

Simple-RF: Regularizing Sparse Input Radiance Fields with Simpler Solutions

NAGABHUSHAN SOMRAJ, Indian Institute of Science, India
 SAI HARSHA MUPPARAJU, Indian Institute of Science, India
 ADITHYAN KARANAYIL, Indian Institute of Science, India
 RAJIV SOUNDARARAJAN, Indian Institute of Science, India

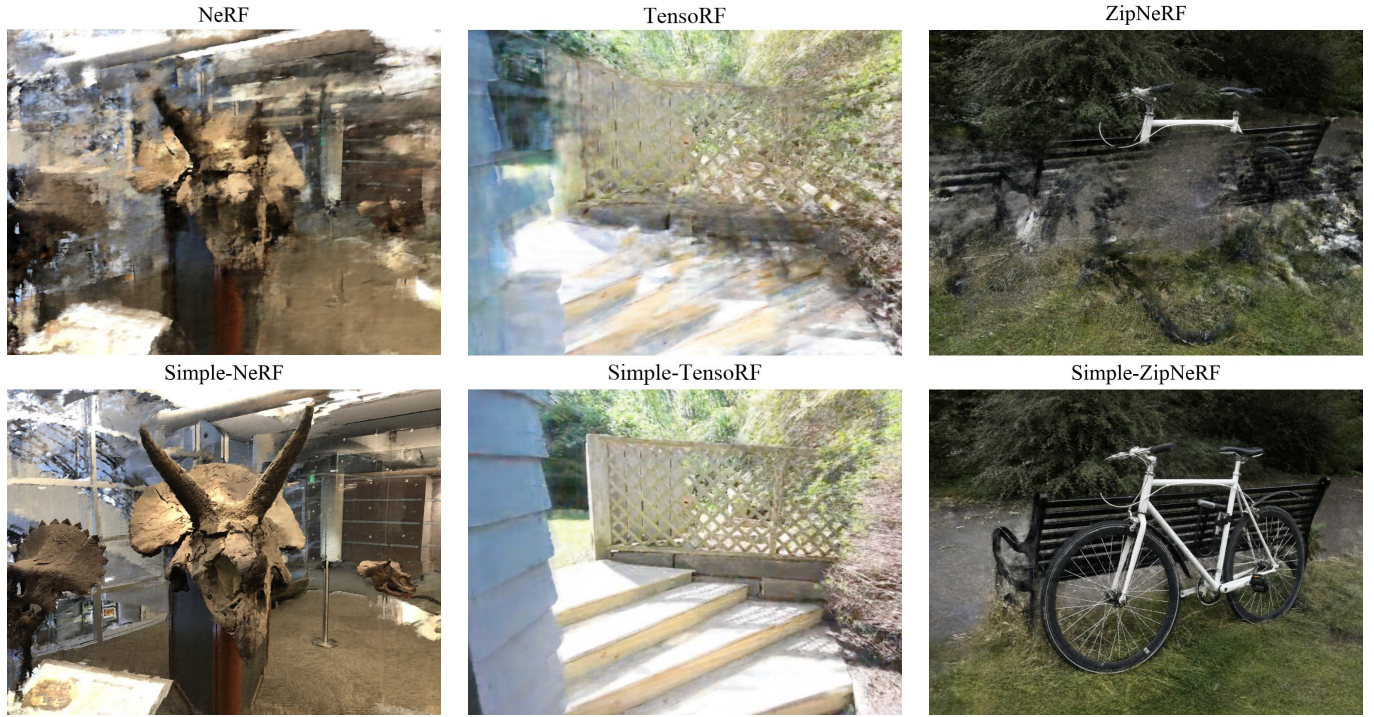


Fig. 1. We show the improvements achieved by our regularizations on the NeRF, TensorRF and ZipNeRF models on NeRF-LLFF, RealEstate-10K and MipNeRF360 datasets respectively. We observe that the vanilla radiance fields suffer from various distortions. Regularizing the radiance fields with simpler solutions leads to significantly better reconstructions with all the three radiance fields.

Neural Radiance Fields (NeRF) show impressive performances in photo-realistic free-view rendering of scenes. Recent improvements such as TensorRF and ZipNeRF employ explicit models for faster optimization and rendering. However, all these radiance fields require a dense sampling of images in the given scene for effective training. Their performances degrade significantly when only a sparse set of views is available. Existing depth priors used to supervise the radiance fields are either sparse or suffer from generalization

Authors' addresses: Nagabhushan Somraj, Indian Institute of Science, Bengaluru, Karnataka, 560012, India, nagabhushans@iisc.ac.in; Sai Harsha Mupparaju, Indian Institute of Science, Bengaluru, Karnataka, India, saiharsham@iisc.ac.in; Adithyan Karanayil, Indian Institute of Science, Bengaluru, Karnataka, India, adithyanv@iisc.ac.in; Rajiv Soundararajan, Indian Institute of Science, Bengaluru, Karnataka, India, rajivs@iisc.ac.in.

© 2024 Association for Computing Machinery.

This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *ACM Transactions on Graphics*, <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>.

issues. We seek to learn scene-specific dense depth priors to regularize the radiance fields. Further, we desire a framework of regularizations that can work across different radiance field models. We observe that certain features of the radiance fields, such as positional encoding, number of decomposed tensor components or size of the hash table, cause overfitting in the sparse-input scenario. We design augmented models by reducing the capacity of these features and train them along with the main radiance field. These augmented models learn simpler solutions, which estimate better depth in certain regions. By supervising the main radiance field with such depths, we significantly improve the performance of the radiance fields on popular forward-facing and 360° datasets by employing the above regularization.

