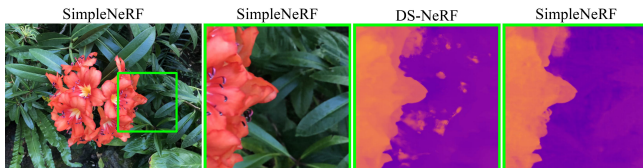


SimpleNeRF: Regularizing Sparse Input Neural Radiance Fields with Simpler Solutions

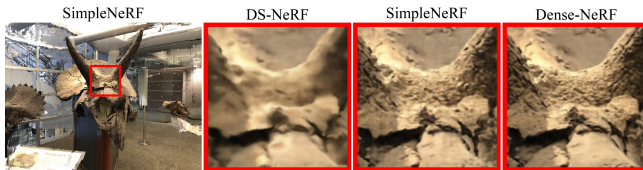
Nagabhushan Somraj
Indian Institute of Science
Bengaluru, Karnataka, India
nagabhushans@iisc.ac.in

Adithyan Karanayil
Indian Institute of Science
Bengaluru, India
adithyanv@iisc.ac.in

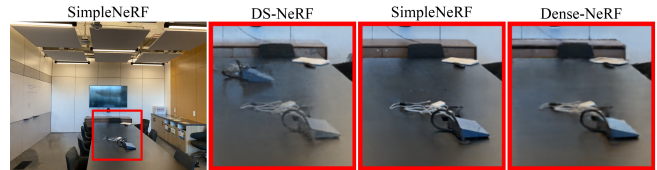
Rajiv Soundararajan
Indian Institute of Science
Bengaluru, India
rajivs@iisc.ac.in



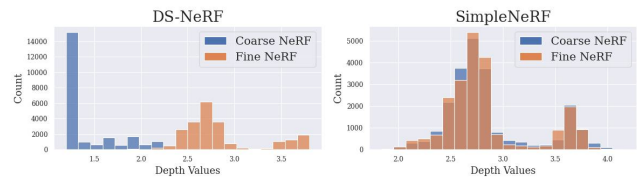
(a) Mitigation of floaters: For a frame from the LLFF flower scene, we show the depths predicted by DS-NeRF and SimpleNeRF. We show the complete frame for reference and focus on a small region to better observe the floaters.



(c) Improving Sharpness: The above images correspond to the LLFF horns scene. We enlarge a small region of the frame to better observe the improvement in sharpness.



(b) Reducing shape-radiance ambiguity: This is an example from the LLFF room scene. For reference, we show the prediction of Vanilla NeRF trained with dense input views.



(d) Improving Sharpness: Histogram of depth values predicted by the coarse and fine NeRF models for the image patch shown in Fig (c).

Figure 1: We show three shortcomings of DS-NeRF when trained with two input views on the LLFF dataset. SimpleNeRF introduces regularizations on DS-NeRF to mitigate these distortions. Notice the floaters in DS-NeRF predictions shown by small orange regions in Fig. (a) which are cleaned by SimpleNeRF. In Fig. (b), we find DS-NeRF suffers from shape-radiance ambiguity and ends up changing the color of the table in different viewpoints to match the observed images. DS-NeRF blends the two input images based on the viewpoint instead of learning correct geometry. This leads to ghosting artifacts in novel viewpoints, which is removed by our model. In Fig (c), we observe SimpleNeRF producing sharper reconstructions than DS-NeRF. Fig (d) shows a possible reason, where the coarse and fine NeRF models in DS-NeRF converge to different depth estimates. This leads to ineffective hierarchical sampling resulting in blurry predictions. We find that SimpleNeRF mitigates this by predicting consistent depth estimates.

ABSTRACT

Neural Radiance Fields (NeRF) show impressive performance for the photo-realistic free-view rendering of scenes. However, NeRFs require dense sampling of images in the given scene, and their performance degrades significantly when only a sparse set of views are available. Researchers have found that supervising the depth estimated by the NeRF helps train it effectively with fewer views. The depth supervision is obtained either using classical approaches or neural networks pre-trained on a large dataset. While the former may provide only sparse supervision, the latter may suffer from generalization issues. As opposed to the earlier approaches, we seek

to learn the depth supervision by designing augmented models and training them along with the NeRF. We design augmented models that encourage simpler solutions by exploring the role of positional encoding and view-dependent radiance in training the few-shot NeRF. The depth estimated by these simpler models is used to supervise the NeRF depth estimates. Since the augmented models can be inaccurate in certain regions, we design a mechanism to choose only reliable depth estimates for supervision. Finally, we add a consistency loss between the coarse and fine multi-layer perceptrons of the NeRF to ensure better utilization of hierarchical sampling. We achieve state-of-the-art view-synthesis performance on two popular datasets by employing the above regularizations. The source code for our model can be found on our project page: <https://nagabhushansn95.github.io/publications/2023/SimpleNeRF.html>

CCS CONCEPTS

• **Computing methodologies** → **Rendering; Volumetric models**; *Computer vision; Computational photography*; 3D imaging; Reconstruction.

KEYWORDS

neural rendering, novel view synthesis, sparse input NeRF

ACM Reference Format:

Nagabhushan Somraj, Adithyan Karanayil, and Rajiv Soundararajan. 2023. SimpleNeRF: Regularizing Sparse Input Neural Radiance Fields with Simpler Solutions. In *SIGGRAPH Asia 2023 Conference Papers (SA Conference Papers '23)*, December 12–15, 2023, Sydney, NSW, Australia. ACM, New York, NY, USA, 25 pages. <https://doi.org/10.1145/3610548.3618188>

1 INTRODUCTION

Neural Radiance Fields (NeRFs) [Mildenhall et al. 2020] show unprecedented levels of performance in synthesizing novel views of a scene by learning a volumetric representation implicitly within the weights of multi-layer perceptrons (MLP). Although NeRFs are very promising for view synthesis, there is a need to improve their design in a wide array of scenarios. For example, NeRFs have been enhanced to learn on unbounded scenes [Barron et al. 2022], handle scenes with highly specular objects [Verbin et al. 2022], and deal with noisy camera poses [Bian et al. 2023]. Yet NeRFs require tens to hundreds of images per scene to learn the scene geometry accurately, and their quality deteriorates significantly when only a few training images are available [Jain et al. 2021]. In this work, we focus on training the NeRF with a sparse set of input images and aim to design novel regularizations for effective training.

Prior work on sparse input NeRFs can be classified into generalized models and regularization based models. Generalized models [Johari et al. 2022] use convolutional neural networks to obtain a latent representation of the scene and use it to condition the NeRF. However, these models require a large dataset for pre-training and may suffer from generalization issues when used to render a novel scene [Niemeyer et al. 2022]. The other thread of work on sparse input NeRFs follows the original NeRF paradigm of training scene-specific NeRFs, and designs novel regularizations to assist NeRFs in converging to a better scene geometry [Zhang et al. 2021]. One popular approach among such models is to supervise the depth estimated by the NeRF. RegNeRF [Niemeyer et al. 2022], DS-NeRF [Deng et al. 2022b] and ViP-NeRF [Somraj and Soundararajan 2023] use simple priors such as depth smoothness, sparse depth or relative depth respectively obtained through classical approaches. On the other hand, DDP-NeRF [Roessle et al. 2022] and SCADE [Uy et al. 2023] pre-train convolutional neural networks (CNN) on a large dataset of scenes to learn a prior on the depth or its probability distribution. These approaches may also suffer from similar issues as the generalized models. This raises the question of whether we can instead learn the depth supervision in-situ without employing any pre-training. We aim to regularize NeRFs by learning augmented models for depth supervision in tandem with the NeRF training.

We follow the popular Occam’s razor principle and regularize the NeRF by designing augmented models to choose simpler solutions over complex ones, wherever possible. Some of the key components

of the NeRF such as positional encoding, view-dependent radiance, and hierarchical sampling provide powerful capabilities to the NeRF and are designed for the case with dense input views. Existing implementations of these components may be sub-optimal with fewer input views, perhaps due to the highly under-constrained system, causing several distortions. Fig. 1 shows three common distortions, namely, floaters [Roessle et al. 2022], shape-radiance ambiguity [Zhang et al. 2020], and blurred renders for the NeRF model in the few-shot setting. We believe that by simplifying some of the capabilities, one can obtain simpler augmented models that may provide better depth supervision for training the NeRF.

We simplify the capability of the augmented NeRF models with respect to positional encoding and view-dependent radiance. Positional encoding maps two nearby points in 3D space to far-apart points in the encoded space. This allows the NeRF to learn sharp depth discontinuities in 3D space via a smooth function in the encoded space. However, in an under-constrained system, the NeRF learns many undesirable depth discontinuities that can still explain the observed images. This leads to floater artifacts as seen in Fig. 1a. Restricting the ability of the NeRF to learn these sharp discontinuities encourages it to mitigate the floater artifacts. Similarly, the ability of the NeRF to predict view-dependent radiance leads to shape-radiance ambiguity, which we find to be more pronounced in the case of few-shot NeRF. The NeRF can learn to explain the observed images by modifying the color of 3D points in accordance with the input images as seen in Fig. 1b. Restricting the NeRF to predict view-independent radiance only encourages the NeRF to explain the observed images using simpler Lambertian surfaces and can help avoid shape-radiance ambiguity.

We use the depth estimated by the simpler augmented models to supervise the depth estimated by the NeRF model. Note that we use the simplified models as augmentations for depth supervision and not as the main NeRF model since naively reducing the capacity of the NeRF may lead to suboptimal solutions [Jain et al. 2021]. For example, the model that predicts view-independent radiance fails if the scene contains specular objects. Further, the simpler augmented models need to be used for supervision only if they explain the observed images accurately. We gauge the reliability of the depths by reprojecting pixels using the estimated depths onto a different nearest train view and comparing them with the corresponding images. We use DS-NeRF [Deng et al. 2022b] as our baseline and design our regularizations on top of it. Our framework can thus be seen as a semi-supervised learning model by considering the sparse depth from a Structure from Motion (SfM) module as providing limited depth labels and the remaining pixels as the unlabeled data. Our approach of using augmented models in tandem with the main NeRF model is perhaps closest to the Dual-Student architecture [Ke et al. 2019] that trains another identical model in tandem with the main model and imposes consistency regularization between the predictions of the two models. In contrast, our augmented models have complementary abilities as compared to the main NeRF model.

Finally, we observe that the coarse and fine NeRFs may converge to different depth estimates when trained with fewer images, as seen in Fig. 1d. This essentially renders the hierarchical sampling ineffective. The resulting model is similar to the one with under-sampled points along the rays. We avoid such degenerate cases by imposing a consistency loss between the depths estimated by the

coarse and fine MLPs. We achieve state-of-the-art view synthesis performance on two popular datasets with the above three regularizations. Further, we show that our model learns significantly improved geometry as compared to the prior art. We refer to our model as SimpleNeRF.

We list the main contributions of our work in the following.

- We design two augmented NeRF models that are biased towards simpler solutions and use the depth estimates from these models to regularize the NeRF training. Through the two augmentations, we mitigate incorrect depth discontinuities and shape-radiance ambiguity.
- We design a mechanism to determine whether the depths estimated by the augmented models are accurate and utilize only the accurate estimates to supervise the NeRF.
- We improve the effectiveness of hierarchical sampling by introducing a consistency constraint between the coarse and fine NeRFs. This generates sharper frames.
- We achieve the state-of-the-art performance of few shot NeRFs on two popular datasets.

2 RELATED WORK

Novel view synthesis is a classic problem traditionally solved broadly using image based rendering [Chen and Williams 1993] or light fields [Gortler et al. 1996; Levoy and Hanrahan 1996]. The seminal work by Mildenhall et al. [2020] started a new pathway in neural view synthesis and led to NeRF based models being employed in a wide variety of applications such as 3D editing of scenes [Yuan et al. 2022], gaming [Menapace et al. 2023], extended reality [Deng et al. 2022a] and image reconstruction [Ma et al. 2022; Mildenhall et al. 2022; Pearl et al. 2022], among others. However, many of the above models require dense sampling of input views for a faithful 3D geometry generation. With fewer views, the quality of rendered novel views and the learned 3D geometry degrade significantly introducing severe distortions. This can limit the widespread usage of NeRFs in multiple applications and hence addressing this limitation is of considerable interest.

2.1 Generalized Sparse Input NeRF

Prior work involves various approaches to learning a 3D neural representation with fewer input views. One line of approach attempts to train a generalized model on a large dataset of scenes such that the model can utilize the learned prior to generate a 3D scene representation from the few input images [Chen et al. 2021; Lee et al. 2023; Tancik et al. 2021]. Early pieces of work such as PixelNeRF [Yu et al. 2021], GRF [Trevithick and Yang 2021], and IBNet [Wang et al. 2021] obtain convolutional features of the input images and additionally condition the NeRF by projecting the 3D points onto the feature grids. MVSNeRF [Chen et al. 2021] and SRF [Chibane et al. 2021] incorporate cross-view knowledge into the features by constructing a cost volume and pair-wise post-processing of the individual frame features respectively. GeoNeRF [Johari et al. 2022] further improves the performance by employing a transformer to effectively reason about the occlusions in the scene. More recent work such as GARF [Shi et al. 2022] and MatchNeRF [Chen et al. 2023] try to provide explicit knowledge about the scene geometry through depth maps and similarity of the projected features

respectively. This approach of conditioning the NeRF on learned features is also popular among single image NeRF models [Lin et al. 2023], which can be considered an extreme case of sparse input NeRF. However, the need for pre-training on a large dataset of scenes with multi-view images and issues due to domain shift have motivated researchers to adopt regularization based approaches.

2.2 Regularization based Sparse Input NeRF

A popular approach among regularization based models is to regularize the depth estimated by the NeRF. DS-NeRF [Deng et al. 2022b] uses sparse depth provided by a SfM module to supervise the NeRF estimated depth at sparse keypoints. RegNeRF [Niemeyer et al. 2022] imposes the depth smoothness constraint on the rendered depth maps. ViP-NeRF [Somraj and Soundararajan 2023] instead attempts to regularize the relative depth of objects by obtaining a prior on the visibility of objects. Unlike the above, few other works exploit the advances in depth-estimation using deep neural networks. DDP-NeRF [Roessle et al. 2022] extends DS-NeRF by employing a CNN to convert the sparse depth into dense depth for more supervision. SCADE [Uy et al. 2023] and SparseNeRF [Wang et al. 2023] use the depth map output by single image depth models to constrain the absolute and the relative order of pixel depths respectively. DiffusioNeRF [Wynn and Turmukhambetov 2023] learns the joint distribution of RGBD patches using denoising diffusion models (DDM) and utilizes the gradient of the distribution provided by the DDM to regularize NeRF rendered RGBD patches. Different from the above, our work obtains depth supervision by harnessing the power of learning through augmented models, but without the need for pre-training on a large dataset.

Another line of regularization based approaches hallucinate new viewpoints and regularize the NeRF on different aspects such as semantic consistency [Jain et al. 2021], depth smoothness [Niemeyer et al. 2022], sparsity of mass [Kim et al. 2022] and depth based reprojection consistency [Bortolon et al. 2022; Chen et al. 2022a; Kwak et al. 2023; Xu et al. 2022]. More recent works have explored other forms of regularizations. FreeNeRF [Yang et al. 2023] anneals the frequency range of positional encoded NeRF inputs as the training progresses. MixNeRF [Seo et al. 2023] regularizes the NeRF by modeling the volume density along a ray as a mixture of Laplacian distributions. A recent work VDN-NeRF [Zhu et al. 2023], aims to resolve shape-radiance ambiguity, but is designed for training the NeRF with dense input views. However, our regularization is aimed at sparse input NeRF with as few as two views.

3 NERF PRELIMINARIES

We first provide a brief recap of the NeRF and describe the notation required for further sections. To render a pixel \mathbf{q} , the NeRF shoots a corresponding ray into the scene and samples N 3D points $\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_N$, where \mathbf{p}_1 and \mathbf{p}_N are the closest to and farthest from the camera, respectively. Two MLPs $\mathcal{F}_1, \mathcal{F}_2$ then predict view-independent volume density σ_i and view-dependent color \mathbf{c}_i as

$$\sigma_i, \mathbf{h}_i = \mathcal{F}_1(\gamma(\mathbf{p}_i, 0, l_p)); \quad \mathbf{c}_i = \mathcal{F}_2(\mathbf{h}_i, \gamma(\mathbf{v}, 0, l_v)), \quad (1)$$

where \mathbf{v} is the viewing direction, \mathbf{h}_i is a latent feature of \mathbf{p}_i and

$$\gamma(x, d_1, d_2) = [x, \sin(2^{d_1}x), \cos(2^{d_1}x), \dots, \sin(2^{d_2-1}x), \cos(2^{d_2-1}x)]$$

is the positional encoding. l_p and l_v are the highest positional encoding frequencies for \mathbf{p}_i and \mathbf{v} respectively. Volume rendering is then applied along every ray to obtain the color for each pixel as $\mathbf{c} = \sum_{i=1}^N w_i \mathbf{c}_i$, where the weights w_i are computed as

$$w_i = \exp\left(-\sum_{j=1}^{i-1} \delta_j \sigma_j\right) \cdot (1 - \exp(-\delta_i \sigma_i)), \quad (2)$$

and δ_i is the distance between \mathbf{p}_i and \mathbf{p}_{i+1} . The expected ray termination length is computed as $z = \sum_{i=1}^N w_i z_i$, where z_i is the depth of \mathbf{p}_i . z is typically also used as the depth of the pixel \mathbf{q} [Deng et al. 2022b]. \mathcal{F}_1 and \mathcal{F}_2 are trained using the mean squared error loss (MSE) as $\mathcal{L}_{\text{color}} = \|\mathbf{c} - \hat{\mathbf{c}}\|^2$, where $\hat{\mathbf{c}}$ is the true color of \mathbf{q} .

