# IE-NeRF: Inpainting Enhanced Neural Radiance Fields in the Wild

**Shuaixian Wang[1], Haoran Xu[1,2], Yaokun Li[1], Jiwei Chen[1], Guang Tan[1*]**

[1]**Sun Yat-sen University, Guangdong, China**
[2]**Pengcheng Laboratory, Shenzhen, China**
{wangshx29, xuhr9, liyk58, chenjw269}@mail2.sysu.edu.cn, tanguang@mail.sysu.edu.cn

## Abstract

We present a novel approach for synthesizing realistic novel views using Neural Radiance Fields (NeRF) with uncontrolled photos in the wild. While NeRF has shown impressive results in controlled settings, it struggles with transient objects commonly found in dynamic and time-varying scenes. Our framework called *Inpainting Enhanced NeRF*, or IE-NeRF, enhances the conventional NeRF by drawing inspiration from the technique of image inpainting. Specifically, our approach extends the Multi-Layer Perceptrons (MLP) of NeRF, enabling it to simultaneously generate intrinsic properties (static color, density) and extrinsic transient masks. We introduce an inpainting module that leverages the transient masks to effectively exclude occlusions, resulting in improved volume rendering quality. Additionally, we propose a new training strategy with frequency regularization to address the sparsity issue of low-frequency transient components. We evaluate our approach on internet photo collections of landmarks, demonstrating its ability to generate high-quality novel views and achieve state-of-the-art performance.

## 1 Introduction

Synthesizing novel views of a scene from limited captured images is a long-standing problem in computer vision, which is fundamental for applications in mixed reality [7], 3D reconstruction [8]. Canonical view synthesizing techniques [14; 31] based on structure-from-motion and image rendering have encountered challenges in maintaining consistency across views, as well as addressing occlusion and distortion. Recently, with the development of implicit scene representation and neural rendering, Neural Radiance Fields (NeRF) [23] have achieved excellent performance in novel view synthesis (NVS).

NeRF employs neural networks to encode radiance properties within a continuous spatial domain, enabling intricate scene reconstructions by learning from multiple viewpoints. While this approach has achieved success in various fields such as computer graphics, computer vision, and immersive
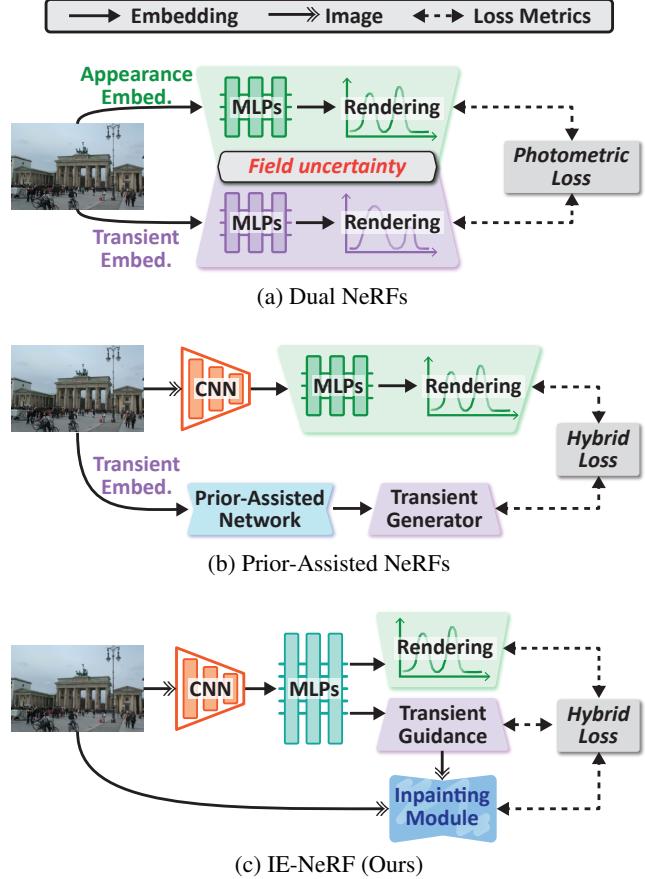


Figure 1: Comparison of different NeRF pipelines that mitigate transient occlusion. (a) Dual NeRFs, which extract transient components by introducing an additional NeRF branch; (b) Prior-Assisted NeRFs, which leverage prior knowledge to assist in separating transient objects; (c) IE-NeRF (Ours), which integrates the inpainting module to enhance NeRF.

technologies [32], conventional NeRFs often operate under controlled settings with static scenes and consistent lighting conditions [21]. However, in real-world scenarios characterized by time-varying and transien t occlusions, NeRF encounters significant performance degradation.

Existing solutions to this problem can be roughly catego-

rized into two approaches **(i) Dual NeRFs**. As depicted in Figure 1a, this approach extracts transient components by introducing additional NeRF pipelines. Specifically, NeRF-W [20] and subsequent work NRW [29] optimize appearance and transient embeddings through individual NeRF modules, rendering static fields and transient fields respectively.**(ii) Prior-Assisted NeRFs**. As illustrated in Figure 1b, this approach leverages prior knowledge to assist in separating transient objects from the background. In particular, SF-NeRF [15] introduced an occlusion filtering module to remove transient objects by a pre-trained semantic segmentation model. While Ha-NeRF [5] eliminated transient components pixel-wise by an anti-occlusion module image dependently. Despite showcasing promising results, these approaches still face the problem of inaccurate transient decomposition from the complex scene and entanglement reconstruction of static appearances and occlusion.