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

Additional Key Words and Phrases: neural rendering, novel view synthesis, sparse input radiance fields

ACM Reference Format:

Nagabhushan Somraj, Sai Harsha Mupparaju, Adithyan Karanayil, and Rajiv Soundararajan. 2024. Simple-RF: Regularizing Sparse Input Radiance Fields with Simpler Solutions. *ACM Trans. Graph.* 1, 1 (May 2024), 39 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

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 optimize and render quickly [Müller et al. 2022], reduce aliasing artifacts [Barron et al. 2021], and learn on unbounded scenes [Barron et al. 2022]. Yet all these models 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 both implicit radiance fields such as NeRF and explicit radiance fields such as TensorRF [Chen et al. 2022b] and ZipNeRF [Barron et al. 2023] with a sparse set of input images and aim to design novel regularizations for effective training.

Researchers have extensively studied the training of NeRFs with sparse input views. One approach to training NeRFs with sparse input views is to use generalized NeRFs, where the NeRF is additionally conditioned on a latent scene representation obtained using a convolutional neural network. However, these models require a large multi-view 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 [Guo et al. 2024; Ni et al. 2024; Zhang et al. 2021b]. 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. 2022] and ViP-NeRF [Somraj and Soundararajan 2023] use simple priors such as depth smoothness, sparse depth or relative depth respectively obtained through classical methods. 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 dense depth prior. These approaches may also suffer from issues similar to those of the generalized models. This raises the question of whether we can instead learn the dense depth supervision in-situ without employing any pre-training.

Recently, there is also a growing interest in sparse input explicit radiance fields owing to their fast optimization and rendering times. However, the regularizations designed for these models are limited to a specific explicit radiance field and do not generalize to more recent models. For example, while the regularizations designed for ZeroRF [Shi et al. 2024] can be applied to TensorRF based models only, other regularizations are applicable to implicit models only [Yang et al. 2023; Zhu et al. 2024]. It is desirable to design regularizations that are relevant to different radiance field models through a single framework. While there exists a plethora of implicit and explicit radiance field models, we consider the NeRF as the representative

model for implicit radiance fields and consider two explicit radiance fields, namely, TensorRF [Chen et al. 2022b] and ZipNeRF [Barron et al. 2023]. We note that although the NeRF based models are slow in optimization and rendering, NeRFs are versatile in learning different properties of the scenes [Bi et al. 2020; Srinivasan et al. 2021; Verbin et al. 2022; Zhang et al. 2021a] and may be easier to optimize as compared to the explicit models. Hence, we believe that designing regularizations for sparse-input NeRF is also of considerable interest.

We first observe that the radiance field models often exploit their high capability to learn unnecessary complex solutions when training with sparse input views. While these solutions perfectly explain the observed images, they can cause severe distortions in novel views. For example, some of the key components of the radiance fields, such as positional encoding in the NeRF or vector-matrix decomposition employed in TensorRF, provide powerful capabilities to the radiance field and are designed for training the model with dense input views. Existing implementations of these components may be sub-optimal with fewer input views due to the highly under-constrained system, causing several distortions. Figs. 4, 7 and 8 show common distortions observed with NeRF, TensorRF and ZipNeRF in the few-shot setting respectively. We follow the popular Occam's razor principle and regularize the radiance fields to choose simpler solutions over complex ones, wherever possible. In particular, we design augmented models by reducing the capabilities of the radiance fields and use the depth estimated by these models to supervise the main radiance field.

We design different augmentations for NeRF, TensorRF and ZipNeRF based on different shortcomings and architectures of these models. The high positional encoding degree used in the NeRF leads to undesired depth discontinuities, creating floaters. Further, the view-dependent radiance feature leads to shape-radiance ambiguity, creating duplication artifacts. We design augmentations for the NeRF by reducing the positional encoding degree and disabling the view-dependent radiance feature. On the other hand, the large number of high-resolution factorized components in TensorRF and the large hash table in ZipNeRF cause floaters in these models in the few-shot setting. Thus, we design augmentations to constrain the model with respect to such components to learn simpler solutions.

We use the simplified models as augmentations for depth supervision and not as the main NeRF model since naively reducing the capacity of the radiance fields may lead to sub-optimal solutions in certain regions [Jain et al. 2021]. For example, the model that can learn only smooth depth transitions may fail to learn sharp depth discontinuities at object boundaries. Further, the 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 refer to our family of regularized models as Simple Radiance Fields (Simple-RF) since we regularize the models to choose simple solutions over complex ones, wherever feasible. We refer to the individual models as Simple-NeRF, Simple-TensorRF and Simple-ZipNeRF respectively. We evaluate our models on four popular datasets that include forward-facing scenes (NeRF-LLFF), unbounded forward-facing scenes (RealEstate-10K), unbounded 360° scenes (MipNeRF360)

and bounded 360° scenes (NeRF-Synthetic) and show that our models achieve significant improvement in performance on all the datasets. Fig. 1 qualitatively shows the improvements achieved through our regularizations on NeRF, TensorRF and ZipNeRF on three different datasets. Further, we show that our model learns geometry significantly better than prior art.

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

- We find that the high positional encoding degree and view-dependent radiance of the NeRF cause floaters and duplication artifacts when training with sparse inputs. We design augmented models on both these fronts to supervise the main NeRF and mitigate both artifacts.
- We observe that the large number of high-resolution decomposed components in TensorRF leads to floater artifacts with sparse inputs. Thus, the augmented model is obtained by reducing the number and resolutions of the decomposed components.
- We find that the large hash table in ZipNeRF causes floaters when training with sparse inputs. The augmented model is designed by reducing the size of the hash table.
- 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 main radiance field.
- We show that our regularization achieves substantial improvements on different radiance fields and on four different datasets.

