# ZS-SRT: An Efficient Zero-Shot Super-Resolution Training Method for Neural Radiance Fields

Xiang Feng[a],[**], Yongbo He[a],[**], Yubo Wang[a], Chengkai Wang[a], Zhenzhong Kuang[a],[*], Jiajun Ding[a], Feiwei Qin[a] and Jianping Fan[b]

[a]School of Computer Science and Technology, Hangzhou Dianzi University, Hangzhou, 310018, China
[b]AI Lab at Lenovo Research, Beijing, China

**ABSTRACT**

Neural Radiance Fields (NeRF) have achieved great success in the task of synthesizing novel views that preserve the same resolution as the training views. However, it is challenging for NeRF to synthesize high-quality high-resolution novel views with low-resolution training data. To solve this problem, we propose a zero-shot super-resolution training framework for NeRF. This framework aims to guide the NeRF model to synthesize high-resolution novel views via single-scene internal learning rather than requiring any external high-resolution training data. Our approach consists of two stages. First, we learn a scene-specific degradation mapping by performing internal learning on a pretrained low-resolution coarse NeRF. Second, we optimize a super-resolution fine NeRF by conducting inverse rendering with our mapping function so as to backpropagate the gradients from low-resolution 2D space into the super-resolution 3D sampling space. Then, we further introduce a temporal ensemble strategy in the inference phase to compensate for the scene estimation errors. Our method is featured on two points: (1) it does not consume high-resolution views or additional scene data to train super-resolution NeRF; (2) it can speed up the training process by adopting a coarse-to-fine strategy. By conducting extensive experiments on public datasets, we have qualitatively and quantitatively demonstrated the effectiveness of our method.

## 1. Introduction

Multi-view reconstruction is a long-standing problem in computer vision, with applications spanning virtual reality, augmented reality, telepresence, etc. Its goal is to estimate both scene appearance and geometry by harnessing information from multiple views. Over the past few decades, many traditional methods [1, 2] were proposed for this purpose. However, they often encounter multiple challenges, such as geometric structure estimation and texture recovery.

Recently, neural radiance fields (NeRF) [3] have attracted lots of attention for scene reconstruction based on neural inverse rendering and implicit representation. It's important to note that NeRF can render images at any desired resolution. However, NeRF struggles to generate super-resolution novel views effectively. This is because when it tries to render at a higher resolution than the training views, it encounters artifacts such as blurry due to the interpolation characteristics and the gap in sampling between training and testing. A straightforward way to solve this problem is to use high-resolution scene data to train a high-resolution NeRF. However, it is usually hard to obtain high-resolution scene data in many practical scenarios, such as the task of novel view synthesis and even text-to-3D content creation [4–6]. Therefore, synthesizing high-resolution novel views without high-resolution ground truth is still challenging.

One possible solution to the above problem is to up-sample the synthesized low-resolution novel views by using a well-trained single-image super-resolution (SISR) model

[7]. While the SISR model performs well on single images, it struggles to maintain multi-view consistency when processing multi-view images. In [8], Bahat et al. proposed a neural volume super-resolution (NVSR) method to synthesize multi-view consistent high-resolution novel views, but it is time-consuming to train a generalizable model, and it is also difficult to apply NVSR due to the lack of scene data. In [9], Wang et al. proposed NeRF-SR by directly optimizing dense radiance fields through a super-sampling strategy to ensure view consistency. Although no additional scene data are required, NeRF-SR may lead to perceptual blurring and aliasing by enforcing the average value of high-resolution sub-pixels to match the value of low-resolution pixels.

In this paper, we propose a zero-shot super-resolution framework to improve the efficiency of reconstruction and generate reliable high-frequency details without consuming additional scene data. This two-stage framework integrates our idea of single-scene internal learning. First, we learn a scene-specific degradation mapping model on top of a pretrained coarse low-resolution NeRF. Second, we train a super-resolution NeRF in a coarse-to-fine manner to synthesize super-resolution novel views by conducting inverse rendering with the above well-trained degradation mapping model so as to backpropagate the gradients from low-resolution 2D space into the super-resolution 3D sampling space. Then, we further introduce a temporal ensemble strategy in the inference phase to compensate for the scene estimation errors. It is worth mentioning that our approach does not consume any high-resolution inputs and additional scene data during the training process and can ensure cross-view consistency. Besides, the employment of the coarse-to-fine strategy would help to speed up the training process. As

---

*Corresponding Author: Zhenzhong Kuang.
**Equal Contribution.
✉ zzkuang@hdu.edu.cn (Z. Kuang)

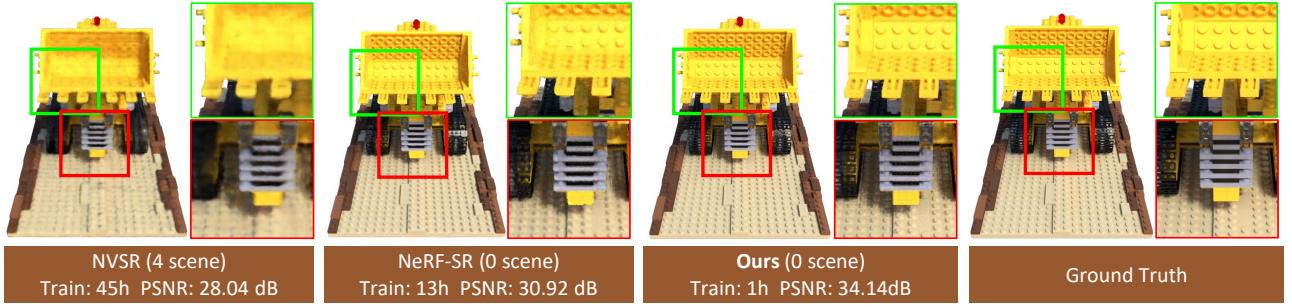| NVSR (4 scene) Train: 45h PSNR: 28.04 dB | NeRF-SR (0 scene) Train: 13h PSNR: 30.92 dB | **Ours** (0 scene) Train: 1h PSNR: 34.14dB | Ground Truth |

**Figure 1:** Comparing our method with existing super-resolution novel view synthesis methods. We have achieved optimal performance in three aspects: model training time, additional scene data consumed by model training, and the quality of synthesized novel views.

shown in Figure 1, it is easy to observe that our approach outperforms NVSR [8] and NeRF-SR [9] in terms of the training time, the amount of data required, and the synthetic quality. In summary, our contributions are as follows:

- We propose a zero-shot and coarse-to-fine super-resolution training framework that only utilizes low-resolution ground truth to optimize the high-resolution 3D sampling space of radiance fields.

- We propose an internal learning method to obtain a scene-specific degradation mapping model which is further used in inverse rendering to optimize our fine NeRF in super-resolution training without consuming additional scene data.

- We propose a temporal ensemble strategy to compensate for the scene estimation errors.

- We qualitatively and quantitatively evaluate our method on a couple of public datasets, which has demonstrates the effectiveness of it.

## 2. Related work

### 2.1. Neural Radiance Fields

Since implicit neural representation (INR) has obvious advantage in parameterizing various types of signals, it has been successfully applied in many recent works [10–14] to address problems that have plagued explicit representation for a long time. Unlike conventional discrete signal representation, INR relies on continuous function to describe various kinds of signals (e.g., image, audio, and 3D shape) by using neural networks to approximate the function for flexible and expressive representations. A straightforward advantage of employing continuous function lies in that it makes the signal representation no longer coupled with spatial resolution. As a result, the memory required for signal parameterization is independent of the spatial resolution and only scales with the complexity of the signal, indicating the infinite resolution of INR.

