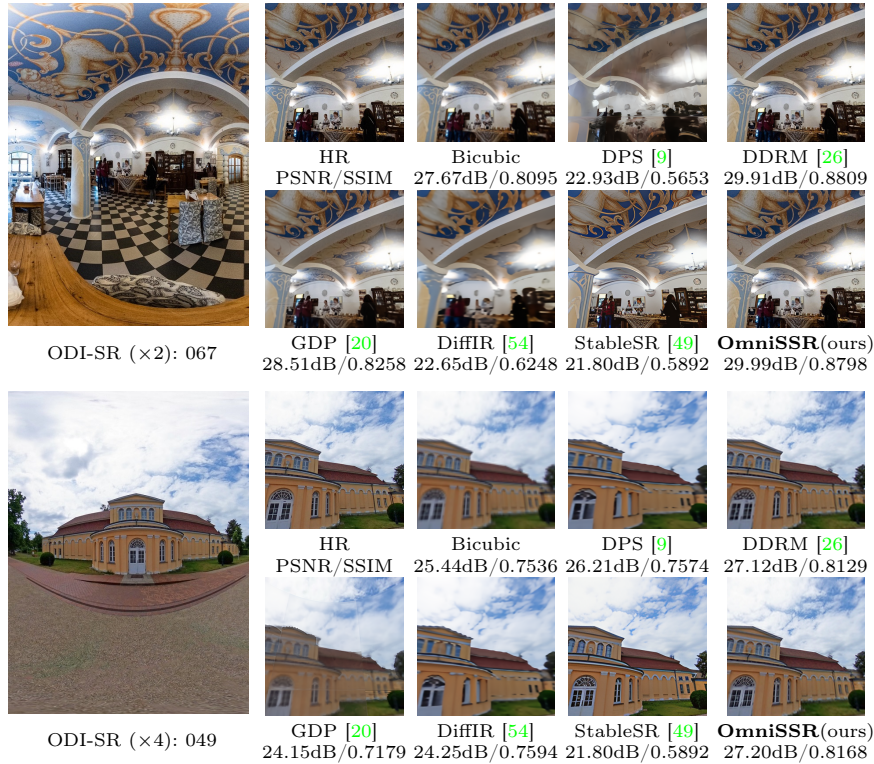


Fig. 5: Visualized comparison of $\times 2$ and $\times 4$ SR results on ODI-SR test set. 067 and 049 are the id numbers in test set filenames. We also calculate the PSNR and SSIM between ground truth and each SR result as well as downsampled image.

4.4 Ablation Studies

We first sequentially validate the performance improvement of the proposed strategy in OmniSSR including input image type, OTII and GD correction, on the ODI-SR test-set with $\times 2$ SR task, thereby demonstrating the significance of these strategies. The details are demonstrated as follows:

- 1) we do not use any proposed strategy in the SR task, which is equivalent to the vanilla StableSR baseline;
- 2) we transform the degraded ERP image to TP images and feed them separately into StableSR pipeline, instead of directly inputting ERP images;
- 3) based on 2), we add OTII strategy during the denoising process of SD (Algo. 2 line 7);
- 4) based on 2), we add GD correction at the *post-processing* stage (Algo. 1 line 12) of the overall pipeline;
- 5) based on 3) and 4), we add GD correction at *every step* and *post-processing* stage of sampling, to improve the consistency of the restored result.

Table 2: Comparison on $\times 4$ SR task with supervised methods trained on ODI-SR dataset, including SwinIR and OSRT. The best results are shown in **Bold**.

Method	Dataset	WS-PSNR \uparrow	WS-SSIM \uparrow	FID \downarrow	LPIPS \downarrow	NIQE \downarrow	DISTS \downarrow
SwinIR [31]	ODI-SR	26.76	0.7620	27.94	0.3321	5.3961	0.1710
OSRT [61]		26.89	0.7646	27.39	0.3258	5.4364	0.1695
OmniSSR		25.77	0.7279	30.97	0.2977	5.2891	0.1541
SwinIR [31]	SUN 360	26.02	0.7692	39.90	0.3419	5.2440	0.1325
OSRT [61]		26.33	0.7766	39.22	0.3364	5.2984	0.1312
OmniSSR		26.01	0.7481	34.58	0.2963	5.1329	0.1299

Note that the execution of GD correction requires the execution of OTII in the denoising process simultaneously, there is no scenario where only GD correction is executed without the execution of OTII in the denoising process.

Table 3: Ablation studies of OmniSSR on input type, OTII, and GD correction, on the test set of the ODI-SR dataset. Best results are shown in **Bold**.

Input type	OTII	GD Correction	WS-PSNR \uparrow	WS-SSIM \uparrow	FID \downarrow	LPIPS \downarrow
ERP	\times	\times	22.69	0.6458	44.87	0.3039
TP	\times	\times	23.53	0.6849	43.91	0.3113
TP	\checkmark	\times	23.74	0.6847	65.35	0.3748
TP	\times	\checkmark (in post-process only)	26.77	0.8192	15.41	0.1691
TP	\checkmark	\checkmark	28.58	0.8540	13.01	0.1575

Table 4: Results of pre-upsampling strategy on different scales, where (x,y) denotes bicubic-based upsampling at $x\times$ scale to ERP before ERP \rightarrow TP, and $y\times$ scale to TP before TP \rightarrow ERP transformation. Best results are shown in **Bold**.