In this paper, we address this problem from a fresh perspective by drawing inspiration from recent advances in image inpainting [4; 36], which aims to remove unwanted objects and make imaginative restoration. The key insight of our work is that eliminating transient objects in NeRF is, in essence, an inpainting process during reconstruction. In addressing the mentioned challenges of NeRF, our objective is to perceive and separate undesired, blurry foreground areas in rendered images and produce plausible and consistent background scenes, which is exactly the expertise of image inpainting.

To this, we propose Inpainting Enhanced NeRF, or **IE-NeRF**, a novel approach that utilizes inpainting to separate transient content within NeRF. As illustrated in Figure 1c, our model comprises three modules: the regular NeRF for static scene image rendering, the transient mask generator, and the inpainting module for removing transient components and repairing static images. Given an image, the model first encodes it into a high-dimensional vector using a CNN. Unlike the conventional NeRF, the MLP network in IE-NeRF predicts not only the color and voxel density, but also the masks for the transient components. The original image along with the masks are fed into the inpainting module to generate the restored static image. Finally, the rendered static image, obtained through volume rendering, is optimized by minimizing the photometric loss with the restored static image. Comprehensive experiments validate the performance of our method.

Our contributions can be summarized as follows:

- We propose to enhance NeRF in uncontrolled environments by incorporating an image inpainting module. Drawing on the success of the inpainting technique, this represents a novel approach compared to prior efforts.

- We enhance the MLPs of NeRF to generate both static elements and transient masks for an image. The transient masks enable the removal of transient elements and contribute to the optimization of static rendering.

- We introduce a training strategy that adopts frequency regularization with integrated positional encoding. Our strategy facilitates faster inference and early separation of transient components during training.

## 2 Related Work

### 2.1 Novel View Synthesis

Novel View Synthesis is a task that aims to generate new views of a scene from existing images [27]. NVS usually involves geometry-based image reprojection [3; 12] and volumetric scene representations [9; 19]. The former applies techniques such as Structure-from-Motion [11] and bundle adjustment [34] to construct a point cloud or triangle mesh to represent the scene from multiple images, while the latter focuses on unifying reconstruction and rendering in an end-to-end learning fashion.

Inspired by the layered depth images, explicit scene representations such as multi-plane images [42; 35] and multiple sphere images [10; 1] have also been explored. They use an alpha-compositing [26] technique or learning compositing along rays to render novel views. In contrast, implicit representation learning techniques like NeRF [27; 23] exhibit remarkable capability of rendering novel views from limited sampled data. While NVS has made significant progress, challenges persist, especially in addressing occlusions and enhancing the efficiency of rendering complex scenes.

### 2.2 Neural Rendering

Neural rendering techniques are now increasingly employed for synthesizing images and reconstructing geometry from real-world observations in scene reconstruction [32]. Various approaches utilize image translation networks and different learning components, such as learned latent textures [33], meshes [13], deep voxel [28], 3D point clouds [6], occupancy fields [22], and signed distance functions [25], to enhance realistic content re-rendering and reconstruction.

The prominent NeRF model utilizes a Multi-Layer Perceptron (MLP) to restore a radiance field. Subsequent research has aimed at extending NeRF's capabilities to dynamic scenes [17], achieving more efficient rendering [37], refining pose estimation [41], and exploring few-shot view synthesis [35]. Noteworthy studies [20; 5] addressed the problem of view synthesis using internet photo collections, which often include transient occlusions and varying illumination. Additional work [15] focused on NeRF training in a few-shot setting. In contrast to these approaches, our work seeks to enhance view synthesis in the wild by separating dynamic and static scenes via image inpainting.

### 2.3 Image Inpainting

Image inpainting is a technique that fills in missing or damaged regions in an image, widely used for removing unwanted objects from images, and reconstructing deteriorated images [39]. There have also been researches looking into integrating inpainting within NeRF. Liu *et al.* [18] removed unwanted objects or retouched undesired regions in a 3D scene represented by a trained NeRF. Weder *et al.* [38] proposed a method that utilizes neural radiance fields for plausible removal of objects from output renderings, utilizing a confidence-based view selection scheme for multi-view consistency. Mirzaei *et al.* [24] utilized image inpainting to guide both the geometry and appearance, performing inpainting in
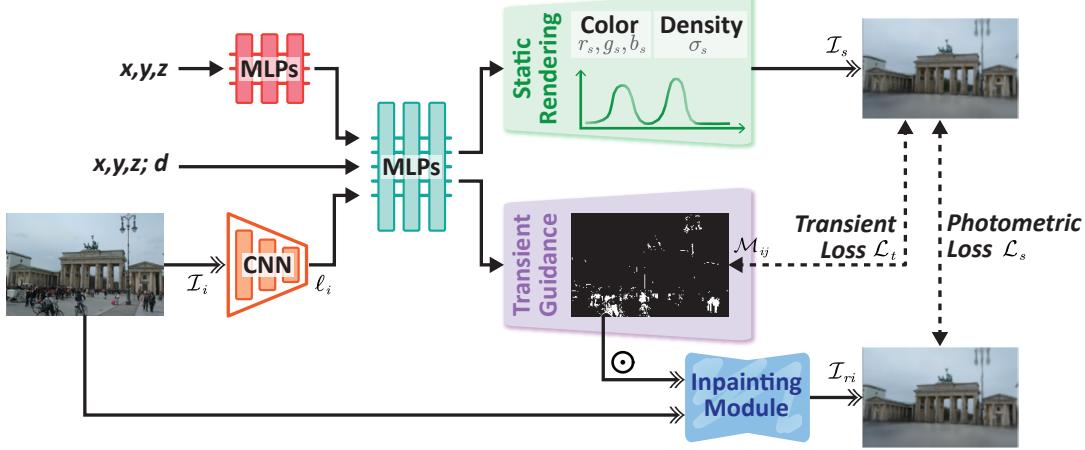
Figure 2: IE-NeRF framework: Given an image $\mathcal{I}_i$, a CNN is used to generate a feature embedding $\ell_i$. Then this embedding, along with the sample location $(x, y, z)$ and the view direction $d$ of the camera ray, is fed into MLPs to produce static color elements $r_s, g_s, b_s$ and radiance intensity $\sigma_s$, as well as the transient mask $\mathcal{M}_{ij}$. The former is utilized to generate the new static scene image $\mathcal{I}_s$ through volumetric rendering, while the latter guides the Inpainting Module for the restoration of the static image $\mathcal{I}_{ri}$. We finally optimize the static image $\mathcal{I}_s$ by minimizing the photometric loss with $\mathcal{I}_{ri}$.