NeRF circumvents the need for the dense sampling of 3D points by employing two sets of MLPs, a coarse NeRF and a fine NeRF, both trained using $\mathcal{L}_{\text{color}}$. The coarse NeRF is trained with a coarse stratified sampling, and the fine NeRF with dense sampling around object surfaces, where object surfaces are coarsely localized based on the predictions of the coarse NeRF.

4 METHOD

Our key idea is to employ simpler NeRF models to obtain better depth supervision in certain regions of the scene. We describe our regularizations based on the simpler solutions in Sec. 4.1. We explain our approach to selecting reliable depth estimates for supervision in Sec. 4.2. Finally, Sec. 4.3 describes our solution to mitigate the sub-optimal utilization of hierarchical sampling. Fig. 2 shows the overall architecture of our model.

4.1 Regularization with Simpler Solutions

Our regularization consists of simplifying the NeRF model with respect to the positional encoding and view-dependent radiance capabilities to obtain better depth supervision. Both positional encoding and view-dependent radiance are elements designed to increase the capability of the NeRF to explain complex phenomena. For example, the former helps in learning thin, repetitive objects against a farther background and the latter for specular objects. However, when training with sparse views, the fewer constraints coupled with the higher capacity of the NeRF lead to solutions that overfit the observed images and learn implausible scene geometries.

We rethink ways of employing positional encoding and view-dependent radiance such that the NeRF utilizes the higher capability only when needed. The challenge here is that it is not known a priori where one needs to employ the higher capability of the NeRF. Our solution here is to use the higher capability NeRF as the main model and employ lower capability NeRFs as augmentations to provide guidance on where to use simpler solutions. Since the scene geometry is mainly learned by the coarse NeRF, we add the augmentations only to the coarse NeRF. We employ two augmentations, one each for regularizing positional encoding and view-dependent radiance, which we describe in the following subsections. In the remainder of this paper, we refer to the two augmentations as points (Sec. 4.1.1) and views (Sec. 4.1.2) augmentations respectively.

4.1.1 Undesirable Depth Discontinuities. The positional encoding maps two nearby points in \mathbb{R}^3 to two farther away points in $\mathbb{R}^{3(2l_p+1)}$ allowing the NeRF to learn sharp discontinuities between the two

points in \mathbb{R}^3 as a smooth function in $\mathbb{R}^{3(2l_p+1)}$. However, this is mainly required at depth edges, and most regions of natural scenes are typically smooth in depth. With sparse input views, the positional encoding causes the NeRF to learn depth discontinuities even in smooth regions due to incorrect matching of pixels between the input views. This gives rise to “floaters” [Barron et al. 2022] where a part of an object is broken away from it and “floats” freely in space. We reduce the depth discontinuities by reducing the highest positional encoding frequency for \mathbf{p}_i to $l_p^{\text{ap}} < l_p$ as

$$\sigma_i, \mathbf{h}_i = \mathcal{F}_1^{\text{ap}}(\gamma(\mathbf{p}_i, 0, l_p^{\text{ap}})), \quad (3)$$

where $\mathcal{F}_1^{\text{ap}}$ is the MLP of the augmented model. The main model is more accurate where depth discontinuities are required and the augmented model is more accurate where discontinuities are not required. We determine the above using a ternary mask m_{ap} as we explain in Sec. 4.2. We supervise the depth predicted by the main NeRF using that of the augmented model, for those pixels for which the depth estimated by the augmented model is reliable, by setting $m_{\text{ap}} = 1$. Similarly, we also determine the pixels for which the depth estimated by the main model is more accurate than that of the augmented model, where we set $m_{\text{ap}} = -1$. In such pixels, we supervise the augmented model with the depth estimated by the main model, which helps the augmented model to improve further and provide better depth supervision for the main model. For pixels where the depth estimated by both the main and augmented models are unreliable, we set $m_{\text{ap}} = 0$. If z_c and z_{ap} are the depths estimated by the coarse NeRF of the main and the augmented models respectively, we impose the depth supervision as

$$\mathcal{L}_{\text{ap}} = \mathbb{1}_{\{m_{\text{ap}}=1\}} \odot \|z_c - \mathcal{T}(z_{\text{ap}})\|^2 + \mathbb{1}_{\{m_{\text{ap}}=-1\}} \odot \|\mathcal{T}(z_c) - z_{\text{ap}}\|^2, \quad (4)$$

where \odot denotes element-wise product, $\mathbb{1}$ is the indicator function and \mathcal{T} is the stop-gradient operator.

Since color tends to have more discontinuities than depth in regions such as textures, we include the remaining high-frequency positional encoding components of \mathbf{p}_i in the input for \mathcal{F}_2 as

$$\mathbf{c}_i = \mathcal{F}_2^{\text{ap}}(\mathbf{h}_i, \gamma(\mathbf{p}_i, l_p^{\text{ap}}, l_p), \gamma(\mathbf{v}_i, 0, l_v)). \quad (5)$$

Note that \mathbf{h}_i already includes the low-frequency positional encoding components of \mathbf{p}_i .

4.1.2 Shape-Radiance Ambiguity. The ability of the NeRF to predict view-dependent radiance helps it learn non-Lambertian surfaces. With fewer images, the NeRF can simply learn any random geometry and change the color of 3D points in accordance with the input viewpoint to explain away the observed images [Zhang et al. 2020]. To bias the NeRF against this, we disable the view-dependent radiance in the second augmented NeRF model to output color based on \mathbf{p}_i alone. If z_{av} is the depth estimated by the augmented model, we impose the depth supervision as

$$\mathcal{L}_{\text{av}} = \mathbb{1}_{\{m_{\text{av}}=1\}} \odot \|z_c - \mathcal{T}(z_{\text{av}})\|^2 + \mathbb{1}_{\{m_{\text{av}}=-1\}} \odot \|\mathcal{T}(z_c) - z_{\text{av}}\|^2, \quad (6)$$

where m_{av} is a ternary mask indicating where the depths estimated by the augmented and the main model are reliable. We note that while the augmented model is more accurate in Lambertian regions, the main model is better equipped to handle specular objects.

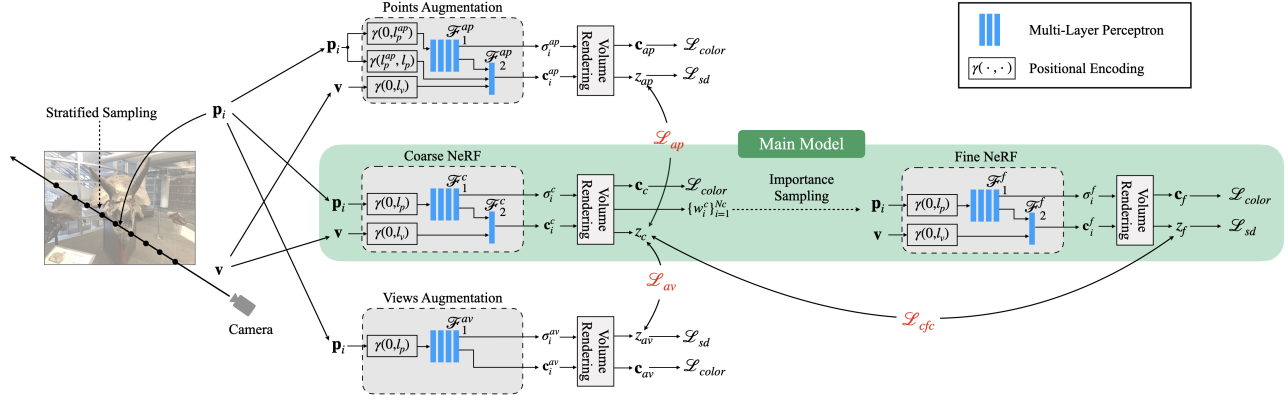


Figure 2: Architecture of SimpleNeRF. We train two augmented NeRF models in tandem with the NeRF to obtain depth supervision. In points augmentation, we reduce the positional encoding frequencies input to \mathcal{F}_1 and concatenate them to the input of \mathcal{F}_2 . For views augmentation, we ask \mathcal{F}_1 to output both volume density and color based on position alone. We add depth supervision losses \mathcal{L}_{ap} and \mathcal{L}_{av} between the coarse NeRFs of the main and augmented models and a consistency loss \mathcal{L}_{cfc} between the coarse and fine NeRFs of the main model. During inference, only the Main Model is employed.

Table 1: Quantitative results on LLFF dataset.

Model	2 views			3 views			4 views		
	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
RegNeRF	0.3056	0.5712	18.52	0.2908	0.6334	20.22	0.2794	0.6645	21.32
FreeNeRF	0.2638	0.6322	19.52	0.2754	0.6583	20.93	0.2848	0.6764	21.91
DS-NeRF	0.3106	0.5862	18.24	0.3031	0.6321	20.20	0.2979	0.6582	21.23
DDP-NeRF	0.2851	0.6218	18.73	0.3250	0.6152	18.73	0.3042	0.6558	20.17
ViP-NeRF	0.2768	0.6225	18.61	0.2798	0.6548	20.54	0.2854	0.6675	20.75
SimpleNeRF	0.2688	0.6501	19.57	0.2559	0.6940	21.37	0.2633	0.7016	21.99

4.2 Determining Reliable Depth Estimates

We follow similar procedures to determine the masks m_{ap} and m_{av} and hence explain the mask computation using a generic variable m_a to denote either of the masks. Given the depths z_c and z_a estimated by the main and augmented models respectively for pixel \mathbf{q} , we reproject a $k \times k$ patch around \mathbf{q} to the nearest training view using both z_c and z_a . We compute the MSE in intensities between the reprojected patch and the corresponding patch in the training image and choose the depth corresponding to lower MSE as the reliable depth. To filter out the cases where both the main and augmented models predict incorrect depth, we define a threshold e_τ and mark the depth to be reliable if its corresponding MSE is also less than e_τ . If e_c and e_a are the reprojection MSE corresponding to z_c and z_a respectively, we compute the mask as

$$m_a = \begin{cases} 1 & \text{if } (e_a \leq e_c) \text{ and } (e_a \leq e_\tau) \\ -1 & \text{if } (e_c < e_a) \text{ and } (e_c \leq e_\tau) \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

For specular regions, the reprojected patches may not match in intensities leading to m_a being zero. This implies supervision for fewer pixels and not supervision with incorrect depth estimates. We note that although both the views augmentation and the depth

reliability estimation work only in the non-specular regions, there is no redundancy in the model. The views augmentation model may still make errors in the non-specular regions, and hence it is necessary to determine the reliability of its depth estimates.

4.3 Hierarchical Sampling

Since multiple solutions can explain the observed images in the few-shot setting, the coarse and fine MLP may converge to different depth estimates for a given pixel as shown in Fig. 1c. Thus, dense sampling may not be employed around the region where fine NeRF predicts the object surface, which is equivalent to using only the coarse sampling for the fine NeRF. This can lead to blur in rendered images as seen in Fig. 1d. To prevent such inconsistencies, we drive the two NeRFs to be consistent in their solutions by imposing an MSE loss between the depths predicted by the two NeRFs. If z_c and z_f are the depths estimated by the coarse and fine NeRFs respectively, we define the coarse-fine consistency loss as

$$\mathcal{L}_{cfc} = \mathbb{1}_{\{m_{cfc}=1\}} \odot \|z_c - \mathcal{T}(z_f)\|^2 + \mathbb{1}_{\{m_{cfc}=-1\}} \odot \|\mathcal{T}(z_c) - z_f\|^2, \quad (8)$$

where the mask m_{cfc} is determined as described in Sec. 4.2.

Table 2: Quantitative results on RealEstate-10K dataset.

Model	2 views			3 views			4 views		
	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
RegNeRF	0.4129	0.5916	17.14	0.4171	0.6132	17.86	0.4316	0.6257	18.34
FreeNeRF	0.5036	0.5354	14.70	0.4635	0.5708	15.26	0.5226	0.6027	16.31
DS-NeRF	0.2709	0.7983	26.26	0.2893	0.8004	26.50	0.3103	0.7999	26.65
DDP-NeRF	0.1290	0.8640	27.79	0.1518	0.8587	26.67	0.1563	0.8617	27.07
ViP-NeRF	0.0687	0.8889	32.32	0.0758	0.8967	31.93	0.0892	0.8968	31.95
SimpleNeRF	0.0635	0.8942	33.10	0.0726	0.8984	33.21	0.0847	0.8987	32.88

Table 3: Evaluation of depth estimated by different models with two input views. The reference depth is obtained using NeRF with dense input views.

model	LLFF		RealEstate-10K	
	MAE ↓	SROCC ↑	MAE ↓	SROCC ↑
DS-NeRF	0.2074	0.7230	0.7164	0.6660
DDP-NeRF	0.2048	0.7480	0.4831	0.7921
ViP-NeRF	0.1999	0.7344	0.3856	0.8446
SimpleNeRF	0.1420	0.8480	0.3269	0.9215

4.4 Overall Loss

We impose the pixel color reconstruction loss on the main model and the augmented models as

$$\mathcal{L}_{\text{color}} = \|c_c - \hat{c}\|^2 + \|c_f - \hat{c}\|^2 + \|c_{\text{ap}} - \hat{c}\|^2 + \|c_{\text{av}} - \hat{c}\|^2, \quad (9)$$

where the subscripts $c, f, \text{'ap'},$ and 'av' denote the outputs of the coarse NeRF, fine NeRF, points augmentation model, and the views augmentation model respectively. We also include the sparse depth loss on the models as,

$$\mathcal{L}_{\text{sd}} = \|z_f - \hat{z}\|^2 + \|z_{\text{ap}} - \hat{z}\|^2 + \|z_{\text{av}} - \hat{z}\|^2, \quad (10)$$

where \hat{z} is the sparse depth given by the SfM model. We do not impose \mathcal{L}_{sd} on the coarse NeRF of the main model following Deng et al. [2022b]. Our final loss is a combination of all the losses as

$$\mathcal{L} = \lambda_1 \mathcal{L}_{\text{color}} + \lambda_2 \mathcal{L}_{\text{sd}} + \lambda_3 \mathcal{L}_{\text{ap}} + \lambda_4 \mathcal{L}_{\text{av}} + \lambda_5 \mathcal{L}_{\text{cfc}} \quad (11)$$

5 EXPERIMENTS

5.1 Evaluation Setup

We compare the performance of sparse input NeRF models on LLFF [Mildenhall et al. 2019] and RealEstate-10K [Zhou et al. 2018] datasets with 2, 3, and 4 input views. We assume the camera parameters are known for the input images, since in applications such as robotics or extended reality, external sensors or pre-calibrated set of cameras may provide the camera poses. We follow prior work [Somraj and Soundararajan 2023] to choose the train and test images. We provide more details in the supplementary.

We quantitatively evaluate the predicted frames from various models using peak signal to noise ratio (PSNR), structural similarity (SSIM) [Wang et al. 2004] and LPIPS [Zhang et al. 2018] measures. We employ depth mean absolute error (MAE) and spearman rank order correlation coefficient (SROCC) to evaluate the models on

their ability to predict absolute and relative depth in novel views. We train a NeRF model with dense input views and use its depth predictions as pseudo ground truth. On the LLFF dataset, we normalize the predicted depths by the median ground truth depth, since the scenes have different depth ranges. Following RegNeRF [Niemeyer et al. 2022], we evaluate the predictions only in the regions of interest. We mask out the regions of test frames which are not visible in the train frames. To determine such regions, we use the depth estimated by a NeRF trained with dense input views and compute the visible region mask through reprojection error in depth. We provide more details on the mask computation in the supplementary along with the unmasked evaluation scores for reference.

5.2 Comparisons

We evaluate the performance of our model against various sparse input NeRF models on both datasets. We compare with DS-NeRF [Deng et al. 2022b], DDP-NeRF [Roessle et al. 2022] and RegNeRF [Niemeyer et al. 2022] which regularize the depth estimated by the NeRF. Further, we include two recent models, FreeNeRF [Yang et al. 2023] and ViP-NeRF [Somraj and Soundararajan 2023], among the comparisons. We train the models on both datasets using the codes provided by the respective authors. Implementation details of our model are provided in the supplement.

5.3 Results

Tabs. 1 and 2 show the view-synthesis performance of SimpleNeRF and other prior art on LLFF and RealEstate-10K datasets. We find that SimpleNeRF achieves state-of-the-art performance on both datasets in most cases. The higher performance of all the models on the RealEstate-10K dataset is perhaps due to the scenes being simpler. Hence, the performance improvement is also smaller as compared to the LLFF dataset. On RealEstate-10K, we observe that all the models struggle on one of the five scenes as compared to the other scenes. Excluding this scene, with two input views, SimpleNeRF improves SSIM over ViP-NeRF from 0.9596 to 0.9685, which we believe is a significant improvement at such high quality regime. In Tab. 2, we show the average performance on all five scenes and show the per-scene performance of various models in the supplementary.

Fig. 3 shows predictions of various models on an example scene from the RealEstate-10K dataset, where we observe that SimpleNeRF is the best in reconstructing the novel view. Figs. 5 to 11 show more comparisons on both datasets. Further, SimpleNeRF improves significantly in estimating the depth of the scene as seen in Tab. 4

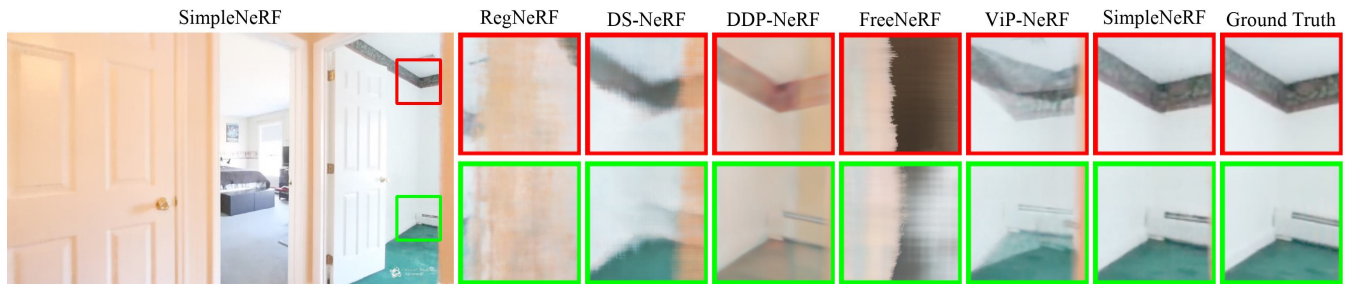


Figure 3: Qualitative examples on RealEstate-10K dataset with three input views. SimpleNeRF predictions are closest to the ground truth among all the models. In particular, DDP-NeRF predictions have a different shade of color and ViP-NeRF suffers from shape-radiance ambiguity creating ghosting artifacts.