2 RELATED WORK

Chen and Williams [1993] introduce the problem of novel view synthesis and propose an image-based rendering (IBR) approach to synthesize novel views. The follow-up approaches introduce the geometry of the scene for synthesizing novel views through approximate representations such as light fields [Levoy and Hanrahan 1996], lumigraphs [Gortler et al. 1996], plenoptic functions [McMillan and Bishop 1995] and layered depth images [Shade et al. 1998]. Chai et al. [2000] study the minimum sampling needed for light field rendering and also show that depth information enables better view synthesis with sparse viewpoints. McMillan Jr [1997] and Mark [1999] introduce depth image based rendering (DIBR) to synthesize new views. Multiple variants of DIBR [Chaurasia et al. 2013; Kanchana et al. 2022; Sun et al. 2010] find use in various applications such as 3D-TV [Fehn 2004] and free-viewpoint video [Carranza et al. 2003; Collet et al. 2015; Smolic et al. 2006]. Ramamoorthi [2023] conducts a detailed survey on classical work for novel view synthesis.

With the advent of deep learning, volumetric models utilize the power of learning by training the model on a large dataset of multi-view images. While the early approaches predict volumetric representations in each of the target views [Flynn et al. 2016; Kalantari et al. 2016], latter approaches predict a single volumetric representation and warp the representation to the target view while rendering [Mildenhall et al. 2019; Penner and Zhang 2017; Shih et al. 2020; Srinivasan et al. 2019; Zhou et al. 2018]. However, these approaches employ discrete depth planes and hence suffer from discretization

artifacts. The seminal work by Mildenhall et al. [2020] employ a continuous representation using multi-layer perceptrons (MLP). This started a new pathway in neural view synthesis. However, these models suffer from two major limitations, namely, the need for the dense sampling of input views and the large time required to render novel views from the given input views. The prior work that address these limitations can be broadly classified into three categories. In Sec. 2.1, we review various approaches in the literature to regularize the NeRF when training with sparse input views. We review the explicit radiance fields that aim at fast optimization and rendering in Sec. 2.2, and also review the recent work on regularizing explicit models for the few-shot setting. Finally, in Sec. 2.3, we review the generalized NeRFs that address both issues jointly.

2.1 Implicit Radiance Fields

There exists extensive literature on regularizing scene-specific NeRFs when training with sparse inputs. Hence, we further group these models based on their approaches.

Hand-Crafted Depth Priors: The prior work on sparse input NeRFs explore a plethora of hand-crafted priors on the NeRF rendered depth. RegNeRF [Niemeyer et al. 2022] imposes a smoothness constraint on the rendered depth maps. DS-NeRF [Deng et al. 2022] uses sparse depth provided by a Structure from Motion (SfM) module to supervise the NeRF estimated depth at sparse keypoints. ViP-NeRF [Somraj and Soundararajan 2023] augments the sparse-depth regularization of DS-NeRF with a regularization on the relative depth of objects by obtaining a prior on the visibility of objects. HG3-NeRF [Gao et al. 2024] uses sparse depth given by colmap to guide the sampling 3D points instead of supervising the NeRF rendered depth. While these priors are more robust across different scenes, they do not exploit the power of learning.

Deep Learning Based Depth Priors: There exist multiple models that utilize the advances in dense depth estimation using deep neural networks. DDP-NeRF [Roessle et al. 2022] extends DS-NeRF by employing a CNN to complete 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. However, the deep-learning based priors require pre-training on a large dataset and may suffer from generalization issues when obtaining the prior on unseen test scenes. Our work obtains depth supervision by harnessing the power of learning without suffering from generalization issues by employing and training augmented models on the given scene alone.

View Hallucination based Methods: Another line of regularization based approaches simulate dense sampling by hallucinating new viewpoints and regularizing 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]. Instead of sampling new viewpoints

randomly, FlipNeRF [Seo et al. 2023a] utilizes ray reflections to determine new viewpoints. Deceptive-NeRF [Liu et al. 2023] and ReconFusion [Wu et al. 2024] employ a diffusion model to generate images in hallucinated views and use the generated views in addition to the input views to train the NeRF. However, supervision with generative models could lead to content hallucinations, leading to poor fidelity [Lee and Lee 2024].

Other regularizations: A few models also explore regularizations other than depth supervision and view hallucinations. FreeNeRF [Yang et al. 2023] and MI-MLP-NeRF [Zhu et al. 2024] regularize the NeRF by modifying the inputs. Specifically, FreeNeRF anneals the frequency range of positional encoded NeRF inputs as the training progresses, and MI-MLP-NeRF adds the 5D inputs to every layer of the NeRF MLP. MixNeRF [Seo et al. 2023b] models the volume density along a ray as a mixture of Laplacian distributions. Philip and Deschaintre [2023] scale the gradients corresponding to 3D points close to the camera when sampling the 3D points in inverse depth to reduce floaters close to the camera. VDN-NeRF [Zhu et al. 2023b] on the other hand, aims to resolve shape-radiance ambiguity in the case of dense input views. However, these approaches are designed for specific cases and are either sub-optimal or do not extend to more recent radiance field models.

2.2 Explicit Radiance Fields

The NeRF takes a long time to optimize and render novel views due to the need to query the NeRF MLP hundreds of times to render a single pixel. Hence, a common approach to fast optimization and rendering is to reduce the time taken per query. Early works such as PlenOctress [Yu et al. 2021a] and KiloNeRF [Reiser et al. 2021] focus on improving only the rendering time by baking the trained NeRF into an explicit structure such as Octrees or thousands of tiny MLPs. PlenOxels [Fridovich-Keil et al. 2022] and DVGO [Sun et al. 2022] reduce the optimization time by directly optimizing voxel grids, but at the cost of large memory requirements to store the voxel grids. TensorRF [Chen et al. 2022b] and K-Planes [Fridovich-Keil et al. 2023] reduce the memory consumption using factorized tensors that exploit the spatial correlation of the radiance field. Alternately, iNGP [Müller et al. 2022] and ZipNeRF [Barron et al. 2023] employ multi-resolution hash-grids to reduce the memory consumption. Recently, 3DGS [Kerbl et al. 2023] propose an alternative volumetric model for real-time rendering of novel views. Specifically, 3DGS employs 3D Gaussians to represent the scene and renders a view by splatting the Gaussians onto the corresponding image plane. While the above methods enable fast optimization and rendering, their performance still reduces significantly with fewer input views.