an inherently 3D manner of NeRFs. However, we investigate the synergistic effect of the inpainting module in joint learning within the NeRF pipeline.

## 3 Methodology

### 3.1 NeRF Preliminary

NeRF models a continuous scene using a 5D vector-valued volumetric function $F(\theta)$ on $\mathbb{R}^3 \times \mathbb{S}^2$, implemented as an MLP. This function takes a 3D location $(x, y, z) \in \mathbb{R}^3$ and a 2D viewing direction $d = (\theta, \varphi) \in \mathbb{S}^2$ as input, produces an emitted color $(r, g, b)$ and volume density $(\sigma)$ as outputs. NeRF represents the volumetric density $\sigma(t)$ and color $c(t)$ at point of camera ray $r(t)$ using MLPs with ReLU activation functions. Formally:

$$[\sigma(t), z(t)] = \text{MLP}_{\theta_1}\left[\gamma_x(r(t))\right], \quad (1)$$

$$c(t) = \text{MLP}_{\theta_2}\left[\gamma_d(d), z(t)\right], \quad (2)$$

where $\theta = [\theta_1, \theta_2]$ represents the collection of learnable weights and biases of MLPs. The functions $\gamma_x$ and $\gamma_d$ are predefined encoding functions used for the spatial position and viewing direction, respectively. NeRF models the neural network using two distinct MLPs. The output of the second MLPs is conditioned on $z(t)$, one of the outputs from the first MLPs. This emphasizes the point that the volume density $\sigma(t)$ is not influenced by the viewing direction $d$. To calculate the color of a single pixel, NeRF uses numerical quadrature to approximate the volume rendering integral along the camera ray. The camera ray is represented as $r(t) = o + td$, where $o$ is the center of the projection and $d$ is the direction vector. NeRF's estimation of the expected color $\hat{C}_r$ for that pixel is obtained by evaluating the integral along the ray:

$$\hat{C}(r) = \mathcal{R}(r, c, \sigma) = \sum_{k=1}^{K} T(t_k) \cdot \alpha(\sigma(t_k)\delta_k) \cdot c(t_k), \quad (3)$$

$$T(t_k) = \exp\left(-\sum_{k'=1}^{k-1} \sigma(t_{k'}) \cdot \delta_{k'}\right), \quad (4)$$

where $\alpha(x) = 1 - \exp(-x)$, $(r, c, \sigma)$ signifies the volume radiance field along the ray $r(t)$, while $\sigma(t)$ indicates the density value and $c(t)$ represents the color at each point along the ray. $\delta_k = t_{k+1} - t_k$ denotes the distance between two integration points. For enhanced sampling efficiency, NeRF employs a dual MLP strategy: it comprises coarse and fine networks sharing the same architecture. The optimization of both models' parameters is achieved by minimizing the following loss function:

$$\sum_{ij} \left\| C(r_{ij}) - \hat{C}_c(r_{ij}) \right\|_2^2 + \left\| C(r_{ij}) - \hat{C}_f(r_{ij}) \right\|_2^2, \quad (5)$$

where $C(r_{ij})$ represents the observed color along the ray $j$ in the image $\mathcal{I}_i$, and $\hat{C}_c, \hat{C}_f$ are the predictions of the coarse and fine models, respectively.

### 3.2 Pipeline Overview

The pipeline of our model is illustrated in Figure 2. We use the core model of NeRF that consists of two MLP modules, namely $\text{MLP}_{\theta_1}$ and $\text{MLP}_{\theta_2}$. At the initial stage, the input image $I_i$ is processed through CNN producing a high-dimensional vector $\ell_i$. While the 3D position $(x, y, z)$ of the image $I_i$ is processed through an $\text{MLP}_{\theta_1}$ with the output of $z_t$. Then, the $\ell_i$ and $z_t$ combined with the original position $(x, y, z)$ and directional $d$ are fed into the second network $\text{MLP}_{\theta_2}$. The two networks are described as:

$$[\sigma(t), z(t)] = \text{MLP}_{\theta_1}\left[\gamma_x(r(t))\right], \quad (6)$$

$$c_i(t) = \text{MLP}_{\theta_2}\left[\gamma_d(d), z_i(t)\right]. \quad (7)$$

$\text{MLP}_{\theta_2}$ outputs two components. The first part includes the static color $(r_s, g_s, b_s)$ and radiance intensity $\sigma_s$. These elements are used to generate a new static scene image $I_s$

through volumetric rendering. We construct the volumetric rendering formula based on Eq. (3):

$$\hat{C}_i(r) = \mathcal{R}(r, c, \sigma) = \sum_{k=1}^{K} T(t_k) \cdot \alpha(\sigma(t_k)\delta_k) \cdot c_i(t_k), \quad (8)$$

where $\sigma(t)$ is the static density and $c_i(t)$ is the color radiance. The second part of the output of $\text{MLP}_{\theta_2}$ is the transient mask $M_{ij}$, which is used to guide the inpainting module for restoring the static scene during training. After the original image $I_i$ is inpainted with the transient mask $M_{ij}$, we obtain the restored static image $\mathcal{I}_{ri}$, which acts as supervisory information for the optimization of the rendered static image $\mathcal{I}_s$.