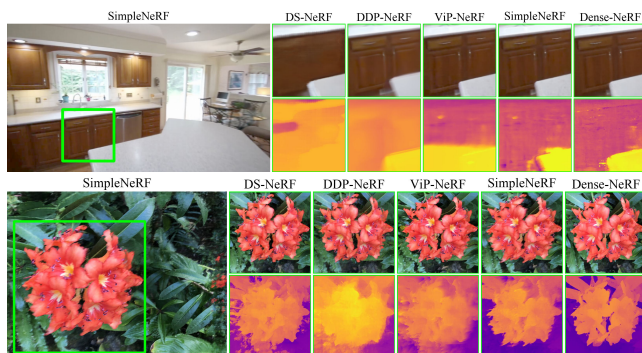


Figure 4: Estimated depth maps on RealEstate-10K and LLLF datasets with two input views. In both examples, the two rows show the predicted images and the depths respectively. We find that SimpleNeRF is significantly better at estimating the scene depth. Also, DDP-NeRF synthesizes the left table edge at a different angle due to incorrect depth estimation.

and Fig. 4. SimpleNeRF also performs significantly better in estimating relative depth, even outperforming ViP-NeRF which uses a prior based on relative depth. Estimating better geometry may be more crucial in downstream applications such as 3D scene editing. We provide video comparisons in the supplementary.

5.3.1 Ablations. We test the importance of each of the components of our model components by disabling them one at a time. We disable the points and views augmentations and coarse-fine consistency loss individually. When disabling \mathcal{L}_{cfc} , we additionally add augmentations to the fine NeRF since the knowledge learned by coarse NeRF may not propagate to the fine NeRF. We also analyze the need to supervise with only the reliable depth estimates by disabling the mask and stop-gradients in \mathcal{L}_{ap} , \mathcal{L}_{av} , and \mathcal{L}_{cfc} . Tab. 4 shows a quantitative comparison between the ablated models.

We observe that each of the components is crucial and disabling any of them leads to a drop in the performance. Further, using all the depths for supervision instead of only the reliable depths leads to a significant drop in performance. Finally, disabling \mathcal{L}_{cfc} also leads to a drop in performance in addition to increasing the training time by almost 2 \times due to the inclusion of augmentations for the fine NeRF.

Table 4: Ablation experiments on both datasets with two input views.

model	RealEstate-10K		LLFF	
	LPIPS ↓	MAE ↓	LPIPS ↓	MAE ↓
SimpleNeRF	0.0635	0.33	0.2688	0.14
w/o points augmentation	0.0752	0.38	0.2832	0.15
w/o views augmentation	0.0790	0.39	0.2834	0.15
w/o coarse-fine consistency	0.0740	0.42	0.3002	0.19
w/o reliable depth	0.0687	0.45	0.3020	0.22
w/ identical augmentations	0.0777	0.40	0.2849	0.15
w/ smaller n/w as points aug	0.0740	0.38	0.2849	0.15

We conduct two additional experiments to further analyze the impact of the augmentation we design. Firstly, we analyze if there is a need to design simpler augmentations by replacing our novel augmentations with identical replicas of the NeRF as augmentations. In the second experiment, we use a smaller network, by reducing the number of layers from eight to four, for the points augmentation instead of using fewer positional encoding frequencies. Tab. 4 shows that in both the above cases, the performance drops to that of SimpleNeRF without the points augmentation or lower. In particular, reducing positional encoding frequencies has a higher impact than reducing the number of network layers perhaps because the network with fewer layers may still be capable of learning floaters on account of using all the positional encoding frequencies. In the supplement, we analyze how the performance of SimpleNeRF varies as l_p^{ap} varies.

5.4 Limitations

Our approach to determining reliable depth estimates for supervision depends on the reprojection error, which may be high for specular objects even if the depth estimates are correct. It may be helpful to explore approaches to determine the reliability of depth estimates without employing the reprojection error. The shape of the reprojected patches to determine reliable depth estimates may change significantly if the input viewpoints are diverse. This can lead to incorrect mask estimation. Further, our model requires accurate camera poses of the sparse input images. Finally, the use of augmented models adds to computational and memory overhead during training by about 1.5 times.

6 CONCLUSION

We address the problem of few-shot NeRF by obtaining depth supervision through simpler augmented models that are trained in tandem with the NeRF. We design two augmentations that learn simpler solutions and help the main NeRF model mitigate floater artifacts and shape-radiance ambiguity. By imposing a consistency loss between the coarse and fine NeRFs, we ensure better application of hierarchical sampling leading to sharper predictions. SimpleNeRF achieves state-of-the-art performance on two commonly used datasets in synthesizing novel views as well as estimating better scene geometry. In future, we plan to explore the role of simpler augmentations for newer models such as instant-ngp [Müller et al. 2022] and related applications such as neural surfaces [Yariv et al. 2021] to train them with sparse input views.

ACKNOWLEDGMENTS

This work was supported in part by a grant from Qualcomm. The first author was supported by the Prime Minister's Research Fellowship awarded by the Ministry of Education, Government of India. The authors would also like to thank Shankhanil Mitra for valuable discussions.

REFERENCES

- Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. 2023. NoPe-NeRF: Optimising Neural Radiance Field with No Pose Prior. (June 2023).
- Matteo Bortolon, Alessio Del Bue, and Fabio Poiesi. 2022. Data augmentation for NeRF: a geometric consistent solution based on view morphing. *arXiv e-prints*, Article arXiv:2210.04214 (2022), arXiv:2210.04214 pages. arXiv:2210.04214
- Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. 2022b. TensorRF: Tensorial Radiance Fields. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. 2021. MVSNeRF: Fast Generalizable Radiance Field Reconstruction from Multi-View Stereo. *arXiv e-prints*, Article arXiv:2103.15595 (March 2021), arXiv:2103.15595 pages. arXiv:2103.15595
- Di Chen, Yu Liu, Lianghua Huang, Bin Wang, and Pan Pan. 2022a. GeoAug: Data Augmentation for Few-Shot NeRF with Geometry Constraints. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Shenchang Eric Chen and Lance Williams. 1993. View Interpolation for Image Synthesis. In *Proceedings of the Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/166117.166153>
- Yuedong Chen, Haofei Xu, Qianyi Wu, Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. 2023. Explicit Correspondence Matching for Generalizable Neural Radiance Fields. *arXiv e-prints*, Article arXiv:2304.12294 (2023), arXiv:2304.12294 pages. arXiv:2304.12294
- Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. 2021. Stereo Radiance Fields (SRF): Learning View Synthesis for Sparse Views of Novel Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jea-Hyung Cho, Wonseok Song, Hyuk Choi, and Taejeong Kim. 2017. Hole Filling Method for Depth Image Based Rendering Based on Boundary Decision. *IEEE Signal Processing Letters (SPL)* 24, 3 (2017), 329–333.
- Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. 2022b. Depth-Supervised NeRF: Fewer Views and Faster Training for Free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Nianchen Deng, Zhenyi He, Jiannan Ye, Budmonde Duinkharjav, Praneeth Chakravarthula, Xubo Yang, and Qi Sun. 2022a. FoV-NeRF: Foveated Neural Radiance Fields for Virtual Reality. *IEEE Transactions on Visualization and Computer Graphics (TVCG)* 28, 11 (2022), 3854–3864. <https://doi.org/10.1109/TVCG.2022.3203102>
- Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. 1996. The Lumigraph. In *Proceedings of the Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/237170.237200>
- Ajay Jain, Matthew Tancik, and Pieter Abbeel. 2021. Putting NeRF on a Diet: Semantically Consistent Few-Shot View Synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. 2022. GeoNeRF: Generalizing NeRF With Geometry Priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Vijayalakshmi Kanchana, Nagabhushan Somraj, Suraj Yadwad, and Rajiv Soundararajan. 2022. Revealing Disocclusions in Temporal View Synthesis through Infilling Vector Prediction. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*.
- Zhanghan Ke, Daoye Wang, Qiong Yan, Jimmy Ren, and Rynson W.H. Lau. 2019. Dual Student: Breaking the Limits of the Teacher in Semi-Supervised Learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Mijeong Kim, Seonguk Seo, and Bohyung Han. 2022. InfoNeRF: Ray Entropy Minimization for Few-Shot Neural Volume Rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Minseop Kwak, Jiuh Song, and Seungryong Kim. 2023. GeCoNeRF: Few-shot Neural Radiance Fields via Geometric Consistency. *arXiv e-prints*, Article arXiv:2301.10941 (2023), arXiv:2301.10941 pages. arXiv:2301.10941
- Seokyeong Lee, JunYong Choi, Seungryong Kim, Ig-Jae Kim, and Junghyun Cho. 2023. ExtremeNeRF: Few-shot Neural Radiance Fields Under Unconstrained Illumination. *arXiv e-prints*, Article arXiv:2303.11728 (2023), arXiv:2303.11728 pages. arXiv:2303.11728
- Marc Levoy and Pat Hanrahan. 1996. Light Field Rendering. In *Proceedings of the Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/237170.237199>
- Kai-En Lin, Yen-Chen Lin, Wei-Sheng Lai, Tsung-Yi Lin, Yi-Chang Shih, and Ravi Ramamoorthi. 2023. Vision Transformer for NeRF-Based View Synthesis From a Single Input Image. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*.
- Li Ma, Xiaoyu Li, Jing Liao, Qi Zhang, Xuan Wang, Jue Wang, and Pedro V. Sander. 2022. Deblur-NeRF: Neural Radiance Fields From Blurry Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Willi Menapace, Aliaksandr Siarohin, Stéphane Lathuilière, Panos Achlioptas, Vladislav Golyanik, Elisa Ricci, and Sergey Tulyakov. 2023. Plotting Behind the Scenes: Towards Learnable Game Engines. *arXiv e-prints*, Article arXiv:2303.13472 (2023), arXiv:2303.13472 pages. arXiv:2303.13472
- Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P. Srinivasan, and Jonathan T. Barron. 2022. NeRF in the Dark: High Dynamic Range View Synthesis From Noisy Raw Images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. 2019. Local Light Field Fusion: Practical View Synthesis with Prescriptive Sampling Guidelines. *ACM Transactions on Graphics (TOG)* 38, 4 (July 2019), 1–14. <https://doi.org/10.1145/3306346.3322980>
- Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Transactions on Graphics (TOG)* 41, 4 (2022), 1–15.
- Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. 2022. RegNeRF: Regularizing Neural Radiance Fields for View Synthesis From Sparse Inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Naama Pearl, Tali Treibitz, and Simon Korman. 2022. NAN: Noise-Aware NeRFs for Burst-Denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. 2022. Dense Depth Priors for Neural Radiance Fields From Sparse Input Views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Seunghyeon Seo, Donghoon Han, Yeonjin Chang, and Nojun Kwak. 2023. MixNeRF: Modeling a Ray with Mixture Density for Novel View Synthesis from Sparse Inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yue Shi, Dingyi Rong, Bingbing Ni, Chang Chen, and Wenjun Zhang. 2022. GARF: Geometry-Aware Generalized Neural Radiance Field. *arXiv e-prints*, Article arXiv:2212.02280 (2022), arXiv:2212.02280 pages. arXiv:2212.02280
- Nagabhushan Somraj and Rajiv Soundararajan. 2023. ViP-NeRF: Visibility Prior for Sparse Input Neural Radiance Fields. In *Proceedings of the ACM Special Interest Group on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/3588432.3591539>
- Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P. Srinivasan, Jonathan T. Barron, and Ren Ng. 2021. Learned Initializations for Optimizing Coordinate-Based Neural Representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.



Figure 5: Qualitative examples on the RealEstate-10K dataset with two input views. While DDP-NeRF predictions contain blurred regions, ViP-NeRF predictions are color-saturated in certain regions of the door. SimpleNeRF does not suffer from these distortions and synthesizes a clean frame. For reference, we also show the input images.

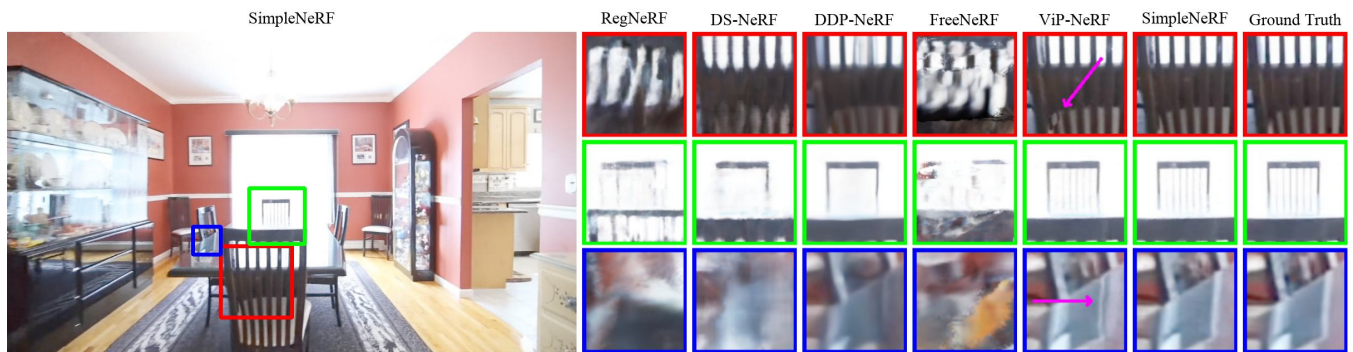


Figure 6: Qualitative examples on the RealEstate-10K dataset with four input views. We find that SimpleNeRF and ViP-NeRF perform the best among all the models. However, ViP-NeRF predictions contain minor distortions as pointed out by the magenta arrow, which is rectified by SimpleNeRF.

Alex Trevithick and Bo Yang. 2021. GRF: Learning a General Radiance Field for 3D Representation and Rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.

Mikaela Angelina Uy, Ricardo Martin-Brualla, Leonidas Guibas, and Ke Li. 2023. SCADE: NeRFs from Space Carving with Ambiguity-Aware Depth Estimates. (June 2023).

Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T. Barron, and Pratul P. Srinivasan. 2022. Ref-NeRF: Structured View-Dependent Appearance for Neural Radiance Fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. <https://doi.org/10.1109/CVPR52688.2022.00541>

Guangcong Wang, Zhaoxi Chen, Chen Change Loy, and Ziwei Liu. 2023. SparseNeRF: Distilling Depth Ranking for Few-shot Novel View Synthesis. *arXiv e-prints*, Article arXiv:2303.16196 (2023), arXiv:2303.16196 pages. arXiv:2303.16196

Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. 2021. IBRNet: Learning Multi-View Image-Based Rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing (TIP)* 13, 4 (2004), 600–612.

Jamie Wynn and Daniyar Turmukhambetov. 2023. DiffusioNeRF: Regularizing Neural Radiance Fields with Denoising Diffusion Models. *arXiv e-prints*, Article arXiv:2302.12231 (2023), arXiv:2302.12231 pages. arXiv:2302.12231

Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. 2022. SinNeRF: Training Neural Radiance Fields on Complex Scenes from a Single Image. In *Proceedings of the European Conference on Computer Vision (ECCV)*.

Jiawei Yang, Marco Pavone, and Yue Wang. 2023. FreeNeRF: Improving Few-shot Neural Rendering with Free Frequency Regularization. (June 2023).

Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. 2021. Volume Rendering of Neural Implicit Surfaces. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*.

Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021. pixelNeRF: Neural Radiance Fields From One or Few Images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yewen Ma, Rongfei Jia, and Lin Gao. 2022. NeRF-Editing: Geometry Editing of Neural Radiance Fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. 2021. NeRS: Neural Reflectance Surfaces for Sparse-view 3D Reconstruction in the Wild. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*.

Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. 2020. NeRF++: Analyzing and Improving Neural Radiance Fields. *arXiv e-prints*, Article arXiv:2010.07492 (2020), arXiv:2010.07492 pages. arXiv:2010.07492

Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. 2018. Stereo Magnification: Learning View Synthesis Using Multiplane Images. *ACM Transactions on Graphics (TOG)* 37, 4 (July 2018).