As a representative of INR, Neural Radiance Fields (NeRF) [3] receives increasing attention due to its excellent modeling ability on three-dimensional scene structure and surface characteristics (i.e. complex lighting and geometric effects) as well as the consistency across various viewpoints and lighting conditions. Initially, many variants of NeRF are proposed to address its inherent limitations. For example, some approaches [15–18] use voxel grids instead of MLPs to speed up training, and several other methods [19, 20] reconstruct NeRF without camera poses. Subsequently, some works focus on generalizing NeRF to resolve challenges that have plagued other fields effectively. For example, NAF [21], SNAF [22], and NeXF [23] explore new practical solutions to reconstruct medical objects by using sparse CBCT images. To further boost rendering quality, TensoRF [24] proposes to factorize the 4D scene tensor into compact vector and matrix factors. In terms of multi-scale representation, Mip-NeRF [25] and Zip-NeRF [26] are proposed based on tapered sampling and adaptive position encoding to solve the aliasing problem caused by low-sampling rendering. Although NeRF and its variants can support the synthesis of novel views at any resolution, they suffer from difficulty in rendering high-quality, high-resolution views.

### 2.2. Image Super-Resolution

Single Image Super-Resolution (SISR) is a classic problem in computer vision that aims at recovering the lost details [27, 28]. In literature, there exist many supervised super-resolution methods, such as dense connections [29–31], recursive structures [32, 33], back-projection [34], and self-attention mechanisms [35–37]. Besides, the recent generative model is also exploited for image super-resolution by learning a mapping function from low resolution (LR) image to high resolution (HR) image [38–41].

The supervised super-resolution methods usually require paired training data, but this may not always be available in reality. Therefore, many researchers turn to explore self-supervised or unsupervised approaches. Most recently, Shocher et al. [7] presents a zero-shot super-resolution (ZSSR) approach to produce high-quality super-resolution results by exploiting the internal recurrence of information inside a single image, where ZSSR does not require additional HR images. Although existing methods have achieved some success, they may confront the problem of maintaining multi-view consistency when dealing with multi-view images.

### 2.3. Super-Resolution Novel View Synthesis

To synthesize super-resolution novel views, many works focus on integrating the merits of both NeRF and SISR.
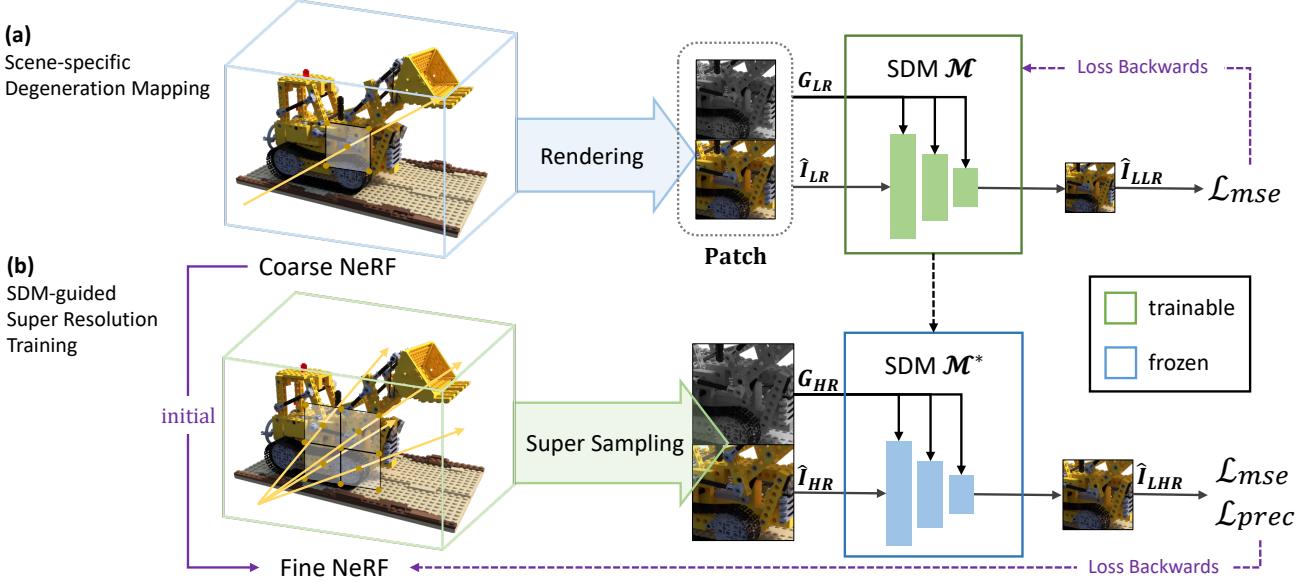
**Figure 2:** Overview of the proposed two-stage ZS-SRT framework for obtaining the high-resolution radiance field from low-resolution ground truth in a coarse-to-fine manner. In the first stage (a), a scene-specific degeneration mapping (SDM) model is learned on top of a coarsely NeRF. G denotes the gradient of the synthetic view, used to guide the training of the mapping relationship. In the second stage (b), a SDM-guided super-resolution training method is proposed to obtain a super-resolution fine NeRF so that our model can produce high-quality views with trusted, high-frequency details.

In [42], CGO-RF first relies on SISR to process multi-view images in advance and then trains a super-resolution update model with a large amount of high and low-resolution scene data to correct the spatial ambiguity caused by SISR. Similarly, NVSR [8] trains an EDSR model to super-resolve tri-planes. Since both of them need to learn generalization modules, they require large amounts of scene data pairs in the training process, which is also time-consuming. NeRF-SR [9] attempts to recover high-frequency information in a sub-pixel way. However, it may confront the problem of blurry details by using the simple average pooling operation during the supervised training process of the model. In this paper, we address this problem by learning a degenerate mapping (from high-resolution radiance fields to low-resolution ground truth) on top of our proposed single-scene internal learning method so that we can generate super-resolution results with richer high-frequency details.

## 3. Preliminary: NeRF

**Scene expression.** A radiance field is represented as a continuous function, $f$, which aims at mapping a 3D coordinate, $\mathbf{x} \in \mathbb{R}^3$, and a directional view unit vector, $\mathbf{d} \in \mathbb{S}^2$, to a volume density, $\sigma$, and $RGB$ values, $\mathbf{c}$. Typically, NeRF uses the following Multilayer Perception (MLP) to parameterize this function:

$$f_\theta : (\mathbf{x}, \mathbf{d}) \mapsto (\sigma, \mathbf{c}) \qquad (1)$$

To improve the efficiency, Tensorial Radiance Fields (TensoRF) [24] is adopted to realize NeRF, and its factorized

definition is as follows:

$$\sigma, c = \sum_r \sum_m \mathcal{A}^m_{\sigma,r}(\mathbf{x}), S\left(\mathbf{B}\left(\oplus\left[\mathcal{A}^m_{c,r}(\mathbf{x})\right]_{m,r}\right), \mathbf{d}\right) \qquad (2)$$

Here, $\mathcal{A}^m_{\sigma,r}$ and $\mathcal{A}^m_{c,r}$ represent the factorized components for density and color, respectively. The matrix $\mathbf{B}$ is used as a global appearance dictionary that abstracts the appearance commonalities across the entire scene. Function $S$ processes the transformed color representations $\mathbf{B}(\oplus[\mathcal{A}^m_{c,r}(\mathbf{x})]_{m,r})$ alongside the viewing direction $\mathbf{d}$ to produce the final color $c$.