ERP \rightarrow TP \rightarrow ERP	(1, 1)	(1, 4)	(4, 1)	(4, 2)	(2, 4)	(4, 4)
WS-PSNR \uparrow	28.98	38.11	28.99	33.91	38.05	38.18
WS-SSIM \uparrow	0.8859	0.9838	0.8862	0.9626	0.9837	0.9841

Quantitative results of ablation studies are shown in Tab. 3. From the result shown below, we could come to the claim that the OTII helps improve the performance on the domain level, and the transformation between ERP and TP images provides information fusion among adjacent TP images. Our proposal of Gradient Decomposition corrects such restoration result, improving fidelity and realness significantly at the same time, and it would be better if it is applied at each step of the overall denoising pipeline. Tab. 4 shows the effect of mitigating information loss via proposed pre-upsampling strategy.

For γ in the GD correction technique, we use grid search to obtain better results on ODI-SR dataset and $\times 4$ SR task. Fig. 6 shows performance on different choices of γ_p in Algo. 1 line 12, γ_e in Algo. 2 line 8, and γ_l in Algo. 2 line 11. The entire ablation of γ_p , γ_e and γ_l , with WS-PSNR, WS-SSIM, FID and LPIPS score all calculated and compared, will be provided in Supplementary Materials.

To evaluate the generalizability of our proposed modules, including Pre-Upsampling, OTII, and GD correction, we further conducted ablation studies on two super-resolution backbones, StableSR and SwinIR. The results underscore substantial performance enhancements facilitated by our modules across both backbones, which is provided in Supplementary Materials.

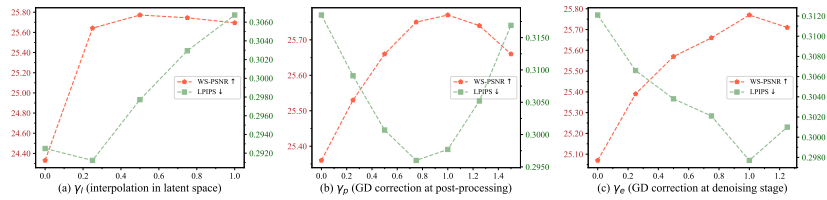


Fig. 6: Ablation of choices on γ_p , γ_e and γ_l . For better readability, WS-PSNR and LPIPS are chosen as evaluation metrics for fidelity and visual quality, respectively, to demonstrate the performance under different choices of the gamma parameter. We illustrate the results of (a) γ_p and γ_e fixed, while adjusting γ_l ; (b) γ_e and γ_l fixed, while adjusting γ_p ; (c) γ_p and γ_l fixed, while adjusting γ_e . It can be observed that when $\gamma_p = 1$, $\gamma_e = 1$, and $\gamma_l = 0.5$, OmniSSR achieves the relatively best performance.

5 Limitation and Discussion

Although OmniSSR bridges the gap between omnidirectional and planar images, achieving competitive performance and better visual results in ODISR, it still exhibits the following limitations: (1) The inference of the diffusion model requires a considerable amount of time, approximately 14 minutes per ERP-formatted omnidirectional image to be super-resolved into size 1024×2048 , making real-time super-resolution challenging; (2) Multiple conversions between ERP and TP are required in the pipeline, leading to improved performance but consuming additional inference time; (3) Further exploration of the convex optimization properties of GD correction is warranted, such as designing gradient term coefficients adaptive to reconstruction results and degradation types.

This study explores the application of image generation models to ODISR tasks. In future work, the framework behind OmniSSR can be extended beyond the confines of image super-resolution in a single scenario and venture into more complex ODI-based real-world scenarios. These include ODI editing, ODI inpainting, enhancing the quality of 3D Gaussian Splatting scenes [27, 43] obtained after super-resolving ERP images, as well as enhancing the quality of omnidirectional videos [50].

6 Conclusion

This paper leverages the image prior of Stable Diffusion (SD) and employs the Octadecaplex Tangent Information Interaction (OTII) to achieve *zero-shot* omnidirectional image super-resolution. Additionally, we propose the Gradient Decomposition (GD) correction based on convex optimization algorithms to refine the initial super-resolution results, enhancing the fidelity and realism of the restored images. The superior performance of our proposed method, OmniSSR, is demonstrated on benchmark datasets. By bridging the gap between omnidirectional and planar images, we establish a training-free approach, mitigating the data demand and over-fitting associated with end-to-end training. The ap-

plication scope of our method can be further extended to various applications, presenting potential value across multiple visual tasks.