## 3.3 Transient Masks and Inpainting Module

We take advantage of the MLPs to generate the transient mask representation and the pre-trained inpainting model for the inpainting task. The mask is designed to capture dynamic elements in a scene, such as moving objects or changing conditions. Rather than relying on a 3D transient field to reconstruct transient elements specific to each image, as adopted in NeRF-W [21], we use an image-dependent 2D mask map to eliminate the transient effects. The transient mask is generated along with the regular output of density and color radiance from the same MLPs as $\mathcal{M}_{ij}$. The mask maps a 2D-pixel location $l = (u, v)$ and an image-specific transient image embedding $\ell_i$ to a 2D pixel-wise probability representation as follows:

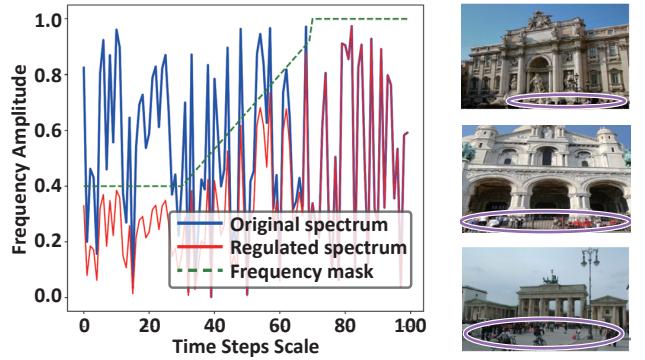$$\mathcal{M}_{ij} = f(l_{ij}, \ell_i). \quad (9)$$

Once the transient mask is generated, it serves as guidance for the inpainting module. Here we use the inpainting model LaMa [30], which is the state-of-the-art single-stage image inpainting system and is robust to large masks with less trainable parameters and inference time. LaMa performs inpainting on a color image $x$ that has been masked by a binary mask $m$, denoted as $x \odot m$. The input to LaMa is a four-channel tensor $x' = stack(x \odot m, m)$, where the mask $m$ is stacked with the masked image $x \odot m$. Taking $x'$, the inpainting network processes the input using a feed-forward network $f_\theta$ in a fully convolutional manner and produces an inpainted image $\hat{x} = f_\theta(x')$. In our pipeline, the original input image $I_i$, along with the thresholded transient mask $m_{ij}$, is fed into the pre-trained inpainting module. With the guidance of $M_{ij}$, LaMa accurately separates the transient components from the original image, then restores and repairs the features of the static scene. This process results in the repaired static scene image $I_{ri}$, which acts as the ground truth for the static rendered image.

## 3.4 Loss Function and Optimization

The rendered static scene image $I_s$ and the repaired static scene image $I_{ri}$ generated by the image inpainting module LaMa are used to calculate the photometric loss as follows:

$$\mathcal{L}_s = \left\| C_s(r_{ij}) - \hat{C}_{ci}(r_{ij}) \right\|_2^2 + \left\| C_s(r_{ij}) - \hat{C}_{fi}(r_{ij}) \right\|_2^2, \quad (10)$$

where $C_s(r_{ij})$ represents the true color of ray $j$ for the rendered static scene image $I_s$ in image $I_i$. $\hat{C}_{ci}(r_{ij})$ and $\hat{C}_{fi}(r_{ij})$



(a) Frequency amplitude during training time steps.

(b) Transient phenomena.

Figure 3: Training strategy: Frequency regularization with integrated positional encoding. We use a step piecewise linearly increasing frequency mask (marked as the dotted green line) to regulate the frequency spectrum based on the training time steps.

represent the color estimate derived from the coarse and the fine model, respectively. In addition to the photometric loss of the static image, we also consider the transient components. The transient image is derived from the original image using a transient mask probability map indicating the visibility of rays originating from the static scene. To separate static and transient components, we optimize the mask map during the training process in an unsupervised manner. Thus, we provide the transient loss as follows:

$$\mathcal{L}_t = (1 - \mathcal{M}_{ij}) \left\| C(r_{ij}) - \hat{C}_s(r_{ij}) \right\|_2^2$$
$$+ \lambda \mathcal{M}_{ij} \left\| C_s(r_{ij}) - \hat{C}_s(r_{ij}) \right\|_2^2. \quad (11)$$

Specifically, the first term tackles the occlusion error by taking into account transient visibility when comparing the rendered static image with the original image. In Eq. (10), we use the repaired image color $C_s(r_{ij})$ as ground truth, while here we rely on the original input color $C_{(r_{ij})}$ to deal with the existence of transient, and $\hat{C}_s(r_{ij})$ is the rendered static scene color optimized by Eq. (10). The second term addresses the reconstruction error of static components between the rendered and repaired ground truth colors, under the assumption that the value of $\mathcal{M}_{ij}$ belongs to the static phenomena. The parameter $\lambda$ is used to adjust the balance between the transient and static components, helping to avoid the neglect of either phenomenon. Then, we can obtain the final optimizing function by combining the loss terms with weight $\beta$:

$$\mathcal{L} = \sum_{ij} \mathcal{L}_s + \beta \sum_{ij} \mathcal{L}_t. \quad (12)$$