**Volume rendering.** For volume rendering, we employ differentiable volume rendering. For each pixel, we trace a ray $\mathbf{r} = \mathbf{o} + t\mathbf{d}$ by sampling a set of points and computing the pixel color as follows:

$$\hat{C} = \sum_{i=1}^N \tau_i \left(1 - \exp\left(-\sigma_i \Delta_i\right)\right) c_i \qquad (3)$$

$$\tau_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \Delta_j\right) \qquad (4)$$

where $\sigma_i$ and $c_i$ are the density and color computed at sampled locations $x_i$; $\Delta_i$ is the ray step size, and $\tau_i$ denotes transmittance.

**Reconstruction.** To reconstruct a scene, given a set of multi-view input images with known camera poses, the tensorial radiance field per scene is optimized by minimizing the following photometric loss:

$$\mathcal{L} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|C(r) - \hat{C}(r)\|_2^2 \qquad (5)$$

where $\mathcal{R}$ denotes the set of rays randomly sampled in each minibatch, $C(r)$ represents the ground truth pixel colors, and $\hat{C}(r)$ is the color computed through our model.

# 4. Method

Figure 2 illustrates our pipeline. The first stage involves learning a scene-specific degradation mapping (SDM) network to obtain low-resolution results from high-resolution input by performing internal learning on a coarsely trained NeRF (i.e. coarse NeRF). In the second stage, we learn a fine super-resolution NeRF (i.e. fine NeRF) by using the low-resolution scene data (i.e. ground truth), where the learned SDM model is used to guide the training of the fine NeRF by conducting inverse rendering. Next, we introduce the details in different subsections. In Section 4.1, we present how to obtain the SDM model by using the internal learning method. In Section 4.2, we present how to obtain the fine NeRF model by using our super-resolution training method. In Section 4.3, we combine radiance fields from multiple time steps in a single training process to alleviate the variance in super-resolution scene estimation.

## 4.1. Internal Learning of Radiance Fields

As shown in Figure 2 (a), we establish a SDM model $\mathcal{M}$ based on deep convolutional neural network (CNN) to downsample a rendered view $\hat{I}_{LR}$ and produce the following low-resolution result $\hat{I}_{LLR}$:

$$\hat{I}_{LLR} = \mathcal{M}\left(G_{LR}, \hat{I}_{LR}\right) \tag{6}$$

where $G_{LR}$ denotes the corresponding **gradient view** of $\hat{I}_{LR}$. To train our **SDM network**, we introduce a single-scene internal learning method in the radiance fields of NeRF. First, we pretrain a **coarse NeRF** $C^*_{coarse}$ by using sparse ray sampling and then employ it to render a novel view $\hat{I}_{LR}$. Second, we use the down-sampled result $I_{LLR}$ (H/s x W/s) of the ground truth image $I_{LR}$ (H x W) to supervise the internal learning, where $s$ denotes the scale factor. The goal is to optimize the following loss function by using mean square error (MSE):

$$\mathcal{L}_{MSE} = \left\|\hat{I}_{LLR} - I_{LLR}\right\|_2^2. \tag{7}$$

Note that the above internal learning method can enable us to (a) avoid the consumption of additional scene data, and (b) simulate a realistic and complex down-sampling effect.

**SDM network.** We propose to exploit a lightweight CNN to realize our SDM model by performing layer-wise down-sample. To avoid bringing convolution to a higher dimension, we use Pixel Adaptive Convolution [43] (PAC) as the down-sample layer, which generates adaptive convolution weights for each pixel based on the provided guidance information (i.e. $G_{LR}$) to fit complex mapping relationships:

$$v'_i = \sum_{j \in \Omega(i)} K\left(f_i, f_j\right) W\left[p_i - p_j\right] v_j + b \tag{8}$$

where $\Omega(i)$ denotes the neighbors of pixel $i$, $K$: $\mathbb{R}^{c' \times c \times s \times s}$ denotes the kernel function, $W$ denotes the weight and $b$ denotes the bias. Another advantage of using PAC lies in that it can effectively establish and transfer the mapping

relationship across scales, which can enable us to train the SDM model by using low-resolution data and generalize it to process the high-resolution data. As a result, our internal learning method can effectively model the internal recurrence of information across the different scales. Different from most existing works [8, 9] that lack in-depth mining of the information hidden in a single scene, SDM takes the scene structure and surrounding pixel information into consideration to model the transformation between different data scales.

**Coarse NeRF.** To facilitate the down-sampling process, we initially train a coarse NeRF by using patches of low-resolution rays and low-resolution ground truth images $I_{LR}$ (HxW). A novel view can be rendered by using:

$$\hat{I}_{LR} = C^*_{\text{coarse}}\left(P_{LR}\right) \tag{9}$$

where $P_{LR}$ denotes the low-resolution patches of rays.

**Gradient View $G$** is used to regularize the possible solutions of SDM model. To obtain $G$, We first utilize the Sobel operator to extract the horizontal and vertical gradients of the rendered low-resolution image $\hat{I}_{LR}$:

$$D_u = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} * \hat{I}_{LR} \tag{10}$$

$$D_v = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 2 \end{bmatrix} * \hat{I}_{LR} \tag{11}$$

Then, we calculate the following gradient magnitude as $G_{LR}$ to summarize the edge information in $\hat{I}_{LR}$:

$$G_{LR} = \sqrt{D_u{}^2 + D_v{}^2} \tag{12}$$

## 4.2. SDM-Guided Super-Resolution Training

Initially, the coarse NeRF is trained by using an available low-resolution ground truth dataset of HxW, which can render high-fidelity novel views at the same resolution. However, when performing super-resolution rendering at a higher scale factor of $s$, it may inevitably confront the blurring effect because there exist many sampling points (or sampling gaps) whose values are usually undefined due to the lack of detailed information at higher resolutions. The distortion would become more severe for larger $s$.

To address the problem, we propose a SDM-guided super-resolution training method to obtain a fine NeRF, $C_{fine}$, by sampling a grid of $s^2$ rays within each pixel (i.e. super-sampling) and filling in the high-frequency details that are undefined in the low-resolution NeRF. The main difficulty lies in the sub-pixel supervision of the super-resolution (i.e. sHxsW) rendering results by using a low-resolution (i.e. HxW) ground truth view. As shown in Figure 2 (b), we adopt an inverse rendering strategy by generalizing the well-learned SDM model to assist the training of the fine NeRF, where the resolution of the output of the fine NeRF is sHxsW, the input and output resolutions of SDM model are sHxsW and HxW, respectively.

**Inverse Rendering.** Instead of using the single ray-based optimization method, we advocate using the patch-wise optimization to capture the correlated ray information within the spatial neighborhood, which would favor the internal learning of the scene information. A novel super-resolution view $\hat{I}_{HR}$ can be rendered by using the patches of high-resolution rays $P_{HR}$:

$$\hat{I}_{HR} = C_{fine}\left(P_{HR}\right) \tag{13}$$

where the fine NeRF $C_{fine}$ is updated in the training process. Then, the obtained high-resolution patch $\hat{I}_{HR}$ is mapped to its corresponding low-resolution version $\hat{I}_{LHR}$ by performing inverse rendering using our pretrained SDM model $\mathcal{M}^*$:

$$\hat{I}_{LHR} = \mathcal{M}^*\left(G_{HR}, \hat{I}_{HR}\right) \tag{14}$$

where $G_{HR}$ denotes the gradient view of $\hat{I}_{HR}$. The resolution of $\hat{I}_{LHR}$ is HxW, which is the same as that of the ground truth training data.