Recently, there is increasing interest in regularizing explicit models to learn with sparse inputs [Li et al. 2024; Yang et al. 2024]. However, the regularizations designed in these models are limited to a specific explicit radiance field and do not generalize to other explicit models. For example, ZeroRF [Shi et al. 2024] imposes a deep image prior [Ulyanov et al. 2018] on the components of the TensorRF [Chen et al. 2022b] model. FSGS [Zhu et al. 2023a] and SparseGS [Xiong et al. 2023] improve the performance of 3DGS [Kerbl et al. 2023] in the sparse input case by improving the initialization of the 3D Gaussian point cloud and pruning Gaussians responsible for floaters

respectively. On the other hand, our approach of regularizing with simpler solutions is applicable to various explicit models, such as TensorRF, iNGP and ZipNeRF as well as to implicit models such as NeRF.

Despite the recent work on sparse input 3DGS models, we do not explore designing augmentations for 3DGS. As noted in the recent literature, 3DGS mainly suffers from poor initialization with few input views [Chen et al. 2024]. We believe 3DGS requires a combination of good initialization and supervision from augmentations to learn from few input views. This necessitates a separate study on designing better initializations for 3DGS, which is beyond the scope of this work.

2.3 Generalized Sparse Input NeRF

Obtaining a volumetric model of a scene by optimizing the NeRF is a time-consuming process. In order to reduce the time required to obtain a volumetric model of a scene and learn with fewer input views, generalized NeRF models train a neural network on a large dataset of multi-view scenes that can be directly applied to a test scene without any optimization [Chen et al. 2021; Lee et al. 2023; Tancik et al. 2021]. Early pieces of work such as PixelNeRF [Yu et al. 2021b], GRF [Trevithick and Yang 2021], and IBRNet [Wang et al. 2021b] 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] incorporates cross-view knowledge into the features by constructing a 3D cost volume. However, the resolution of the 3D cost volume is limited by the available memory size, which limits the performance of MVS-NeRF [Lin et al. 2023]. On the other hand, SRF [Chibane et al. 2021] processes individual frame features in a pair-wise manner, and GNT [Wang et al. 2022] employs a transformer to efficiently incorporate cross-view knowledge.

NeuRay [Liu et al. 2022] and GeoNeRF [Johari et al. 2022] further improve the performance by employing visibility priors and a transformer respectively to effectively reason about the occlusions in the scene. More recent work such as GARF [Shi et al. 2022], DINER [Prinzler et al. 2023] and MatchNeRF [Chen et al. 2023] try to provide explicit knowledge about the scene geometry through depth maps and similarity of the projected features. 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 as an extreme case of the sparse input NeRF. However, the need for pre-training on a large dataset of scenes with multi-view images and generalization issues due to domain shift have motivated researchers to continue to be interested in regularizing scene-specific radiance fields.

3 RADIANCE FIELDS AND VOLUME RENDERING PRELIMINARIES

We first provide a brief recap of the radiance fields and volume rendering. We also describe the notation required for further sections. To render a pixel \mathbf{q} , we shoot the corresponding ray into the scene and sample 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. At every 3D point \mathbf{p}_i , the radiance field $\mathcal{F} = \mathcal{F}_1 \circ \mathcal{F}_2$ is queried to obtain a

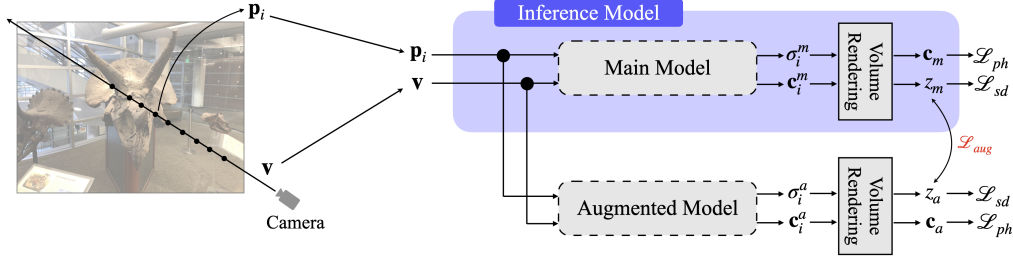


Fig. 2. Architecture of Simple-RF family of models. We train the augmented model that only learns simpler solutions in tandem with the main model. The augmented models learn better depth in certain regions, which is propagated to the main model through the depth supervision loss \mathcal{L}_{aug} . During inference, only the Main Model is employed.

view-independent volume density σ_i and a view-dependent color c_i as

$$\sigma_i, \mathbf{h}_i = \mathcal{F}_1(\mathbf{p}_i), \quad c_i = \mathcal{F}_2(\mathbf{h}_i, \mathbf{v}), \quad (1)$$

where \mathbf{v} is the viewing direction and \mathbf{h}_i is a latent feature of \mathbf{p}_i . Volume rendering is then applied along every ray to obtain the color for each pixel as $\mathbf{c} = \sum_{i=1}^N w_i c_i$, 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. 2022]. \mathcal{F}_1 and \mathcal{F}_2 are modelled differently for NeRF, TensorRF and ZipNeRF, and are trained using the photometric loss $\mathcal{L}_{ph} = \|\mathbf{c} - \hat{\mathbf{c}}\|^2$, where $\hat{\mathbf{c}}$ is the true color of \mathbf{q} .