## 3.5 Training Strategy with Frequency Regularization

To optimize the training process, we design a frequency regularization scheme with position integral encoding. Instead of employing a single ray per pixel, as done in NeRF, our

| | Brandenburg Gate | | | Sacre Coeur | | | Trevi Fountain | | | Taj Mahal | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| NeRF | 18.90 | 0.816 | 0.232 | 15.60 | 0.716 | 0.292 | 16.14 | 0.601 | 0.366 | 15.77 | 0.697 | 0.427 |
| NeRF-W | 24.17 | 0.891 | 0.167 | 19.20 | 0.807 | 0.192 | 18.97 | 0.698 | 0.265 | **26.36** | **0.904** | 0.207 |
| Ha-NeRF | 24.04 | 0.877 | **0.139** | 20.02 | 0.801 | 0.171 | 20.18 | 0.691 | 0.223 | 19.82 | 0.829 | 0.243 |
| SF-NeRF (30-fews) | 23.23 | 0.846 | 0.178 | 19.64 | 0.757 | 0.186 | 20.24 | 0.657 | 0.243 | 20.86 | 0.820 | 0.218 |
| IE-NeRF (Ours) | **25.33** | **0.898** | 0.158 | **20.37** | **0.861** | **0.169** | **20.76** | **0.719** | **0.217** | 25.86 | 0.889 | **0.196** |

Table 1: Comparison of PSNR, SSIM, and LPIPS for IE-NeRF and other models: NeRF, NeRF-W, Ha-NeRF on Phototourism datasets across specific scenes (Brandenburg Gate, Sacre Coeur, Trevi Fountain, and Taj Mahal).

approach utilizes mip-NeRF [2], casting a cone whose radius adapts to variations in image resolution. This alteration transforms the positional encoding scheme from encoding an infinitesimally small point to integrating within the conical frustum (Integrated Positional Encoding) for each segment of the ray. This encoding method not only enables NeRF to learn multiscale representations but also demonstrates a performance where the participation of high-frequency signals in the encoding gradually increases as the training progresses [40].

For unrestricted scenarios, the presence of transient elements often indicates a low density of low-frequency signals in the training set. It is noteworthy that the frequency of input can be regulated by position encoding. Therefore, we design the training strategy that initiates the process with raw inputs devoid of positional encoding and incrementally boosts the frequency amplitude in each training iteration with a visible mask. This frequency regularization strategy mitigates the instability and vulnerability associated with high-frequency signals, resulting in the separation of transient components during the initial stages of training. Furthermore, the early transient separation mask facilitates adjustments in the subsequent stages of training progress, gradually enhancing NeRF with high-frequency information and preventing both over-smoothing and interference with transient phenomena in the static scene reconstruction. The training strategy is illustrated in Figure 3.

## 4 Experiments

### 4.1 Experiments settings

**Datasets:** Our approach is evaluated on unconstrained internet photo collections highlighting cultural landmarks from the Phototourism dataset. We reconstruct four training datasets based on scenes from Brandenburg Gate, Sacre Coeur, Trevi Fountain, and Taj Mahal. We perform the train set and test set split using the same approach as employed by HA-NeRF [5]. Additionally, we downsample the original images by a ratio of 2 during training, consistent with the approach taken by NeRF-W and HA-NeRF.

**Implementation Details:** Our implementation of the NeRF in the wild network is structured as follows: The entire neural radiance field consists of eight fully connected layers with 256 channels each, followed by two different activation tasks, one is ReLU activations to generate $\sigma$, the other is sigmoid activation following another 128 channels connected layer,

to generate the transient mask possibility $M_{ij}$. Additionally, there is one more fully connected layer with 128 channels and a sigmoid activation, responsible for outputting the static RGB color $c$. For the image inpainting module, we utilize the LaMa model [30], a pre-trained repairer employing a ResNet-like architecture with 3 downsampling blocks, 12 residual blocks using Fast Fourier Convolution (FFC), and 3 upsampling blocks.

**Baselines:** We evaluate our proposed method against several state-of-the-art NeRF models in the Wild, including NeRF, NeRF-W, HA-NeRF, and two variations of IE-NeRF. To ensure a fair comparison, we maintain consistency in the main NeRF architecture across all models. This architecture comprises 8 layers with 256 hidden units for generating density $\sigma$, along with an additional layer of 128 hidden units for color $c$.

**Evaluation:** The performance of our IE-NeRF and baselines are assessed by utilizing a held-out image and its associated camera parameters. After rendering an image from the matching pose, we evaluate its similarity with the ground truth. We provide a set of standard image quality metrics to assess the performance of the models. These metrics include Peak Signal Noise Ratio (PSNR), measuring the fidelity of the reconstructed image; Structural Similarity Index Measure (SSIM), evaluating the structural similarity between the generated and ground truth images; and Learned Perceptual Image Patch Similarity (LPIPS), which leverages perceptual image similarity through insights derived from learned features.