**Objective Function.** The overall objective of super-resolution training is to optimize the following loss function:

$$\mathcal{L} = \mathcal{L}_{MSE} + \lambda \mathcal{L}_{perc} \tag{15}$$

where $\lambda$ denotes the hyper-parameter. $\mathcal{L}_{MSE}$ is the MSE loss used to optimize the parameters of the fine NeRF:

$$\mathcal{L}_{MSE} = \left\| \hat{I}_{LHR} - I_{LR} \right\|_2^2. \tag{16}$$

$\mathcal{L}_{perc}$ is the patch-based perceptual loss used to improve high-frequency details by estimating the similarity between predicted patches $\hat{I}_{LHR}$ and ground-truth $I_{LR}$ in the feature space via a pretrained 19-layer VGG network $\varphi$:

$$\mathcal{L}_{perc} = \left\| \varphi(\hat{I}_{LHR}) - \varphi(I_{LR}) \right\|_2^2. \tag{17}$$

**Coarse-to-Fine Optimization.** To effectively train our model, we present a coarse-to-fine strategy to speed up the convergence. We use the first $N_1$ epochs to train the coarse NeRF. Then, we use $N_2$ epochs to train the fine NeRF. The employed pretrained SDM model is optimized separately before training the fine NeRF by using $N_3$ epochs. Note that by integrating the information from multiple views during the training process, our inverse rendering strategy can bridge more information to fill in the undefined gaps in super-resolution rendering, which can help to handle the complex rendering challenges.

### 4.3. Temporal Ensemble

The super-resolution training of 3D NeRF is analogous to that of 2D images. Since it is an inverse problem, there is more than one solution for the high-quality NeRF for the corresponding low-resolution view. In fact, from the perspective of ill-posed problems, the constraints on a gradient are equivalent to the regularization, allowing the gradient
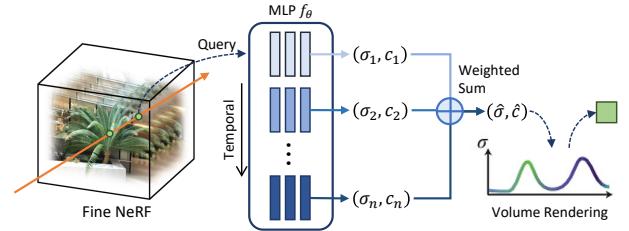


**Figure 3:** Illustration of the temporal ensemble strategy. The radiance fields at multiple time steps are ensembled with respect to color and density.

descent method to maintain a relatively ideal optimization direction. Although regularization can narrow the solution space, local degradation still exists, which may affect the robustness of the results.

To deal with the problem, we further propose a temporal ensemble strategy to narrow the solution space and obtain a more smooth estimation of density and color corresponding to 5D coordinates. Figure 3 illustrates the records of the radiance fields at multiple time steps. The estimated density and color information are averaged by using

$$\hat{\sigma}, \hat{c} = \frac{1}{N} \sum_i^N \sigma_i, \frac{1}{N} \sum_i^N c_i \tag{18}$$

where $N$ denotes the number of radiance fields to be ensembled, $\sigma_i$ and $c_i$ denote the estimated volume density and view-dependent color. We need to point out that our temporal ensemble strategy can compensate for the errors of the estimated super-resolution scene, effectively alleviating the inescapable local degradation issue associated with super-resolution training.

## 5. Experiment

In this section, we conduct quantitative and qualitative studies to show the performance of the proposed method.

### 5.1. Datasets

**Blender Dataset** is a popular public dataset that contains eight synthetic objects. Each scene is captured from virtual cameras positioned in a hemisphere layout, focusing inwards. The dataset includes 100 training images per scene and reserves 200 images for testing. All images are in a high-resolution format of 800×800.

**LLFF Dataset** is a popular public dataset that comprises eight real-world scenes, primarily composed of forward-facing images. These scenes are captured using 20 to 62 images each. For testing purposes, 1/8 of these images are used as the test set. The resolution of all images in this dataset is 1008×756.

### 5.2. Implementation Details

**Basic Settings.** TensoRF-VM-192 is used to implement our NeRF model, and the configurations, such as tensor resolution and number, remain unchanged. We use the Adam optimizer with initial learning rates of 0.02 for tensor factors
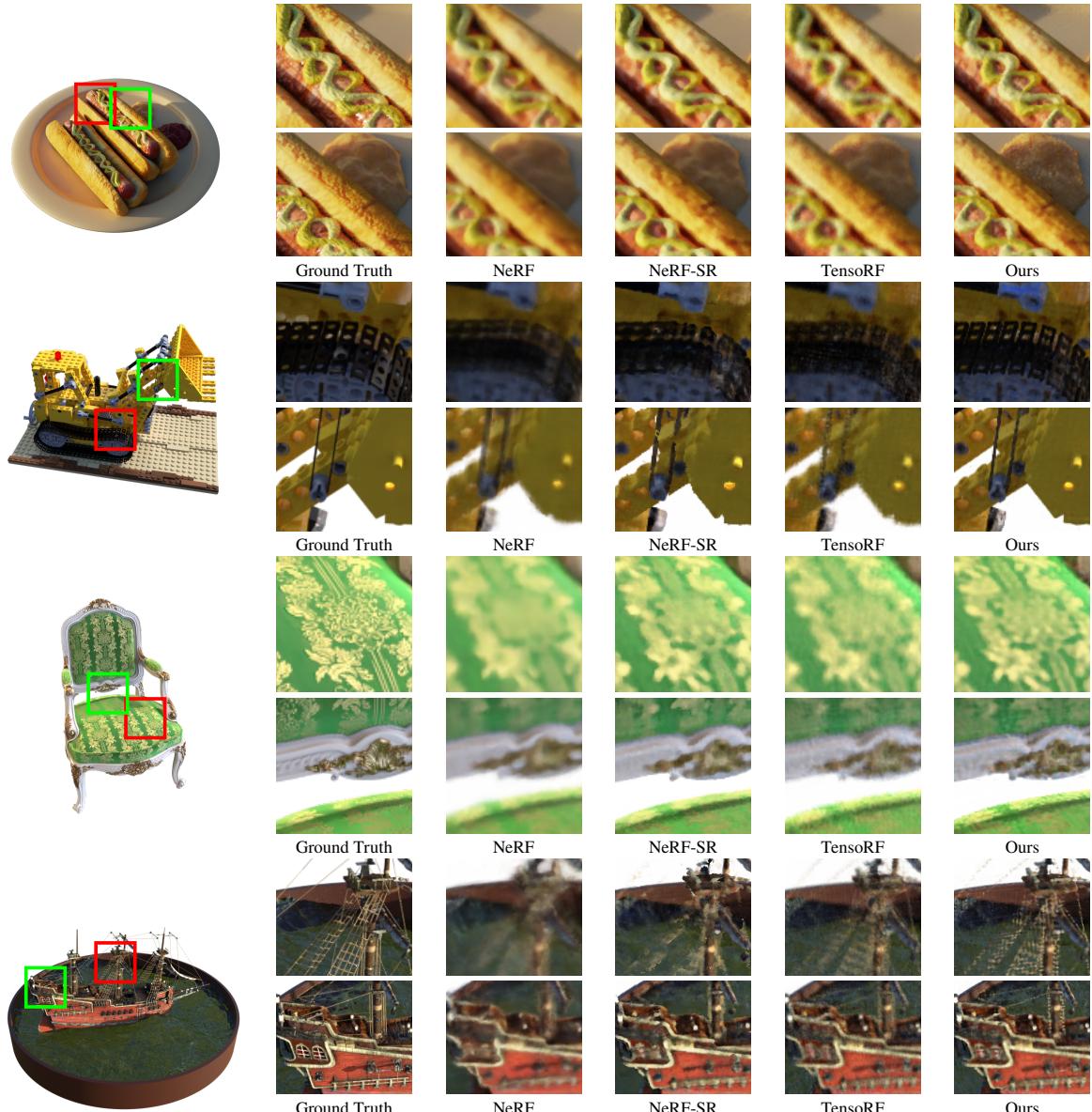
**Figure 4:** Comparison of the synthesized super-resolution views (x4) of different methods on the Blender dataset. For each of the views in the first column, the results in the first line show the larger version of the red box region, and the results in the second line show the larger version of the green box region. Our results show clearer details than NeRF, NeRF-SR, and TensoRF.

and 0.001 for the MLP decoder. All our experiments are done on a single RTX 3090 GPU.

