

# Towards Degradation-Robust Reconstruction in Generalizable NeRF

Chan Ho Park<sup>1</sup>, Ka Leong Cheng<sup>1</sup>, Zhicheng Wang<sup>2</sup>, Qifeng Chen<sup>1</sup>

<sup>1</sup>HKUST <sup>2</sup>UCSD

## Abstract

*Generalizable Neural Radiance Field (GNeRF) across scenes has been proven to be an effective way to avoid per-scene optimization by representing a scene with deep image features of source images. However, despite its potential for real-world applications, there has been limited research on the robustness of GNeRFs to different types of degradation present in the source images. The lack of such research is primarily attributed to the absence of a large-scale dataset fit for training a degradation-robust generalizable NeRF model. To address this gap and facilitate investigations into the degradation robustness of 3D reconstruction tasks, we construct the Objaverse Blur Dataset, comprising 50,000 images from over 1000 settings featuring multiple levels of blur degradation. In addition, we design a simple and model-agnostic module for enhancing the degradation robustness of GNeRFs. Specifically, by extracting 3D-aware features through a lightweight depth estimator and denoiser, the proposed module shows improvement on different popular methods in GNeRFs in terms of both quantitative and visual quality over varying degradation types and levels. Our dataset and code will be made publicly available.*

## 1. Introduction

Robust 3D scene reconstruction under various lighting conditions is a practical task for vision and graphics applications. However, in practical image-capturing environments, image quality often degrades with noise and blur (from long exposure time, low-light environment, unstable capture, etc.), which introduces a substantial challenge to ensure the resilience of 3D rendering models to different types of degradation. Inferring the 3D geometry of a scene from a collection of 2D degraded images remains a complex and challenging task due to low image quality and the small number of captured images.

Recently, the advent of neural radiance fields (NeRF) [46] has significantly transformed the landscape of 3D reconstruction, making it more applicable to real-world scenarios. The limitation of vanilla NeRF has prompted exten-

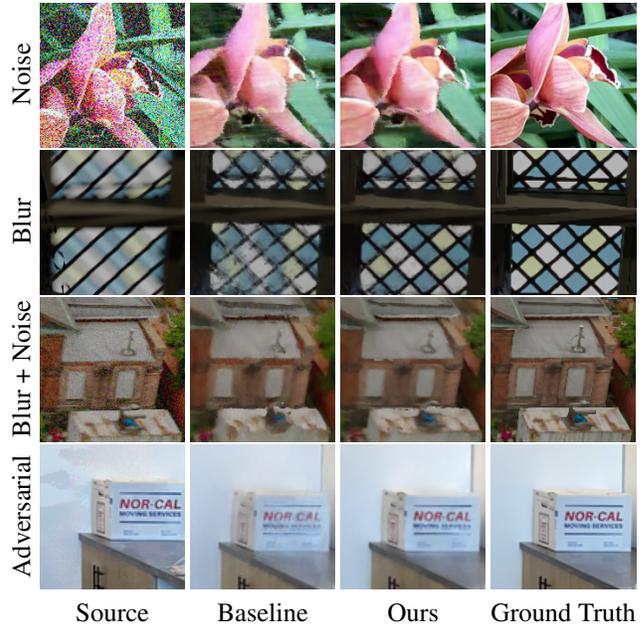


Figure 1. Visual comparison of the different degradation settings against baseline. Note that there is misalignment between source and ground truth because source columns contains the source views nearby to the target view.

sive research focused on the optimization of both rendering [3, 18, 23, 50, 61] and training speed [15, 48, 52]. In this process, researchers have explored various scene representations, including both explicit and implicit ones, and advanced the rendering methods in pursuit of these enhancements. From that regard, generalized NeRF (GNeRF) models [4, 28, 51, 53, 59, 67] are a practical alternative solution to circumvent the need for scene-specific optimization. By aggregating deep image features from source images, GNeRF models construct scene representations through a single forward pass.

Aiming to enhance the robustness of neural 3D reconstruction, several NeRF-based methods adopt source image data with haze [7], blur degradation [34, 42, 58], low light source images [24, 29, 47], and underwater captures [35, 66]. Such NeRF-based models typically modify the rendering

method based on the real physical image formation process by designing a particular renderer or module aiming to learn the latent variables for the underlying degradation. To leverage per-scene optimization methods for training in specific degradation types, prior knowledge regarding the existing degradation in the source image is necessary.

Meanwhile, there has been relatively limited exploration of GNeRF models under image degradation. Handling degradation for 3D reconstruction in a GNeRF setting [16, 49] poses a greater challenge compared to NeRF-based methods [13, 47, 56, 60]. While methods such as [49, 54] discuss 3D reconstruction under signal-dependent synthetic degradation with image-based rendering and multi-plane features, a more general strategy to improve robustness against degraded source images can be devised. From that regard, we propose a 3D-aware feature extraction module to enhance robustness against multiple types of degradation. The design is simple and model-agnostic, which makes it easily adaptable to various GNeRF models, broadening its potential applications in various scenarios.

Furthermore, to address the challenge of lacking large-scale 3D reconstruction datasets for training a robust GNeRF model against blur, we construct a novel 3D reconstruction dataset with motion blur, named *Objaverse Blur Dataset* based on Objaverse [12]. Specifically, we design an algorithm to generate the degradation using 3D models for 3D consistency, which is verified in [30] to be crucial for ensuring compatibility with real-world applications and domain adaptation.

To this end, our contribution can be summarized as follows:

- We construct the *Objaverse Blur Dataset*, the first large-scale 3D reconstruction dataset with blur degradation at multiple levels.
- We propose a lightweight *3D-aware feature extractor*, with depth estimation and differentiable warping as a plugin module for GNeRF.
- Experimental results show that various GNeRF models with our plugin module can achieve *consistent improvement* under *multiple types of degradations*.

## 2. Related Work

**Generalizable neural radiance field.** Although Neural Radiance Fields (NeRF) have proven to be a powerful implicit optimization method for modeling complex 3D scenes, the lack of generalizability in NeRF makes it impractical in real-world scenarios. Consequently, a parallel line of research in 3D reconstruction with neural rendering focuses on developing models that can generalize to new scenes without the need for fine-tuning. Generalizable NeRF (GNeRF) [4, 28, 40, 51, 53, 59, 67], as inference models, map the extracted deep image features from the source images to the parameterized radiance field of the scene for volume

rendering. The key variations among these methods are in their image feature aggregation, coordinate system design and rendering mechanism.

**Degradation Robustness in NeRF.** To address the presence of degradation in the source images, many degradation-robust NeRF approaches adopt a renderer that incorporate the physical image formation process. For example, approaches [8, 43, 70] tackle occlusion and transient objects in the source images by incorporating additional outputs capturing the radiance of the transient image. In a low-lighted source image setting, Cui *et al.* [11] incorporates the concept of the concealing field. Similarly, works such as [34, 42, 58] deal with blur degradation through the joint optimization of rays per source image and model the camera trajectory or the focal plane. With regard to GNeRF-based approaches, Pearl *et al.* [49] addresses source images with noise degradation. On a related note, Tanay *et al.* [54] proposes to mitigate burst noise by constructing multi-plane features (MPF) for scene representation. Another approach proposed by RaFE [62] is to introduce Generative Adversarial Networks to refine the scene representation to account for the inconsistencies in source images. Lastly, NeRFool [16] explores the adversarial robustness of GNeRF models by applying the strategy in [20] and also studies the vulnerability pattern.

**Image restoration.** Both single-image and multi-image restoration models have been proven successful under corruptions such as blur [6, 38, 68, 69], noise [1, 14, 44, 63], rain [26], mosaic [31], super-resolution [25], and haze [2]. Initially, CNN-based methods pioneered this domain. Subsequently, advanced architectures [6, 38, 68] and generative adversarial learning [32, 33] have demonstrated effectiveness in restoring realistic visual results. Some works [5, 9, 37, 41] devise a way to learn an adequate representation of the degraded image through a Transformer and self-supervised pre-training by low-level pretext task. For multi-frame image restoration models [1, 14, 44, 63], one of the main challenges is to align the different frames for better restoration. This is done by predicting a blending kernel or introducing modules like Transformer for alignment. While our module shares the same purpose with the methods above, which is to represent degraded images effectively, the main difference is that the GNeRF setting makes use of camera pose information.

## 3. Preliminaries

NeRF is an implicit field that maps a query xyz coordinate  $\mathbf{x} \in \mathbb{R}^3$  and a viewing direction  $\mathbf{d} \in \mathbb{R}^3$  into a radiance value  $\mathbf{c} \in \mathbb{R}^3$  and density  $\sigma \in \mathbb{R}$ , *i.e.*  $(\mathbf{c}(\mathbf{x}, \mathbf{d}), \sigma(\mathbf{x})) = f(\mathbf{x}, \mathbf{d})$ . Each pixel in an image is seen as a ray  $\mathbf{r}$  starting from the camera origin  $\mathbf{o}$  with a direction  $\mathbf{d} \in \mathbb{R}^3$  which is dependent on the camera pose  $\mathbf{R}|\mathbf{t}$  and intrinsic  $\mathbf{K}$ . By accumulating the radiance value from the camera origin up to a predefined