Bingfan Zhu, Yanchao Yang, Xulong Wang, Youyi Zheng, and Leonidas Guibas. 2023. VDN-NeRF: Resolving Shape-Radiance Ambiguity via View-Dependence Normalization. *arXiv e-prints*, Article arXiv:2303.17968 (2023), arXiv:2303.17968 pages. arXiv:2303.17968

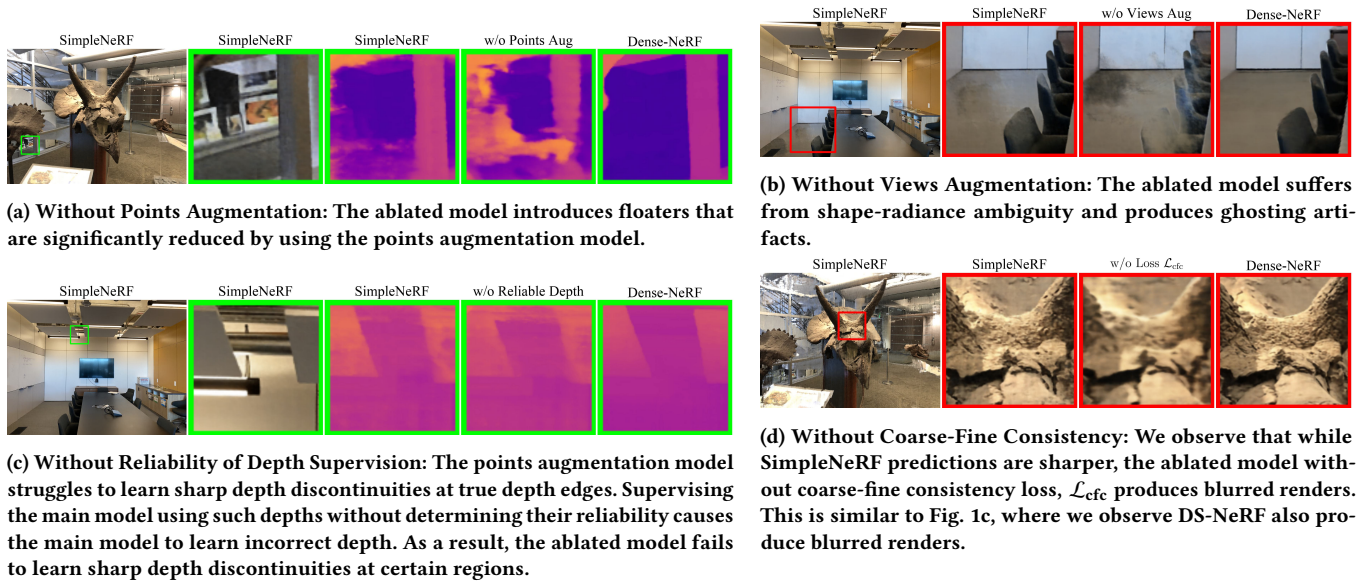


Figure 7: Qualitative examples for ablated models on the LLFF dataset with two input views. We also show the outputs of the dense-input NeRF for reference.

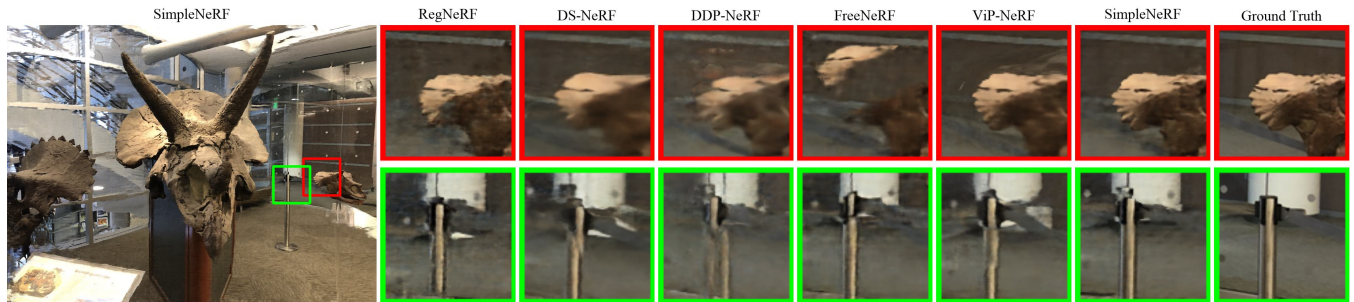


Figure 8: Qualitative examples on the LLFF dataset with two input views. DDP-NeRF and ViP-NeRF synthesize frames with broken objects in the second row and FreeNeRF breaks the object in the first row, due to incorrect depth estimations. SimpleNeRF produces sharper frames devoid of such artifacts.

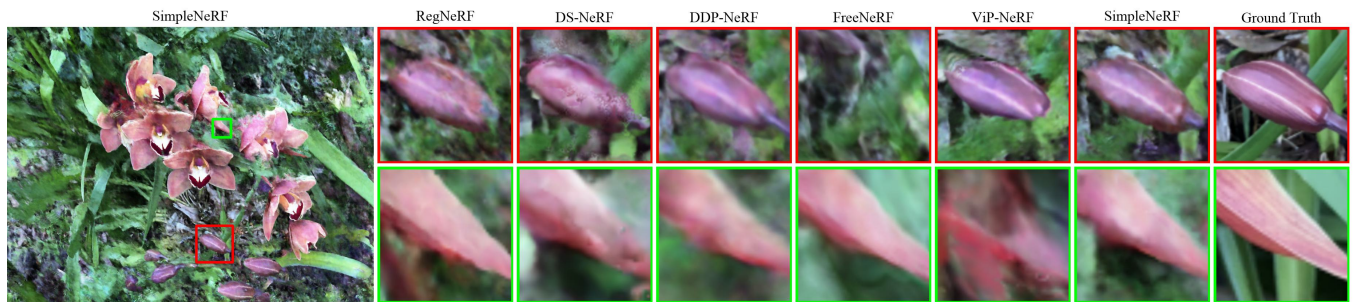


Figure 9: Qualitative examples on the LLFF dataset with three input views. In the first row, the orchid is displaced out of the cropped box in the FreeNeRF prediction, due to incorrect depth estimation. ViP-NeRF and RegNeRF fail to predict the complete orchid accurately and contain distortions at either ends. In the second row, ViP-NeRF prediction contains severe distortions. SimpleNeRF reconstructs the best among all the models in both examples.

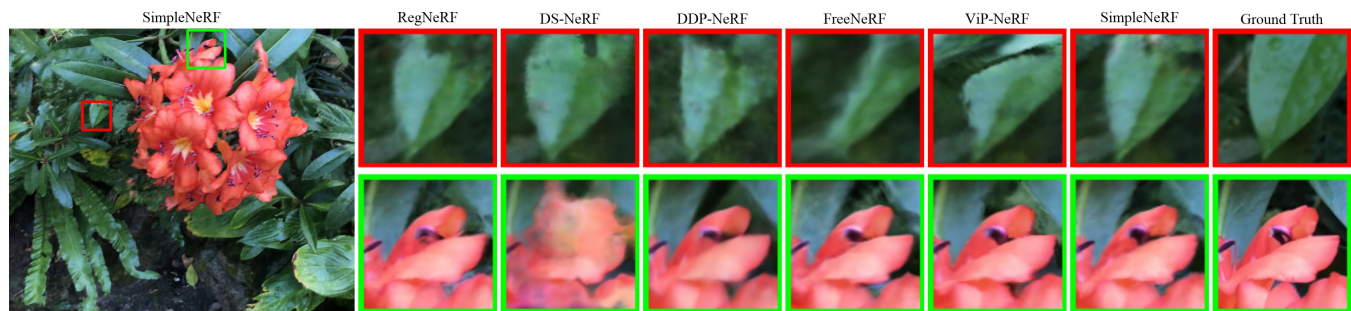


Figure 10: Qualitative examples on the LLFF dataset with four input views. In the first row, we find that ViP-NeRF, FreeNeRF, and DDP-NeRF struggle to reconstruct the shape of the leaf accurately. In the second row, DS-NeRF introduces floaters. SimpleNeRF does not suffer from such artifacts and reconstructs the shapes better.



Figure 11: Qualitative examples on the LLFF dataset with two, three, and four input views. We observe errors in depth estimation with two input views causing a change in the position of the roof. While this is corrected with three input views, there are a few shape distortions in the metal rods. With four input views, even such distortions are corrected.

Supplement

Table 5: Train and test frame numbers of RealEstate-10K dataset used in the three different settings.

No. of i/p frames	Train frame nos.	Test frame nos.
2	10, 20	5–9, 11–19, 21–25
3	10, 20, 30	5–9, 11–19, 21–29, 31–35
4	0, 10, 20, 30	1–9, 11–19, 21–29, 31–35

The contents of this supplement include

- Details on databases, evaluation, and implementations.
- Video examples on LLFF and RealEstate-10K datasets.
- Discussions on FreeNeRF, RegNeRF and ViP-NeRF.
- Additional comparisons.
- Additional analysis - positional encoding frequency, depth reliability masks and comparison with vanilla NeRF.
- Extensive quantitative evaluation reports.

A DATABASE, EVALUATION, AND IMPLEMENTATION DETAILS

A.1 Database Details

We compare the performance of sparse input NeRF models on LLFF and RealEstate-10K datasets with $n = 2, 3, 4$ input views.

LLFF [Mildenhall et al. 2019] dataset contains 8 forward-facing unbounded scenes. Each scene contains a variable number of frames at a spatial resolution of 1008×756 . Following prior art [Niemeyer et al. 2022], we use every 8th frame for testing in each scene. We uniformly sample n training views from the remaining frames.

RealEstate-10K [Zhou et al. 2018] is a dataset of videos with camera motion popularly used for novel view synthesis. From its large test set, we choose five scenes following ViP-NeRF [Somraj and Soundararajan 2023]. Each scene contains 50 temporally continuous frames at a spatial resolution of 1024×576 . We reserve frames 10, 20, 30, 0, and 40 for training and the remaining 45 frames for testing in every scene. We use the first n frames from the above list as training views for n -view NeRF. For testing, we choose all the frames between the training views that correspond to interpolation and five frames on either side that correspond to extrapolation. Tab. 5 provides exact details on train and test frames.

A.2 Evaluation Details

As mentioned in the main paper, we evaluate the model predictions only in the regions visible in the training images. We now explain our reasoning behind masking the model predictions for evaluation and then provide the details of how we compute the masks.

Recall that NeRFs are designed to memorize a scene and are therefore not equipped to predict unseen regions by design. Further, many regularization based sparse input NeRF models do not employ pre-trained prior. As a result, NeRFs are also ill-equipped to predict the depth of objects seen in only one of the input views, again by design. Thus, NeRFs require the objects to be visible in at least two views to estimate their 3D geometry accurately. Hence, we generate a mask that denotes the pixels visible in at least two input views, and evaluate the NeRF predictions in such regions only.

We generate the mask by using the depths predicted by the dense input NeRF model as follows. For every train view, we warp its depth predicted by Dense-NeRF to every other test view and compare it with the Dense-NeRF predicted depth of the test view. Intuitively, if a pixel in a test view is visible in the considered train view, then the two depths should be close to each other. Thus, we threshold the depth difference to obtain the mask. That is, if the difference in the two depths is less than a threshold, then we mark the pixel in the test view as visible in the considered train view. Our final mask for every test view is generated by marking pixels as visible if they are visible in at least two input views. We warp the depth maps using depth based reprojection similar to Cho et al. [2017]; Kanchana et al. [2022] and set the threshold to 0.05 times the maximum depth of the train view. By computing the mask using depths and not color, our approach is robust to the presence of specular objects, as long as the depth estimated by Dense-NeRF is accurate. We will release the code used to generate the masks and the masks along with the main code release.

Nonetheless, we report the performance without masking the unseen regions in Tabs. 6 to 11.

A.3 Implementation Details

We train the competing models using the official code releases by the respective authors using the configurations recommended in the code releases. Our code is developed in PyTorch and on top of DS-NeRF [Deng et al. 2022b]. We employ Adam Optimizer with an initial learning rate of $5e-4$ and exponentially decay it to $5e-6$. We adjust the weights for the different losses such that their magnitudes after scaling are of a similar order. For the first 10k iterations of the training, we impose \mathcal{L}_{color} and \mathcal{L}_{sd} only. \mathcal{L}_{ap} , \mathcal{L}_{av} and \mathcal{L}_{cfc} are imposed after 10k iterations. We set the hyper-parameters as follows: $l_p = 10$, $l_v = 4$, $l_p^{ap} = 3$, $k = 5$, $e_r = 0.1$, $\lambda_1 = 1$, and $\lambda_2 = \lambda_3 = \lambda_4 = \lambda_5 = 0.1$. The network architecture is exactly the same as DS-NeRF. For the augmented models, we only change the input dimension of the MLPs \mathcal{F}_1 and \mathcal{F}_2 appropriately. The augmented models are employed only during training, and the network is exactly the same as Vanilla NeRF for inference. We train the models on a single NVIDIA RTX A4000 16GB GPU for 100k iterations. Our model takes about 22 hours to train per scene.

B VIDEO COMPARISONS

We compare various models by rendering videos along a continuous trajectory. For the LLFF dataset, we render the videos along the spiral trajectory that is commonly used in the literature. Since RealEstate-10K is a dataset of videos, we combine the train and test frames to get the continuous trajectories.

We divide the video comparisons into two sets. In the first set, we show how our regularizations reduce the floaters and shape-radiance ambiguity by comparing the videos rendered by our model with those of the competing models and the ablated models. In the second set, we compare the videos rendered by our model with those of the competing models. The videos are available on our project website <https://nagabhushansn95.github.io/publications/2023/SimpleNeRF.html>

C DISCUSSIONS ON COMPETING MODELS

C.1 FreeNeRF

FreeNeRF [Yang et al. 2023] is a concurrent work to ours and is also designed to train NeRFs with sparse input views. FreeNeRF also employs the idea of using fewer positional encoding frequencies. However, we developed our model independently, and is also sufficiently different from FreeNeRF. Firstly, our idea is based on using simpler solutions and is more general than FreeNeRF. For example, we employ view-independent color predicting NeRF as an augmentation, unlike FreeNeRF, which uses fewer positional encoding frequencies with respect to the viewing direction too. Our approach is a stronger regularization than reducing positional encoding frequencies and is also applicable in models that may not use positional encoding [Chen et al. 2022b; Müller et al. 2022]. Secondly, FreeNeRF starts with a lower positional encoding degree and anneals it to a higher value as the training progresses. In contrast, we employ the simpler NeRF models (with fewer positional encoding frequencies or with view-independent radiance) as an augmentation and use it to obtain reliable depth supervision. Further, we employ reprojection error to filter out incorrect depth estimations by the augmented models and thus supervise using only the reliable depth predictions. As a result, in the later training stages of FreeNeRF, there is nothing constraining the model to not use the high frequency components. Our model does not suffer from this issue. Our points augmentation model is richer compared to FreeNeRF, since we include the residual positional encoded values to predict the color. This is highly useful in textured regions which are smooth in terms of depth, but contain sharp discontinuities in color as a function of position. Finally, we also include the coarse-fine consistency loss, \mathcal{L}_{cfc} , which is not employed in FreeNeRF.

We also note that FreeNeRF uses an additional regularization term to penalize the NeRF if it predicts any objects close to the camera. We believe this is a dataset bias and may not work in general. For example, many scenes from the LLFF dataset satisfy the above constraint, but this is not satisfied on the RealEstate-10K dataset. This may also be contributing to the poor performance of FreeNeRF on the RealEstate-10K dataset. An exception to the above assumption in the LLFF dataset is the fortress scene, which contains a table that extends toward the camera. Fig. 12 shows the depth map predicted by FreeNeRF for the fortress scene, which is unsurprisingly incorrect near the camera. This is further validated by the poor performance of FreeNeRF on the fortress scene as witnessed in Tab. 13. SimpleNeRF achieves a SSIM of 0.84 while FreeNeRF only achieves 0.62, with three input views on the fortress scene.

C.2 RegNeRF

RegNeRF [Niemeyer et al. 2022] is a popular few-shot NeRF model that regularizes NeRF by imposing depth smoothness in the input and various hallucinated viewpoints. Our approach of using fewer positional encoding frequencies in the points augmentation model also promotes depth smoothness. However, the depth smoothness imposed by RegNeRF is uniformly applied across all pixels which probably tries to place all the objects in a single plane. The depth smoothness promoted by SimpleNeRF is different in terms of two

aspects. Firstly, we use another network to impose smoothness which can still learn to place the objects at different depths and learn the desirable depth discontinuities. Secondly, we determine whether the depth estimated by the augmented model is accurate before using it to supervise the main NeRF model. These two reasons perhaps contribute for the superior performance of SimpleNeRF over RegNeRF.

C.3 ViP-NeRF

We note that the quantitative results in Tabs. 1 and 2 in the main paper differ from the values reported in ViP-NeRF [Somraj and Soundararajan 2023]. In the following, we describe the reasons for this difference.

- (1) In ViP-NeRF, the quality evaluation metrics are computed on full frames. However, we exclude the regions not seen in the input views. We explain the reasoning behind this in the main paper Sec. 5 and the details in Appendix A.2.
- (2) On the RealEstate-10K dataset, while we use the same train set as that of ViP-NeRF, we modify the test set as shown in Tab. 5. We explain the reasoning behind reducing the test set in Appendix A.1. More specifically, the test views that are very far away from the train views may contain large unobserved regions. Since the NeRF is not trained to reproduce such large unobserved regions, we exclude such extreme viewpoints.

D ADDITIONAL COMPARISONS

We compare with two other sparse-input NeRF models, namely, DietNeRF [Jain et al. 2021] and InfoNeRF [Kim et al. 2022]. We also evaluate the performance of SimpleNeRF without using the residual high frequency positional encodings $\gamma(\mathbf{p}_i, l_p^{\text{ap}}, l_p)$ while predicting the color in the points augmentation model. In Tabs. 6 to 11, we report the performance of all the comparisons on both datasets with 2, 3, and 4 input views in terms of all the five evaluation measures. The values in parenthesis show the unmasked evaluation scores.

E ADDITIONAL ANALYSIS

E.1 Positional Encoding Frequency

We analyze the impact of the highest positional encoding frequency l_p^{ap} , used in the points augmentation. We vary l_p^{ap} from 1 to 6 and test the performance of SimpleNeRF on the horns scene of the LLFF dataset. We show the quantitative performance in terms of LPIPS in Fig. 13. We observe only small variations in the performance as l_p^{ap} is varied. Thus, our framework is robust to the choice of l_p^{ap} . We note that using $l_p^{\text{ap}} = l_p = 10$ is equivalent to using an identical augmentation.