**Patch-wise Ray Sampling.** To address the issue of empty space in the blender dataset, we construct a mask for each patch to speed up training using the density estimation of the coarse NeRF. In the Super-Resolution Training stage, we set patch size to 16x16 (x2) and 32x32 (x4), and the corresponding batch sizes are set to 32 and 8.

**Coarse NeRF.** The configuration of training the coarse NeRF is $N_1 = 5,000$ for the blender dataset, and $N_1 = 10,000$ for the LLFF dataset, and the batch size is set to 4096 for the randomly sampled ray.

**Fine NeRF.** The configuration of training the fine NeRF is $\lambda = 0.03$, $N_2 = 25,000$ for the blender dataset and $N_2 = 20,000$ for the LLFF dataset, and the batch size is set to 8192.

**SDM Model.** The configuration of training the SDM model is $N_3 = 10,000$, patch size 16 (x2), 32 (x4) and the corresponding batch sizes are set to 32 and 8.

### 5.3. Evaluation

In this section, we rely on comparative study to show the performance of our approach by using the following baseline and state-of-the-art methods:

- **NeRF [3].** Leveraging NeRF's implicit function representation, which enables new view synthesis at any resolution, we utilize publicly available NeRF code. This allows us to implement LR inputs and directly render HR images.
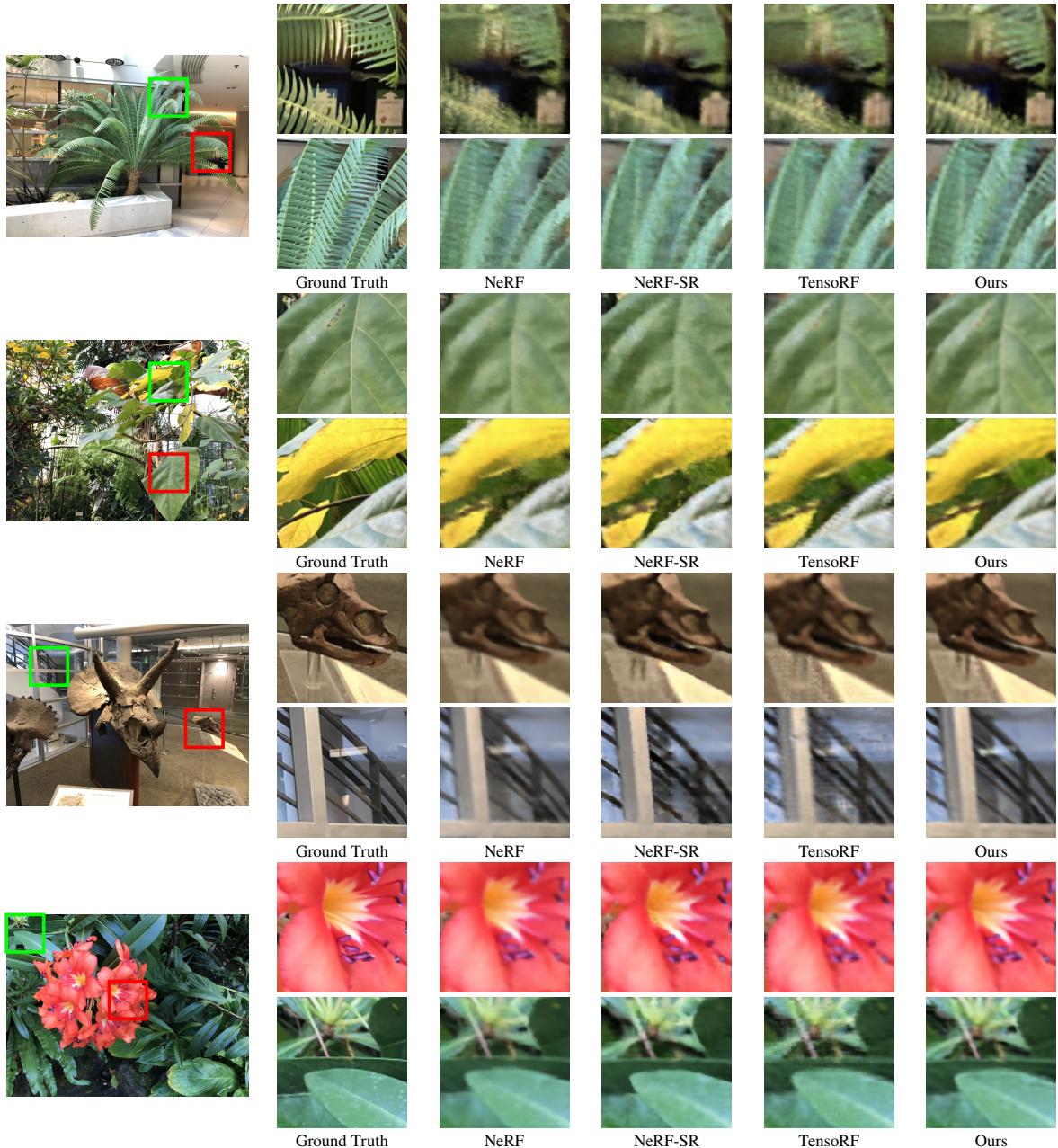
**Figure 5:** Comparison of the synthesized super-resolution views (x4) of different methods on the LLFF dataset. Our results show clearer details than NeRF, NeRF-SR, and TensoRF, such as the textures of leaves and flowers.

- **TensoRF [24].** We compare with TensoRF implementation of NeRF by using the official source code, which acts as the baseline of our method.

- **NeRF-SR [9].** NeRF-SR is designed for synthesizing HR novel views from LR image inputs. Our experimental data is derived from both the NeRF-SR source code and related literature.

- **NVSR [8].** As a contemporary work to NeRF-SR, NVSR synthesizes high-resolution views through super-resolution tri-planes. To save training time, We use four scenes from the Blender dataset to train NVSR and the remaining four for testing. The official source code is used.

- **SISR based methods [7, 35].** We also compare our method with SISR-based methods by performing pre-processing or post-processing. Two representative SISR models are used: ZSSR [7] and SWINIR [35]. For pre-processing: SISR is used to preprocess the LR ground truth images to obtain HR images for training Super-Resolution TensoRF: ZSSR-TensorRF, SWINIR-TensorRF. For post-processing: SISR is used to post-process the rendered LR views of TensoRF: TensoRF-ZSSR and ZSSR-TensorRF.

We qualitatively evaluate the synthesized view quality against the ground truth under identical poses. We quantitatively evaluate our approach with the following three metrics: Peak Signal Noise Ratio (PSNR), Structural Similarity

**Table 1**

Results of novel view synthesis on blender and LLFF datasets for scale factors ×2 and ×4 on four input resolutions: blender x2 (400x400→800x800), blender x4 (200x200→800x800), LLFF x2 (504x378→1008x756), LLFF x4 (252x189→1008x756).