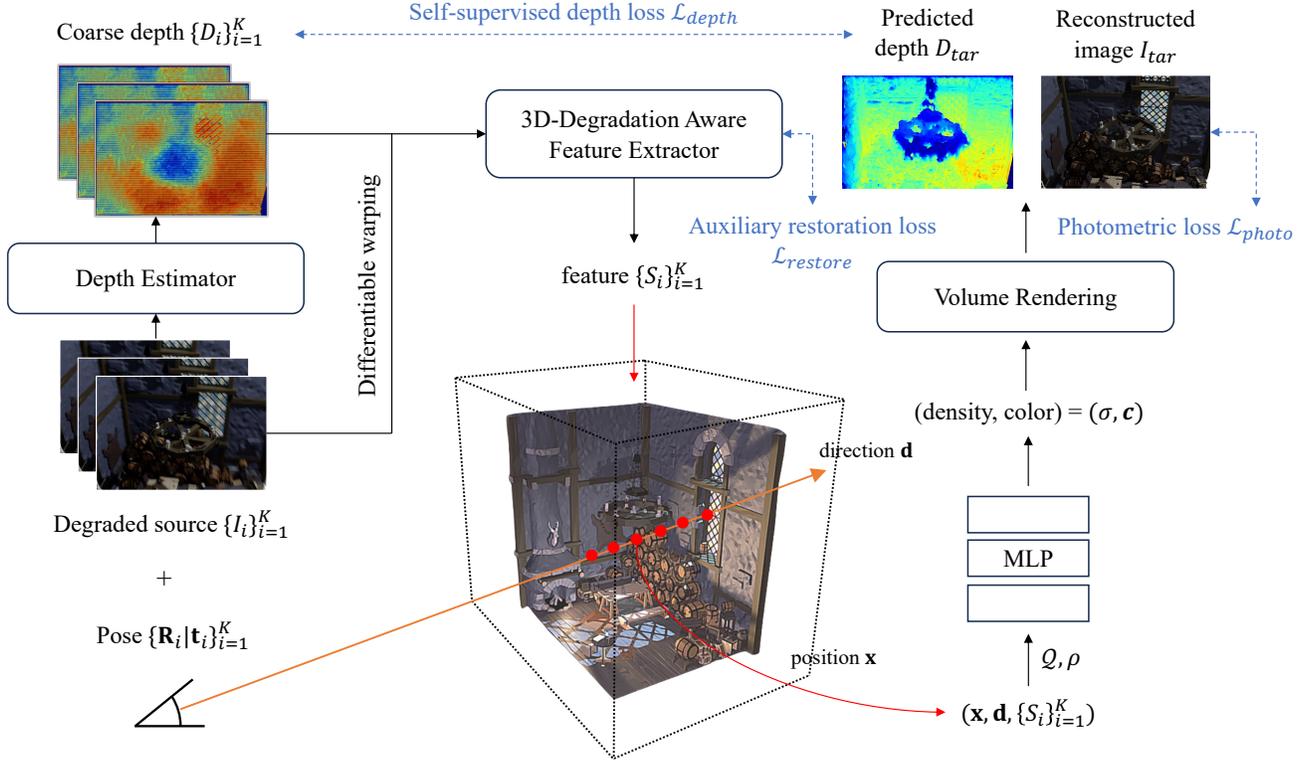


Figure 2. Training pipeline for degradation-robust GNeRF. This diagram illustrates the process for training a GNeRF that is robust to image degradation. The procedure involves a 3D-degradation-aware feature extractor  $\mathcal{F}$ , consisting of two steps. Initially, grouped degraded input images  $\{I_i\}_{i=1}^N$  predict the corresponding depths  $\{D_i\}_{i=1}^N$  through the depth estimator  $\mathcal{D}$ , followed by a differentiable warping. Subsequently, the 3D-degradation-aware feature extractor  $\mathcal{F}$  extracts feature  $\{S_i\}_{i=1}^N$  to parameterize the scene representation for the volume rendering stage.

far bound  $t_f$ , NeRF predicts a pixel value  $\hat{C}(\mathbf{r})$ :

$$\hat{C}(\mathbf{r}) = \int_0^{t_f} T(t) \sigma(\mathbf{o} + t\mathbf{d}) \mathbf{c}(\mathbf{o} + t\mathbf{d}, \mathbf{d}) dt, \quad (1)$$

where  $T(t) = \exp(-\int_0^t \sigma(\mathbf{o} + s\mathbf{d}) ds)$  is the accumulated density, also referred as transmittance. NeRF is typically optimized by an  $\mathcal{L}_{photo}$  MSE loss between the ground truth RGB pixel  $C(\mathbf{r})$  and the prediction  $\hat{C}(\mathbf{r})$ . Additionally, by replacing  $\mathbf{c}$  in Eq. (1) with the depth  $t$ , we can obtain the predicted depth.

Under this setting, GNeRF projects the query 3D coordinate  $\mathbf{x}$  to the input source images  $\{I_i\}_{i=1}^N$  through the corresponding camera pose  $\{\mathbf{R}_i | \mathbf{t}_i\}_{i=1}^N$  and intrinsics  $\{\mathbf{K}_i\}_{i=1}^N$  for aggregating source image information. Denoting the projection from 3D space to the  $i$ -th image  $I_i$  as  $\rho_i(\cdot) : \mathbb{R}^3 \rightarrow \mathbb{R}^2$  and a shared CNN feature extraction module as  $\mathcal{F}$ . For a query point  $\mathbf{x}$ , the extracted feature  $\mathbf{s}_i$  from the  $i$ -th image  $I_i$  is

$$\mathbf{s}_i = \text{GridSample}(\mathcal{F}(I_i), \rho_i(\mathbf{x})), \quad (2)$$

yielding a set of image features  $\mathcal{S}_{\mathbf{x}} = \{\mathbf{s}_i\}_{i=1}^N$ , where  $F_i$  is the extracted feature for image  $I_i$  through  $\mathcal{F}$ . Subsequently,

GNeRF models aggregate the collected feature through an aggregation function  $Q$  and decode  $\mathcal{S}_{\mathbf{x}}$  through an implicit renderer MLP to obtain the radiance  $\mathbf{c}$  and density value  $\sigma$ , *i.e.*  $(\mathbf{c}, \sigma(\mathbf{x})) = \text{MLP}(Q(\mathcal{S}_{\mathbf{x}}), \mathbf{d})$ , where  $\mathbf{d}$  is the viewing direction of the query point  $\mathbf{x}$ . Then, GNeRF models can be optimized using the same objective function as vanilla NeRF models. By training over a large number of scenes, the renderer optimizes to map the query 3D point to radiance and density value conditioned on the aggregated scene-variant feature  $Q(\mathcal{S}_{\mathbf{x}})$ , allowing generalization to new scenes.

## 4. Method

Given that there is a broad spectrum of potential degradation in the real-world application of GNeRF, source images can be degraded in both the training and inference stages. Since GNeRF models are implemented by conditioning on the projected source image features  $\mathcal{S}_{\mathbf{x}}$ , as shown in [16, 49], GNeRF models can be vulnerable without noise-aware components. To this end, we introduce both the 3D and degradation awareness to the feature extractor  $\mathcal{F}$ .

Specifically, the proposed module comprises two main

components: (i) a self-supervised depth estimator  $\mathcal{D}$ ; (ii) a 3D-degradation-aware feature extractor  $\mathcal{F}$  with restoration head  $\mathcal{R}$ . The main insight is from hypothesizing that the performance drop of GNeRF with degraded source images originates from the inaccurate geometry estimation caused by the high variance amongst the learned feature in  $\mathcal{S}_x$  when not specifically handling the degradation. The objective of the proposed module is two-fold. Firstly, it aims to extract visual features that are independent of degradation, ensuring accurate depth estimation. Secondly, it aims to extract features inherent in natural-looking images, facilitating effective image reconstruction. As explored in NeRFool [16], the perturbation in depth is the leading cause of a significant drop in reconstruction accuracy compared to the inaccurate texture or radiance. Therefore, the initial stage of our 3D-aware feature extractor  $\mathcal{F}$  involves estimating the depth of each source image by leveraging information from nearby images. Subsequently, these nearby images are aligned and stacked together to form a composite representation, which encapsulates degradation-robust and inherent features  $\mathcal{S}_x$ .

With the two-step module design, our proposed feature extractor aims not only to replace traditional convolution-based feature extractors but also to maintain comparable inference speeds to the original GNeRF model, as detailed in Sec. 6. Additionally, it is worth noting that individual components of the 3D-degradation-aware feature extractor  $\mathcal{F}$  are modular, allowing for substitution with alternative designs to accommodate variations in computational resource constraints.

#### 4.1. Depth Estimation and Differentiable Warping

Methods like [49] have demonstrated that denoising techniques that incorporate camera poses significantly outperform multi-frame denoising methods without the use of pose information. This is mainly through the improved alignment of the input images with larger displacement. To further leverage the camera pose information to enhance the feature extraction and denoising process, we propose a depth estimator  $\mathcal{D}$  to estimate a rough depth as the first step. Specifically, we group the  $K$  closest source images for each source image as in Eq. (3). Denoting the set of  $K$  nearby source images for the  $i$ -th input image  $I_i$  as  $\mathcal{I}_i^K$ :

$$\mathcal{I}_i^K = \{I_k : I_k \in \text{NearestK}(I_i)\}. \quad (3)$$

We can then compute the relative camera pose  $\mathbf{R}_i^k | \mathbf{t}_i^k$  for image  $I_i$  against its nearby image  $I_k \in \mathcal{I}_i^K$ :

$$[\mathbf{R}_i^k | \mathbf{t}_i^k] = [\mathbf{R}_k | \mathbf{t}_k] \cdot [\mathbf{R}_i | \mathbf{t}_i]^{-1}, \quad (4)$$

where  $\mathbf{R}_i | \mathbf{t}_i$  and  $\mathbf{R}_k | \mathbf{t}_k$  denote the camera pose of  $I_i$  and  $I_k$ , respectively.

Subsequently, inspired by multi-view stereo (MVS) [21, 57, 65], we can create the group-correlation-based matching cost volumes using the relative camera poses, which

in turn estimate the dense geometric structure of the scene. Following the efficient approach of PatchMatchNet [57], our framework incorporates a multi-scale PatchMatch algorithm for depth estimation, where we obtain the estimated depths  $\{D_i\}_{i=1}^N$  for the  $N$  source images. Together with the relative camera poses, our system aligns the nearby  $K$  views  $\mathcal{I}_i^K$  onto the source image  $I_i$ . Specifically, we modify the homographic warping as differentiable warping to account for the predicted depth. The mapping  $\pi_i^k(\cdot) : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  from a homogeneous pixel  $\mathbf{p}$  in source view  $I_i$  to the projected pixel in the nearby view  $I_k$  is defined as follows:

$$\pi_i^k(\mathbf{p}) = \mathbf{K}_k \cdot (\mathbf{R}_i^k \cdot (\mathbf{K}_i^{-1} \cdot \mathbf{p} \cdot D_i[\mathbf{p}]) + \mathbf{t}_i^k), \quad (5)$$

where  $D_i[\mathbf{p}]$  denotes the depth value of pixel  $\mathbf{p}$ .