References

1. An, H., Zhang, X.: Perception-oriented omnidirectional image super-resolution based on transformer network. In: Proceedings of the IEEE International Conference on Image Processing (ICIP) (2023)
2. Arican, Z., Frossard, P.: Joint registration and super-resolution with omnidirectional images. *IEEE Transactions on Image Processing (TIP)* (2011)
3. Cao, M., Mou, C., Yu, F., Wang, X., Zheng, Y., Zhang, J., Dong, C., Li, G., Shan, Y., Timofte, R., et al.: Ntire 2023 challenge on 360deg omnidirectional image and video super-resolution: Datasets, methods and results. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) (2023)
4. Chan, K.C., Xu, X., Wang, X., Gu, J., Loy, C.C.: Glean: Generative latent bank for image super-resolution and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2022)
5. Chen, Y., Liu, S., Wang, X.: Learning continuous image representation with local implicit image function. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2021)
6. Chen, Z., Zhang, Y., Gu, J., Kong, L., Yang, X., Yu, F.: Dual aggregation transformer for image super-resolution. In: Proceedings of the IEEE/CVF international conference on computer vision (ICCV) (2023)
7. Cheng, M., Ma, H., Ma, Q., Sun, X., Li, W., Zhang, Z., Sheng, X., Zhao, S., Li, J., Zhang, L.: Hybrid transformer and cnn attention network for stereo image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
8. Chong, M., Yanze, W., Xintao, W., Chao, D., Jian, Z., Ying, S.: Metric learning based interactive modulation for real-world super-resolution. In: Proceedings of the European Conference on Computer Vision (ECCV) (2022)
9. Chung, H., Kim, J., Mccann, M.T., Klasky, M.L., Ye, J.C.: Diffusion posterior sampling for general noisy inverse problems. *arXiv preprint arXiv:2209.14687* (2022)
10. Chung, H., Sim, B., Ye, J.C.: Improving diffusion models for inverse problems using manifold constraints. In: Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) (2022)
11. Chung, H., Ye, J., Milanfar, P., Delbracio, M.: Prompt-tuning latent diffusion models for inverse problems. *arXiv preprint arXiv:2310.01110* (2023)
12. Coxeter, H.S.M.: Introduction to geometry. John Wiley & Sons, Inc. (1961)
13. Daras, G., Dean, J., Jalal, A., Dimakis, A.: Intermediate layer optimization for inverse problems using deep generative models. In: Proceedings of the International Conference on Machine Learning (ICML) (2021)
14. Deng, X., Wang, H., Xu, M., Guo, Y., Song, Y., Yang, L.: Lau-net: Latitude adaptive upscaling network for omnidirectional image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2021)
15. Deng, X., Wang, H., Xu, M., Li, L., Wang, Z.: Omnidirectional image super-resolution via latitude adaptive network. *IEEE Transactions on Multimedia (TMM)* (2022)

16. Ding, K., Ma, K., Wang, S., Simoncelli, E.P.: Image quality assessment: Unifying structure and texture similarity. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2020)
17. Dong, C., Loy, C.C., He, K., Tang, X.: Image super-resolution using deep convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2015)
18. Duan, H., Zhai, G., Min, X., Zhu, Y., Fang, Y., Yang, X.: Perceptual quality assessment of omnidirectional images. In: *Proceedings of the IEEE International Symposium on Circuits and Systems (ISCAS)* (2018)
19. Fakour-Sevom, V., Guldogan, E., Kämäräinen, J.K.: 360 panorama super-resolution using deep convolutional networks. In: *Proceedings of the Int. Conf. on Computer Vision Theory and Applications (VISAPP)* (2018)
20. Fei, B., Lyu, Z., Pan, L., Zhang, J., Yang, W., Luo, T., Zhang, B., Dai, B.: Generative diffusion prior for unified image restoration and enhancement. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2023)
21. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)* (2014)
22. Guo, L., Tao, T., Cai, X., Zhu, Z., Huang, J., Zhu, L., Gu, Z., Tang, H., Zhou, R., Han, S., et al.: Cas-diffcom: Cascaded diffusion model for infant longitudinal super-resolution 3d medical image completion. *arXiv preprint arXiv:2402.13776* (2024)
23. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. In: *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)* (2017)
24. Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. In: *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)* (2020)
25. Jiang, Y., Chan, K.C., Wang, X., Loy, C.C., Liu, Z.: Reference-based image and video super-resolution via c^2 -matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2022)
26. Kawar, B., Elad, M., Ermon, S., Song, J.: Denoising diffusion restoration models. In: *Proceedings of the ICLR Workshop on Deep Generative Models for Highly Structured Data (ICLRW)* (2022)
27. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)* (2023)
28. Kim, J., Park, G.Y., Chung, H., Ye, J.C.: Regularization by texts for latent diffusion inverse solvers. *arXiv preprint arXiv:2311.15658* (2023)
29. Li, W., Chen, B., Zhang, J.: D3c2-net: Dual-domain deep convolutional coding network for compressive sensing. *arXiv preprint arXiv:2207.13560* (2022)
30. Li, Y., Guo, Y., Yan, Z., Huang, X., Duan, Y., Ren, L.: Omnifusion: 360 monocular depth estimation via geometry-aware fusion. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2022)
31. Liang, J., Cao, J., Sun, G., Zhang, K., Van Gool, L., Timofte, R.: Swinir: Image restoration using swin transformer. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)* (2021)
32. Liu, J., Wang, Q., Fan, H., Wang, Y., Tang, Y., Qu, L.: Residual denoising diffusion models. *arXiv preprint arXiv:2308.13712* (2023)
33. Lu, Z., Li, J., Liu, H., Huang, C., Zhang, L., Zeng, T.: Transformer for single image super-resolution. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)* (2022)