E.2 Depth Reliability Masks Visualization

In Fig. 14, we present visualizations that motivate the design of our augmentations, namely the points and views augmentations. We train our model without augmentations and the individual augmentations separately with only $\mathcal{L}_{\text{color}}$ and \mathcal{L}_{sd} for 100k iterations. Using the depth maps predicted by the models for an input training view, we determine the mask that indicates which depth estimates are more accurate, as explained in Sec. 4.2. For two scenes from the

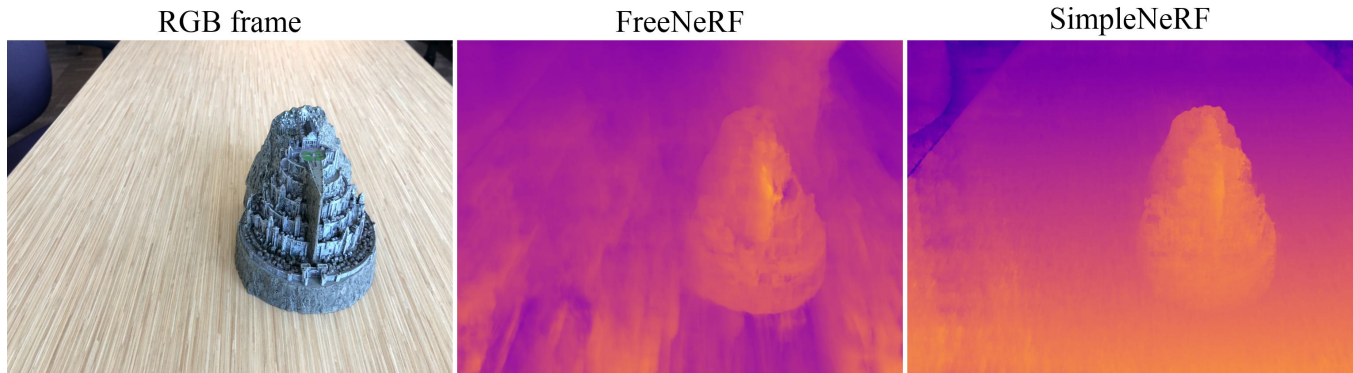


Figure 12: Comparison of depth maps predicted by FreeNeRF and SimpleNeRF on fortress scene of the LLFF dataset. Notice the distortions in the depth map estimated by FreeNeRF in the regions near the camera. SimpleNeRF does not suffer from such issues.

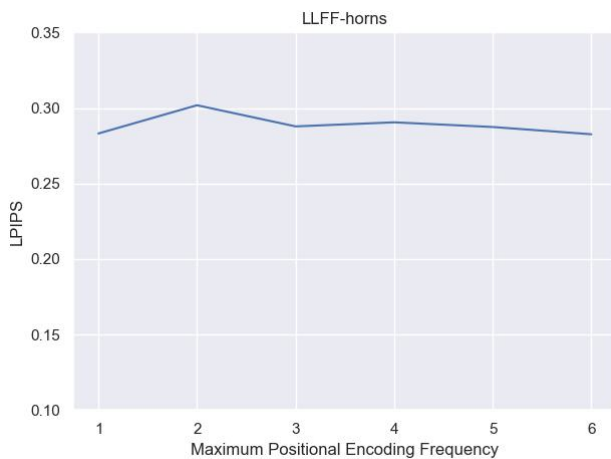


Figure 13: LPIPS scores on the horns scene as l_p^{ap} is varied.

LLFF dataset, we show an input training view and focus on a small region to visualize the corresponding masks.

We observe that the points augmented model is determined to have estimated better depths in smooth regions. At edges, the depth estimated by the main model is more accurate. Similarly, the views augmented model estimates better depth in Lambertian regions, and the main model estimates better depth in specular regions. We note that the masks shown are not the masks obtained by our final model. Since the masks are computed at every iteration, and the training of the main and augmented models are coupled, it is not possible to determine the exact locations where the augmented models help the main model learn better.

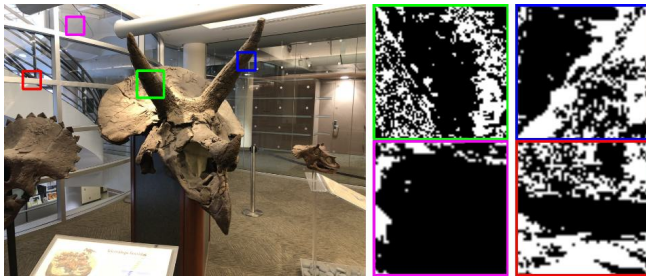
E.3 Comparison with Vanilla NeRF

We analyze the impact of our regularizations when more input views are available. Specifically, we trained SimpleNeRF with all the available input views on the room scene of the LLFF dataset. We observed that the performance dropped a little lower than the

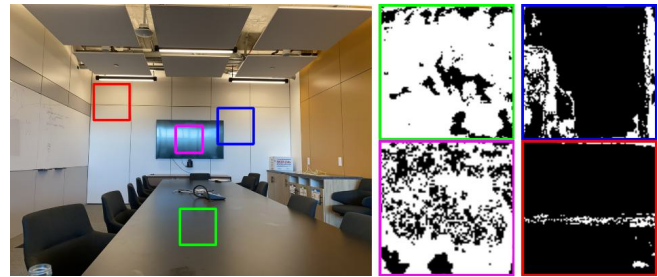
Vanilla NeRF trained with the same input views, perhaps due to the following reason. During the initial part of the training, the augmented models may be more accurate and hence may influence the learning of the main model. As the training progresses, the augmented models do not learn the scene accurately due to their respective constraints. Since the augmented models constrain the main model in terms of estimated depths, it may lead to constraining the main model from freely learning the scene. This is perhaps the reason for a drop in the performance of SimpleNeRF as compared to the Vanilla NeRF model.

F PERFORMANCE ON INDIVIDUAL SCENES

For the benefit of follow-up work, where researchers may want to analyze the performance of different models or compare the models on individual scenes, we provide the performance of various models on individual scenes in Tabs. 12 to 19.



(a) Points augmentation: The green and blue boxes focus on the two horns, where we observe that the augmented model depth is preferred in the depth-wise smooth regions on horns, and the main model depth is preferred at the edges. The magenta box focuses on a completely smooth region, so the augmented model depth is preferred for most pixels. In the red box, augmented model depth is preferred along the horizontal bar. The main model depth is preferred on either side of the bar that contains multiple depth discontinuities.



(b) Views augmentation: The green and magenta boxes focus on the TV and the table, respectively, which are highly specular in this scene (please view the supplementary videos of the room scene to observe the specularity of these objects). In these regions, the main model depth is determined to be more accurate since the main model can handle specular regions. The red and blue boxes focus on Lambertian regions of the scene where the depth estimated by the augmented model is preferred.

Figure 14: Visualizations of depth reliability mask for the two augmentations. White pixels in the mask indicate that the main model depth is determined to be more accurate at the corresponding locations. Black pixels indicate that the augmented model depth is determined to be more accurate.

Table 6: Quantitative Results on LLFF dataset with two input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
InfoNeRF	0.6024(0.7561)	0.2219(0.2095)	9.16(9.23)	0.9797(1.1000)	-0.0188(-0.0092)
DietNeRF	0.5465(0.7265)	0.3283(0.3209)	11.94(11.89)	0.8886(1.0105)	-0.0099(-0.0045)
RegNeRF	0.3056(0.4297)	0.5712(0.4885)	18.52(16.88)	-	0.7141(0.6513)
DS-NeRF	0.3106(0.4176)	0.5862(0.5074)	18.24(16.93)	0.2074(0.3372)	0.7230(0.5787)
DDP-NeRF	0.2851(0.3920)	0.6218(0.5424)	18.73(17.19)	0.2048(0.3494)	0.7480(0.6109)
FreeNeRF	0.2638(0.3760)	0.6322(0.5432)	19.52(17.55)	-	0.8066(0.7137)
ViP-NeRF	0.2768(0.3725)	0.6225(0.5230)	18.61(16.66)	0.1999(0.3413)	0.7344(0.6221)
SimpleNeRF	0.2688(0.3899)	0.6501(0.5529)	19.57(17.57)	0.1420(0.2777)	0.8480(0.7531)
SimpleNeRF w/o points aug	0.2832(0.4100)	0.6402(0.5448)	19.33(17.38)	0.1505(0.2805)	0.8334(0.7465)
SimpleNeRF w/o views aug	0.2834(0.4115)	0.6396(0.5438)	19.27(17.37)	0.1529(0.2760)	0.8306(0.7481)
SimpleNeRF w/o coarse-fine cons	0.3002(0.4240)	0.6068(0.5226)	19.02(17.25)	0.1864(0.3308)	0.8028(0.7031)
SimpleNeRF w/o reliable depth	0.3020(0.4212)	0.6012(0.5115)	18.41(16.68)	0.2186(0.3644)	0.7564(0.6724)
SimpleNeRF w/o residual pos enc	0.2837(0.4114)	0.6397(0.5436)	19.20(17.27)	0.1588(0.2890)	0.8160(0.7231)
SimpleNeRF w/ identical augs	0.2849(0.4110)	0.6379(0.5433)	19.27(17.38)	0.1509(0.2803)	0.8354(0.7466)

Table 7: Quantitative Results on LLFF dataset with three input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
InfoNeRF	0.6732(0.7679)	0.1953(0.1859)	8.38(8.52)	1.0012(1.1149)	-0.0144(-0.0176)
DietNeRF	0.6120(0.7254)	0.3405(0.3297)	11.76(11.77)	0.9093(1.0242)	-0.0598(-0.0471)
RegNeRF	0.2908(0.3602)	0.6334(0.5677)	20.22(18.65)	-	0.8238(0.7589)
DS-NeRF	0.3031(0.3641)	0.6321(0.5774)	20.20(18.97)	0.1787(0.2699)	0.7852(0.7173)
DDP-NeRF	0.3250(0.3869)	0.6152(0.5628)	18.73(17.71)	0.1941(0.3032)	0.7433(0.6707)
FreeNeRF	0.2754(0.3415)	0.6583(0.5960)	20.93(19.30)	-	0.8379(0.7656)
ViP-NeRF	0.2798(0.3365)	0.6548(0.5907)	20.54(18.89)	0.1721(0.2795)	0.7891(0.7082)
SimpleNeRF	0.2559(0.3259)	0.6940(0.6222)	21.37(19.47)	0.1199(0.2201)	0.8935(0.8153)

Table 8: Quantitative Results on LLFF dataset with four input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
InfoNeRF	0.6985(0.7701)	0.2270(0.2188)	9.18(9.25)	1.0411(1.1119)	-0.0394(-0.0390)
DietNeRF	0.6506(0.7396)	0.3496(0.3404)	11.86(11.84)	0.9546(1.0259)	-0.0368(-0.0249)
RegNeRF	0.2794(0.3227)	0.6645(0.6159)	21.32(19.89)	-	0.8933(0.8528)
DS-NeRF	0.2979(0.3376)	0.6582(0.6135)	21.23(20.07)	0.1451(0.2097)	0.8506(0.8130)
DDP-NeRF	0.3042(0.3467)	0.6558(0.6121)	20.17(19.19)	0.1704(0.2487)	0.8322(0.7664)
FreeNeRF	0.2848(0.3280)	0.6764(0.6303)	21.91(20.45)	-	0.9091(0.8626)
ViP-NeRF	0.2854(0.3203)	0.6675(0.6182)	20.75(19.34)	0.1555(0.2316)	0.8622(0.8070)
SimpleNeRF	0.2633(0.3083)	0.7016(0.6521)	21.99(20.44)	0.1110(0.1741)	0.9355(0.8952)

Table 9: Quantitative Results on RealEstate-10K dataset with two input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
InfoNeRF	0.5924(0.6384)	0.4342(0.4343)	12.27(12.17)	2.1793(2.2314)	0.0912(0.0942)
DietNeRF	0.4381(0.4862)	0.6534(0.6520)	18.06(17.83)	2.0247(2.0807)	0.2339(0.2452)
RegNeRF	0.4129(0.4483)	0.5916(0.5864)	17.14(16.87)	–	0.1118(0.1117)
DS-NeRF	0.2709(0.3171)	0.7983(0.7859)	26.26(25.44)	0.7164(0.8323)	0.6660(0.6433)
DDP-NeRF	0.1290(0.1481)	0.8640(0.8502)	27.79(26.15)	0.4831(0.5944)	0.7921(0.7663)
FreeNeRF	0.5036(0.5471)	0.5354(0.5336)	14.70(14.50)	–	-0.1937(-0.1813)
ViP-NeRF	0.0687(0.0783)	0.8889(0.8717)	32.32(29.55)	0.3856(0.5337)	0.8446(0.7851)
SimpleNeRF	0.0635(0.0745)	0.8942(0.8783)	33.10(30.30)	0.3269(0.4584)	0.9215(0.8781)
SimpleNeRF w/o points aug	0.0752(0.0889)	0.8886(0.8722)	32.54(29.85)	0.3795(0.5109)	0.8973(0.8487)
SimpleNeRF w/o views aug	0.0790(0.0925)	0.8884(0.8714)	32.13(29.62)	0.3870(0.5110)	0.8837(0.8428)
SimpleNeRF w/o coarse-fine cons	0.0740(0.0865)	0.8869(0.8693)	31.86(29.47)	0.4223(0.5526)	0.8576(0.8140)
SimpleNeRF w/o reliable depth	0.0687(0.0802)	0.8913(0.8754)	32.63(30.33)	0.4485(0.5778)	0.8729(0.8404)
SimpleNeRF w/o residual pos enc	0.0790(0.0909)	0.8875(0.8715)	32.00(29.91)	0.4040(0.5255)	0.8838(0.8392)
SimpleNeRF w/ identical augs	0.0777(0.0916)	0.8875(0.8713)	32.31(29.85)	0.4037(0.5269)	0.8949(0.8453)

Table 10: Quantitative Results on RealEstate-10K dataset with three input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
InfoNeRF	0.6561(0.6846)	0.3792(0.3780)	10.57(10.57)	2.2198(2.2830)	0.1929(0.1994)
DietNeRF	0.4636(0.4886)	0.6456(0.6445)	18.01(17.89)	2.0355(2.1023)	0.0240(0.0438)
RegNeRF	0.4171(0.4362)	0.6132(0.6078)	17.86(17.73)	–	0.0574(0.0475)
DS-NeRF	0.2893(0.3211)	0.8004(0.7905)	26.50(25.94)	0.5400(0.6524)	0.8106(0.7910)
DDP-NeRF	0.1518(0.1601)	0.8587(0.8518)	26.67(25.92)	0.4139(0.5222)	0.8612(0.8331)
FreeNeRF	0.5146(0.5414)	0.5708(0.5675)	15.26(15.12)	–	-0.2590(-0.2445)
ViP-NeRF	0.0758(0.0832)	0.8967(0.8852)	31.93(30.27)	0.3365(0.4683)	0.9009(0.8558)
SimpleNeRF	0.0726(0.0829)	0.8984(0.8879)	33.21(31.40)	0.2770(0.3885)	0.9266(0.8931)

Table 11: Quantitative Results on RealEstate-10K dataset with four input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
InfoNeRF	0.6651(0.6721)	0.3843(0.3830)	10.62(10.59)	2.1874(2.2742)	0.2549(0.2594)
DietNeRF	0.4853(0.4954)	0.6503(0.6475)	18.01(17.89)	2.0398(2.1273)	0.0990(0.1011)
RegNeRF	0.4316(0.4383)	0.6257(0.6198)	18.34(18.25)	–	0.1422(0.1396)
DS-NeRF	0.3103(0.3287)	0.7999(0.7920)	26.65(26.28)	0.5154(0.6171)	0.8145(0.8018)
DDP-NeRF	0.1563(0.1584)	0.8617(0.8557)	27.07(26.48)	0.3832(0.4813)	0.8739(0.8605)
FreeNeRF	0.5226(0.5323)	0.6027(0.5989)	16.31(16.25)	–	-0.2152(-0.2162)
ViP-NeRF	0.0892(0.0909)	0.8968(0.8894)	31.95(30.83)	0.3658(0.4761)	0.8414(0.8080)
SimpleNeRF	0.0847(0.0891)	0.8987(0.8917)	32.88(31.73)	0.2692(0.3565)	0.9209(0.9035)