#### 4.2. 3D-Degradation-Aware Feature Extractor

As discussed in literature such as DIP [55] and follow-up works [17, 22, 36], deep image prior obtained through an auxiliary restoration module can have better capture the local-level statistics of a single natural image. Motivated by this, we introduce an auxiliary 3D-aware restoration head  $\mathcal{R}$  to support the learning of our 3D-degradation-aware feature extractor  $\mathcal{F}$  for better feature matching and extraction concerning multi-view consistency.

Given the source view of image  $I_i$  and the nearby image  $I_k$  with estimated depth map  $D_k$ , we can warp both  $I_k$  and  $D_k$  to the source view denoted as  $\tilde{I}_k$  and  $\tilde{D}_k$ , respectively, based on Eq. (5). We concatenate all the nearby warped images  $\{\tilde{I}_k\}_{k=1}^K$  and warped depths  $\{\tilde{D}_k\}_{k=1}^K$  with the source image  $I_i$  and its estimated depth  $D_i$  to form the aligned feature  $J_i$ :

$$J_i = \text{Concat}(I_i, D_i, \tilde{I}_1, \dots, \tilde{I}_K, \tilde{D}_1, \dots, \tilde{D}_K). \quad (6)$$

Inspired by multi-frame degraded-image restoration methods [1, 14, 44, 63] where source images have relatively small displacement, we leverage on the depth-aligned degraded source images ( $J_i$ ) through image restoration module. The intent of aligning and stacking degraded source images is to provide the image restoration module with more information from other views; the main rationale for including the predicted depth  $D_i$  and the warped depths  $\{\tilde{D}_i\}_{k=1}^K$  as part of the input for is that the levels of degradation within the source images such as motion blur are usually dependent on the depths [30].

The aligned feature  $J_i$  is then sent to the feature extractor  $\mathcal{F}$  as input, which outputs a 3D-degradation-aware feature  $S_i \in \mathbb{R}^{H/4 \times W/4}$ . The auxiliary restoration head  $\mathcal{R}$  operates in a residual manner to restore a degradation-free reconstruction  $\hat{I}_i$ , with an auxiliary restoration loss against the corresponding clean source image  $I_i^{\text{GT}}$ . The features  $S_i$  are used in the subsequent stage of constructing the scene representation through the aggregation function  $Q$  as in Eq. (2).

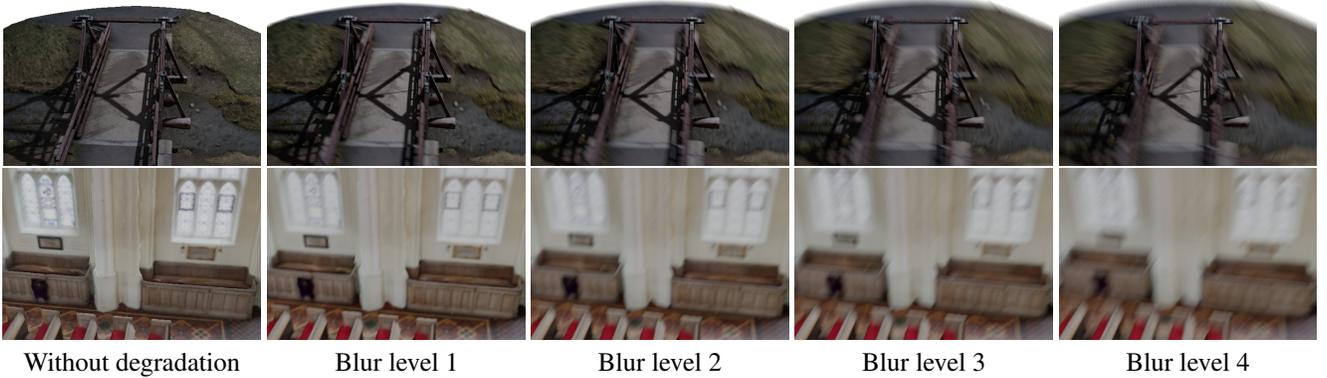


Figure 3. Examples of Objaverse Blur Dataset at different blur levels.

As the image feature shape is identical to that of the original CNN-based image feature, the proceeding steps of the volume rendering process and training strategy remain unchanged. Through such a model-agnostic feature extractor design, the proposed module can be easily plugged into any existing GNeRF under different degradation settings.

### 4.3. Optimization and Loss

During the training stage, the depth estimator  $\mathcal{D}$  predicts the depth of the target view  $[\mathbf{R}_{\text{tar}}|\mathbf{t}_{\text{tar}}]$ . This prediction is based on the degraded target image  $I_{\text{tar}}$  and its nearby  $K$  source images  $\mathcal{I}_{\text{tar}}^K$ . As we assume the absence of ground-truth depth, the depth estimator is supervised by the target view fine depth prediction  $\hat{D}_{\text{tar}}$  of the GNeRF as the pseudo ground truth. With this pseudo ground truth, we detach the depth estimator from the computational graph of the main 3D reconstruction. This ensures self-supervised training of the depth estimator without requiring ground-truth depth data. We employ smooth L1 function [19] to the predicted depth prediction  $D_{\text{tar}}$ :

$$\mathcal{L}_{\text{depth}} = \text{SmoothL1}(\hat{D}_{\text{tar}}, D_{\text{tar}}). \quad (7)$$

Similar to [28], During the training stage, the output reconstructions of the 3D-aware restoration head  $\mathcal{R}$  are supervised by the clean source images as ground truth signals are available as a supervision of the 3D reconstruction task. Formally, this loss is termed as the auxiliary restoration loss:

$$\mathcal{L}_{\text{restore}} = \sum_{i=1}^K \left\| \hat{I}_i - I_i^{\text{GT}} \right\|_1, \quad (8)$$

where  $I_i^{\text{GT}}$  is the  $i$ -th clean image.

Combined with the basic 3D reconstruction photometric loss  $\mathcal{L}_{\text{photo}}$ , the total optimization loss can be formulated as:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{photo}} + \lambda_1 \mathcal{L}_{\text{depth}} + \lambda_2 \mathcal{L}_{\text{restore}} \quad (9)$$

## 5. Objaverse Blur Dataset

In the evaluation of robustness against degradation in GNeRF models, Pearl *et al.* [49] introduces a variant of the real-world forward-facing dataset called LLFF-N. This dataset comprises 43 scenes and is specifically designed to test the model performance in the presence of signal-dependent synthetic degradation, which is applied to the linear image space through inverse gamma correction.

In practical settings, image blur is another significant and common degradation that leads to the loss of sharp details in captured images, which is widely observed in various data capture scenarios. While Ma *et al.* [42] introduced five sets of blurry images constructed using Blender [10], the size of the dataset is insufficient to train a NeRF model with generalizable performance.

Motivated by these considerations, we introduce the first 3D reconstruction dataset with motion blur called the *Objaverse Blur Dataset*. Leveraging the publicly available *Objaverse* [12], a large-scale 3D object dataset with diverse models, we simulate the real-world effects of camera motion during the image capture process. To render a single image with a specific blur level  $l$  for a given 3D model, the following steps are performed:

---

#### Algorithm 1 Render a blurry image from 3D model

---

**Input:** a 3D scene  $u$  from Objaverse

**Output:** an image with blur corruption

- 1: Sample a camera position  $\mathbf{p} = (r, \phi, \theta) \in \mathbb{R}^3$  distant enough from the scene or object.
  - 2: Sample trajectory direction  $\delta = (\Delta r, \Delta \phi, \Delta \theta) \in \mathbb{R}^3$ .
  - 3: Sample a trajectory weight  $w_l$  at blur level  $l$  for controlling the blur strength.
  - 4: Uniformly sample  $m$  positions  $\mathbf{p}_i \sim \mathcal{U}(\mathbf{p}, \mathbf{p} + w_l \delta)$ .
  - 5: Render latent images  $\mathbf{x}_i$  at each position  $\mathbf{p}_i$ .
  - 6: Obtain the synthesized image  $\mathbf{x} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$ .
- 

We can repeat Algorithm 1 to render multi-frame blurry

images of a given scene with 3D consistency. Building upon the steps above, we create the Objaverse Blur Dataset, which comprises an extensive collection of over 50,000 images derived from more than 1000 unique settings obtained from a diverse set of 250 3D models.

To discuss the details of the rendering, as the 3D models contained in Objaverse dataset [12] are diverse, the main challenge is to find the rendering configuration that aligns with the reasonable rendering in most 3D models. After tuning, the final rendering configuration is set to be as follows:

1. **Camera position  $\mathbf{p}$ .** We uniformly sample the camera’s azimuthal and polar angles in  $\phi \in (\phi_0 - \phi_1, \phi_0 + \phi_1)$  and  $\theta \in (\theta_0 - \theta_1, \theta_0 + \theta_1)$ , respectively. The parameters are set to  $\phi_0 = 60^\circ$ ,  $\phi_1 = 7.5^\circ$ , and  $\theta_0$  to each of  $0^\circ$ ,  $90^\circ$ ,  $180^\circ$ , and  $270^\circ$  with  $\theta_1 = 7.5^\circ$ .
2. **Camera trajectory  $\delta$ .** Spherical coordinate  $\Delta r, \Delta \phi, \Delta \theta \in \mathbb{R}^3$ , is sampled uniformly within the range of  $[\frac{1}{2}\delta_0, \delta_0]$ . The parameter is set to  $\delta_0 = 2.5$ .
3. **Scene dependent blur weight  $w_u$ .** To account for the geometric attributes of a give 3D scene  $u$  at the camera position  $\mathbf{p}$ , such as near  $n_u$ , far  $f_u$  and 3D bounding box  $d_u$ , the blur weight  $w_u$  of a given 3D model is dependent on three factors:

(a) Flatness  $F_u = \frac{f_u - n_u}{\max(d_u)}$ ,

(b) Depth Range  $R_u = \frac{f_u}{n_u}$ ,

(c) Orientation  $O_u = \frac{\max(d_u)}{\min(d_u)}$ .