34. Lugmayr, A., Danelljan, M., Timofte, R.: Ntire 2020 challenge on real-world image super-resolution: Methods and results. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) (2020)
35. Menon, S., Damian, A., Hu, S., Ravi, N., Rudin, C.: Pulse: Self-supervised photo upsampling via latent space exploration of generative models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
36. Mittal, A., Soundararajan, R., Bovik, A.C.: Making a “completely blind” image quality analyzer. *IEEE Signal Processing Letters (SPL)* (2013)
37. Nishiyama, A., Ikehata, S., Aizawa, K.: 360° single image super resolution via distortion-aware network and distorted perspective images. In: Proceedings of the IEEE International Conference on Image Processing (ICIP) (2021)
38. Ozcinar, C., Rana, A., Smolic, A.: Super-resolution of omnidirectional images using adversarial learning. In: Proceedings of the IEEE International Workshop on Multimedia Signal Processing (MMSPW) (2019)
39. Pan, X., Zhan, X., Dai, B., Lin, D., Loy, C.C., Luo, P.: Exploiting deep generative prior for versatile image restoration and manipulation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2021)
40. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022)
41. Rout, L., Raouf, N., Daras, G., Caramanis, C., Dimakis, A., Shakkottai, S.: Solving linear inverse problems provably via posterior sampling with latent diffusion models. In: Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) (2023)
42. Saharia, C., Ho, J., Chan, W., Salimans, T., Fleet, D.J., Norouzi, M.: Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* (2022)
43. Schönbein, M., Geiger, A.: Omnidirectional 3d reconstruction in augmented manhattan worlds. In: Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (2014)
44. Song, J., Zhang, Q., Yin, H., Mardani, M., Liu, M.Y., Kautz, J., Chen, Y., Vahdat, A.: Loss-guided diffusion models for plug-and-play controllable generation. In: Proceedings of the International Conference on Machine Learning (ICML) (2023)
45. Song, Y., Sohl-Dickstein, J., Kingma, D.P., Kumar, A., Ermon, S., Poole, B.: Score-based generative modeling through stochastic differential equations. In: Proceedings of the International Conference on Learning Representations (ICLR) (2020)
46. Sun, X., Li, W., Zhang, Z., Ma, Q., Sheng, X., Cheng, M., Ma, H., Zhao, S., Zhang, J., Li, J., et al.: Opdn: Omnidirectional position-aware deformable network for omnidirectional image super-resolution. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW) (2023)
47. Sun, Y., Lu, A., Yu, L.: Weighted-to-spherically-uniform quality evaluation for omnidirectional video. *IEEE Signal Processing Letters (SPL)* (2017)
48. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) (2017)
49. Wang, J., Yue, Z., Zhou, S., Chan, K., Loy, C.: Exploiting diffusion prior for real-world image super-resolution. *arXiv preprint arXiv:2305.07015* (2023)
50. Wang, Q., Li, W., Mou, C., Cheng, X., Zhang, J.: 360dvd: Controllable panorama video generation with 360-degree video diffusion model. In: Proceedings of the

- IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2024)
51. Wang, X., Xie, L., Dong, C., Shan, Y.: Real-esrgan: Training real-world blind super-resolution with pure synthetic data. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2021)
 52. Wang, X., Yu, K., Wu, S., Gu, J., Liu, Y., Dong, C., Qiao, Y., Change Loy, C.: Esrgan: Enhanced super-resolution generative adversarial networks. In: Proceedings of the European Conference on Computer Vision Workshops (ECCVW) (2018)
 53. Wang, Y., Yu, J., Zhang, J.: Zero-shot image restoration using denoising diffusion null-space model. In: Proceedings of the International Conference on Learning Representations (ICLR) (2022)
 54. Xia, B., Zhang, Y., Wang, S., Wang, Y., Wu, X., Tian, Y., Yang, W., Van Gool, L.: Diffir: Efficient diffusion model for image restoration. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2023)
 55. Xiao, J., Ehinger, K.A., Oliva, A., Torralba, A.: Recognizing scene viewpoint using panoramic place representation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2012)
 56. Yagi, Y.: Omnidirectional sensing and its applications. IEICE Transactions on Information and Systems (TOIS) (1999)
 57. Yamazawa, K., Yagi, Y., Yachida, M.: Omnidirectional imaging with hyperboloidal projection. In: Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (1993)
 58. Yang, S., Zhou, Y., Liu, Z., Loy, C.C.: Rerender a video: Zero-shot text-guided video-to-video translation. In: Proceedings of the SIGGRAPH Asia 2023 Conference Papers (2023)
 59. Yinhuai, W., Yujie, H., Jiwen, Y., Jian, Z.: Gan prior based null-space learning for consistent super-resolution. In: Proceedings of the AAAI Conference on Artificial Intelligence (AAAI) (2023)
 60. Yoon, Y., Chung, I., Wang, L., Yoon, K.J.: Spheresr: 360deg image super-resolution with arbitrary projection via continuous spherical image representation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022)
 61. Yu, F., Wang, X., Cao, M., Li, G., Shan, Y., Dong, C.: Osrt: Omnidirectional image super-resolution with distortion-aware transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
 62. Yu, J., Zhang, X., Xu, Y., Zhang, J.: Cross: Diffusion model makes controllable, robust and secure image steganography. In: Proceedings of the Advances in Neural Information Processing Systems (NeurIPS) (2023)
 63. Yue, Z., Wang, J., Loy, C.C.: Resshift: Efficient diffusion model for image super-resolution by residual shifting. In: Advances in Neural Information Processing Systems (NeurIPS) (2023)
 64. Zhang, K., Liang, J., Van Gool, L., Timofte, R.: Designing a practical degradation model for deep blind image super-resolution. In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) (2021)
 65. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
 66. Zhang, W., Li, X., Shi, G., Chen, X., Qiao, Y., Zhang, X., Wu, X.M., Dong, C.: Real-world image super-resolution as multi-task learning. In: Advances in Neural Information Processing Systems (NeurIPS) (2023)