### 4.2 Results Comparision

**Quantitative results:** We conduct experiments to compare the performance of our proposed method with existing baselines on the Phototourism dataset, specifically focusing on scenes of Brandenburg Gate, Sacre Coeur, Trevi Fountain, and Taj Mahal. The four scenes present unique challenges and variations in lighting conditions with transient phenomena. The quantitative results are summarized in Table 1, where we report PSNR/SSIM (higher is better) and LPIPS (lower is better). As depicted in Table 1, our proposed model, IE-NeRF, mostly demonstrates superior performance compared to the baselines. The reported PSNR and SSIM values, indicative of image fidelity and structural similarity, show that our model achieves higher scores, highlighting its effectiveness in reconstructing static image quality. Additionally, the lower LPIPS scores further emphasize the superiority of IE-NeRF in minimizing perceptual differences compared to the
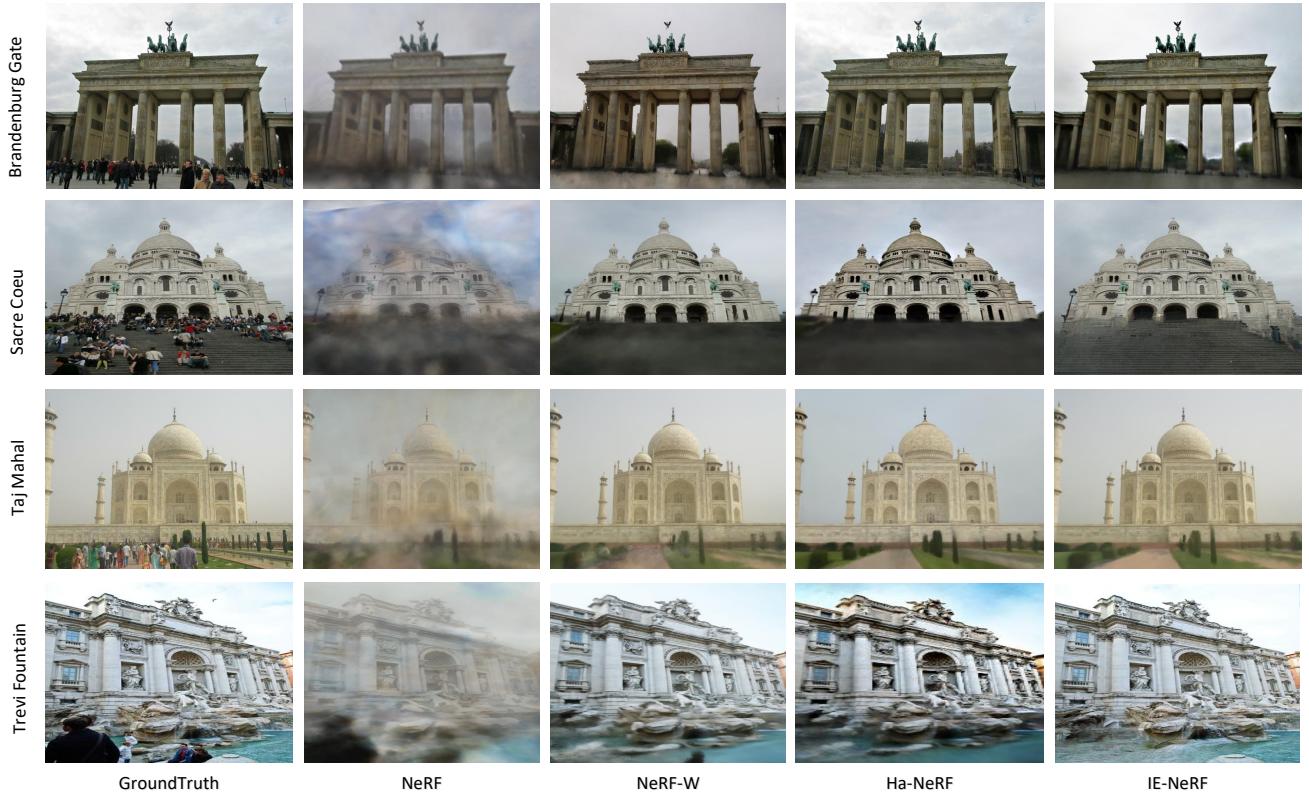
Figure 4: Visual comparison results on four scenes of phototourism constructed datasets. IE-NeRF can remove the transient occlusions more naturally and render a consistent 3D scene geometry with finer details than existing methods.
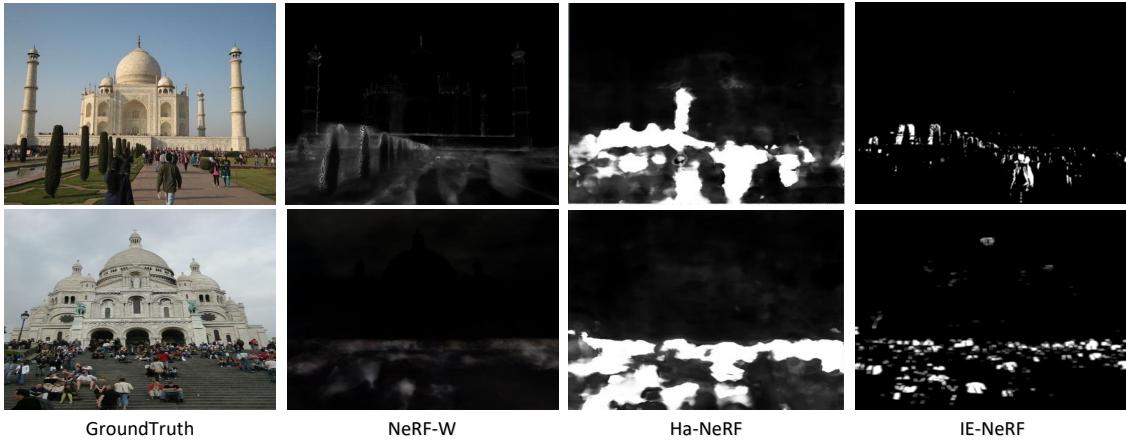


Figure 5: Comparisons of transient components predicted by IE-NeRF and baselines. The predicted transient components of NeRF-W are rendered with the 3D transient field, while Ha-NeRF predicts the transient visibility map with the pre-trained MLP.

baselines.