4 METHOD

Learning a radiance field with sparse input views leads to overfitting on the input views with severe distortions in novel views. Our key observation is that most of the distortions are due to the sub-optimal use of the high capabilities of the radiance field model. Further, we find that reducing the capability of the radiance field helps constrain the model to learn only simpler solutions, which can provide better depth supervision in certain regions of the scene. However, the lower capability models are not optimal either since they cannot learn complex solutions where necessary. Our solution here is to use the higher capability model as the main model and employ the lower capability models as augmentations to provide guidance on where to use simpler solutions. The challenge is that, it is not known apriori where one needs to employ supervision from the augmented model. We determine the more accurate model among the main and augmented models in terms of the estimated depth and use the more reliable depth to supervise the other. We note that the augmented models are employed only during the learning phase and not during inference. Thus, there is no additional overhead during inference. The augmented models are similar to the main model, but we modify their parameters to reduce their capability, and train them in tandem with the main model.

To design the augmented models, we first analyze the shortcomings of the radiance field with sparse input views. Specifically, we determine the components of the model that cause overfitting with fewer input views, causing distortions in novel views. We then design the augmented models by reducing the model capability with respect to such components. Thus, designing the augmented models is non-trivial, and the design may need to be different for different radiance fields based on the architecture of the radiance fields and the distortions observed. Nonetheless, the core idea of designing augmentations by reducing their capability to learn simpler solutions is common across all radiance fields.

We discuss the design of augmentations for NeRF, TensorRF and ZipNeRF in Secs. 4.1 to 4.3 respectively. We describe our approach to determining the reliability of the depth estimates in Sec. 4.1.4. Finally, Sec. 4.4 summarizes all the loss functions used to train our full model. Fig. 2 shows the architecture of our family of simple radiance fields.

4.1 Simple-NeRF

We start by discussing the specific details of the NeRF that are relevant for the design of augmentations in Sec. 4.1.1, then analyze the shortcomings of the NeRF with sparse input views in Sec. 4.1.2 and finally discuss the design of augmentations in Sec. 4.1.3. Fig. 3 shows the detailed architecture of Simple-NeRF.

4.1.1 NeRF Preliminaries. The NeRF learns the radiance field \mathcal{F} using two neural networks $\mathcal{N}_1, \mathcal{N}_2$ and predicts the view-independent volume density σ_i and view-dependent color c_i as

$$\sigma_i, \mathbf{h}_i = \mathcal{N}_1(\gamma(\mathbf{p}_i, 0, l_p)); \quad c_i = \mathcal{N}_2(\mathbf{h}_i, \gamma(\mathbf{v}, 0, l_v)), \quad (3)$$

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

$$\gamma(x, d_1, d_2) = [\sin(2^{d_1}x), \cos(2^{d_1}x), \dots, \sin(2^{d_2-1}x), \cos(2^{d_2-1}x)] \quad (4)$$

is the positional encoding from frequency d_1 to d_2 . l_p and l_v are the highest positional encoding frequencies for \mathbf{p}_i and \mathbf{v} respectively. When $d_1 = 0$, x is concatenated to the positional encoding features in Eq. (4). NeRF circumvents the need for the dense sampling of 3D points along a ray by employing two sets of MLPs, a coarse NeRF and a fine NeRF, both trained using \mathcal{L}_{ph} . 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

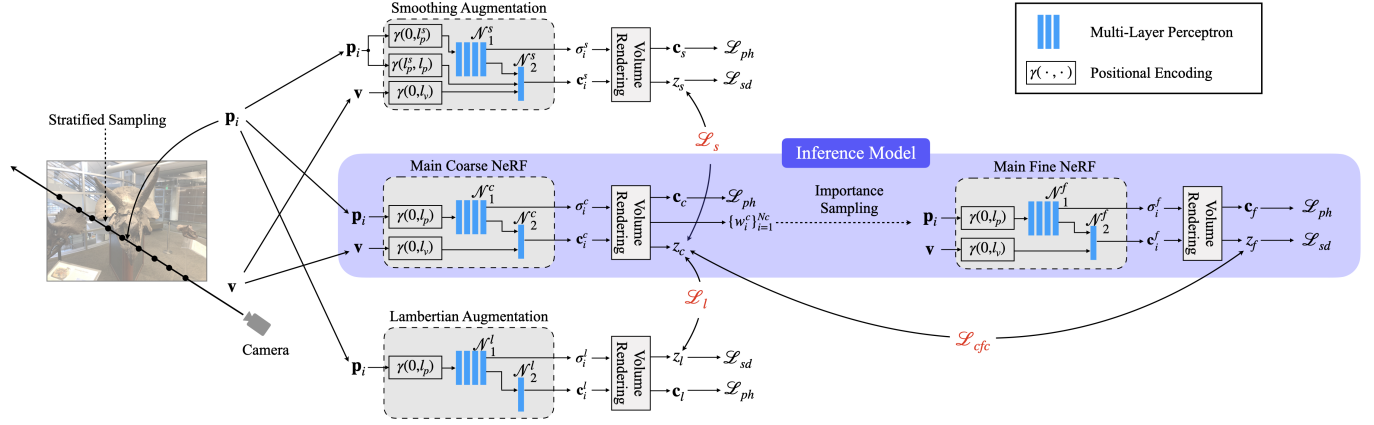
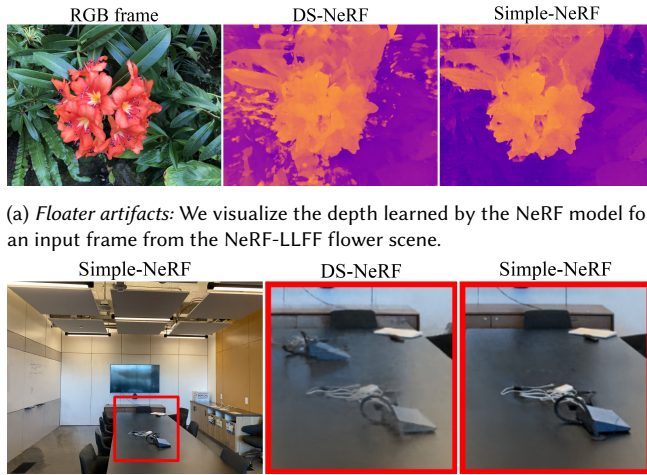