Composing all the factors, the scene dependent weight  $w_u$  is then calculated as  $w_u = (F_u \cdot R_u \cdot O_u)^{1/3}$ .

4. **Blur level  $l$ .** The final blur weight  $w_l$  at blur level  $l$  is sampled from the range of  $[0.9w_u l, 1.1w_u l]$ .

Following the dataset convention of Deblur-NeRF [42], we render  $n = 34$  images at each blur level for every viewpoint sampled in the 3D model, where each image is obtained by averaging  $m = 34$  latent images along the simulated camera trajectory.

## 6. Evaluation

### 6.1. Implementation Details

Training of a degradation-robust GNeRF generally follows the original respective GNeRF training procedures. While employing our proposed module, two losses are added in along with the original photometric loss  $\mathcal{L}_{\text{photo}}$ : the auxiliary restoration loss  $\mathcal{L}_{\text{restore}}$  and the self-supervised depth estimation loss  $\mathcal{L}_{\text{depth}}$ . These additional losses to train the 3D-degradation-aware feature extractor  $\mathcal{F}$  are balanced by the corresponding weight parameters  $\lambda_{\text{depth}}$  and  $\lambda_{\text{rec}}$ , which are empirically set to 1.0 and 0.01, respectively. For our experiments, we have employed multi-Dconv head transposed attention (MDTA), proposed by Restormer [68], as building blocks of the 3D-aware restoration head  $\mathcal{R}$ .

In order to train GNeRFs with our proposed module, we have gradually annealed the auxiliary restoration

loss  $\mathcal{L}_{\text{restore}}$  weight parameter every iteration  $\lambda_{\text{restore}}^{\text{nstep}} = \text{Max}(0.01, \alpha^{\text{nstep}}) \times \lambda_{\text{restore}}$ . Here, we set  $\alpha = 0.99997$  and clip the minimum to 0.01 to keep the clean source image supervision to a moderate level. For training GNT [53] for testing adversarial noise [16], we have trained and inferenced following the same setting as the original work where the number of source is 4.

For our newly constructed *Objaverse Blur Dataset*, we randomly chose 52 settings and selected test views using a similar sampling method as LLFF [45, 59] to sub-sample every 16 viewpoints. As a result, our blur dataset consists of 156 test images. Furthermore, for the blur and noise degradation experimental settings, we follow the construction proposed by Pearl *et al.* [49] to test the 156 test images at gain levels 4,8,16,20. More details can be found in the supplementary materials.

### 6.2. Results

To demonstrate the invariance of our proposed module over the choice of GNeRF, we conduct a series of experiments with three GNeRF models: NAN [49], GeoNeRF [28], and GNT [53]. Particularly, we have selected recently proposed effective GNeRF methods of varying mechanisms. GNT [53] is an attention-based method without an explicit rendering formula; NAN [49] is an image-based rendering method like IBRNet [59]; and GeoNeRF [28] is a method based on an explicit cost volume scene representation like MVSNeRF [4].

Degradation conditions can be classified into two categories: single degradation type, including blur degradation in Tab. 8, noise degradation in Tab. 9, and adversarial degradation in Tab. 10; multiple degradations of blur and noise in Tab. 4. Additionally, to demonstrate the practicality of our proposed module, we include an evaluation of the real-noise image dataset in Sec. 6.5, stability of the depth estimation over varying degradation levels in Tab. 5, and lastly impact of our proposed module on inference speed in Tab. 6. For conciseness, we have only included the PSNR metric and complete details of the results will be included in the supplementary materials.

### 6.3. Source Images with Single Degradation

To assess the robustness against noise degradation, we adhered to the same setting established by NAN with 5 source images, specifically evaluating on LLFF-N dataset. Moreover, for the blur degradation setting, we have tested NAN with 7 source images and GeoNeRF with 5 source images on the Objaverse Blur Dataset while ensuring that all other components remained consistent with NAN’s methodology. Additionally, for blur degradation, we compare our performance against two baselines. The first one is the clean GNeRF inferencing the pre-processed images by the state-of-the-art image restoration model Restormer [68], and the second is the GNeRF model pretrained with degraded images. For

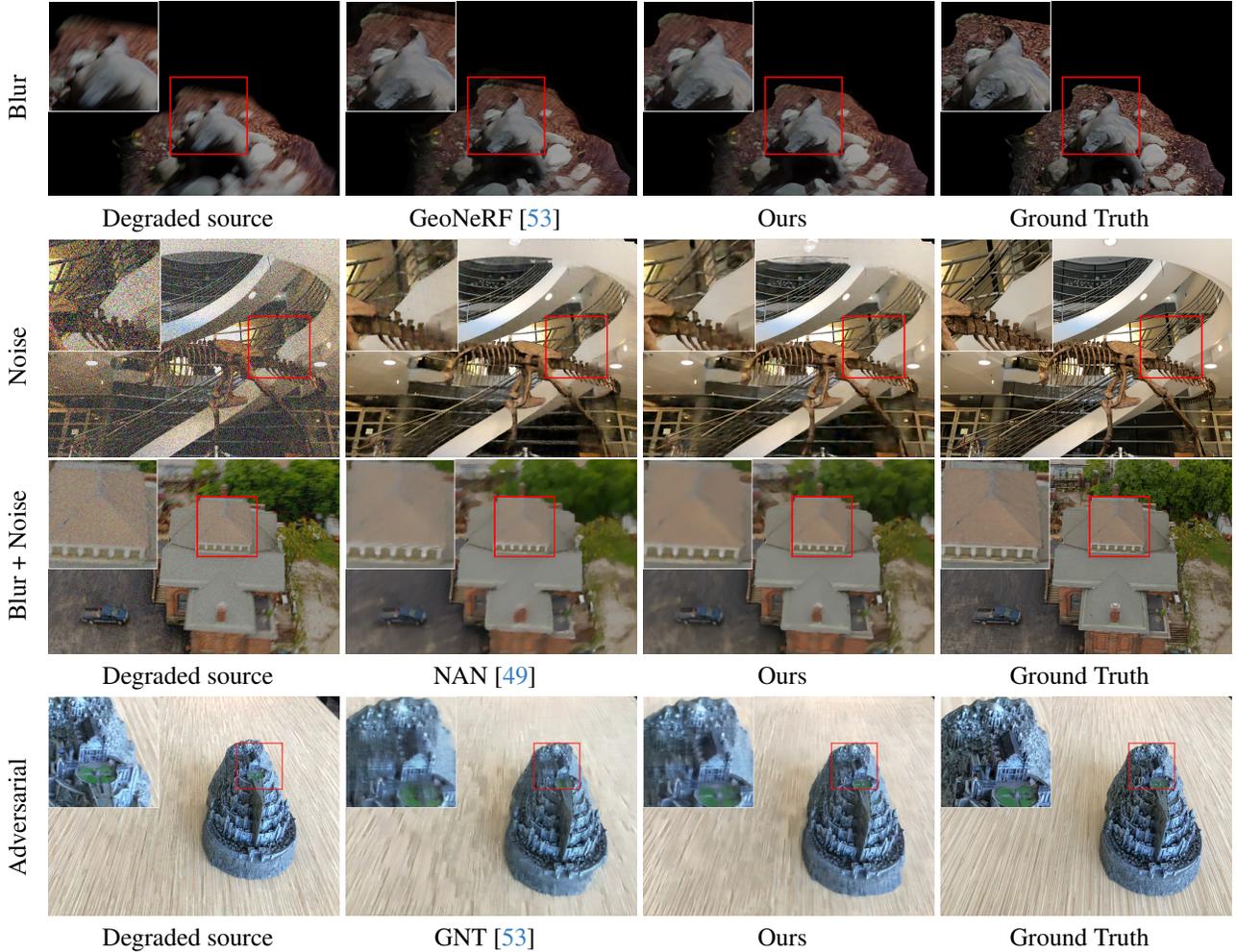


Figure 4. Qualitative comparison of GNeRFs [28, 49, 53] on different degradations and datasets [45, 49]. Each row contains the nearest degraded source image, rendered RGB with and without our proposed module, and ground-truth target view.

fairness, we train Restormer with the given dataset.

Similarly, for noise degradation, we compare the performance against the original NAN model, as other baselines have been tested in the original paper. Several results are shown in Fig. 4 and Tabs. 8 and 9 summarize the quantitative results, demonstrating consistent improvements of our model at varying blur and gain levels. Although our depth prediction seems to be smoothed compared to the original prediction, we can spot the wrong original GNeRF depth predictions, especially in terms of predicting the far surfaces.

To evaluate the robustness of GNeRF models with adversarial degradation in the source images, NeRFool [16] compared the performance of GNT [53] with and without adversarial training. Following the same scheme, we have trained 4 different pretrain GNT models with the two factors of whether to use adversarial training and whether to incorporate our proposed module. Based on these four models, Tab. 10 shows that the proposed module is effective

across all training and test combinations in enhancing the robustness of the GNeRF model for both with or without adversarial perturbation. Notably, the clean image training and clean image evaluation setups (rows 1-3) also benefit from the module. This implies that the module’s integration does not merely prevent performance drop in the presence of perturbation but actively contributes to the model’s ability to render scenes more accurately, even without degradation in source images. The clean image training and adversarial image evaluation settings (rows 4-6) demonstrate significant improvements. This shows that the proposed module enhances the inherent adversarial robustness of GNeRF models. Such an enhancement is crucial as it suggests that the module provides robustness against adversarial or random perturbation.

Table 1. Novel view reconstruction results of GeoNeRF [28] and NAN [49] on the Objaverse Blur Dataset across different blur levels. Methods with \* indicate the inference results of the pre-trained models trained on clean source images using pre-processed source images by Restormer [68].