**Qualitative results:** We obtain qualitative comparison results based on different scenes of our model and the baselines. The rendered static images are shown in Figure 4, showing that the rendering process with NeRF is challenged by the persistence of transient phenomena, leading to global color deviations and shadowing effects. While NeRF-W and HA-NeRF demonstrate the capability to model diverse photometric effects, facilitated by the incorporation of appearance embeddings, it is important to note that they can not avoid the rendering of ghosting artifacts in Tajcompari Mahal and Brandenburg Gate (seriously on NeRF-W) and blurry arti-

| Method | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| IE-NeRF(IM) | 22.27 | 0.802 | 0.191 |
| IE-NeRF(SM) | 24.65 | 0.879 | 0.187 |
| IE-NeRF(Ours) | **25.33** | **0.898** | **0.158** |

Table 2: Ablation comparison on metrics of PSNR, SSIM, and LPIPS results of different transient mask generation methods on Phototourism datasets in the scene Brandenburg Gate.

| Training strategy | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| IPE | 23.58 | 0.864 | 0.191 |
| WT-IPE | 19.86 | 0.723 | 0.289 |
| RegFre-IPE | **25.33** | **0.898** | **0.158** |

Table 3: Ablation comparisons of PSNR, SSIM, and LPIPS for IT-NeRF with different training strategies on Phototourism datasets in the scene of Brandenburg Gate.
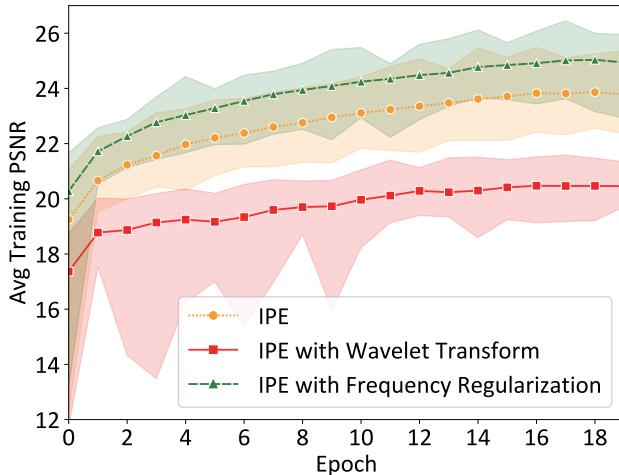


Figure 6: Comparisons of Training PSNR across epochs were conducted with different training strategies, including Integrated Positional Encoding (IPE), Integrated Positional Encoding with wavelet transform (WT-IPE), and Integrated Positional Encoding with frequency regularization (RegFre-IPE).

facts in Trevi Fountain and Sacre Coeur. On the contrary, IE-NeRF disentangles transient elements from the static scene consistently which proves the effectiveness of the transient mask-guided inpainting module. In addition, we present the transient components of NeRF-W, Ha-NeRF, and IE-NeRF in Figure 5. Results from IE-NeRF reveal more detailed texture in the pixel mask. Transient performance further supports the capability of our IE-NeRF method.

### 4.3 Ablations Studies

In this section, we conduct ablation studies exploring diverse transient mask generation approaches with the same pipeline of static image reconstruction on the Phototourism datasets. We compare three methods: (1) IE-NeRF(ours), our primary approach utilizing a transient mask generated by the NeRF MLP network. (2) IE-NeRF(IM), which generates the transient mask with an independent MLP network, following a way similar to the Ha-NeRF [5]. (3) IE-NeRF(SM), where we leverage a pre-trained semantic model, specifically the object instance segmentation model MaskDINO [16], in which we predefine the transient objects, including but not limited to people, cars, bicycles, flags, and slogans. We evaluate the ablation models on Phototourism datasets in the Brandenburg Gate scene of PSNR, Sthe SIM (SSIM), and LPIPS metrics and provide the results in Table 2. Our primary approach outperforms the other two methods.

We also conduct ablation studies to analyze the effectiveness of our training strategy. We evaluate the performance on the scene of Brandenburg Gate datasets in the following ways: First, we employed the conventional approach using integrated positional encoding. Second, we implemented our proposed strategy, which incorporates regularization frequency into integrated positional encoding (RegFre-IPE). Third, we experimented with the wavelet transformer during the embedding process, considering the non-stationarity of input signals caused by transient phenomena. Comparisons of training PSNR across epochs with the three training approaches are shown in Figure 6. The results indicate that RegFre-IPE achieves the highest PSNR in the early epochs and maintains a lead throughout the entire training period. However, integrated positional encoding with wavelet transform (WT-IPE) has negative impacts on performance. The metric results are summarized in Table 3. Our proposed training strategy with RegFre-IPE consistently demonstrates the best performance across all three considered metrics under the same training epochs.

## 5 Conclusion

This paper proposes a novel method to enhance the NeRF in the wild. Our approach improves traditional NeRF by integrating an inpainting module that helps restore the static scene image, guided by the transient mask generated from the MLPs of the NeRF network. Results from both qualitative and quantitative experiments on Phototourism datasets show the effectiveness of our method. Currently, our proposed approach still encounters challenges on small datasets or under sparse inputs, as the transient mask lacks sufficient information to guide the inpainting process. For further optimization, we plan to explore an approach to learning the consistent appearance of the static scene with varying transient phenomena under the few-shot setting.

## References

[1] Benjamin Attal, Selena Ling, Aaron Gokaslan, Christian Richardt, and James Tompkin. Matryodshka: Real-time 6dof video view synthesis using multi-sphere images. In *European Conference on Computer Vision (ECCV)*, pages 441–459. Springer, 2020.

[2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (CVPR)*, pages 5855–5864, 2021.

[3] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 497–504. 2023.

[4] Chenjie Cao, Qiaole Dong, and Yanwei Fu. Zits++: Image inpainting by improving the incremental transformer on structural priors. *IEEE Transactions on Pattern Analysis and Machine Intelligence (T-PAMI)*, 45(10):12667–12684, 2023.