Fig. 3. Architecture of Simple-NeRF. We train two augmented NeRF models in tandem with the main NeRF. In smoothing augmentation, we reduce the positional encoding frequencies that are input to \mathcal{N}_1^s and concatenate the remaining frequencies to the input of \mathcal{N}_2^s . For Lambertian augmentation, we ask \mathcal{N}_2^l to output the color based on position alone, independent of the viewing direction. We add depth supervision losses \mathcal{L}_s and \mathcal{L}_l between the coarse 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.



(a) *Floater artifacts*: We visualize the depth learned by the NeRF model for an input frame from the NeRF-LLFF flower scene.

(b) *Duplication artifacts*: To visualize the duplication artifacts that arise due to the shape-radiance ambiguity in sparse-input NeRF, we render an input frame by only changing the viewing direction. This is an example from the NeRF-LLFF room scene.

Fig. 4. **Failure of sparse-input NeRF**: We show two shortcomings of the NeRF when trained with two input views on the NeRF-LLFF dataset. In Fig (a), we observe the floaters as small orange regions in the depth map. In Fig (b), we observe the duplication of the object on the table caused by the NeRF trying to blend the input images. Simple-NeRF introduces regularizations to mitigate these distortions as seen in both figures. We note that the models used to synthesize the above images include the sparse depth supervision (Sec. 4.4).

localized based on the predictions of the coarse NeRF. Since the scene geometry is mainly learned by the coarse NeRF, we add the augmentations only to the coarse NeRF.

4.1.2 Analysing Sparse Input NeRF. With sparse input views, we find that two components of the NeRF, namely positional-encoding and view-dependent radiance, can cause overfitting, leading to distortions in novel views. Both positional encoding and view-dependent radiance are elements designed to increase the capability of the NeRF to explain different complex phenomena. For example, the former helps in learning thin objects against a farther background, and the latter helps in learning 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. Specifically, the high positional encoding degree leads to undesired depth discontinuities in smooth-depth regions resulting in floater artifacts, where a part of an object is broken away from it and floats freely in space [Barron et al. 2022], as shown in Fig. 4a. The view-dependent radiance causes shape-radiance ambiguity, leading to duplication artifacts in the novel views as shown in Fig. 4b. With sparse input views, the NeRF explains different objects by varying the color of the same 3D points based on the viewing direction, thereby giving us an illusion of the object without learning the correct geometry of the object. This is, in a way, similar to the illusion created by lenticular images and can be observed better in the supplementary video. Our augmentations consist of reducing the capability of the NeRF model with respect to the positional encoding and view-dependent radiance to obtain better depth supervision.

4.1.3 Design of Augmentations. We employ two augmentations, one each for regularizing positional encoding and view-dependent radiance, which we describe in the following. We refer to the two augmentations as smoothing and Lambertian augmentations, respectively.

Smoothing augmentation: 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 in volume density between the

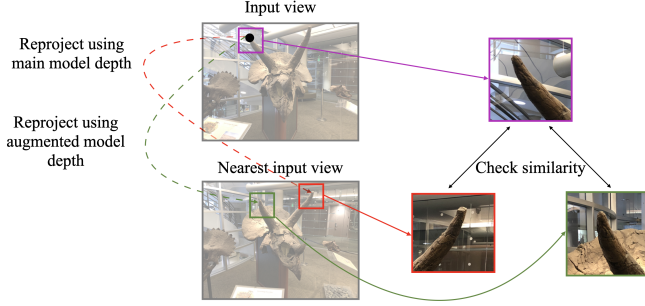


Fig. 5. **Determining the reliability of depths for supervision:** We choose the depth that has higher similarity, with respect to the patches reprojected to the nearest input view, to supervise the other model (Sec. 4.1.4). The patches are only representative and are not to scale.

two points in \mathbb{R}^3 as a smooth function in $\mathbb{R}^{3(2l_p+1)}$. We reduce the depth discontinuities, which are caused by discontinuities in the volume density, by reducing the highest positional encoding frequency for \mathbf{p}_i to $l_p^s < l_p$ as

$$\sigma_i, \mathbf{h}_i = \mathcal{N}_1^s(\gamma(\mathbf{p}_i, 0, l_p^s)), \quad (5)$$

where \mathcal{N}_1^s 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 respective locations as binary masks and use only the reliable depth estimates from one model to supervise the other model, as we explain in Sec. 4.1.4.

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{N}_2 as

$$\mathbf{c}_i = \mathcal{N}_2^s(\mathbf{h}_i, \gamma(\mathbf{p}_i, l_p^s, l_p), \gamma(\mathbf{v}_i, 0, l_v)). \quad (6)$$

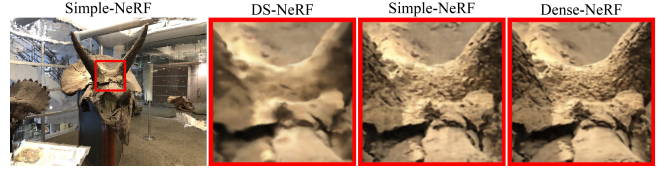
Note that \mathbf{h}_i already includes the low-frequency positional encoding components of \mathbf{p}_i .