Table 12: Per-scene performance of various models with two input views on LLFF dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	Trex	Average
InfoNeRF	0.66(0.78)	0.57(0.69)	0.65(0.82)	0.58(0.76)	0.45(0.64)	0.60(0.69)	0.68(0.80)	0.61(0.79)	0.60(0.76)
	0.26(0.26)	0.21(0.19)	0.18(0.18)	0.20(0.20)	0.15(0.10)	0.17(0.11)	0.36(0.38)	0.23(0.21)	0.22(0.21)
	10.7(11.0)	10.8(11.0)	6.6(6.5)	8.9(8.9)	9.7(9.5)	9.1(9.4)	10.4(10.8)	8.5(8.4)	9.2(9.2)
	0.96(1.05)	0.76(0.89)	1.18(1.41)	1.08(1.26)	0.85(0.92)	0.87(1.12)	1.05(1.06)	0.94(0.95)	0.98(1.10)
	-0.00(-0.04)	-0.25(-0.18)	0.09(0.06)	0.07(0.08)	-0.14(-0.12)	0.00(0.01)	-0.01(0.00)	-0.00(0.01)	-0.02(-0.01)
DietNeRF	0.63(0.77)	0.54(0.69)	0.50(0.66)	0.54(0.76)	0.43(0.67)	0.57(0.73)	0.62(0.77)	0.54(0.75)	0.55(0.73)
	0.30(0.29)	0.29(0.26)	0.44(0.44)	0.27(0.28)	0.18(0.12)	0.21(0.15)	0.49(0.52)	0.35(0.37)	0.33(0.32)
	12.1(12.3)	12.1(12.2)	15.2(14.2)	10.5(10.7)	10.8(10.6)	10.5(10.6)	13.0(13.1)	11.1(11.3)	11.9(11.9)
	0.83(0.92)	0.70(0.84)	0.94(1.17)	1.02(1.20)	0.83(0.90)	0.78(1.03)	0.95(0.97)	0.89(0.90)	0.89(1.01)
	-0.02(-0.05)	-0.02(0.00)	0.13(0.13)	-0.01(-0.01)	-0.13(-0.10)	0.02(-0.01)	0.02(0.03)	-0.09(-0.07)	-0.01(-0.00)
RegNeRF	0.42(0.51)	0.28(0.43)	0.28(0.37)	0.34(0.51)	0.19(0.35)	0.38(0.45)	0.30(0.38)	0.30(0.42)	0.31(0.43)
	0.50(0.45)	0.62(0.51)	0.47(0.46)	0.49(0.42)	0.54(0.37)	0.45(0.30)	0.81(0.74)	0.64(0.54)	0.57(0.49)
	16.1(15.8)	19.9(17.0)	21.2(20.6)	17.5(15.9)	17.4(14.5)	15.1(13.9)	21.1(18.7)	17.8(16.7)	18.5(16.9)
	-	-	-	-	-	-	-	-	-
	0.70(0.64)	0.88(0.68)	0.58(0.61)	0.67(0.67)	0.82(0.60)	0.74(0.60)	0.76(0.69)	0.65(0.68)	0.71(0.65)
DS-NeRF	0.41(0.50)	0.32(0.43)	0.21(0.30)	0.34(0.49)	0.29(0.47)	0.38(0.43)	0.27(0.35)	0.31(0.41)	0.31(0.42)
	0.50(0.46)	0.52(0.44)	0.68(0.65)	0.56(0.49)	0.35(0.24)	0.45(0.32)	0.82(0.76)	0.64(0.53)	0.59(0.51)
	16.7(16.4)	17.3(16.1)	23.6(23.0)	18.1(16.6)	13.3(12.4)	14.7(13.7)	21.3(18.9)	17.4(15.7)	18.2(16.9)
	0.18(0.28)	0.30(0.43)	0.04(0.19)	0.21(0.35)	0.54(0.60)	0.21(0.38)	0.14(0.24)	0.15(0.31)	0.21(0.34)
	0.56(0.37)	0.71(0.55)	0.99(0.99)	0.76(0.68)	0.33(0.18)	0.75(0.63)	0.76(0.73)	0.71(0.29)	0.72(0.58)
DDP-NeRF	0.35(0.44)	0.33(0.46)	0.12(0.17)	0.30(0.46)	0.31(0.52)	0.33(0.41)	0.25(0.30)	0.33(0.43)	0.29(0.39)
	0.55(0.49)	0.53(0.45)	0.80(0.77)	0.60(0.52)	0.34(0.23)	0.53(0.38)	0.82(0.76)	0.64(0.54)	0.62(0.54)
	17.8(17.2)	17.3(16.2)	23.4(22.7)	19.3(17.1)	13.5(12.6)	16.6(15.1)	21.6(18.7)	17.2(15.7)	18.7(17.2)
	0.13(0.23)	0.33(0.46)	0.06(0.27)	0.21(0.39)	0.58(0.64)	0.15(0.32)	0.10(0.19)	0.18(0.32)	0.20(0.35)
	0.72(0.56)	0.60(0.44)	0.99(0.98)	0.85(0.71)	0.25(0.08)	0.85(0.74)	0.94(0.92)	0.60(0.29)	0.75(0.61)
FreeNeRF	0.36(0.46)	0.24(0.38)	0.25(0.33)	0.27(0.43)	0.19(0.36)	0.35(0.42)	0.26(0.34)	0.24(0.33)	0.26(0.38)
	0.55(0.49)	0.66(0.55)	0.55(0.53)	0.62(0.53)	0.56(0.38)	0.51(0.35)	0.81(0.76)	0.70(0.60)	0.63(0.54)
	17.7(17.1)	20.6(17.6)	22.0(21.3)	19.5(17.1)	17.9(14.4)	15.6(14.1)	20.4(18.3)	19.8(18.1)	19.5(17.6)
	-	-	-	-	-	-	-	-	-
	0.64(0.62)	0.97(0.80)	0.71(0.71)	0.86(0.74)	0.82(0.55)	0.76(0.61)	0.83(0.80)	0.78(0.74)	0.81(0.71)
ViP-NeRF	0.37(0.45)	0.31(0.42)	0.15(0.21)	0.25(0.39)	0.29(0.46)	0.34(0.40)	0.29(0.36)	0.29(0.38)	0.28(0.37)
	0.51(0.45)	0.53(0.43)	0.77(0.71)	0.65(0.54)	0.33(0.21)	0.53(0.36)	0.80(0.72)	0.66(0.54)	0.62(0.52)
	16.7(16.2)	16.2(14.9)	24.6(22.6)	19.9(17.1)	12.5(11.7)	15.8(14.2)	21.0(17.7)	17.5(15.9)	18.6(16.7)
	0.17(0.26)	0.27(0.42)	0.04(0.24)	0.19(0.37)	0.61(0.68)	0.16(0.35)	0.14(0.25)	0.15(0.25)	0.20(0.34)
	0.51(0.41)	0.70(0.45)	0.99(0.99)	0.86(0.73)	0.21(0.07)	0.81(0.70)	0.84(0.76)	0.65(0.54)	0.73(0.62)
SimpleNeRF	0.39(0.51)	0.28(0.43)	0.16(0.25)	0.26(0.42)	0.25(0.44)	0.33(0.41)	0.27(0.35)	0.28(0.39)	0.27(0.39)
	0.54(0.50)	0.65(0.53)	0.72(0.67)	0.64(0.54)	0.45(0.30)	0.52(0.37)	0.83(0.77)	0.68(0.58)	0.65(0.55)
	17.4(17.0)	20.1(16.9)	23.8(22.5)	19.6(17.1)	15.2(13.5)	16.3(14.7)	22.5(19.5)	18.3(16.8)	19.6(17.6)
	0.14(0.21)	0.10(0.25)	0.05(0.25)	0.18(0.35)	0.31(0.44)	0.15(0.34)	0.12(0.21)	0.13(0.19)	0.14(0.28)
	0.72(0.68)	0.94(0.72)	0.99(0.99)	0.88(0.76)	0.64(0.38)	0.85(0.74)	0.86(0.84)	0.78(0.73)	0.85(0.75)

Table 13: Per-scene performance of various models with three input views on LLFF dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	Trex	Average
InfoNeRF	0.74(0.86)	0.62(0.69)	0.76(0.83)	0.63(0.76)	0.50(0.65)	0.68(0.71)	0.75(0.83)	0.69(0.79)	0.67(0.77)
	0.22(0.21)	0.24(0.23)	0.16(0.17)	0.18(0.18)	0.16(0.11)	0.12(0.08)	0.29(0.29)	0.19(0.18)	0.20(0.19)
	7.5(7.5)	10.3(10.8)	5.1(5.2)	8.9(8.9)	10.2(9.9)	7.9(8.3)	8.7(8.9)	8.6(8.7)	8.4(8.5)
	0.99(1.08)	0.84(0.94)	1.21(1.42)	1.07(1.26)	0.87(0.92)	0.95(1.15)	1.04(1.07)	0.93(0.96)	1.00(1.11)
	0.01(-0.07)	-0.18(-0.11)	0.17(0.17)	0.04(0.05)	-0.12(-0.10)	0.07(0.08)	-0.11(-0.15)	-0.03(-0.06)	-0.01(-0.02)
DietNeRF	0.68(0.80)	0.58(0.67)	0.63(0.70)	0.60(0.75)	0.49(0.67)	0.65(0.75)	0.68(0.76)	0.60(0.72)	0.61(0.73)
	0.33(0.32)	0.33(0.31)	0.40(0.41)	0.28(0.29)	0.18(0.13)	0.19(0.15)	0.54(0.55)	0.36(0.35)	0.34(0.33)
	11.9(12.0)	13.2(13.1)	12.6(12.4)	10.9(11.1)	10.8(10.6)	10.0(10.1)	12.6(12.6)	11.8(11.9)	11.8(11.8)
	0.86(0.95)	0.73(0.84)	1.00(1.22)	1.03(1.21)	0.85(0.90)	0.86(1.06)	0.95(0.97)	0.88(0.91)	0.91(1.02)
	-0.06(-0.03)	-0.16(-0.15)	-0.08(-0.09)	0.03(0.06)	-0.03(-0.03)	-0.07(-0.08)	0.03(0.09)	-0.16(-0.17)	-0.06(-0.05)
RegNeRF	0.40(0.47)	0.22(0.27)	0.26(0.31)	0.33(0.44)	0.26(0.39)	0.42(0.44)	0.22(0.25)	0.28(0.36)	0.29(0.36)
	0.53(0.48)	0.67(0.58)	0.66(0.64)	0.59(0.53)	0.48(0.37)	0.41(0.31)	0.84(0.81)	0.70(0.63)	0.63(0.57)
	18.3(17.9)	21.5(19.6)	24.6(22.7)	20.2(18.2)	16.6(14.6)	15.1(14.2)	22.1(21.0)	19.7(18.4)	20.2(18.7)
	-	-	-	-	-	-	-	-	-
	0.59(0.59)	0.80(0.64)	0.92(0.90)	0.84(0.81)	0.79(0.66)	0.70(0.61)	0.89(0.86)	0.87(0.79)	0.82(0.76)
DS-NeRF	0.40(0.47)	0.22(0.25)	0.21(0.25)	0.37(0.47)	0.37(0.50)	0.43(0.45)	0.19(0.22)	0.31(0.37)	0.30(0.36)
	0.56(0.52)	0.72(0.66)	0.74(0.72)	0.58(0.52)	0.33(0.25)	0.43(0.33)	0.87(0.84)	0.65(0.59)	0.63(0.58)
	19.0(18.5)	22.8(21.3)	26.3(24.8)	19.0(17.5)	13.5(12.6)	14.9(14.1)	24.5(23.0)	18.1(17.1)	20.2(19.0)
	0.15(0.22)	0.08(0.17)	0.07(0.20)	0.19(0.31)	0.52(0.57)	0.22(0.34)	0.10(0.17)	0.20(0.25)	0.18(0.27)
	0.58(0.51)	0.93(0.81)	1.00(1.00)	0.87(0.77)	0.29(0.17)	0.76(0.70)	0.92(0.91)	0.67(0.60)	0.79(0.72)
DDP-NeRF	0.38(0.47)	0.25(0.29)	0.18(0.20)	0.38(0.48)	0.37(0.52)	0.41(0.45)	0.31(0.32)	0.36(0.42)	0.33(0.39)
	0.58(0.53)	0.70(0.63)	0.78(0.75)	0.58(0.53)	0.32(0.24)	0.46(0.35)	0.78(0.76)	0.59(0.54)	0.62(0.56)
	19.4(18.5)	21.3(20.2)	22.7(22.1)	19.3(17.4)	13.5(12.8)	16.1(15.1)	19.3(18.3)	16.5(16.0)	18.7(17.7)
	0.13(0.21)	0.10(0.19)	0.09(0.27)	0.18(0.36)	0.59(0.63)	0.17(0.32)	0.12(0.19)	0.24(0.30)	0.19(0.30)
	0.67(0.54)	0.94(0.82)	0.99(0.99)	0.88(0.76)	0.07(0.01)	0.83(0.75)	0.92(0.92)	0.45(0.36)	0.74(0.67)
FreeNeRF	0.34(0.40)	0.24(0.28)	0.28(0.32)	0.30(0.41)	0.26(0.40)	0.39(0.41)	0.19(0.22)	0.26(0.33)	0.28(0.34)
	0.59(0.54)	0.68(0.61)	0.62(0.60)	0.65(0.58)	0.51(0.40)	0.48(0.37)	0.88(0.85)	0.71(0.64)	0.66(0.60)
	19.6(18.9)	22.2(20.7)	23.4(22.0)	21.3(18.7)	17.3(15.0)	15.9(14.7)	24.2(22.6)	20.2(19.0)	20.9(19.3)
	-	-	-	-	-	-	-	-	-
	0.76(0.72)	0.87(0.75)	0.67(0.64)	0.91(0.80)	0.84(0.70)	0.81(0.73)	0.88(0.88)	0.89(0.83)	0.84(0.77)
ViP-NeRF	0.43(0.51)	0.21(0.24)	0.16(0.19)	0.32(0.42)	0.33(0.44)	0.40(0.41)	0.24(0.27)	0.26(0.32)	0.28(0.34)
	0.53(0.49)	0.73(0.65)	0.80(0.76)	0.64(0.57)	0.34(0.25)	0.45(0.34)	0.84(0.81)	0.69(0.62)	0.65(0.59)
	17.9(17.3)	22.6(20.8)	27.0(24.5)	20.7(18.2)	13.5(12.4)	15.2(14.2)	23.1(21.7)	19.4(18.1)	20.5(18.9)
	0.21(0.30)	0.10(0.20)	0.06(0.22)	0.16(0.32)	0.54(0.59)	0.20(0.34)	0.12(0.20)	0.14(0.19)	0.17(0.28)
	0.31(0.22)	0.93(0.79)	1.00(0.99)	0.89(0.78)	0.20(0.11)	0.79(0.73)	0.91(0.89)	0.83(0.71)	0.79(0.71)
SimpleNeRF	0.37(0.43)	0.19(0.24)	0.12(0.17)	0.31(0.42)	0.28(0.42)	0.37(0.39)	0.22(0.26)	0.26(0.34)	0.26(0.33)
	0.57(0.52)	0.75(0.66)	0.82(0.78)	0.65(0.57)	0.50(0.38)	0.50(0.38)	0.86(0.83)	0.74(0.66)	0.69(0.62)
	18.8(18.2)	23.1(20.7)	27.5(24.7)	20.6(18.4)	17.0(14.8)	16.2(15.0)	23.6(22.0)	20.4(18.9)	21.4(19.5)
	0.14(0.22)	0.09(0.19)	0.05(0.19)	0.16(0.31)	0.20(0.27)	0.15(0.29)	0.11(0.17)	0.09(0.14)	0.12(0.22)
	0.68(0.62)	0.94(0.79)	1.00(0.99)	0.89(0.78)	0.85(0.71)	0.87(0.81)	0.88(0.83)	0.93(0.86)	0.89(0.82)

Table 14: Per-scene performance of various models with four input views on LLFF dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	Trex	Average
InfoNeRF	0.71(0.80)	0.63(0.68)	0.84(0.87)	0.68(0.77)	0.58(0.66)	0.72(0.76)	0.72(0.78)	0.69(0.80)	0.70(0.77)
	0.20(0.19)	0.24(0.23)	0.16(0.16)	0.19(0.19)	0.12(0.09)	0.11(0.09)	0.45(0.45)	0.26(0.26)	0.23(0.22)
	9.6(9.8)	11.4(11.5)	4.7(4.7)	8.9(8.8)	9.4(9.3)	7.8(8.1)	11.9(11.9)	9.9(10.1)	9.2(9.2)
	0.95(1.03)	0.83(0.91)	1.29(1.41)	1.18(1.27)	0.89(0.92)	1.00(1.18)	1.05(1.06)	0.97(0.97)	1.04(1.11)
	0.02(0.00)	-0.27(-0.21)	0.09(0.08)	0.08(0.09)	-0.12(-0.11)	0.03(0.04)	-0.06(-0.06)	-0.13(-0.17)	-0.04(-0.04)
DietNeRF	0.66(0.80)	0.62(0.69)	0.67(0.70)	0.66(0.78)	0.57(0.70)	0.69(0.77)	0.68(0.73)	0.64(0.75)	0.65(0.74)
	0.36(0.35)	0.31(0.29)	0.42(0.42)	0.30(0.30)	0.18(0.14)	0.20(0.16)	0.59(0.59)	0.34(0.35)	0.35(0.34)
	12.7(12.9)	12.6(12.6)	12.8(12.6)	10.9(10.8)	11.0(10.8)	10.2(10.1)	14.0(13.9)	10.9(11.2)	11.9(11.8)
	0.85(0.93)	0.76(0.84)	1.07(1.19)	1.13(1.22)	0.87(0.90)	0.88(1.06)	0.98(0.99)	0.91(0.92)	0.95(1.03)
	-0.11(-0.08)	-0.05(-0.06)	-0.04(-0.04)	-0.00(0.02)	-0.12(-0.10)	-0.02(-0.02)	0.08(0.12)	-0.10(-0.10)	-0.04(-0.02)
RegNeRF	0.28(0.35)	0.24(0.29)	0.35(0.37)	0.27(0.34)	0.26(0.32)	0.42(0.43)	0.17(0.19)	0.28(0.32)	0.28(0.32)
	0.67(0.63)	0.71(0.64)	0.56(0.55)	0.69(0.64)	0.52(0.44)	0.43(0.34)	0.89(0.87)	0.71(0.66)	0.66(0.62)
	21.6(20.8)	22.3(19.8)	23.0(22.4)	22.3(20.1)	17.3(15.9)	15.6(14.8)	25.6(23.9)	19.7(18.9)	21.3(19.9)
	-	-	-	-	-	-	-	-	-
	0.90(0.87)	0.94(0.83)	0.94(0.94)	0.93(0.91)	0.83(0.75)	0.78(0.71)	0.93(0.90)	0.85(0.82)	0.89(0.85)
DS-NeRF	0.28(0.35)	0.24(0.28)	0.29(0.31)	0.35(0.41)	0.34(0.41)	0.41(0.41)	0.14(0.16)	0.34(0.39)	0.30(0.34)
	0.67(0.63)	0.70(0.64)	0.67(0.66)	0.64(0.59)	0.47(0.39)	0.47(0.38)	0.91(0.89)	0.63(0.59)	0.66(0.61)
	21.5(20.9)	22.6(20.6)	24.9(24.1)	21.0(19.5)	16.9(15.8)	16.0(15.2)	27.2(25.6)	17.6(17.1)	21.2(20.1)
	0.07(0.12)	0.08(0.15)	0.12(0.21)	0.15(0.21)	0.25(0.30)	0.16(0.28)	0.09(0.14)	0.21(0.25)	0.15(0.21)
	0.88(0.83)	0.95(0.86)	1.00(0.99)	0.93(0.90)	0.77(0.69)	0.86(0.82)	0.96(0.96)	0.49(0.46)	0.85(0.81)
DDP-NeRF	0.32(0.40)	0.26(0.30)	0.18(0.18)	0.36(0.42)	0.35(0.45)	0.41(0.42)	0.25(0.26)	0.35(0.39)	0.30(0.35)
	0.65(0.60)	0.69(0.63)	0.74(0.73)	0.63(0.59)	0.44(0.37)	0.50(0.41)	0.84(0.82)	0.64(0.60)	0.66(0.61)
	21.1(20.1)	21.5(20.0)	23.9(23.4)	20.6(19.3)	16.1(15.1)	16.8(15.8)	22.0(20.8)	17.9(17.3)	20.2(19.2)
	0.11(0.19)	0.09(0.17)	0.16(0.26)	0.19(0.27)	0.29(0.34)	0.16(0.30)	0.14(0.19)	0.21(0.26)	0.17(0.25)
	0.77(0.58)	0.94(0.84)	0.99(0.99)	0.93(0.88)	0.66(0.57)	0.85(0.79)	0.92(0.91)	0.55(0.45)	0.83(0.77)
FreeNeRF	0.29(0.37)	0.25(0.30)	0.34(0.35)	0.30(0.37)	0.27(0.35)	0.41(0.42)	0.16(0.19)	0.28(0.31)	0.28(0.33)
	0.68(0.64)	0.70(0.64)	0.60(0.60)	0.68(0.63)	0.54(0.47)	0.46(0.37)	0.90(0.88)	0.72(0.68)	0.68(0.63)
	21.8(21.1)	22.9(20.5)	23.8(23.2)	22.4(20.4)	18.1(16.6)	15.7(14.9)	27.0(24.8)	20.5(19.6)	21.9(20.5)
	-	-	-	-	-	-	-	-	-
	0.89(0.86)	0.95(0.85)	0.95(0.95)	0.95(0.89)	0.87(0.80)	0.77(0.69)	0.91(0.89)	0.90(0.87)	0.91(0.86)
ViP-NeRF	0.33(0.39)	0.24(0.27)	0.24(0.25)	0.33(0.38)	0.30(0.36)	0.40(0.40)	0.21(0.23)	0.29(0.32)	0.29(0.32)
	0.63(0.58)	0.70(0.63)	0.71(0.70)	0.65(0.60)	0.47(0.40)	0.48(0.39)	0.87(0.85)	0.69(0.64)	0.67(0.62)
	19.3(18.2)	21.8(19.5)	24.3(23.3)	20.7(19.0)	16.0(14.8)	15.8(14.8)	24.9(23.2)	19.5(18.6)	20.7(19.3)
	0.12(0.19)	0.09(0.18)	0.16(0.27)	0.20(0.28)	0.28(0.33)	0.18(0.32)	0.11(0.14)	0.12(0.16)	0.16(0.23)
	0.70(0.60)	0.93(0.81)	0.99(0.98)	0.90(0.85)	0.65(0.55)	0.82(0.77)	0.90(0.90)	0.85(0.79)	0.86(0.81)
SimpleNeRF	0.27(0.33)	0.21(0.27)	0.26(0.28)	0.31(0.38)	0.27(0.35)	0.36(0.36)	0.17(0.19)	0.27(0.32)	0.26(0.31)
	0.69(0.65)	0.74(0.67)	0.70(0.69)	0.68(0.63)	0.54(0.46)	0.51(0.42)	0.90(0.88)	0.74(0.68)	0.70(0.65)
	21.9(21.1)	23.2(20.8)	25.4(24.3)	21.7(19.7)	17.7(16.3)	16.8(15.7)	26.3(24.3)	20.4(19.3)	22.0(20.4)
	0.07(0.12)	0.07(0.13)	0.10(0.18)	0.16(0.24)	0.17(0.21)	0.14(0.27)	0.07(0.10)	0.10(0.13)	0.11(0.17)
	0.90(0.87)	0.97(0.89)	1.00(0.99)	0.94(0.89)	0.85(0.78)	0.90(0.85)	0.97(0.96)	0.91(0.88)	0.94(0.90)