[5] Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Ying Feng, Xuan Wang, and Jue Wang. Hallucinated neural radiance fields in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12943–12952, 2022.

[6] Peng Dai, Yinda Zhang, Zhuwen Li, Shuaicheng Liu, and Bing Zeng. Neural point cloud rendering via multiplane projection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7830–7839, 2020.

[7] Nianchen Deng, Zhenyi He, Jiannan Ye, Budmonde Duinkharjav, Praneeth Chakravarthula, Xubo Yang, and Qi Sun. Fov-nerf: Foveated neural radiance fields for virtual reality. *IEEE Transactions on Visualization and Computer Graphics*, 28(11):3854–3864, 2022.

[8] Anis Farshian, Markus Götz, Gabriele Cavallaro, Charlotte Debus, Matthias Nießner, Jón Atli Benediktsson, and Achim Streit. Deep-learning-based 3-d surface reconstruction—a survey. *Proceedings of the IEEE*, 2023.

[9] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world's imagery. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5515–5524, 2016.

[10] Tewodros Habtegebrial, Christiano Gava, Marcel Rogge, Didier Stricker, and Varun Jampani. Somsi: Spherical novel view synthesis with soft occlusion multi-sphere images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15725–15734, 2022.

[11] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.

[12] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)*, 37(6):1–15, 2018.

[13] Dominic Jack, Jhony K Pontes, Sridha Sridharan, Clinton Fookes, Sareh Shirazi, Frederic Maire, and Anders Eriksson. Learning free-form deformations for 3d object reconstruction. In *Computer Vision–ACCV 2018: 14th Asian Conference on Computer Vision, Perth, Australia, December 2–6, 2018, Revised Selected Papers, Part II 14*, pages 317–333. Springer, 2019.

[14] Sebastian Knorr and Thomas Sikora. An image-based rendering (ibr) approach for realistic stereo view synthesis of tv broadcast based on structure from motion. In *2007 IEEE International Conference on Image Processing*, volume 6, pages VI–572. IEEE, 2007.

[15] Jaewon Lee, Injae Kim, Hwan Heo, and Hyunwoo J Kim. Semantic-aware occlusion filtering neural radiance fields in the wild. *arXiv preprint arXiv:2303.03966*, 2023.

[16] Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3041–3050, 2023.

[17] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6498–6508, 2021.

[18] Hao-Kang Liu, I Shen, Bing-Yu Chen, et al. Nerf-in: Free-form nerf inpainting with rgb-d priors. *arXiv preprint arXiv:2206.04901*, 2022.

[19] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751*, 2019.

[20] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7210–7219, 2021.

[21] Moustafa Meshry, Dan B Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla. Neural rerendering in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6878–6887, 2019.

[22] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (CVPR)*, pages 4743–4752, 2019.

[23] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.

[24] Ashkan Mirzaei, Tristan Aumentado-Armstrong, Marcus A Brubaker, Jonathan Kelly, Alex Levinshtein, Konstantinos G Derpanis, and Igor Gilitschenski.

Reference-guided controllable inpainting of neural radiance fields. *arXiv preprint arXiv:2304.09677*, 2023.

[25] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 165–174, 2019.

[26] Thomas Porter and Tom Duff. Compositing digital images. In *Proceedings of the 11th annual conference on Computer graphics and interactive techniques*, pages 253–259, 1984.

[27] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIX 16*, pages 623–640. Springer, 2020.

[28] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2437–2446, 2019.

[29] Jiaming Sun, Xi Chen, Qianqian Wang, Zhengqi Li, Hadar Averbuch-Elor, Xiaowei Zhou, and Noah Snavely. Neural 3d reconstruction in the wild. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–9, 2022.

[30] Roman Suvorov, Elizaveta Logacheva, Anton Mashikhin, Anastasia Remizova, Arsenii Ashukha, Aleksei Silvestrov, Naejin Kong, Harshith Goka, Kiwoong Park, and Victor Lempitsky. Resolution-robust large mask inpainting with fourier convolutions. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2149–2159, 2022.

[31] Richard Szeliski and Richard Szeliski. Image-based rendering. *Computer Vision: Algorithms and Applications*, pages 681–722, 2022.

[32] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. State of the art on neural rendering. In *Computer Graphics Forum*, volume 39, pages 701–727. Wiley Online Library, 2020.

[33] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)*, 38(4):1–12, 2019.

[34] Bill Triggs, Philip F McLauchlan, Richard I Hartley, and Andrew W Fitzgibbon. Bundle adjustment—a modern synthesis. In *Vision Algorithms: Theory and Practice: International Workshop on Vision Algorithms Corfu, Greece, September 21–22, 1999 Proceedings*, pages 298–372. Springer, 2000.

[35] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 551–560, 2020.

[36] Juan Wang, Chunfeng Yuan, Bing Li, Ying Deng, Weiming Hu, and Stephen Maybank. Self-prior guided pixel adversarial networks for blind image inpainting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(10):12377–12393, 2023.

[37] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and Wenping Wang. F2-nerf: Fast neural radiance field training with free camera trajectories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4150–4159, 2023.

[38] Silvan Weder, Guillermo Garcia-Hernando, Aron Monszpart, Marc Pollefeys, Gabriel J Brostow, Michael Firman, and Sara Vicente. Removing objects from neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16528–16538, 2023.

[39] Hanyu Xiang, Qin Zou, Muhammad Ali Nawaz, Xianfeng Huang, Fan Zhang, and Hongkai Yu. Deep learning for image inpainting: A survey. *Pattern Recognit.*, 134:109046, 2023.

[40] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8254–8263, 2023.

[41] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. inerf: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1323–1330. IEEE, 2021.

[42] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018.