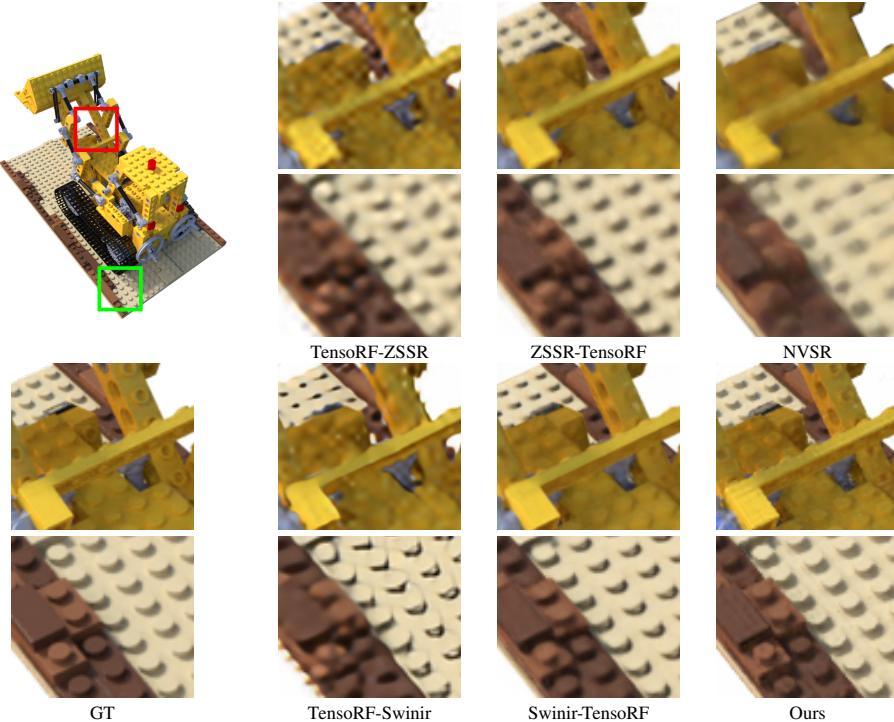| Method | Factor | Blender | | | LLFF | | |
|--------|--------|---------|--------|--------|---------|--------|--------|
| | | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
| NeRF | | 28.06 | 0.923 | 0.052 | 24.57 | 0.750 | 0.192 |
| TensoRF | x2 | 31.45 | 0.952 | 0.049 | 25.88 | 0.810 | 0.175 |
| NeRF-SR | | 30.08 | 0.939 | 0.050 | 25.26 | 0.755 | 0.220 |
| Ours | | **31.93** | **0.954** | **0.042** | **26.49** | **0.837** | **0.155** |
| NeRF | | 27.47 | 0.910 | 0.128 | 21.69 | 0.626 | 0.313 |
| TensoRF | x4 | 28.01 | 0.910 | 0.113 | 23.82 | 0.694 | 0.358 |
| NeRF-SR | | 28.46 | 0.921 | 0.076 | 23.51 | 0.693 | 0.297 |
| Ours | | **29.69** | **0.929** | **0.069** | **24.80** | **0.749** | **0.283** |



**Figure 6:** Visual comparison of our results with that of NVSR and SISR-based methods on Lego.

Index Measure (SSIM) [44], and Learned Perceptual Image Patch Similarity (LPIPS) [45].

## 5.4. Main Results

**Qualitative evaluation.** Figure 4 and 5 visually compare super-resolution rendering results of different methods on both the Blender and LLFF datasets, where the input resolution is 200, and the upsampling factor is 4. NeRF and TensoRF suffer from obvious blurring effects due to the existence of a sampling gap. NeRF-SR can obtain better visual quality because the employment of the supersampling strategy can make up the sampling gap. In contrast, our results show more precise details that are closer to the ground truth. This can be contributed to the advantage of our parameterized SDM model which is superior than the linear interpolation method (i.e. average pooling) used in NeRF-SR. In Figure 6, we visually compare our method with NVSR and the SISR-based methods. The results of NVSR, TensoRF-ZSSR and ZSSR-TensoRF suffer from the

problems of blurry and structure distortion. Although the results of TensoRF-Swinir and Swinir-TensoRF perform better than those of ZSSR, they still suffer from the problem of structure distortion. Compared with pre-processing, the post-processing may easily confront the problem of structure distortions in image details. Instead, our results exhibit better visual quality because the optimization of our model fully happens in the 3D space supervised by our SDM model. Besides, our method also demonstrates its advantages on the training speed, the consumption of scene data and the quality of the rendered view in Figure 1. The speedup can be contributed to the employment of our coarse-to-fine optimization strategy, where the training cost for the lightweight SDM model is very low and can be neglected compared to training NeRF.

**Quantitative evaluation.** Table 1 outlines the quantitative evaluation results. For the Blender dataset, we train the models on resolutions of $400 \times 400$ and $200 \times 200$ while

## TensoRF-SwinIR ( inconsistency )

**Multi-views**

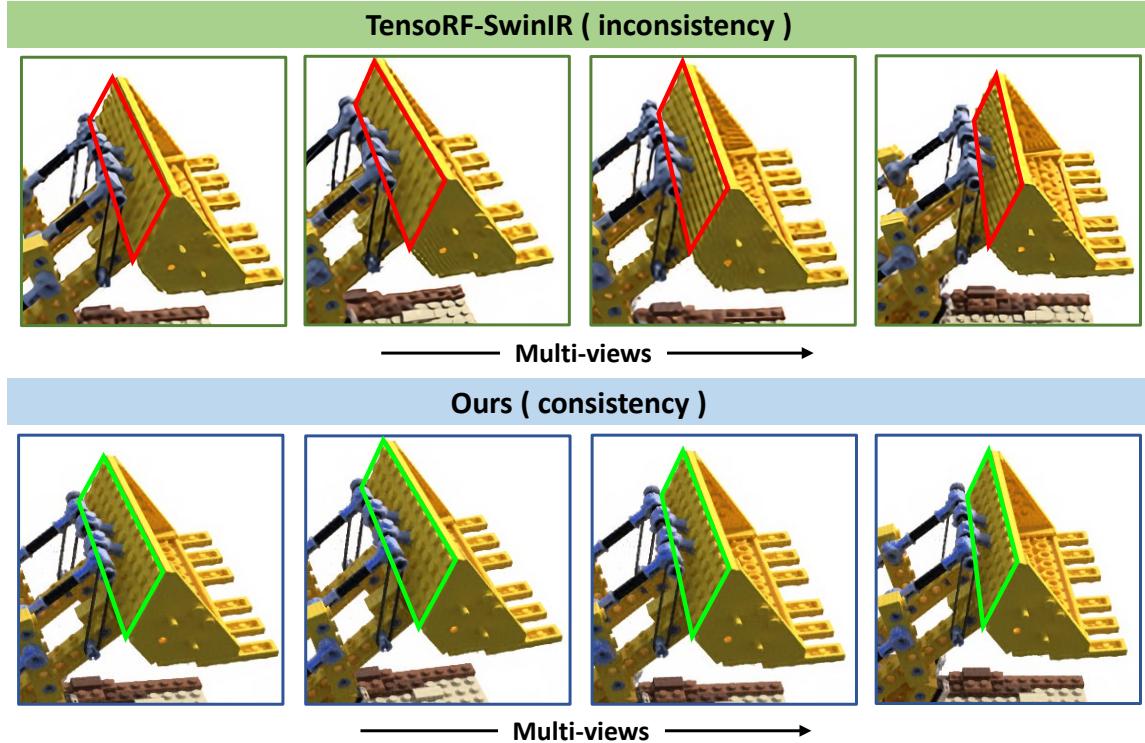## Ours ( consistency )

**Multi-views**

**Figure 7:** Qualitative Comparison of Multi-view Consistency. The synthetic views of TensoRF-SwinIR suffer from inconsistent structure or texture (see the red boxes in the first row), whereas our model maintains remarkable multi-view consistency.