Blur	Method	PSNR $\uparrow$	Method	PSNR $\uparrow$
1	GeoNeRF*	26.28	NAN*	23.73
	GeoNeRF	27.96	NAN	25.59
	+ Proposed	<b>28.47 (+0.51)</b>	+ Proposed	<b>26.51 (+0.92)</b>
2	GeoNeRF*	24.84	NAN*	22.11
	GeoNeRF	26.17	NAN	24.19
	+ Proposed	<b>26.73 (+0.56)</b>	+ Proposed	<b>24.96 (+0.77)</b>
3	GeoNeRF*	23.78	NAN*	21.14
	GeoNeRF	25.12	NAN	23.28
	+ Proposed	<b>25.69 (+0.57)</b>	+ Proposed	<b>23.93 (+0.65)</b>
4	GeoNeRF*	22.91	NAN*	20.39
	GeoNeRF	24.42	NAN	22.63
	+ Proposed	<b>24.98 (+0.56)</b>	+ Proposed	<b>23.16 (+0.53)</b>

Table 2. Novel view reconstruction results of NAN [49] on noise degradation (*LLFF-N* dataset) across different gain levels.

Gain Level	4	8	16	20
NAN	23.85	23.33	22.36	21.94
<b>+ Proposed</b>	<b>23.90</b>	<b>23.57</b>	<b>22.89</b>	<b>22.57</b>
<b>Difference</b>	<b>+0.05</b>	<b>+0.24</b>	<b>+0.53</b>	<b>+0.63</b>

Table 3. Novel view reconstruction results of GNT [53] on different training and test methods across the *LLFF* Dataset under adversarial degradation.

Training	Test	Proposed	Average PSNR $\uparrow$
Clean	Clean	×	23.50
		✓	24.14
		<b>Difference</b>	<b>+0.64</b>
Clean	Adversarial	×	17.68
		✓	18.90
		<b>Difference</b>	<b>+1.22</b>
Adversarial	Clean	×	24.26
		✓	24.40
		<b>Difference</b>	<b>+0.14</b>
Adversarial	Adversarial	×	20.05
		✓	20.22
		<b>Difference</b>	<b>+0.17</b>

## 6.4. Source Images with Multiple Degradations

Table 4 illustrates an evaluation of the NAN’s robustness when given 5 source images with multiple degradations, specifically blur and noise. The improvement in the performance metric with the inclusion of the proposed module is consistent across all levels of blur and noise. Although

Table 4. Novel view reconstruction results of NAN [49] on the *Objaverse Blur Dataset* across different blur levels and burst gain levels.

Blur	Proposed	Gain 1	Gain 2	Gain 4
1	×	23.59	23.47	23.07
	✓	<b>24.28 (+0.69)</b>	<b>24.21 (+0.74)</b>	<b>23.97 (+0.90)</b>
2	×	22.26	22.20	21.91
	✓	<b>22.70 (+0.44)</b>	<b>22.65 (+0.45)</b>	<b>22.39 (+0.48)</b>
3	×	21.48	21.43	21.20
	✓	<b>21.80 (+0.32)</b>	<b>21.73 (+0.30)</b>	<b>21.41 (+0.21)</b>
4	×	20.93	20.88	20.66
	✓	<b>21.17 (+0.24)</b>	<b>21.09 (+0.21)</b>	<b>20.74 (+0.08)</b>

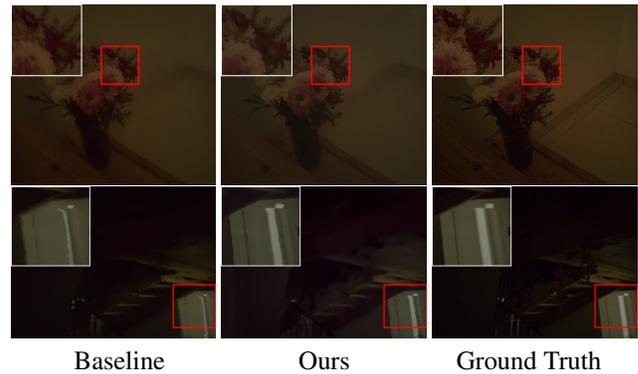


Figure 5. Qualitative evaluation on real low-light noise captured with Google Pixel 4 from NAN [49].

the improvement weakens when the degradation strength is too large, by enhancing the model’s capability across a spectrum of degradation levels, the proposed module potentially broadens the applicability and utility of GNeRF models, making them more viable for real-world, practical applications where robustness is desirable.

## 6.5. Real Degraded Images

To assess the effectiveness of our proposed module in real applications, we qualitatively evaluate our proposed modules on captured images with real noise in the low-light environment with Google Pixel 4 provided by [49]. We use models trained on Objaverse Blur Dataset with noise as shown in Tab. 8. For camera pose information, we referred to instructions described by [49] and used COLMAP on noisy images without any color adjustments. Fig. 5 compares the rendering results with and without our proposed module.

## 6.6. Analysis

As demonstrated in NeRFool [16], stable depth prediction is crucial for robustness against degradation. In other words, given degraded source images, the model’s depth prediction ideally should not change. Therefore, to evaluate the stability

Table 5. Analysis of depth prediction over varying blur degradation compared to the clean source image prediction.

Blur	Proposed	$\Delta$ Abs. Depth $\downarrow$	Diff.	$\Delta$ Rel. Depth $\downarrow$	Diff.
0	×	0.0	-	0.0	-
	✓	0.0		0.0	
1	×	0.394	-28.17%	0.068	-26.47%
	✓	0.283		0.050	
2	×	0.474	-32.91%	0.080	-32.50%
	✓	0.318		0.054	
3	×	0.528	-33.14%	0.085	-34.12%
	✓	0.353		0.056	
4	×	0.562	-32.56%	0.089	-34.83%
	✓	0.379		0.058	

Table 6. Image rendering time (in seconds) of different GNeRFs with and without 3D-degradation-aware feature extractor  $\mathcal{F}$ .

GNeRF	# Source	Original	Proposed	Difference
GeoNeRF	3	62.04	62.27	+0.23
	5	95.79	94.81	-0.98
GNT	3	75.54	76.11	+0.57
	5	90.08	90.38	+0.31
NAN	3	54.34	44.98	-9.36
	5	76.95	68.80	-8.15

of depth prediction over varying blur degradation levels, we have measured the change in absolute and relative depth prediction per pixel. Here, relative depth is the absolute depth divided by the scene range. The results are summarized in Tab. 5. Unsurprisingly, increments in blur degradation result in larger changes in depth per pixel, and the result also indicates that incorporating our module results in less deviation from the clean source image-based prediction.

Additionally, since our approach replaces the conventional CNN-based feature extractor with a novel two-step process, it is essential to ensure that this modification does not adversely affect rendering speed for practicality. Therefore, to assess the impact on inference speed, we conduct tests focusing on rendering images with a resolution of  $756 \times 1008$  pixels using an RTX 3090 graphics card. The results, presented in Tab. 6, detail the performance of our module in comparison to the inference speeds of GNeRF models evaluated in our experiments. Notably, the inference speeds observed with our module are commensurate with those recorded for the tested GNeRF models.

## 7. Conclusion

We present a model-agnostic module that can be easily incorporated into GNeRF training pipelines in order to enhance the degradation robustness. The proposed 3D-degradation-aware feature extractor has demonstrated effectiveness in

enhancing the quality of the rendered images and stability of depth prediction under conditions where source images contain various or multiple types of degradations. Our plugin module consists of a two-step process: depth prediction and latent image reconstruction. Each component of the module is supervised and self-supervised from the clean source images and GNeRF depth prediction, respectively. Our proposed plugin module is modular in a way such that other designs of the two components can be easily explored. Similarly, the capacity of the modules can be adjusted based on the computational overhead of the deploy environment. Additionally, in order to evaluate the robustness against various degradations, we have constructed a novel 3D reconstruction dataset named Objaverse Blur Dataset, which simulates realistic blur image captures at different levels. Lastly, possible follow-up research directions include utilizing a physical image formulation model to enhance the robustness of GNeRF models and making these modules more interpretable through other scene representations.

## References

- [1] G. Bhat, M. Danelljan, F. Yu, Luc Van Gool, and R. Timofte. Deep reparametrization of multi-frame super-resolution and denoising. In *Proceedings of ICCV*, pages 2440–2450. IEEE, 2021. [2](#), [4](#), [13](#)
- [2] Bolun Cai, Xiangmin Xu, Kui Jia, Chunmei Qing, and Dacheng Tao. Dehazenet: An end-to-end system for single image haze removal. *IEEE Trans. Image Process.*, 25(11): 5187–5198, 2016. [2](#)
- [3] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of CVPR*, pages 130–141. IEEE, 2023. [1](#)
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of ICCV*, pages 14104–14113. IEEE, 2021. [1](#), [2](#), [6](#)
- [5] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *Proceedings of CVPR*, pages 12299–12310. Computer Vision Foundation / IEEE, 2021. [2](#)
- [6] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *Proceedings of ECCV*, pages 17–33. Springer, 2022. [2](#)
- [7] Wei-Ting Chen, Yifan Wang, Sy-Yen Kuo, and Gordon Wetstein. Dehazenerf: Multiple image haze removal and 3d shape reconstruction using neural radiance fields. *CoRR*, abs/2303.11364, 2023. [1](#)
- [8] 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 CVPR*, pages 12933–12942. IEEE, 2022. [2](#)
- [9] Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution

- transformer. In *Proceedings of CVPR*, pages 22367–22377. IEEE, 2023. 2
- [10] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. 5
- [11] Ziteng Cui, Lin Gu, Xiao Sun, Yu Qiao, and Tatsuya Harada. Aleth-nerf: Low-light condition view synthesis with concealing fields. *CoRR*, abs/2303.05807, 2023. 2
- [12] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *Proceedings of CVPR*, pages 13142–13153. IEEE, 2023. 2, 5, 6
- [13] Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B. Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *Proceedings of ICCV*, pages 14304–14314. IEEE, 2021. 2
- [14] Akshay Dudhane, Syed Waqas Zamir, Salman Khan, Fahad Shabbaz Khan, and Ming-Hsuan Yang. Burstormer: Burst image restoration and enhancement transformer. In *Proceedings of CVPR*, pages 5703–5712. IEEE, 2023. 2, 4
- [15] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of CVPR*, pages 5491–5500. IEEE, 2022. 1
- [16] Yonggan Fu, Ye Yuan, Souvik Kundu, Shang Wu, Shunyang Zhang, and Yingyan Celine Lin. Nerfool: Uncovering the vulnerability of generalizable neural radiance fields against adversarial perturbations. In *Proceedings of ICML*, 2023. 2, 3, 4, 6, 7, 8, 13
- [17] Yossi Gandelsman, Assaf Shocher, and Michal Irani. "double-dip": Unsupervised image decomposition via coupled deep-image-priors. In *Proceedings of CVPR*, pages 11026–11035. Computer Vision Foundation / IEEE, 2019. 4
- [18] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien P. C. Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of ICCV*, pages 14326–14335. IEEE, 2021. 1
- [19] Ross B. Girshick. Fast R-CNN. In *Proceedings of ICCV*, pages 1440–1448. IEEE Computer Society, 2015. 5
- [20] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *Proceedings of ICLR*, 2015. 2
- [21] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *Proceedings of CVPR*, pages 2492–2501. Computer Vision Foundation / IEEE, 2020. 4
- [22] Reinhard Heckel and Paul Hand. Deep decoder: Concise image representations from untrained non-convolutional networks. In *Proceedings of ICLR*. OpenReview.net, 2019. 4
- [23] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of ICCV*, pages 5855–5864. IEEE, 2021. 1
- [24] Xin Huang, Qi Zhang, Ying Feng, Hongdong Li, Xuan Wang, and Qing Wang. Hdr-nerf: High dynamic range neural radiance fields. In *Proceedings of CVPR*, pages 18377–18387. IEEE, 2022. 1
- [25] Takashi Isobe, Songjiang Li, Xu Jia, Shanxin Yuan, Gregory G. Slabaugh, Chunjing Xu, Ya-Li Li, Shengjin Wang, and Qi Tian. Video super-resolution with temporal group attention. In *Proceedings of CVPR*, pages 8005–8014. Computer Vision Foundation / IEEE, 2020. 2
- [26] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Baojin Huang, Yimin Luo, Jiayi Ma, and Junjun Jiang. Multi-scale progressive fusion network for single image deraining. In *Proceedings of CVPR*, pages 8343–8352. Computer Vision Foundation / IEEE, 2020. 2
- [27] Yifan Jiang, Peter Hedman, Ben Mildenhall, Dejia Xu, Jonathan T. Barron, Zhangyang Wang, and Tianfan Xue. Alignerf: High-fidelity neural radiance fields via alignment-aware training. In *Proceedings of CVPR*, pages 46–55. IEEE, 2023. 13
- [28] Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. Geonerf: Generalizing nerf with geometry priors. In *Proceedings of CVPR*, pages 18344–18347. IEEE, 2022. 1, 2, 5, 6, 7, 8, 13
- [29] Kim Jun-Seong, Kim Yu-Ji, Moon Ye-Bin, and Tae-Hyun Oh. Hdr-plenoxels: Self-calibrating high dynamic range radiance fields. In *Proceedings of ECCV*, pages 384–401. Springer, 2022. 1
- [30] Oguzhan Fatih Kar, Teresa Yeo, Andrei Atanov, and Amir Zamir. 3d common corruptions and data augmentation. In *Proceedings of CVPR*, pages 18941–18952. IEEE, 2022. 2, 4
- [31] Filippos Kokkinos and Stamatios Lefkimmiatis. Iterative residual cnns for burst photography applications. In *Proceedings of CVPR*, pages 5929–5938. Computer Vision Foundation / IEEE, 2019. 2
- [32] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiri Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *Proceedings of CVPR*, pages 8183–8192. Computer Vision Foundation / IEEE Computer Society, 2018. 2
- [33] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *Proceedings of ICCV*, pages 8877–8886. IEEE, 2019. 2
- [34] Dogyoon Lee, Minhyeok Lee, Chajin Shin, and Sangyoun Lee. Dp-nerf: Deblurred neural radiance field with physical scene priors. In *Proceedings of CVPR*, pages 12386–12396. IEEE, 2023. 1, 2
- [35] Deborah Levy, Amit Peleg, Naama Pearl, Dan Rosenbaum, Derya Akkaynak, Simon Korman, and Tali Treibitz. Seathru-nerf: Neural radiance fields in scattering media. In *Proceedings of CVPR*, pages 56–65. IEEE, 2023. 1
- [36] Taihui Li, Hengkang Wang, Zhong Zhuang, and Ju Sun. Deep random projector: Accelerated deep image prior. In *Proceedings of CVPR*, pages 18176–18185. IEEE, 2023. 4
- [37] Wenbo Li, Xin Lu, Shengju Qian, and Jiangbo Lu. On efficient transformer-based image pre-training for low-level vision. In *Proceedings of IJCAI*, pages 1089–1097. ijcai.org, 2023. 2
- [38] Jingyun Liang, Jiezhong Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. Swinir: Image restoration

- using swin transformer. In *Proceedings of ICCV Workshops*, pages 1833–1844. IEEE, 2021. 2
- [39] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. BARF: bundle-adjusting neural radiance fields. In *Proceedings of ICCV*, pages 5721–5731. IEEE, 2021. 13
- [40] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *Proceedings of CVPR*, pages 7814–7823. IEEE, 2022. 2
- [41] Yihao Liu, Jingwen He, Jinjin Gu, Xiangtao Kong, Yu Qiao, and Chao Dong. Degae: A new pretraining paradigm for low-level vision. In *Proceedings of CVPR*, pages 23292–23303. IEEE, 2023. 2
- [42] Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V. Sander. Deblur-nerf: Neural radiance fields from blurry images. In *Proceedings of CVPR*, pages 12851–12860. IEEE, 2022. 1, 2, 5, 6
- [43] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of CVPR*, pages 7210–7219. Computer Vision Foundation / IEEE, 2021. 2
- [44] Ben Mildenhall, Jonathan T. Barron, Jiawen Chen, Dillon Sharlet, Ren Ng, and Robert Carroll. Burst denoising with kernel prediction networks. In *Proceedings of CVPR*, pages 2502–2510. Computer Vision Foundation / IEEE Computer Society, 2018. 2, 4, 13
- [45] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: practical view synthesis with prescriptive sampling guidelines. *ACM Trans. Graph.*, 38(4):29:1–29:14, 2019. 6, 7
- [46] 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. In *Proceedings of ECCV*, pages 405–421. Springer, 2020. 1
- [47] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P. Srinivasan, and Jonathan T. Barron. Nerf in the dark: High dynamic range view synthesis from noisy raw images. In *Proceedings of CVPR*, pages 16169–16178. IEEE, 2022. 1, 2
- [48] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. 1
- [49] Naama Pearl, Tali Treibitz, and Simon Korman. NAN: noise-aware nerfs for burst-denoising. In *Proceedings of CVPR*, pages 12662–12671. IEEE, 2022. 2, 3, 4, 5, 6, 7, 8, 13, 18
- [50] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of ICCV*, pages 14315–14325. IEEE, 2021. 1
- [51] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Generalizable patch-based neural rendering. In *Proceedings of ECCV*, pages 156–174. Springer, 2022. 1, 2
- [52] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of CVPR*, pages 5449–5459. IEEE, 2022. 1
- [53] Mukund Varma T, Peihao Wang, Xuxi Chen, Tianlong Chen, Subhashini Venugopalan, and Zhangyang Wang. Is attention all that nerf needs? In *Proceedings of ICLR*, 2023. 1, 2, 6, 7, 8, 13, 17
- [54] Thomas Tanay, Ales Leonardis, and Matteo Maggioni. Efficient view synthesis and 3d-based multi-frame denoising with multiplane feature representations. In *Proceedings of CVPR*, pages 20898–20907. IEEE, 2023. 2
- [55] Dmitry Ulyanov, Andrea Vedaldi, and Victor S. Lempitsky. Deep image prior. *CoRR*, abs/1711.10925, 2017. 4
- [56] Chen Wang, Angtian Wang, Junbo Li, Alan L. Yuille, and Cihang Xie. Benchmarking robustness in neural radiance fields. *CoRR*, abs/2301.04075, 2023. 2
- [57] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. Patchmatchnet: Learned multi-view patchmatch stereo. In *Proceedings of CVPR*, pages 14194–14203. Computer Vision Foundation / IEEE, 2021. 4
- [58] Peng Wang, Lingzhe Zhao, Ruijie Ma, and Peidong Liu. Bad-nerf: Bundle adjusted deblur neural radiance fields. In *Proceedings of CVPR*, pages 4170–4179. IEEE, 2023. 1, 2
- [59] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas A. Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of CVPR*, pages 4690–4699. Computer Vision Foundation / IEEE, 2021. 1, 2, 6
- [60] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural radiance fields without known camera parameters. *CoRR*, abs/2102.07064, 2021. 2
- [61] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *Proceedings of CVPR*, pages 8534–8543. Computer Vision Foundation / IEEE, 2021. 1
- [62] Zhongkai Wu, Ziyu Wan, Jing Zhang, Jing Liao, and Dong Xu. Rafe: Generative radiance fields restoration, 2024. 2
- [63] Zhihao Xia, Federico Perazzi, Michaël Gharbi, Kalyan Sunkavalli, and Ayan Chakrabarti. Basis prediction networks for effective burst denoising with large kernels. In *Proceedings of CVPR*, pages 11841–11850. Computer Vision Foundation / IEEE, 2020. 2, 4, 13
- [64] Qingsong Yan, Qiang Wang, Kaiyong Zhao, Jie Chen, Bo Li, Xiaowen Chu, and Fei Deng. Cf-nerf: Camera parameter free neural radiance fields with incremental learning. In *Proceedings of AAAI*, pages 6440–6448. AAAI Press, 2024. 13
- [65] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of ECCV*, pages 785–801. Springer, 2018. 4
- [66] Tian Ye, Sixiang Chen, Yun Liu, Yi Ye, Erkang Chen, and Yuche Li. Underwater light field retention: Neural rendering for underwater imaging. In *Proceedings of CVPR*, pages 487–496. IEEE, 2022. 1

- [67] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of CVPR*, pages 4578–4587. Computer Vision Foundation / IEEE, 2021. [1](#), [2](#)
- [68] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *Proceedings of CVPR*, pages 5718–5729. IEEE, 2022. [2](#), [6](#), [8](#)
- [69] Zhihang Zhong, Mingdeng Cao, Xiang Ji, Yinqiang Zheng, and Imari Sato. Blur interpolation transformer for real-world motion from blur. In *Proceedings of CVPR*, pages 5713–5723. IEEE, 2023. [2](#)
- [70] Chengxuan Zhu, Renjie Wan, Yunkai Tang, and Boxin Shi. Occlusion-free scene recovery via neural radiance fields. In *Proceedings of CVPR*, pages 20722–20731. IEEE, 2023. [2](#)

## A. Appendix Section

Fig. 1 includes additional examples of our rendered images. For a better overview of the dataset, Figs. 6 and 7 include the histogram of the scene bound, scene dimension, and blur weight for each blur level. Here, scene range is the ratio between the near and far bound of the scene, and *scene dimension* is defined as the geometric mean of the three dimensions of the scene. As a dataset aiming to train a degradation robust Generalizable NeRF, more extensive variety and complexity should result in better generalizability to various use-cases. Additionally, we also show the distribution of blur weights, which is denoted as the  $w_l$  of Algorithm 1. We can notice that our dataset contains a wide variety of scenes with a well-distributed depth range and dimensions.