67. Zhang, X., Zhang, Y., Xiong, R., Sun, Q., Zhang, J.: Heronet: Hyperspectral explainable reconstruction and optimal sampling deep network for snapshot compressive imaging. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022)
68. Zhou, Y., Yu, M., Ma, H., Shao, H., Jiang, G.: Weighted-to-spherically-uniform ssim objective quality evaluation for panoramic video. In: Proceedings of the IEEE International Conference on Signal Processing (ICSP) (2018)

Supplementary Materials of “OmniSSR: Zero-shot Omnidirectional Image Super-Resolution using Stable Diffusion Mode”

A Extra Experiments

A.1 Ablation Studies

Ablation study of γ on Gradient Decomposition (GD) correction According to the principle of GD correction, the super-resolution (SR) result in equirectangular projection (ERP) format $\mathbf{E}_{0|t}$ generated by StableSR [49] can be further corrected to $\tilde{\mathbf{E}}_{0|t} = \mathbf{E}_{0|t} + \gamma \mathbf{A}^\dagger(\mathbf{E}_{init} - \mathbf{A}\mathbf{E}_{0|t})$, where γ balances realness and fidelity. To improve the convergence of this gradient-based technique, we perform a grid search over different γ values to obtain the best results, presented in Tab. 5. For an overall performance superiority, we choose $\gamma_l = 0.5, \gamma_p = 1, \gamma_e = 1$.

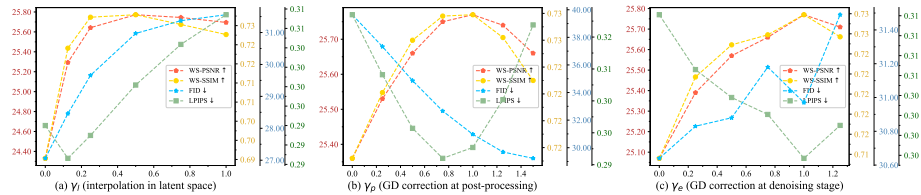


Fig. 7: Visualization of different choices of γ . (a) γ_p and γ_e fixed, while adjusting γ_l ; (b) γ_e and γ_l fixed, while adjusting γ_p ; (c) γ_p and γ_l fixed, while adjusting γ_e .

Ablation study of SR backbone We further conducted ablation studies on the selection of the SR backbone network to justify our choice of StableSR as the backbone and demonstrate the effectiveness of our proposed strategy at the same time. We selected the current state-of-the-art method in super-resolution work, SwinIR [31], to compare its results with StableSR [49], which is shown in Tab. 6.

Compared with SwinIR, StableSR significantly improves the fidelity and realness of reconstruction results. On the other hand, it also validates the effectiveness of our proposed Octadecaplex Tangent Information Interaction (OTII) and GD correction techniques on different backbones. Given its iterative updating and continuous correction nature, StableSR indeed has advantages over SwinIR’s end-to-end reconstruction approach.

Table 5: Ablation studies of hyper-parameter γ in GD correction. γ_p denotes γ in post-processing stage, γ_l denotes γ in post-processing stage, γ_e denotes γ in post-processing stage. The best results are shown in **Bold**.