**Table 2**

Comparison of our results (×4) with that of NVSR and the SISR based methods: ZSSR-TensoRF, Swinir-TensoRF, TensoRF-ZSSR and TensoRF-Swinir.

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|--------|--------|--------|---------|
| TensoRF-ZSSR | 27.93 | 0.894 | 0.124 |
| ZSSR-TensoRF | 29.60 | 0.918 | 0.090 |
| TensoRF-Swinir | 27.27 | 0.897 | 0.100 |
| Swinir-TensoRF | 30.05 | 0.920 | 0.070 |
| NVSR | 27.71 | 0.892 | 0.116 |
| Ours | **30.36** | **0.925** | **0.064** |

testing them at 2x and 4x scales, respectively. Similarly, for the LLFF dataset, we train the models on resolutions of 504 × 378 and 252 × 189, and test them at 2x and 4x scales. It is easy to observe that our results exhibit the best performances on all the evaluation measures across the two datasets. In the x2 super-resolution task, tensor neural representation brings higher benefits than the dense radiation field. In the x4 super-resolution task, the sampling gap is larger, and the dense radiation field brings greater benefits. Due to differences in test scenarios, the results of NVSR are not included in Table 1. In Table 2, we compare our results with NVSR as well as that of the methods based on SISR under the same test configurations. It is obvious that our method outperforms the others with substantial improvements in PSNR, SSIM, and LPIPS metrics. The recent Swinir method outperforms ZSSR in super-resolution training. The post-processing is inferior to pre-processing, which is because that the super-resolution methods may easily lead to view-inconsistency without considering multi-view information.

**Multi-view consistency evaluation.** We demonstrate the view consistency through multi-view analysis. In Figure 7, we compare our results with that of TensoRF-SwinIR. It is interesting to observe that SwinIR has disrupted the multi-view 3D consistency of TensorRF for super-resolution. For example, the structure or texture of the red boxes becomes distorted for consistent views. In contrast, our results can well preserve the multi-view consistency.

### 5.5. Ablation Study

In this section, we rely on ablation study to show the effectiveness of each component of method by adding one component each time starting from the baseline model TensoRF, including the super-resolution training, the SDM model, the gradient view, and the temporal ensemble. (a) + *super-resolution training* denotes performing super-resolution training of TensoRF as with that of NeRF-SR [9]. (b) + *SDM* denotes performing super-resolution training of TensoRF supervised by using our SDM model but without gradient view, i.e., (a)+SDM. (c) + *gradient view* denotes the embedding of the gradient information in the SDM model for (b). (d) +*temporal ensemble* denotes adding temporal ensemble for (c), which equals our entire model.

We first present the visual results in Figure 8. It is easy to find that the super-resolution training can help to improve the image quality compared with the baseline TensoRF. SDM can help to improve the details (e.g. structure or texture) of
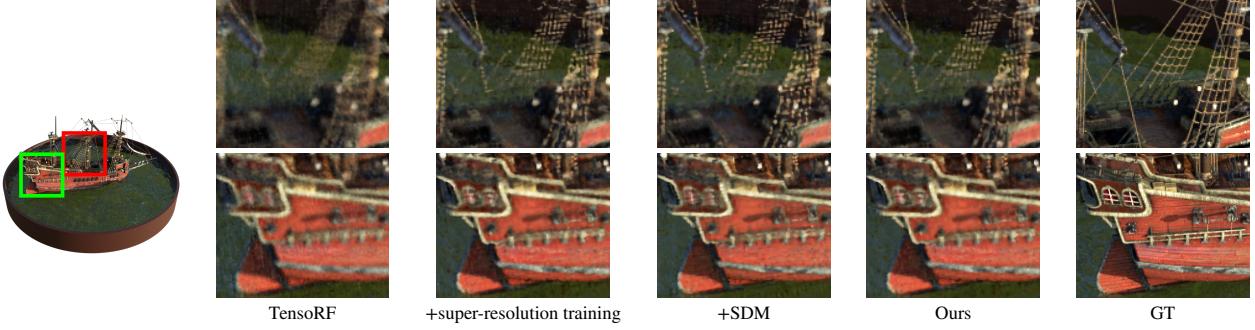
| TensoRF | +super-resolution training | +SDM | Ours | GT |

**Figure 8:** Demonstration of the visual ablation results with respect to super-resolution training, SDM and our final solution. Each of our methods can bring visual improvements to the rendered views.

**Table 3**

Quantitative analysis of ablation study results at a scale factor of x4.

| Method | PSNR ↑ | SSIM ↑ | LPIPS ↓ |
|---|---|---|---|
| Baseline(TensoRF) | 28.01 | 0.910 | 0.113 |
| (a) +super-resolution training | 29.01 | 0.920 | 0.079 |
| (b) +SDM | 29.51 | 0.926 | 0.073 |
| (c) +gradient view | 29.59 | 0.927 | 0.070 |
| (d) +temporal ensemble | **29.69** | **0.929** | **0.069** |

the rendered view. The employment of a temporal ensemble would further improve the smoothness of the result. Then, in Table 3, we present the quantitative ablation results at upscale factor of x4.

**Effectiveness of the super-resolution training.** According to the results listed in Table 3, the super-resolution training can achieve a PSNR gain of 1.0dB, which can be contributed to the employment of the supersampling strategy to make up the sampling gap, where the average pooling mapping is used to perform supervised super-resolution training which is similar to that of NeRF-SR. Super-resolution training can improve the spatial sampling rate of the radiance field, allowing it to obtain high-quality novel views with high-resolution volume rendering.

**Effectiveness of the SDM mapping.** The results of (a) are still insufficient because the average operation for super-resolution training cannot support obtaining the fine details. After replacing the average pooling with our SDM model, our results (b) in Table 3 can achieve a further improvement, such as the PSNR gain of 0.5dB. These improvements indicate that our SDM model can be used as a suitable supervisor for high-frequency super-resolution training.

**Effectiveness of gradient view and temporal ensemble.** As shown in Table 3 (c) and (d), the employment of gradient view and temporal ensemble can also produce positive effects towards lifting the performance of our method. For example, they can totally achieve a 0.18dB improvement on PSNR score. The improvements can contribute to the regularization of gradient view information and the temporal ensemble strategy.

Notice that our final solution has significantly improved the performance of the baseline (i.e. TensoRF) on all the evaluation items. Together with the visual comparison results, we have demonstrated the ability of each component of our approach.

## 6. Conclusion

In this paper, we propose a zero-shot super-resolution training method for neural radiance fields by focusing on addressing the problem of the lack of high-resolution ground truth. By employing the idea of internal learning, we can obtain a single-scene degeneration mapping model, which is further used in inverse rendering to supervise the super-resolution training of the fine NeRF for synthesizing high-quality, high-resolution novel views without consuming any external high-resolution scene data. With the temporal ensemble strategy, our method can compensate for the errors of super-resolution scene estimation. The comparative and ablation studies show that our method can achieve state-of-the-art performances by producing high-quality novel views with more fine details.

## Acknowledgments

## Conflict of interest

The authors declare that they have no conflict of interest.

# References

[1] J. L. Schonberger, J.-M. Frahm, Structure-from-motion revisited, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 4104–4113.

[2] Y. Furukawa, C. Hernández, et al., Multi-view stereo: A tutorial, Foundations and Trends® in Computer Graphics and Vision 9 (2015) 1–148.

[3] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, R. Ng., Nerf: Representing scenes as neural radiance fields for view synthesis, 2020.