Lambertian Augmentation: 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 guard 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 as

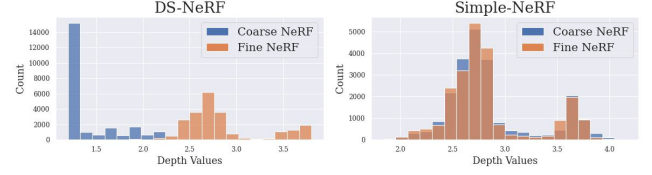
$$\sigma_i, \mathbf{h}_i = \mathcal{N}_1^l(\gamma(\mathbf{p}_i, 0, l_p)); \quad \mathbf{c}_i = \mathcal{N}_2^l(\mathbf{h}_i), \quad (7)$$

We note that while the augmented model is more accurate in Lambertian regions, the main model is better equipped to handle specular objects. We determine the respective locations as we explain in the following and use only the reliable depth estimates for supervision.

4.1.4 Determining Reliable Depth Estimates. Let the depths estimated by the main and augmented models for pixel \mathbf{q} be z_m and z_a respectively. We now seek to determine the more accurate depth among the two. Fig. 5 shows our approach to determining the reliability of the estimated depth. Specifically, we reproject a $k \times k$ patch around \mathbf{q} to the nearest training view using both z_m and z_a .



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



(b) Histogram of depth values predicted by the coarse and fine NeRF models for the image patch shown in Fig (a).

Fig. 6. **Ineffective hierarchical sampling in sparse-input NeRF:** Fig (b) shows that the coarse and fine models in the NeRF converge to different depth estimates when training with sparse input views. This leads to ineffective hierarchical sampling, resulting in blurry predictions in Fig (a). By predicting consistent depth estimates with the help of \mathcal{L}_{fc} , Simple-NeRF predicts consistent depth estimates leading to sharp reconstructions. We note that the models used to synthesize the above images include the sparse depth supervision (Sec. 4.4).

We then compute the similarity of the reprojected patch with the corresponding patch in the first image using the mean squared error (MSE) in intensities. We 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 only if its MSE is also less than e_τ . If e_m and e_a are the reprojection MSE corresponding to z_m and z_a respectively, we compute a mask m_a that indicates where the augmented model is more reliable as

$$m_a = \begin{cases} 1 & \text{if } (e_a \leq e_m) \text{ and } (e_a \leq e_\tau) \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

We compute the reliability mask m_m for the main model similarly. We now impose the depth supervision as

$$\mathcal{L}_{aug} = \mathbb{1}_{\{m_a=1\}} \odot \|z_m - \mathcal{V}(z_a)\|^2 + \mathbb{1}_{\{m_m=1\}} \odot \|\mathcal{V}(z_m) - z_a\|^2, \quad (9)$$

where \odot denotes element-wise product, $\mathbb{1}$ is the indicator function and \mathcal{V} is the stop-gradient operator. We impose two losses, \mathcal{L}_s for the smoothing augmentation and \mathcal{L}_l for the Lambertian augmentation. The final depth supervision loss is the sum of the two losses.

For specular regions, the intensities of the reprojected patches may not match, leading to the masks being zero. This only implies supervision for fewer pixels and not supervision with incorrect depth estimates.

4.1.5 Hierarchical Sampling. Since multiple solutions can explain the observed images in the few-shot setting, the coarse and fine MLP of the NeRF may converge to different depth estimates for a

given pixel as shown in Fig. 6b. Thus, dense sampling may not be employed around the region where the 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. 6a. 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

$$\mathcal{L}_{cfc} = \mathbb{1}_{\{m_f=1\}} \odot \|z_c - \mathcal{V}(z_f)\|^2 + \mathbb{1}_{\{m_c=1\}} \odot \|\mathcal{V}(z_c) - z_f\|^2, \quad (10)$$

where the masks m_c and m_f are determined as we describe in Sec. 4.1.4.

Apart from enforcing consistency between the coarse and fine NeRF models, \mathcal{L}_{cfc} provides two additional benefits. Without \mathcal{L}_{cfc} , the augmentations need to be imposed on the fine NeRF as well, leading to an increase in the training time and memory requirements. Secondly, if one of the coarse or fine NeRFs converges to the correct solution, \mathcal{L}_{cfc} helps quickly convey the knowledge to the other NeRF, thereby facilitating faster convergence.

4.2 Simple-TensorRF

We first provide a brief overview of TensorRF [Chen et al. 2022b] in Sec. 4.2.1 and also describe the notation required to explain the design of our augmentations. We discuss the distortions observed with sparse input TensorRF in Sec. 4.2.2 and then discuss the design of augmentations in Sec. 4.2.3.

4.2.1 TensorRF Preliminaries. TensorRF models the fields \mathcal{F}_1 and \mathcal{F}_2 with a tensor \mathcal{G}_1 and a tiny MLP \mathcal{N}_2 , respectively. The 3D tensor \mathcal{G}_1 is factorized as the sum of outer products of 1D vectors \mathbf{v} and 2D matrices \mathbf{M} . Specifically, \mathcal{G}_1 consists of two 3D tensors, \mathcal{G}_σ to learn the volume density and \mathcal{G}_c to learn the latent features of the color as

$$\mathcal{G}_\sigma = \sum_{r=1}^{R_\sigma} \mathbf{v}_{\sigma,r}^X \circ \mathbf{M}_{\sigma,r}^{YZ} + \sum_{r=1}^{R_\sigma} \mathbf{v}_{\sigma,r}^Y \circ \mathbf{M}_{\sigma,r}^{XZ} + \sum_{r=1}^{R_\sigma} \mathbf{v}_{\sigma,r}^Z \circ \mathbf{M}_{\sigma,r}^{XY}, \quad (11)$$