γ_p	γ_l	γ_e	WS-PSNR \uparrow	WS-SSIM \uparrow	FID \downarrow	LPIPS \downarrow
1	0	1	24.33	0.6903	27.05	0.2925
1	0.25	1	25.64	0.7272	29.66	0.2912
1	0.5	1	25.77	0.7279	30.97	0.2977
1	0.75	1	25.74	0.7253	31.37	0.3029
1	1	1	25.69	0.7227	31.56	0.3067
0	0.5	1	25.37	0.7172	39.64	0.3184
0.25	0.5	1	25.53	0.7221	37.303	0.3090
0.5	0.5	1	25.67	0.7260	34.86	0.3037
0.75	0.5	1	25.75	0.7278	32.66	0.2960
1	0.5	1	25.77	0.7279	30.97	0.2977
1.25	0.5	1	25.74	0.7262	29.69	0.3052
1.5	0.5	1	25.66	0.7230	29.22	0.3169
1	0.5	0	25.07	0.7136	30.64	0.3121
1	0.5	0.25	25.38	0.7217	30.83	0.3066
1	0.5	0.5	25.56	0.7249	30.88	0.3037
1	0.5	0.75	25.66	0.7259	31.18	0.3020
1	0.5	1	25.77	0.7278	30.97	0.2977
1	0.5	1.25	25.71	0.7257	31.49	0.3010

Table 6: Results of our proposed techniques on different backbones, StableSR, and SwinIR. Best results are shown in **Bold**.

Backbone	Whether to use proposed techniques	WS-PSNR \uparrow	WS-SSIM \uparrow	FID \downarrow	LPIPS \downarrow
SwinIR [31]	×	26.11	0.7821	27.11	0.2390
SwinIR [31]	✓	27.89	0.8409	13.33	0.1510
StableSR [49]	✓	28.58	0.8540	13.01	0.1575

A.2 Further Exploration of ERP \leftrightarrow TP Transformation

A simple question arises: can we perform ERP \leftrightarrow TP⁴ transformation in the latent space, thus avoiding the need to transform intermediate results between image and latent space repeatedly? To answer this question, we made two attempts without Stable Diffusion (SD) encoder and decoder during each denoising step. GD correction is also not used in this section.

1) **Projection transformations on latent feature z_0 :** In this experiment, we focus on the impact of projection transformation on image features in the latent space, so here we do not involve the denoising process. Therefore, we first transformed the ground truth ERP image \mathbf{E}_0 to m TP images $\{\mathbf{x}_0^{(i)}\}_{i=1,\dots,m}$ through ERP \rightarrow TP. Then, we sequentially obtain the latent TP image features in the latent space:

$$\mathbf{z}_0^{(i)} = \mathcal{E}(\mathbf{x}_0^{(i)}), i = 1, \dots, m. \quad (7)$$

⁴ TP denotes tangent projection.

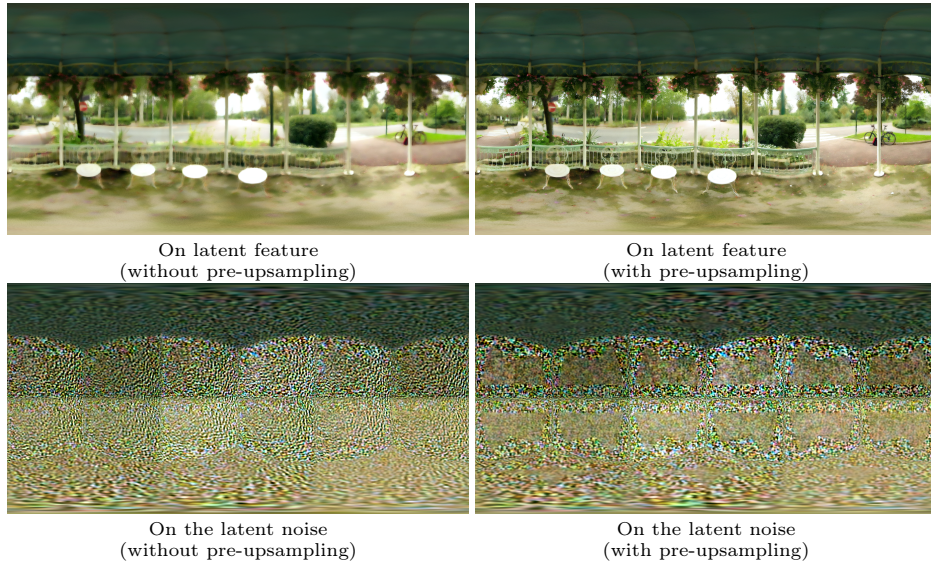


Fig. 8: Visualized comparison of projection transformations on latent image feature and latent noise. Zoom in for details.

Next, we perform $\text{TP} \rightarrow \text{ERP} \rightarrow \text{TP}$ on $\mathbf{z}_0^{(i)}$ to obtain $\hat{\mathbf{z}}_0^{(i)}$ and decode them to TP image as follows:

$$\hat{\mathbf{x}}_0^{(i)} = \mathcal{D}(\hat{\mathbf{z}}_0^{(i)}), i = 1, \dots, m. \quad (8)$$

Finally, the decoded TP image $\hat{\mathbf{x}}_0^{(i)}$ are transformed by $\text{TP} \rightarrow \text{ERP}$ to get $\hat{\mathbf{E}}_0$.