[4] C.-H. Lin, J. Gao, L. Tang, T. Takikawa, X. Zeng, X. Huang, K. Kreis, S. Fidler, M.-Y. Liu, T.-Y. Lin, Magic3d: High-resolution text-to-3d content creation, in: IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2023.

[5] B. Poole, A. Jain, J. T. Barron, B. Mildenhall, Dreamfusion: Text-to-3d using 2d diffusion, arXiv (2022).

[6] J. Tang, T. Wang, B. Zhang, T. Zhang, R. Yi, L. Ma, D. Chen, Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior, arXiv preprint arXiv:2303.14184 (2023).

[7] A. Shocher, N. Cohen, M. Irani, Zero-shot super-resolution using deep internal learning, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 3118–3126.

[8] Y. Bahat, Y. Zhang, H. Sommerhoff, A. Kolb, F. Heide, Neural volume super-resolution, 2022. arXiv:2212.04666.

[9] C. Wang, X. Wu, Y. Guo, S. Zhang, Y. Tai, S. Hu, NeRF-SR: High quality neural radiance fields using supersampling, in: Proceedings of the 30th ACM International Conference on Multimedia, ACM, 2022.

[10] S. Saito, Z. Huang, R. Natsume, S. Morishima, A. Kanazawa, H. Li, Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization, in: The IEEE International Conference on Computer Vision (ICCV), 2019.

[11] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, Y. Sheikh, Neural volumes: Learning dynamic renderable volumes from images, ACM Trans. Graph. 38 (2019) 65:1–65:14.

[12] V. Sitzmann, J. N. Martel, A. W. Bergman, D. B. Lindell, G. Wetzstein, Implicit neural representations with periodic activation functions, in: arXiv, 2020.

[13] J. J. Park, P. Florence, J. Straub, R. Newcombe, S. Lovegrove, Deepsdf: Learning continuous signed distance functions for shape representation, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019.

[14] Y. Shi, R. Yang, Z. Wu, P. Li, C. Liu, H. Zhao, G. Zhou, City-scale continual neural semantic mapping with three-layer sampling and panoptic representation, Knowledge-Based Systems 284 (2024) 111145.

[15] C. Sun, M. Sun, H. Chen, Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction, in: CVPR, 2022.

[16] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, A. Kanazawa, PlenOctrees for real-time rendering of neural radiance fields, in: ICCV, 2021.

[17] L. Liu, J. Gu, K. Z. Lin, T.-S. Chua, C. Theobalt, Neural sparse voxel fields, NeurIPS (2020).

[18] P. Hedman, P. P. Srinivasan, B. Mildenhall, J. T. Barron, P. Debevec, Baking neural radiance fields for real-time view synthesis, ICCV (2021).

[19] W. Bian, Z. Wang, K. Li, J.-W. Bian, V. A. Prisacariu, Nope-nerf: Optimising neural radiance field with no pose prior, 2023. arXiv:2212.07388.

[20] C.-H. Lin, W.-C. Ma, A. Torralba, S. Lucey, Barf: Bundle-adjusting neural radiance fields, in: IEEE International Conference on Computer Vision (ICCV), 2021.

[21] R. Zha, Y. Zhang, H. Li, Naf: Neural attenuation fields for sparse-view cbct reconstruction, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, 2022, pp. 442–452.

[22] Y. Fang, L. Mei, C. Li, Y. Liu, W. Wang, Z. Cui, D. Shen, Snaf: Sparse-view cbct reconstruction with neural attenuation fields, arXiv preprint arXiv:2211.17048 (2022).

[23] W. Song, H. Zheng, J. Yang, C. Liang, L. He, Oral-nexf: 3d oral reconstruction with neural x-ray field from panoramic imaging, arXiv preprint arXiv:2303.12123 (2023).

[24] A. Chen, Z. Xu, A. Geiger, J. Yu, H. Su, Tensorf: Tensorial radiance fields, in: Computer Vision – ECCV 2022, Springer Nature Switzerland, 2022.

[25] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, P. P. Srinivasan, Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields, in: 2021 IEEE/CVF International Conference on Computer Vision (ICCV), 2021, pp. 5835–5844.

[26] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, P. Hedman, Zip-nerf: Anti-aliased grid-based neural radiance fields, ICCV (2023).

[27] H. Feng, L. Wang, Y. Li, A. Du, Lkasr: Large kernel attention for lightweight image super-resolution, Knowledge-Based Systems 252 (2022) 109376.

[28] H. Liu, F. Cao, C. Wen, Q. Zhang, Lightweight multi-scale residual networks with attention for image super-resolution, Knowledge-Based Systems 203 (2020) 106103.

[29] G. Huang, Z. Liu, K. Q. Weinberger, Densely connected convolutional networks, CoRR abs/1608.06993 (2016).

[30] B. Lim, S. Son, H. Kim, S. Nah, K. M. Lee, Enhanced deep residual networks for single image super-resolution, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops, 2017.

[31] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, Y. Fu, Image super-resolution using very deep residual channel attention networks, in: ECCV, 2018.

[32] J. Kim, J. K. Lee, K. M. Lee, Deeply-recursive convolutional network for image super-resolution, CoRR abs/1511.04491 (2015).

[33] Y. Guo, J. Chen, J. Wang, Q. Chen, J. Cao, Z. Deng, Y. Xu, M. Tan, Closed-loop matters: Dual regression networks for single image super-resolution, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2020.

[34] M. Haris, G. Shakhnarovich, N. Ukita, Deep back-projection networks for super-resolution, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 1664–1673.

[35] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, R. Timofte, Swinir: Image restoration using swin transformer, in: 2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW), 2021, pp. 1833–1844.

[36] Z. Chen, Y. Zhang, J. Gu, Y. Zhang, L. Kong, X. Yuan, Cross aggregation transformer for image restoration, in: NeurIPS, 2022.

[37] Z. Chen, Y. Zhang, J. Gu, L. Kong, X. Yang, F. Yu, Dual aggregation transformer for image super-resolution, in: ICCV, 2023.

[38] C. Saharia, J. Ho, W. Chan, T. Salimans, D. J. Fleet, M. Norouzi, Image super-resolution via iterative refinement, arXiv:2104.07636 (2021).

[39] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, C. C. Loy, Esrgan: Enhanced super-resolution generative adversarial networks, in: The European Conference on Computer Vision Workshops (ECCVW), 2018.

[40] H. Li, Y. Yang, M. Chang, S. Chen, H. Feng, Z. Xu, Q. Li, Y. Chen, Srdiff: Single image super-resolution with diffusion probabilistic models, Neurocomputing 479 (2022) 47–59.

[41] K. C. Chan, X. Wang, X. Xu, J. Gu, C. C. Loy, Glean: Generative latent bank for large-factor image super-resolution, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2021.

[42] Y. Yoon, K.-J. Yoon, Cross-guided optimization of radiance fields with multi-view image super-resolution for high-resolution novel view synthesis, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2023, pp. 12428–12438.

[43] H. Su, V. Jampani, D. Sun, O. Gallo, E. Learned-Miller, J. Kautz, Pixel-adaptive convolutional neural networks, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 11158–11167. doi:10.1109/CVPR.2019.01142.

[44] Z. Wang, E. Simoncelli, A. Bovik, Multiscale structural similarity for image quality assessment, in: The Thrity-Seventh Asilomar

Conference on Signals, Systems & Computers, 2003, volume 2, 2003, pp. 1398–1402 Vol.2.

[45] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, O. Wang, The unreasonable effectiveness of deep features as a perceptual metric, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 586–595. doi:10.1109/CVPR.2018.00068.