Table 15: Per-scene performance of ablated models with two input views on LLFF dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	Trex	Average
SimpleNeRF	0.39(0.51)	0.28(0.43)	0.16(0.25)	0.26(0.42)	0.25(0.44)	0.33(0.41)	0.27(0.35)	0.28(0.39)	0.27(0.39)
	0.54(0.50)	0.65(0.53)	0.72(0.67)	0.64(0.54)	0.45(0.30)	0.52(0.37)	0.83(0.77)	0.68(0.58)	0.65(0.55)
	17.4(17.0)	20.1(16.9)	23.8(22.5)	19.6(17.1)	15.2(13.5)	16.3(14.7)	22.5(19.5)	18.3(16.8)	19.6(17.6)
	0.14(0.21)	0.10(0.25)	0.05(0.25)	0.18(0.35)	0.31(0.44)	0.15(0.34)	0.12(0.21)	0.13(0.19)	0.14(0.28)
	0.72(0.68)	0.94(0.72)	0.99(0.99)	0.88(0.76)	0.64(0.38)	0.85(0.74)	0.86(0.84)	0.78(0.73)	0.85(0.75)
SimpleNeRF w/o Points Augmentation	0.40(0.52)	0.27(0.43)	0.17(0.27)	0.29(0.46)	0.25(0.44)	0.36(0.43)	0.31(0.39)	0.29(0.40)	0.28(0.41)
	0.54(0.49)	0.64(0.54)	0.71(0.65)	0.62(0.52)	0.47(0.32)	0.50(0.35)	0.81(0.74)	0.68(0.58)	0.64(0.54)
	17.4(17.0)	20.4(17.6)	22.9(21.5)	19.7(17.2)	16.3(13.9)	15.9(14.5)	21.3(18.5)	18.0(16.7)	19.3(17.4)
	0.14(0.21)	0.11(0.24)	0.06(0.25)	0.20(0.37)	0.28(0.42)	0.16(0.35)	0.14(0.23)	0.14(0.18)	0.15(0.28)
	0.68(0.64)	0.95(0.77)	0.99(0.98)	0.84(0.72)	0.66(0.39)	0.83(0.72)	0.85(0.81)	0.77(0.76)	0.83(0.75)
SimpleNeRF w/o Views Augmentation	0.41(0.52)	0.29(0.45)	0.15(0.24)	0.30(0.47)	0.25(0.44)	0.34(0.41)	0.31(0.41)	0.28(0.40)	0.28(0.41)
	0.54(0.49)	0.63(0.52)	0.74(0.68)	0.60(0.51)	0.46(0.31)	0.52(0.36)	0.80(0.73)	0.68(0.57)	0.64(0.54)
	17.5(17.2)	20.4(17.5)	23.6(22.0)	18.7(16.5)	15.6(13.6)	16.1(14.6)	21.4(18.7)	18.3(16.9)	19.3(17.4)
	0.14(0.20)	0.11(0.24)	0.05(0.24)	0.21(0.37)	0.29(0.43)	0.16(0.34)	0.15(0.22)	0.13(0.19)	0.15(0.28)
	0.69(0.66)	0.95(0.77)	0.99(0.98)	0.82(0.72)	0.70(0.43)	0.84(0.73)	0.79(0.77)	0.78(0.77)	0.83(0.75)
SimpleNeRF w/o Coarse-fine Consistency	0.38(0.50)	0.31(0.47)	0.22(0.31)	0.33(0.52)	0.29(0.49)	0.34(0.41)	0.29(0.36)	0.28(0.38)	0.30(0.42)
	0.56(0.51)	0.57(0.48)	0.69(0.64)	0.57(0.50)	0.34(0.23)	0.54(0.38)	0.79(0.72)	0.66(0.56)	0.61(0.52)
	17.8(17.4)	18.4(16.3)	23.2(21.9)	19.4(17.0)	13.4(12.3)	16.2(14.4)	21.5(18.9)	18.6(17.2)	19.0(17.2)
	0.15(0.23)	0.25(0.40)	0.06(0.30)	0.20(0.39)	0.58(0.65)	0.15(0.34)	0.10(0.21)	0.12(0.20)	0.19(0.33)
	0.70(0.65)	0.78(0.56)	0.99(0.98)	0.87(0.74)	0.28(0.10)	0.85(0.73)	0.93(0.88)	0.78(0.72)	0.80(0.70)
SimpleNeRF w/o reliable depth	0.39(0.50)	0.54(0.69)	0.16(0.25)	0.26(0.42)	0.27(0.46)	0.34(0.42)	0.28(0.35)	0.27(0.37)	0.30(0.42)
	0.55(0.50)	0.28(0.25)	0.72(0.66)	0.65(0.54)	0.38(0.25)	0.52(0.36)	0.82(0.75)	0.69(0.57)	0.60(0.51)
	17.7(17.3)	10.9(10.7)	23.3(22.0)	19.9(17.1)	14.1(12.7)	16.2(14.6)	22.0(19.1)	18.8(17.0)	18.4(16.7)
	0.15(0.23)	0.58(0.71)	0.06(0.31)	0.20(0.38)	0.51(0.60)	0.16(0.34)	0.12(0.22)	0.11(0.20)	0.22(0.36)
	0.71(0.65)	0.16(0.18)	0.99(0.97)	0.87(0.75)	0.52(0.27)	0.83(0.72)	0.91(0.87)	0.83(0.73)	0.76(0.67)
SimpleNeRF w/o Residual Positional Encodings	0.41(0.53)	0.30(0.47)	0.16(0.26)	0.28(0.45)	0.25(0.45)	0.34(0.41)	0.32(0.41)	0.28(0.39)	0.28(0.41)
	0.54(0.49)	0.62(0.52)	0.72(0.65)	0.62(0.53)	0.44(0.30)	0.53(0.37)	0.80(0.73)	0.68(0.58)	0.64(0.54)
	17.5(17.1)	20.4(17.4)	22.9(21.7)	19.2(16.9)	15.3(13.5)	16.1(14.5)	21.3(18.2)	18.1(16.8)	19.2(17.3)
	0.14(0.22)	0.11(0.23)	0.06(0.29)	0.20(0.35)	0.35(0.47)	0.16(0.34)	0.16(0.25)	0.14(0.19)	0.16(0.29)
	0.68(0.61)	0.96(0.79)	0.99(0.97)	0.84(0.72)	0.60(0.35)	0.84(0.74)	0.76(0.68)	0.76(0.75)	0.82(0.72)
SimpleNeRF w/ Identical Augmentations	0.40(0.52)	0.28(0.44)	0.16(0.25)	0.30(0.46)	0.24(0.44)	0.34(0.42)	0.32(0.40)	0.29(0.40)	0.28(0.41)
	0.54(0.49)	0.64(0.53)	0.73(0.67)	0.60(0.51)	0.47(0.32)	0.51(0.36)	0.80(0.73)	0.67(0.57)	0.64(0.54)
	17.4(17.0)	20.5(17.6)	23.6(22.1)	19.0(16.8)	15.7(13.7)	16.2(14.6)	21.4(18.9)	17.8(16.5)	19.3(17.4)
	0.14(0.20)	0.10(0.22)	0.05(0.24)	0.20(0.37)	0.27(0.41)	0.16(0.33)	0.15(0.26)	0.15(0.20)	0.15(0.28)
	0.71(0.69)	0.96(0.79)	0.99(0.99)	0.84(0.74)	0.72(0.45)	0.85(0.74)	0.78(0.72)	0.77(0.75)	0.84(0.75)

Table 16: Per-scene performance of various models with two input views on RealEstate-10K dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	0	1	3	4	6	Average
InfoNeRF	0.5711(0.6035)	0.4592(0.5016)	0.5959(0.6686)	0.7184(0.7605)	0.6175(0.6580)	0.5924(0.6384)
	0.4396(0.4356)	0.6275(0.6339)	0.2549(0.2469)	0.3873(0.3936)	0.4616(0.4614)	0.4342(0.4343)
	11.70(11.68)	16.20(15.61)	10.66(10.64)	10.70(10.83)	12.10(12.09)	12.27(12.17)
	2.2791(2.2552)	2.0442(2.0017)	2.8117(2.9467)	2.3704(2.3495)	1.3912(1.6040)	2.1793(2.2314)
	-0.0335(-0.0449)	0.1286(0.1617)	0.0727(0.0663)	0.0867(0.0712)	0.2014(0.2166)	0.0912(0.0942)
DietNeRF	0.4738(0.5162)	0.3748(0.4149)	0.5748(0.6718)	0.4712(0.4986)	0.2956(0.3296)	0.4381(0.4862)
	0.6335(0.6295)	0.8289(0.8387)	0.3356(0.3284)	0.6972(0.7016)	0.7718(0.7617)	0.6534(0.6520)
	14.80(14.72)	21.04(20.58)	12.93(12.59)	18.47(18.50)	23.05(22.74)	18.06(17.83)
	2.2511(2.2272)	1.9828(1.9476)	2.8672(3.0125)	2.0381(2.0190)	0.9844(1.1973)	2.0247(2.0807)
	0.0047(0.0032)	0.5266(0.5465)	0.2821(0.3065)	-0.1165(-0.0861)	0.4725(0.4559)	0.2339(0.2452)
RegNeRF	0.3357(0.3480)	0.2898(0.3179)	0.4428(0.4944)	0.4969(0.5390)	0.4992(0.5421)	0.4129(0.4483)
	0.6175(0.6010)	0.8226(0.8290)	0.3103(0.2973)	0.6177(0.6148)	0.5902(0.5901)	0.5916(0.5864)
	16.77(16.51)	21.62(21.04)	14.16(13.88)	17.46(17.13)	15.68(15.79)	17.14(16.87)
	-	-	-	-	-	-
	-0.0420(-0.0269)	-0.3059(-0.2176)	0.1796(0.1485)	0.2285(0.1764)	0.4986(0.4781)	0.1118(0.1117)
DS-NeRF	0.2331(0.2588)	0.2340(0.2727)	0.4074(0.5095)	0.2117(0.2369)	0.2685(0.3074)	0.2709(0.3171)
	0.8248(0.8111)	0.9137(0.9085)	0.5264(0.4999)	0.8835(0.8767)	0.8433(0.8332)	0.7983(0.7859)
	25.42(24.68)	29.86(27.93)	19.49(19.24)	29.81(29.18)	26.71(26.18)	26.26(25.44)
	1.0798(1.0705)	0.7518(0.8238)	0.9077(1.1926)	0.2362(0.2607)	0.6065(0.8138)	0.7164(0.8323)
	0.2650(0.2633)	0.7431(0.6798)	0.9088(0.8864)	0.9117(0.8934)	0.5014(0.4934)	0.6660(0.6433)
DDP-NeRF	0.1017(0.1099)	0.0966(0.1167)	0.2827(0.3438)	0.0506(0.0563)	0.1131(0.1140)	0.1290(0.1481)
	0.8973(0.8858)	0.9579(0.9490)	0.5931(0.5589)	0.9419(0.9356)	0.9296(0.9215)	0.8640(0.8502)
	27.01(25.90)	30.45(25.87)	19.59 (18.97)	33.27(32.01)	28.62(28.00)	27.79(26.15)
	0.4971(0.5310)	0.5727(0.6260)	0.7513(1.0264)	0.0855(0.1074)	0.5091(0.6813)	0.4831(0.5944)
	0.7288(0.6851)	0.6975(0.6481)	0.9424(0.9238)	0.9304(0.9070)	0.6616(0.6674)	0.7921(0.7663)
FreeNeRF	0.4362(0.4490)	0.4395(0.4970)	0.5543(0.6376)	0.6310(0.6736)	0.4571(0.4786)	0.5036(0.5471)
	0.5537(0.5446)	0.7697(0.7719)	0.2879(0.2766)	0.4863(0.4914)	0.5795(0.5833)	0.5354(0.5336)
	15.14(15.00)	17.78(17.00)	12.54(12.15)	12.76(12.84)	15.30(15.50)	14.70(14.50)
	-	-	-	-	-	-
	-0.0933(-0.0896)	-0.3832(-0.3058)	0.1582(0.1591)	-0.1832(-0.1963)	-0.4668(-0.4741)	-0.1937(-0.1813)
ViP-NeRF	0.0347(0.0422)	0.0354(0.0497)	0.1793(0.1944)	0.0315(0.0344)	0.0626(0.0708)	0.0687(0.0783)
	0.9578(0.9431)	0.9791(0.9666)	0.6061(0.5675)	0.9572(0.9498)	0.9442(0.9316)	0.8889(0.8717)
	32.93(30.42)	37.48(31.96)	19.36(18.90)	38.05(34.75)	33.76(31.73)	32.32(29.55)
	0.3975(0.4451)	0.2310(0.3632)	0.7525(1.0708)	0.0840(0.1069)	0.4631(0.6827)	0.3856(0.5337)
	0.7263(0.6812)	0.9423(0.8487)	0.9392(0.8270)	0.9470(0.9174)	0.6683(0.6509)	0.8446(0.7851)
SimpleNeRF	0.0276(0.0369)	0.0273(0.0392)	0.1913(0.2130)	0.0300(0.0336)	0.0413(0.0500)	0.0635(0.0745)
	0.9662(0.9532)	0.9823(0.9728)	0.5968(0.5590)	0.9592(0.9511)	0.9663(0.9553)	0.8942(0.8783)
	34.94(31.89)	38.36(33.80)	19.23(18.65)	38.51(34.93)	34.47(32.24)	33.10(30.30)
	0.3299(0.3764)	0.2283(0.3202)	0.7286(1.0194)	0.0730(0.0960)	0.2747(0.4803)	0.3269(0.4584)
	0.8490(0.8007)	0.9733(0.9309)	0.9414(0.8582)	0.9624(0.9363)	0.8812(0.8647)	0.9215(0.8781)