## B. Training Details

In this section, we provide a more comprehensive description of the training and reported results. For training GNeRF models [28, 49] with and without our proposed module, we have used an RTX3090 graphic card. During the pretraining stage of the GNeRF models with noise degradation ( Tabs. 7, 9 and 11), we first apply inverse gamma correction and random white balancing as in [1, 44, 49, 63] to linearize the image, then added noise. The predicted image by the GNeRF renderer is then post-processed by applying the gamma correction and white-balancing for calculating the photometric loss  $\mathcal{L}_{\text{photo}}$ . With such a setting, the models are trained over 150K and 200K iterations respectively during pretraining stage where the batch size is 512 rays.

For the adversarial training [16] of GNT [53], we follow their pretraining configuration which is to train over 250K iterations with batch size 512. We have used RTX3090 when source images are clean during pretraining stage (Row 1-2 of Tab. 10). For source images with adversarially degradation during the pretraining stage (Row 3-4 of Tab. 10), we have used NVIDIA A100 graphics card. The switch was necessary due to the substantial memory demands arising from learning the adversarial perturbation, which is dependent on the gradient of the entire 3D reconstruction process. Especially when depth estimation is integrated into the computation graph, the memory cost increases significantly.

## C. Visual Results

To demonstrate the effectiveness of our proposed 3D-Degradation Aware Feature Extractor, we include visual results of the intermediate outputs. Fig. 8 visualizes the results of differential warping (Eq. 5)) through the coarse depth predicted by the depth estimator  $\mathcal{D}$ . This process aims to align degraded source images with each source view, optimizing the extraction of information from the region. By achieving this alignment across multiple views, our method enhances feature awareness among nearby views, thereby

improving feature matching and overall image reconstruction.

Fig. 9 shows the restored image from auxiliary restoration head  $\mathcal{R}$ , which takes the input of nearby warped image and depths as described in Eq. 6). By engaging in the auxiliary task of restoring natural signals from degraded source images, our feature extractor adeptly captures the statistical properties inherent in natural images, thereby enhancing the overall 3D reconstruction process.

## D. Quantitative Result

Tabs. 8 to 11 include detailed experiment results, including LPIPS and SSIM metrics of 3D reconstruction from degraded source images. Additionally, Tab. 7 is an experiment demonstrating the effect of our proposed module on NAN [49] with a different number of source images. Results indicate that the effectiveness of our module increases with the number of source images. Such empirical result suggests that our module is particularly adept at extracting consistent deep image features from multiple degraded sources.

Next, we detail the methodology and underlying rationale presented in Table 5 of the main paper. Utilizing our novel blur dataset, which encompasses varying levels of a blur for identical viewpoints and scenes, we conducted inferences at the same locations under different blur levels to analyze the stability in depth predictions. The results reported are averages across 13 scenes consisting of the 52 test settings. Our dataset is uniquely constructed by averaging latent images, as specified in Algorithm 1. This approach allows us to emulate realistic blur effects that are depth-dependent, representing a more accurate simulation than the uniform degradation typically achieved with a standard blur kernel.

## E. Limitation

Our proposed module’s current limitation also relates to our future research directions. At present, similar to a predominant number of NeRF-based methods, our approach presupposes the availability of accurate camera poses during both pre-training and inference stages. However, with degraded source images, the reliability of pose estimation methods like COLMAP may be inaccurate. While we have demonstrated our method’s effectiveness on real-dataset in Sec. 6.5, with larger degradation resulting in inaccurate depth estimation and the downstream feature extraction, our method may fail. Recent works, such as those by [27, 39, 64], suggest methodologies for adjusting incorrect camera poses during training. Inspired by these advancements, we plan to enhance our module’s robustness against misalignment and noise in camera information, aiming to improve its performance under more challenging conditions.

Figure 6. Histogram of scene rang eand dimension of Objaverse Dataset.

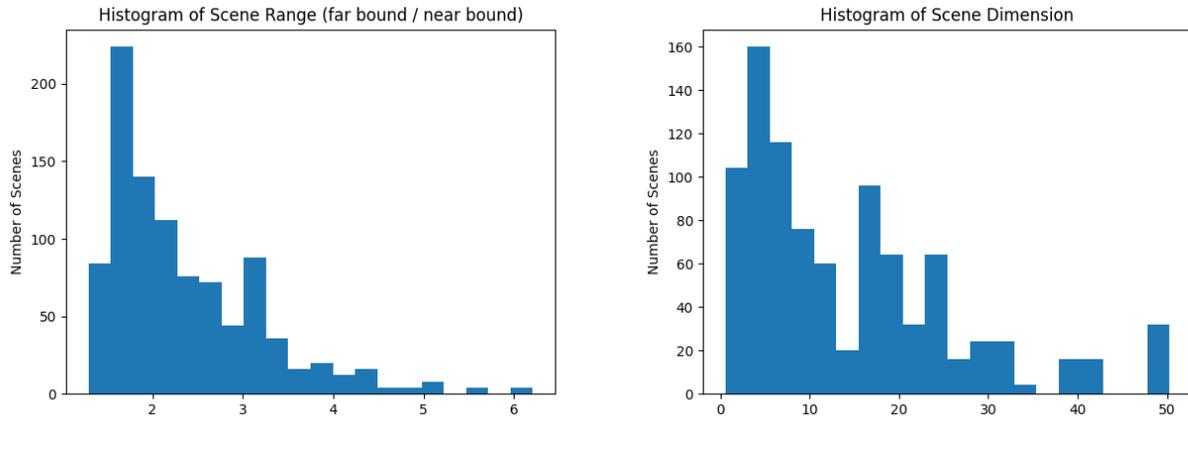


Figure 7. Histogram of blur weights at different levels of *Objaverse Blur* dataset.

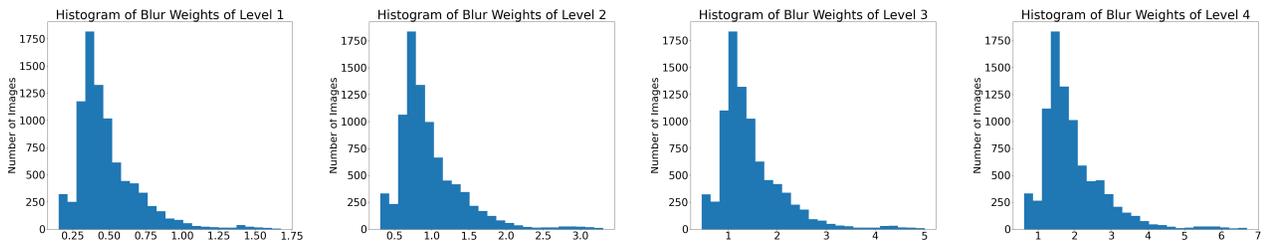


Figure 8. Visual results of 3D Differentiable Warping through Depth Estimator  $\mathcal{D}$

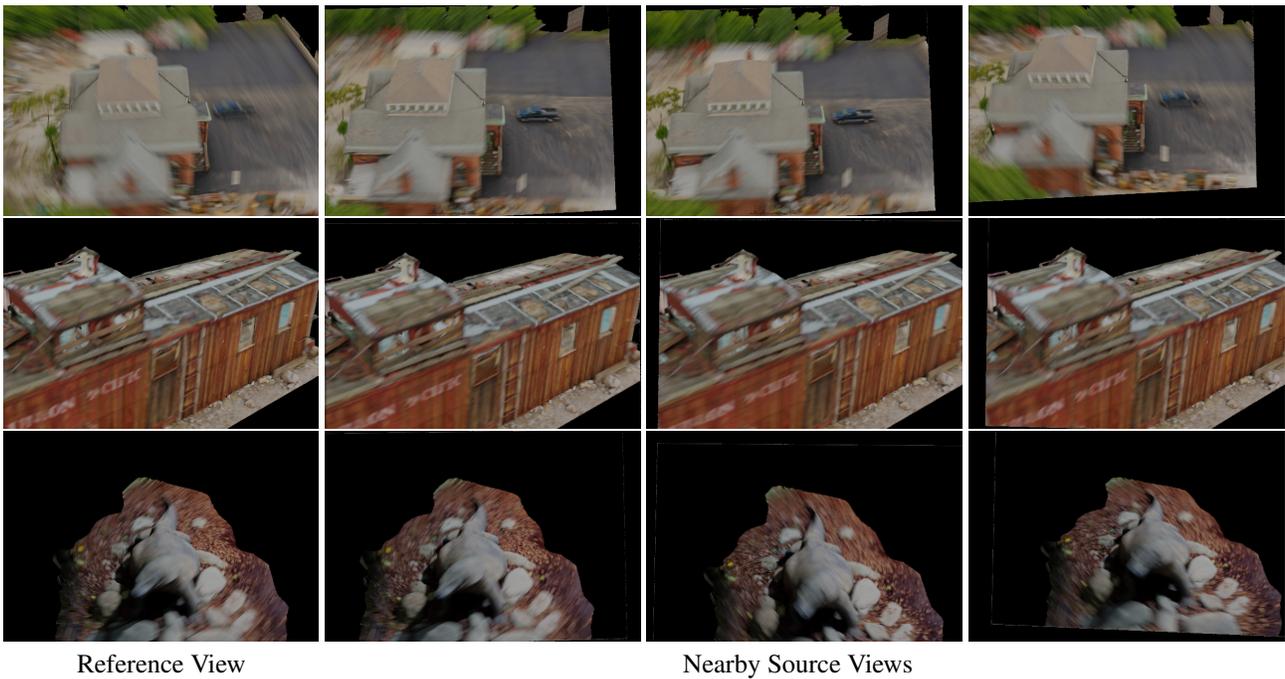


Figure 9. Visual results of 3D-aware Auxiliary Restoration Head  $\mathcal{R}$



Degraded Image

Restored Image

Ground Truth

Table 7. NAN: Novel view results of different numbers of source views on the *Objaverse Blur* dataset across blur degradation.