2) **Projection transformations on latent noise $\epsilon_t^{(i)}$:** In this experiment, we focus on the impact of projection transformation on the noise $\epsilon_t^{(i)}$. We transform the low-resolution ERP image to TP images and feed the latter into StableSR pipeline. At each sampling step, we directly perform $\text{TP} \rightarrow \text{ERP} \rightarrow \text{TP}$ transformation on the predicted noise $\{\epsilon_t^{(i)}\}_{i=1, \dots, m}$ to get $\{\hat{\epsilon}_t^{(i)}\}_{i=1, \dots, m}$, and using $\hat{\epsilon}_t^{(i)}$ for following denoising.

In the two experiments above, we also present the effects of using and not using pre-upsampling in the $\text{TP} \rightarrow \text{ERP} \rightarrow \text{TP}$ transformation process, respectively. We illustrate the visual results of $\hat{\mathbf{E}}_0$, using the 0000.png in image ODI-SR test-set as an example in Fig. 8. When **performing projection transformations on latent feature z_0** , the decoded images exhibit severe blurring. Although using pre-upsampling in the $\text{TP} \rightarrow \text{ERP} \rightarrow \text{TP}$ process can alleviate the blurriness to some extent and present clearer image content in certain areas, the overall image quality remains poor. In the experiment involving **projection transformations on latent noise $\epsilon_t^{(i)}$** , it can be observed that regardless of whether pre-upsampling strategy is used or not, the super-resolved images suffer from significant damage. This may be attributed to the SD encoder’s spatial down-sampling at $\times 8$ scale, compressing image pixels within an 8×8 patch into a single

latent pixel. Projection transformations, on the other hand, operate at the image pixel level with fine granularity. Applying such fine-grained operations directly to latent pixels can greatly disrupt the original image structure. Therefore, projection transformations related to ODIs should be performed in image space rather than in the latent space mapped by the SD Variational Auto Encoder (VAE).

A.3 Exploration of SD Encoder and Decoder

During the ablation study, we observed that OmniSSR, when GD correction is removed while OTII is retained, demonstrates improved fidelity (e.g., WS-PSNR, WS-SSIM) and deteriorated realism (e.g., FID, LPIPS) compared to the original StableSR model. Upon examining the outputs of the ablation model under this configuration, significant color shift issues were identified, as depicted in Fig. 9(a).

We initially suspected that this color shift stemmed from **the utilization of the SD VAE** before and after OTII in each denoising step. To validate this hypothesis, we conducted a visual comparison experiment using image 0006.png from the ODI-SR testset as an example. It can be observed that even when GD correction and OTII are successively removed, as illustrated in Fig. 9(a)(b), the color shift persists. It is only when we eliminate the repeated usage of SD VAE in each denoising step that the color at the boundary of black and white tiles returns to normal, as shown in Fig. 9(c). Ground truth reference can be seen in Fig. 9(d). This phenomenon of color shift indicates the potential problem caused by frequently using SD VAE.

A.4 The Global Continuity of ODIs

The existing ODISR methods directly perform SR on ERP images, resulting in the discontinuity between the left and right sides [3]. Our proposed OTII treats TP images as the direct input for the network. Besides facilitating the transfer use of existing planar image-specific diffusion models, it also effectively considers the omnidirectional characteristics of ODIs. We selected some visualization results of OSRT [61] and OmniSSR, focusing on the continuity near the left and right sides of the ERP. As shown in Fig. 10, OSRT exhibits poor continuity between the left and right sides of the ERP, while OmniSSR naturally inherits the advantage of TP images in seamlessly spanning different areas of the ERP.

A.5 Time Consumption

The inference runtime of different methods are compared as follows. Considering fair comparison, we use the default settings referred to in corresponding papers. The diffusion sampling steps for OmniSSR are 200, DDRM [26] 100, and PSLD [41] 1000.⁵ All experiments are conducted on a single NVIDIA 3090Ti GPU.

⁵ We have tried to use the same sampling accelerate strategy in DDRM, but get bad restored results.