Table 17: Per-scene performance of various models with three input views on RealEstate-10K dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	0	1	3	4	6	Average
InfoNeRF	0.6722(0.6983)	0.6151(0.6280)	0.6187(0.6621)	0.7524(0.7785)	0.6220(0.6562)	0.6561(0.6846)
	0.3622(0.3593)	0.4916(0.4871)	0.1892(0.1848)	0.3749(0.3790)	0.4779(0.4797)	0.3792(0.3780)
	9.74(9.73)	12.08(11.75)	9.08(9.23)	10.19(10.28)	11.74(11.87)	10.57(10.57)
	2.3079(2.2937)	2.1419(2.1368)	2.8991(3.0636)	2.3476(2.3454)	1.4028(1.5757)	2.2198(2.2830)
	0.0640(0.0560)	0.3196(0.3379)	0.0813(0.0884)	0.1210(0.1273)	0.3788(0.3875)	0.1929(0.1994)
DietNeRF	0.4698(0.4893)	0.4979(0.5115)	0.5759(0.6289)	0.4600(0.4719)	0.3142(0.3413)	0.4636(0.4886)
	0.6226(0.6166)	0.8046(0.8158)	0.3235(0.3233)	0.7144(0.7151)	0.7630(0.7514)	0.6456(0.6445)
	16.44(16.36)	18.51(18.29)	13.44(13.30)	18.98(18.95)	22.65(22.54)	18.01(17.89)
	2.1863(2.1710)	2.0866(2.0870)	2.8594(3.0383)	1.9957(1.9918)	1.0496(1.2234)	2.0355(2.1023)
	0.0856(0.0728)	-0.1136(-0.0664)	0.3096(0.3401)	-0.2370(-0.2100)	0.0754(0.0823)	0.0240(0.0438)
RegNeRF	0.3877(0.3952)	0.3179(0.3220)	0.4885(0.5302)	0.5406(0.5635)	0.3506(0.3699)	0.4171(0.4362)
	0.6084(0.5950)	0.8178(0.8237)	0.3008(0.2893)	0.6227(0.6221)	0.7162(0.7087)	0.6132(0.6078)
	16.22(15.99)	21.06(20.89)	14.07(13.87)	17.61(17.60)	20.35(20.28)	17.86(17.73)
	-	-	-	-	-	-
	-0.0854(-0.0849)	0.2210(0.2376)	-0.0298(-0.0517)	0.3273(0.2715)	-0.1459(-0.1349)	0.0574(0.0475)
DS-NeRF	0.2273(0.2436)	0.2404(0.2606)	0.4562(0.5346)	0.2463(0.2602)	0.2761(0.3067)	0.2893(0.3211)
	0.8395(0.8264)	0.9124(0.9099)	0.5121(0.4930)	0.8750(0.8710)	0.8631(0.8521)	0.8004(0.7905)
	26.09(25.24)	29.43(28.68)	19.44(19.14)	29.40(29.08)	28.12(27.58)	26.50(25.94)
	0.4587(0.4890)	0.7153(0.7774)	0.8878(1.1777)	0.2076(0.2207)	0.4305(0.5970)	0.5400(0.6524)
	0.7501(0.7227)	0.7428(0.7184)	0.9158(0.8939)	0.8619(0.8490)	0.7826(0.7711)	0.8106(0.7910)
DDP-NeRF	0.1138(0.1143)	0.1229(0.1069)	0.3283(0.3821)	0.0591(0.0623)	0.1348(0.1349)	0.1518(0.1601)
	0.9046(0.8940)	0.9557(0.9583)	0.5709(0.5548)	0.9389(0.9352)	0.9234(0.9166)	0.8587(0.8518)
	26.17(25.27)	27.99(26.67)	19.11(18.81)	32.68(31.84)	27.40(26.99)	26.67(25.92)
	0.3632(0.4007)	0.4621(0.5204)	0.7746(1.0395)	0.1057(0.1155)	0.3640(0.5347)	0.4139(0.5222)
	0.7805(0.7418)	0.8822(0.8365)	0.9494(0.9272)	0.9022(0.8939)	0.7918(0.7661)	0.8612(0.8331)
FreeNeRF	0.5272(0.5429)	0.4975(0.5138)	0.5800(0.6407)	0.5668(0.5905)	0.4016(0.4193)	0.5146(0.5414)
	0.5337(0.5250)	0.7487(0.7504)	0.2937(0.2863)	0.6111(0.6122)	0.6671(0.6634)	0.5708(0.5675)
	13.87(13.79)	16.07(15.59)	12.60(12.45)	15.74(15.72)	18.01(18.05)	15.26(15.12)
	-	-	-	-	-	-
	-0.1174(-0.1177)	-0.2051(-0.1583)	-0.0762(-0.0521)	-0.3930(-0.3889)	-0.5031(-0.5058)	-0.2590(-0.2445)
ViP-NeRF	0.0405(0.0432)	0.0517(0.0541)	0.1939(0.2170)	0.0351(0.0352)	0.0579(0.0663)	0.0758(0.0832)
	0.9567(0.9450)	0.9715(0.9638)	0.6409(0.6148)	0.9518(0.9484)	0.9624(0.9537)	0.8967(0.8852)
	32.88(30.68)	33.82(31.50)	19.91(19.59)	37.13(35.78)	35.90(33.79)	31.93(30.27)
	0.3021(0.3488)	0.3589(0.4690)	0.6819(0.9977)	0.0850(0.0956)	0.2548(0.4306)	0.3365(0.4683)
	0.8391(0.7931)	0.9002(0.8319)	0.9427(0.8559)	0.9268(0.9173)	0.8956(0.8805)	0.9009(0.8558)
SimpleNeRF	0.0327(0.0379)	0.0263(0.0350)	0.2059(0.2325)	0.0324(0.0328)	0.0660(0.0761)	0.0726(0.0829)
	0.9646(0.9546)	0.9817(0.9765)	0.6347(0.6111)	0.9542(0.9498)	0.9568(0.9477)	0.8984(0.8879)
	34.77(32.23)	39.32(36.44)	19.93(19.65)	37.67(35.85)	34.37(32.81)	33.21(31.40)
	0.2606(0.3062)	0.2189(0.2822)	0.5466(0.8178)	0.0798(0.0989)	0.2791(0.4375)	0.2770(0.3885)
	0.8791(0.8347)	0.9797(0.9687)	0.9462(0.8691)	0.9405(0.9179)	0.8876(0.8747)	0.9266(0.8931)

Table 18: Per-scene performance of various models with four input views on RealEstate-10K dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	0	1	3	4	6	Average
InfoNeRF	0.5948(0.6030)	0.7054(0.6927)	0.6889(0.6976)	0.7918(0.8068)	0.5447(0.5604)	0.6651(0.6721)
	0.5023(0.5002)	0.3621(0.3621)	0.1767(0.1725)	0.3119(0.3176)	0.5683(0.5625)	0.3843(0.3830)
	12.80(12.74)	8.76(8.60)	8.68(8.71)	8.76(8.79)	14.12(14.12)	10.62(10.59)
	2.2648(2.2692)	2.0969(2.1482)	2.8542(3.0882)	2.3608(2.3658)	1.3605(1.4993)	2.1874(2.2742)
	0.2109(0.2053)	0.2990(0.3040)	0.2874(0.2812)	0.1548(0.1696)	0.3227(0.3369)	0.2549(0.2594)
DietNeRF	0.5095(0.5189)	0.4664(0.4628)	0.6182(0.6399)	0.4881(0.4966)	0.3442(0.3589)	0.4853(0.4954)
	0.6163(0.6125)	0.8352(0.8399)	0.3243(0.3198)	0.7157(0.7156)	0.7598(0.7495)	0.6503(0.6475)
	15.78(15.63)	19.91(19.64)	13.29(13.27)	19.06(19.01)	21.99(21.93)	18.01(17.89)
	2.2313(2.2346)	2.0104(2.0647)	2.7825(3.0206)	1.9930(1.9969)	1.1820(1.3199)	2.0398(2.1273)
	0.0064(0.0055)	0.1436(0.1283)	0.3925(0.3897)	-0.1838(-0.1637)	0.1363(0.1458)	0.0990(0.1011)
RegNeRF	0.4215(0.4252)	0.3498(0.3461)	0.5737(0.5903)	0.5450(0.5581)	0.2678(0.2719)	0.4316(0.4383)
	0.5922(0.5824)	0.8249(0.8256)	0.3036(0.2936)	0.6484(0.6479)	0.7595(0.7497)	0.6257(0.6198)
	16.29(16.09)	21.14(20.98)	13.93(13.91)	18.49(18.48)	21.86(21.78)	18.34(18.25)
	-	-	-	-	-	-
	-0.0032(-0.0018)	0.2630(0.2456)	0.2565(0.2652)	0.4133(0.4104)	-0.2187(-0.2214)	0.1422(0.1396)
DS-NeRF	0.2663(0.2746)	0.2513(0.2580)	0.5061(0.5550)	0.2380(0.2461)	0.2899(0.3100)	0.3103(0.3287)
	0.8230(0.8156)	0.9202(0.9176)	0.5184(0.5014)	0.8758(0.8734)	0.8621(0.8520)	0.7999(0.7920)
	25.95(25.40)	29.99(29.40)	19.73(19.64)	29.48(29.26)	28.08(27.69)	26.65(26.28)
	0.4308(0.4611)	0.6204(0.6904)	0.8457(1.1119)	0.2160(0.2265)	0.4644(0.5954)	0.5154(0.6171)
	0.7511(0.7254)	0.7074(0.6988)	0.9340(0.9191)	0.8703(0.8603)	0.8097(0.8055)	0.8145(0.8018)
DDP-NeRF	0.1235(0.1196)	0.0893(0.0797)	0.3667(0.3938)	0.0648(0.0641)	0.1371(0.1346)	0.1563(0.1584)
	0.8989(0.8921)	0.9633(0.9629)	0.5926(0.5783)	0.9329(0.9306)	0.9206(0.9144)	0.8617(0.8557)
	25.67(25.14)	30.10(28.57)	19.64(19.57)	32.44(31.73)	27.52(27.36)	27.07(26.48)
	0.3369(0.3668)	0.3465(0.4096)	0.7186(0.9760)	0.1030(0.1135)	0.4109(0.5405)	0.3832(0.4813)
	0.8081(0.7767)	0.9270(0.9187)	0.9529(0.9422)	0.9067(0.8976)	0.7747(0.7674)	0.8739(0.8605)
FreeNeRF	0.5569(0.5630)	0.4786(0.4750)	0.6252(0.6458)	0.5672(0.5839)	0.3852(0.3938)	0.5226(0.5323)
	0.5398(0.5335)	0.7945(0.7990)	0.3130(0.3063)	0.6635(0.6623)	0.7027(0.6936)	0.6027(0.5989)
	13.91(13.84)	18.05(17.93)	12.71(12.69)	17.37(17.29)	19.52(19.48)	16.31(16.25)
	-	-	-	-	-	-
	-0.0622(-0.0616)	-0.2656(-0.2658)	0.3116(0.3100)	-0.4956(-0.5005)	-0.5640(-0.5630)	-0.2152(-0.2162)
ViP-NeRF	0.0438(0.0437)	0.0719(0.0702)	0.2232(0.2305)	0.0390(0.0387)	0.0681(0.0716)	0.0892(0.0909)
	0.9569(0.9507)	0.9636(0.9591)	0.6563(0.6401)	0.9474(0.9445)	0.9599(0.9527)	0.8968(0.8894)
	33.96(32.32)	33.23(31.95)	20.61(20.42)	36.76(35.65)	35.20(33.81)	31.95(30.83)
	0.2883(0.3216)	0.5206(0.6261)	0.6319(0.8927)	0.0842(0.1009)	0.3039(0.4392)	0.3658(0.4761)
	0.8421(0.8097)	0.6288(0.5661)	0.9222(0.8849)	0.9232(0.9026)	0.8908(0.8766)	0.8414(0.8080)
SimpleNeRF	0.0373(0.0381)	0.0417(0.0459)	0.2308(0.2427)	0.0341(0.0333)	0.0796(0.0857)	0.0847(0.0891)
	0.9616(0.9562)	0.9777(0.9746)	0.6513(0.6355)	0.9509(0.9477)	0.9522(0.9446)	0.8987(0.8917)
	34.52(32.95)	38.01(36.44)	20.68(20.52)	37.29(35.97)	33.90(32.77)	32.88(31.73)
	0.2668(0.2977)	0.2254(0.2825)	0.4750(0.6914)	0.0776(0.0944)	0.3014(0.4162)	0.2692(0.3565)
	0.8631(0.8362)	0.9719(0.9613)	0.9493(0.9278)	0.9417(0.9210)	0.8784(0.8710)	0.9209(0.9035)

Table 19: Per-scene performance of ablated models with two input views on RealEstate-10K dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively. The values within parenthesis show unmasked scores.

Model \ Scene Name	0	1	3	4	6	Average
SimpleNeRF	0.0276(0.0369)	0.0273(0.0392)	0.1913(0.2130)	0.0300(0.0336)	0.0413(0.0500)	0.0635(0.0745)
	0.9662(0.9532)	0.9823(0.9728)	0.5968(0.5590)	0.9592(0.9511)	0.9663(0.9553)	0.8942(0.8783)
	34.94(31.89)	38.36(33.80)	19.23(18.65)	38.51(34.93)	34.47(32.24)	33.10(30.30)
	0.3299(0.3764)	0.2283(0.3202)	0.7286(1.0194)	0.0730(0.0960)	0.2747(0.4803)	0.3269(0.4584)
	0.8490(0.8007)	0.9733(0.9309)	0.9414(0.8582)	0.9624(0.9363)	0.8812(0.8647)	0.9215(0.8781)
SimpleNeRF w/o Points Augmentation	0.0334(0.0426)	0.0387(0.0543)	0.1988(0.2215)	0.0315(0.0365)	0.0736(0.0898)	0.0752(0.0889)
	0.9618(0.9487)	0.9768(0.9672)	0.5927(0.5540)	0.9582(0.9492)	0.9533(0.9418)	0.8886(0.8722)
	34.26(31.31)	37.15(32.80)	19.28(18.76)	38.27(34.50)	33.75(31.87)	32.54(29.85)
	0.3632(0.4080)	0.3145(0.3992)	0.7602(1.0816)	0.0862(0.1146)	0.3733(0.5510)	0.3795(0.5109)
	0.8279(0.7849)	0.9211(0.8817)	0.9074(0.7931)	0.9601(0.9296)	0.8699(0.8540)	0.8973(0.8487)
SimpleNeRF w/o Views Augmentation	0.0351(0.0449)	0.0559(0.0685)	0.2018(0.2258)	0.0308(0.0365)	0.0713(0.0870)	0.0790(0.0925)
	0.9606(0.9461)	0.9717(0.9625)	0.5971(0.5565)	0.9584(0.9493)	0.9541(0.9423)	0.8884(0.8714)
	33.84(30.91)	35.74(32.62)	19.20(18.44)	38.35(34.51)	33.50(31.61)	32.13(29.62)
	0.4028(0.4448)	0.3687(0.4311)	0.7095(1.0231)	0.0907(0.1191)	0.3633(0.5370)	0.3870(0.5110)
	0.7597(0.7208)	0.9366(0.9130)	0.9205(0.8209)	0.9582(0.9264)	0.8437(0.8327)	0.8837(0.8428)
SimpleNeRF w/o Coarse-fine Consistency	0.0407(0.0497)	0.0482(0.0605)	0.2017(0.2288)	0.0309(0.0353)	0.0484(0.0580)	0.0740(0.0865)
	0.9525(0.9371)	0.9733(0.9653)	0.5890(0.5466)	0.9584(0.9486)	0.9612(0.9492)	0.8869(0.8693)
	32.03(29.63)	36.02(33.23)	19.09(18.39)	38.24(34.21)	33.89(31.89)	31.86(29.47)
	0.4844(0.5346)	0.4072(0.4644)	0.8111(1.1278)	0.0826(0.1163)	0.3263(0.5198)	0.4223(0.5526)
	0.6863(0.6411)	0.9236(0.9177)	0.9083(0.7955)	0.9590(0.9206)	0.8108(0.7950)	0.8576(0.8140)
SimpleNeRF w/o reliable depth	0.0304(0.0402)	0.0548(0.0670)	0.1875(0.2104)	0.0296(0.0327)	0.0414(0.0506)	0.0687(0.0802)
	0.9631(0.9519)	0.9685(0.9597)	0.5991(0.5600)	0.9595(0.9516)	0.9661(0.9540)	0.8913(0.8754)
	34.95(32.35)	34.67(32.08)	19.49(19.04)	38.53(35.17)	35.52(33.00)	32.63(30.33)
	0.5129(0.5538)	0.5654(0.6359)	0.7689(1.0549)	0.0737(0.0994)	0.3214(0.5452)	0.4485(0.5778)
	0.6705(0.6365)	0.9246(0.9088)	0.9452(0.8952)	0.9602(0.9304)	0.8641(0.8312)	0.8729(0.8404)
SimpleNeRF w/o Residual Positional Encodings	0.0353(0.0442)	0.0627(0.0726)	0.1967(0.2205)	0.0306(0.0329)	0.0697(0.0844)	0.0790(0.0909)
	0.9605(0.9468)	0.9681(0.9607)	0.5959(0.5552)	0.9586(0.9526)	0.9544(0.9424)	0.8875(0.8715)
	33.64(30.97)	35.19(32.53)	19.16(18.47)	38.44(35.86)	33.57(31.74)	32.00(29.91)
	0.4173(0.4567)	0.3882(0.4344)	0.7512(1.0759)	0.0866(0.1046)	0.3764(0.5560)	0.4040(0.5255)
	0.7540(0.7167)	0.9422(0.9137)	0.9016(0.7851)	0.9595(0.9343)	0.8618(0.8462)	0.8838(0.8392)
SimpleNeRF w/ Identical Augmentations	0.0317(0.0410)	0.0491(0.0623)	0.1996(0.2257)	0.0311(0.0360)	0.0768(0.0930)	0.0777(0.0916)
	0.9621(0.9494)	0.9740(0.9651)	0.5911(0.5514)	0.9585(0.9498)	0.9520(0.9409)	0.8875(0.8713)
	34.36(31.72)	36.22(32.63)	19.13(18.47)	38.36(34.68)	33.49(31.73)	32.31(29.85)
	0.3636(0.4079)	0.3754(0.4348)	0.7880(1.1129)	0.0887(0.1169)	0.4028(0.5621)	0.4037(0.5269)
	0.8239(0.7811)	0.9569(0.9145)	0.8882(0.7710)	0.9583(0.9244)	0.8472(0.8354)	0.8949(0.8453)