Blur Level	Method	# Source	PSNR $\uparrow$	$\Delta$ PSNR	SSIM $\uparrow$	$\Delta$ SSIM	LPIPS $\downarrow$	$\Delta$ LPIPS
1	NAN	3	23.15	-	0.766	-	0.268	-
		5	25.18	-	0.810	-	0.241	-
		7	25.59	-	0.823	-	0.231	-
	NAN+Proposed	3	23.74	0.59	0.797	0.031	0.226	-0.042
		5	25.65	0.47	0.839	0.029	0.199	-0.041
		7	26.51	0.92	0.850	0.027	0.194	-0.037
2	NAN	3	21.79	-	0.705	-	0.328	-
		5	23.71	-	0.751	-	0.299	-
		7	24.19	-	0.766	-	0.288	-
	NAN+Proposed	3	22.28	0.49	0.730	0.025	0.299	-0.029
		5	24.18	0.47	0.774	0.023	0.271	-0.028
		7	24.96	0.77	0.787	0.021	0.264	-0.024
3	NAN	3	21.0	-	0.671	-	0.363	-
		5	22.82	-	0.715	-	0.334	-
		7	23.28	-	0.731	-	0.323	-
	NAN+Proposed	3	21.43	0.43	0.692	0.021	0.339	-0.024
		5	23.26	0.44	0.735	0.020	0.313	-0.021
		7	23.93	0.65	0.748	0.017	0.3038	-0.0192
4	NAN	3	20.43	-	0.649	-	0.386	-
		5	22.17	-	0.691	-	0.357	-
		7	22.63	-	0.707	-	0.347	-
	NAN+Proposed	3	20.84	0.41	0.667	0.018	0.365	-0.021
		5	22.58	0.41	0.708	0.017	0.339	-0.018
		7	23.16	0.53	0.721	0.014	0.330	-0.017

Table 8. NAN and GeoNeRF: novel view results on the *Objaverse Blur* dataset across different blur levels.

Blur	Method	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$	Method	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$
1	GeoNeRF*	26.28	0.279	0.779	NAN*	23.73	0.266	0.771
	GeoNeRF	27.96	0.275	0.786	NAN	25.59	0.231	0.823
	+ Proposed	<b>28.47 (+0.51)</b>	<b>0.253 (-0.022)</b>	<b>0.838 (+0.052)</b>	+ Proposed	<b>26.51 (+0.92)</b>	<b>0.194 (-0.037)</b>	<b>0.850 (+0.027)</b>
2	GeoNeRF*	24.84	0.345	0.721	NAN*	22.11	0.323	0.704
	GeoNeRF	26.17	0.338	0.727	NAN	24.19	0.288	0.766
	+ Proposed	<b>26.73 (+0.56)</b>	<b>0.318 (-0.020)</b>	<b>0.776 (+0.049)</b>	+ Proposed	<b>24.96 (+0.77)</b>	<b>0.264 (-0.024)</b>	<b>0.787 (+0.021)</b>
3	GeoNeRF*	23.79	0.385	0.682	NAN*	21.14	0.353	0.671
	GeoNeRF	25.12	0.376	0.687	NAN	23.28	0.323	0.731
	+ Proposed	<b>25.69 (+0.57)</b>	<b>0.357 (-0.019)</b>	<b>0.734 (+0.047)</b>	+ Proposed	<b>23.93 (+0.65)</b>	<b>0.304 (-0.019)</b>	<b>0.748 (+0.017)</b>
4	GeoNeRF*	22.91	0.415	0.645	NAN*	20.39	0.374	0.650
	GeoNeRF	24.42	0.403	0.657	NAN	22.63	0.347	0.707
	+ Proposed	<b>24.97 (+0.56)</b>	<b>0.384 (-0.019)</b>	<b>0.706 (+0.049)</b>	+ Proposed	<b>23.16 (+0.53)</b>	<b>0.330 (-0.017)</b>	<b>0.721 (+0.014)</b>

Table 9. Novel view reconstruction results of NAN on noise degradation (*LLFF-N* dataset) across different gain levels.

Metric	Method	Gain 4	Gain 8	Gain 16	Gain 20
PSNR $\uparrow$	NAN	23.85	23.33	22.36	21.94
	+ Proposed	23.90	23.58	22.89	22.57
	<b>Difference</b>	<b>+0.05</b>	<b>+0.24</b>	<b>+0.53</b>	<b>+0.63</b>
LPIPS $\downarrow$	NAN	0.318	0.381	0.469	0.501
	+ Proposed	0.299	0.350	0.424	0.453
	<b>Difference</b>	<b>-0.019</b>	<b>-0.031</b>	<b>-0.045</b>	<b>-0.048</b>
SSIM $\uparrow$	NAN	0.754	0.706	0.623	0.586
	+ Proposed	0.774	0.740	0.681	0.655
	<b>Difference</b>	<b>+0.020</b>	<b>+0.034</b>	<b>+0.058</b>	<b>+0.069</b>

Table 10. Novel view reconstruction results of GNT [53] on different training and test methods across the *LLFF* dataset under adversarial degradation.

Training	Test	Proposed	Average	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	Trex
Clean	Clean	$\times$	23.50	22.31	24.62	28.02	24.32	18.47	17.93	27.76	22.53
		$\checkmark$	24.14	22.64	25.79	28.57	24.96	19.16	18.26	28.11	23.70
		<b>Difference</b>	<b>+0.64</b>	<b>+0.33</b>	<b>+1.17</b>	<b>+0.55</b>	<b>+0.64</b>	<b>+0.69</b>	<b>+0.33</b>	<b>+0.35</b>	<b>+1.18</b>
Clean	Adversarial	$\times$	17.68	18.06	18.95	20.78	18.47	15.34	14.38	19.57	15.90
		$\checkmark$	18.90	18.91	20.26	21.63	19.80	16.64	15.02	20.89	17.99
		<b>Difference</b>	<b>+1.22</b>	<b>+0.85</b>	<b>+1.31</b>	<b>+0.85</b>	<b>+1.33</b>	<b>+1.30</b>	<b>+0.64</b>	<b>+1.32</b>	<b>+2.09</b>
Adversarial	Clean	$\times$	24.26	22.44	25.37	27.93	24.52	19.04	18.15	27.87	23.25
		$\checkmark$	24.40	22.56	25.57	27.51	24.69	19.01	18.08	28.00	23.50
		<b>Difference</b>	<b>+0.14</b>	<b>+0.12</b>	<b>+0.20</b>	<b>-0.41</b>	<b>+0.16</b>	<b>-0.04</b>	<b>-0.07</b>	<b>+0.13</b>	<b>+0.25</b>
Adversarial	Adversarial	$\times$	20.05	19.95	21.77	23.20	21.35	17.48	16.32	23.41	20.47
		$\checkmark$	20.22	20.15	22.02	23.89	21.40	17.29	16.22	23.81	20.28
		<b>Difference</b>	<b>+0.17</b>	<b>+0.20</b>	<b>+0.25</b>	<b>+0.69</b>	<b>+0.05</b>	<b>-0.19</b>	<b>-0.10</b>	<b>+0.40</b>	<b>-0.19</b>

Table 11. Novel view reconstruction results of NAN [49] on the *Objaverse Blur* dataset across different blur levels and burst gain levels.

Blur Level	Method	Burst Noise Level	PSNR $\uparrow$	$\Delta$ PSNR	LPIPS $\downarrow$	$\Delta$ LPIPS	SSIM $\uparrow$	$\Delta$ SSIM
1	NAN	1	23.59	-	0.274	-	0.774	-
		2	23.47	-	0.293	-	0.751	-
		4	23.07	-	0.336	-	0.704	-
	NAN+Proposed	1	24.28	+0.69	0.244	-0.029	0.802	+0.028
		2	24.21	+0.74	0.257	-0.036	0.792	+0.041
		4	23.97	+0.90	0.293	-0.043	0.762	+0.059
2	NAN	1	22.26	-	0.328	-	0.712	-
		2	22.20	-	0.342	-	0.695	-
		4	21.91	-	0.376	-	0.650	-
	NAN+Proposed	1	22.70	+0.44	0.308	-0.019	0.733	+0.021
		2	22.65	+0.45	0.316	-0.026	0.725	+0.031
		4	22.39	+0.48	0.343	-0.033	0.698	+0.049
3	NAN	1	21.48	-	0.359	-	0.675	-
		2	21.43	-	0.371	-	0.659	-
		4	21.20	-	0.400	-	0.617	-
	NAN+Proposed	1	21.80	+0.32	0.347	-0.012	0.692	+0.017
		2	21.73	+0.30	0.352	-0.019	0.685	+0.026
		4	21.41	+0.21	0.373	-0.027	0.662	+0.045
4	NAN	1	20.93	-	0.379	-	0.651	-
		2	20.88	-	0.391	-	0.635	-
		4	20.66	-	0.417	-	0.596	-
	NAN+Proposed	1	21.17	+0.24	0.370	-0.009	0.666	+0.015
		2	21.09	+0.21	0.374	-0.017	0.660	+0.025
		4	20.74	+0.08	0.392	-0.025	0.638	+0.042