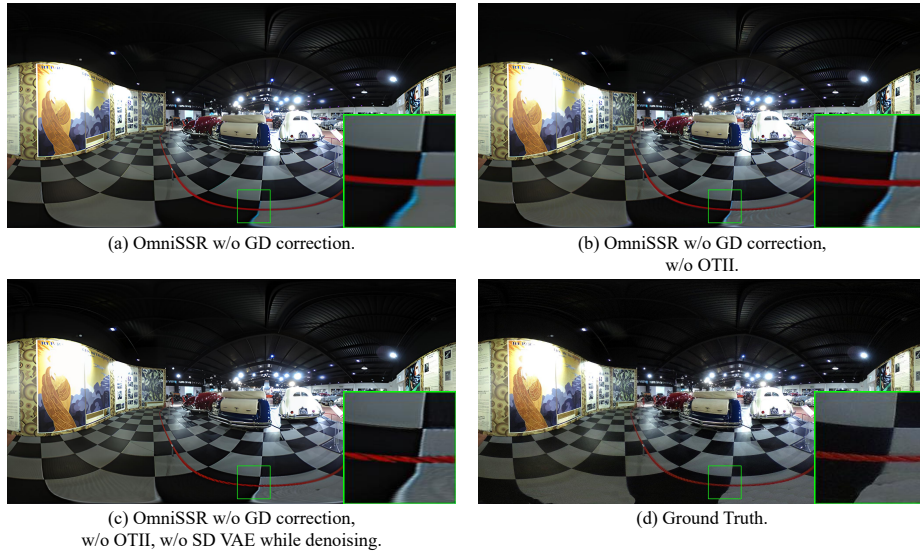


Fig. 9: Phenomenon and causes of color shift: By progressively removing different components of OmniSSR (a)(b)(c), we ultimately discovered that the color shift in the super-resolution results disappears again after removing the SD VAE used in the denoising step. This indicates the potential risk of color shift associated with frequent usage of SD VAE during denoising.

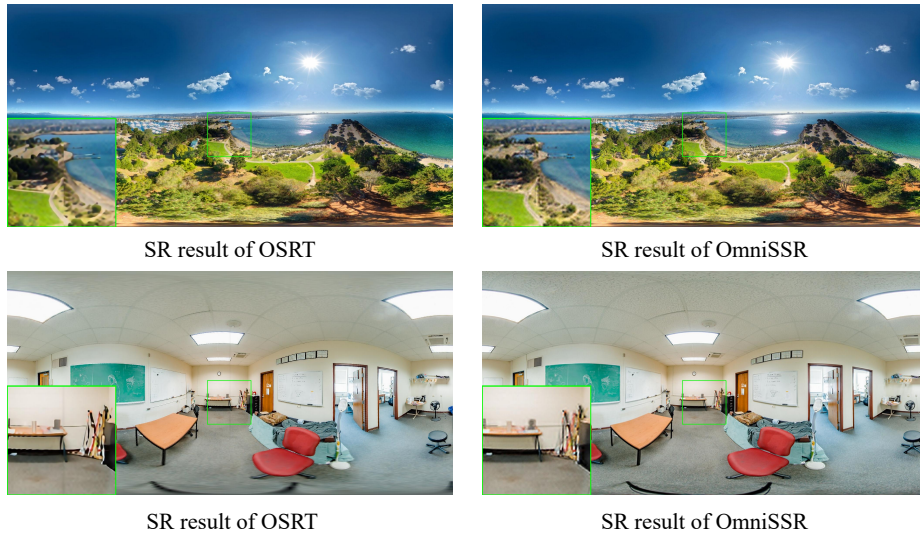


Fig. 10: Continuity of left and right part of SR results on OSRT and our proposed OmniSSR. It is shown that OSRT suffers from serious artifacts and bad continuity. All ERP images have been rotated by 180 degrees to stitch the left and right sides. (Upper image: 0039 of ODI-SR test set, lower image: 0015 of SUN test set.)

Table 7: Time consumption of OmniSSR and other SR methods.

Method	Runtime per ERP image (s)↓
SwinIR [31]	0.87
OSRT [61]	1.44
DDRM	711.95
PSLD	6720.87
OmniSSR (Ours)	726.19

B Theoretical Discussion

In this section, we provide a simple theoretical discussion of our proposed GD correction technique, explaining why a single step of GD would also work and obtain better results.

Take the update step in GD correction as an example, let us first re-examine this step:

$$\tilde{\mathbf{E}}_{0|t} = \mathbf{E}_{0|t} + \gamma_e \mathbf{A}^\dagger (\mathbf{E}_{init} - \mathbf{A}\mathbf{E}_{0|t}), \quad (9)$$

where $\gamma_e \mathbf{A}^\dagger (\mathbf{E}_{init} - \mathbf{A}\mathbf{E}_{0|t})$ is the gradient of fidelity term $\|\mathbf{E}_{init} - \mathbf{A}\mathbf{E}_{0|t}\|_F$, and $\gamma_e = 2 \times \alpha$ (learning rate).

An obvious and direct question is: why did we perform only a single update step rather than multiple steps? Through the following analysis, we will demonstrate that, in this context, multi-step gradient descent and single-step are essentially equivalent, with the number of steps being governed by the coefficient γ_e .

Analysis Suppose we take multiple steps in GD correction and are taking step k to $k - 1$. As $\tilde{\mathbf{E}}_{0|t}^{(k)}$ can be represented via $\tilde{\mathbf{E}}_{0|t}^{(k-1)}$ in linear form, we can use $\tilde{\mathbf{E}}_{0|t}^{(0)}$ to express $\tilde{\mathbf{E}}_{0|t}^{(k)}$, and $\tilde{\mathbf{E}}_{0|t}^{(0)}$ only has linear coefficients composed of γ_e , \mathbf{A} and \mathbf{A}^\dagger . Thus for fixed γ_e , there is no difference between one step and multiple steps of GD correction. For adaptive γ_e , it is also obvious that $\tilde{\mathbf{E}}_{0|t}^{(k)}$ can be represented via $\tilde{\mathbf{E}}_{0|t}^{(0)}$ with linear transforms and different γ_e . Thus for a better trade-off between performance and inference time, we turn to use **one** step of GD correction.