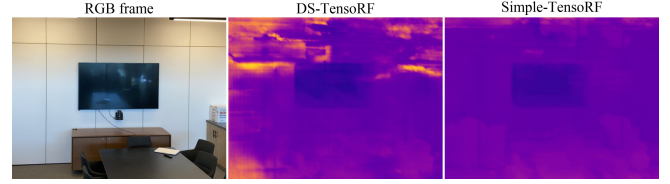
$$\mathcal{G}_c = \sum_{r=1}^{R_c} \mathbf{v}_{c,r}^X \circ \mathbf{M}_{c,r}^{YZ} \circ \mathbf{a}_{3r-2} + \sum_{r=1}^{R_c} \mathbf{v}_{c,r}^Y \circ \mathbf{M}_{c,r}^{XZ} \circ \mathbf{a}_{3r-1} \quad (12)$$

$$+ \sum_{r=1}^{R_c} \mathbf{v}_{c,r}^Z \circ \mathbf{M}_{c,r}^{XY} \circ \mathbf{a}_{3r},$$

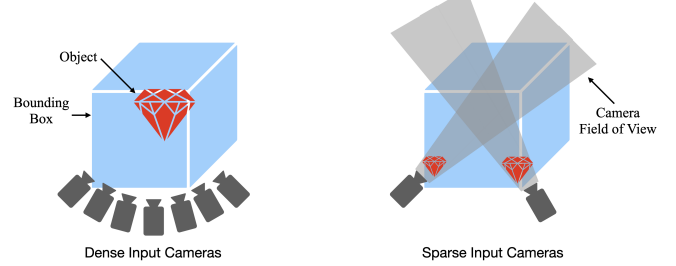
$$\sigma_i = \text{sigmoid}(\mathcal{G}_\sigma(\mathbf{p}_i)); \quad \mathbf{h}_i = \mathcal{G}_c(\mathbf{p}_i), \quad (13)$$

where \circ represents the outer product, and R_σ and R_c represent the number of components in the factorization of sigma and color grids, respectively. $\mathcal{G}(\mathbf{p}_i)$ is obtained by trilinearly interpolating \mathcal{G} at \mathbf{p}_i . $\mathbf{v}^X \in \mathbb{R}^I$ and $\mathbf{M}^{YZ} \in \mathbb{R}^{J \times K}$ represent the vector along the x-axis and the matrix in the yz-plane respectively and so on, where I, J and K represent the resolution of the tensor in the x, y and z dimensions respectively. Thus, the total number of voxels in the tensor is $N_{vox} = I \times J \times K$. Note that \mathcal{G}_c uses an additional vector $\mathbf{a}_r \in \mathbb{R}^D$ to learn appearance as a latent feature of dimension D .

TensorRF assumes that the entire scene is contained within a 3D bounding box \mathbf{b} as shown in Fig. 7b, whose vertices are given



(a) *Floater artifacts*: We visualize the depth learned by the TensorRF model for an input frame from the NeRF-LLFF room scene.



(b) *Objects close to camera*: We illustrate TensorRF incorrectly placing objects close to the cameras using a toy example.

Fig. 7. Failure of sparse-input TensorRF: We show the two shortcomings of TensorRF when trained with few input views. In Fig (a), the orange regions indicate the floaters. For reference, we also show the depth learned by Simple-TensorRF, which is free from floaters. We note that the models used to synthesize these images include the sparse depth supervision (Sec. 4.4). In Fig (b), the image on the left depicts the true scene, which can be accurately learned by the TensorRF model provided with dense input views. The image on the right illustrates how TensorRF can incorrectly place the objects yet perfectly reconstruct the input views, when training with few input views.

by $\{(b_{x_1}, b_{x_2}), (b_{y_1}, b_{y_2}), (b_{z_1}, b_{z_2})\}$. TensorRF handles unbounded forward-facing scenes by transforming the space into normalized device coordinates (ndc) similar to the NeRF. The coarse to fine training is implemented by using lower resolution tensors \mathcal{G}_σ and \mathcal{G}_c during the initial stages of the optimization and gradually increasing the resolution as the training progresses. Finally, the color at \mathbf{p}_i is obtained using the tiny MLP \mathcal{N}_2 as

$$\mathbf{c}_i = \mathcal{N}_2(\mathbf{h}_i, \gamma(\mathbf{v}, 0, l_v)), \quad (14)$$

where γ is the positional encoding described by Eq. (4). Thus, we note that $\mathcal{F}_2 = \gamma \circ \mathcal{N}_2$. The color of the pixel is then obtained through volume rendering using σ_i and \mathbf{c}_i as in Sec. 3. For further details, we refer the readers to TensorRF [Chen et al. 2022b].

4.2.2 Analysing Sparse Input TensorRF. When training a TensorRF model with sparse input views, we find that three of its components cause overfitting, leading to distortions in novel views. Employing a higher resolution tensor \mathcal{G}_σ with a large number of components R_σ allows the TensorRF to learn sharp depth edges, but results in undesired depth discontinuities in smooth regions causing floaters as shown in Fig. 7a. Further, the large bounding box \mathbf{b} allows the TensorRF to handle objects that are truly very close to the camera. On account of large distances between cameras when only a few input views are available, it may be possible to place objects close to one camera such that they are out of the field of view of the other

cameras, even for objects visible in multiple input views. Specifically, TensorRF learns multiple copies of the same object, each visible in only one input view, thereby explaining the observations without learning the geometry of the objects as shown in Fig. 7b. We design the augmentations to reduce the capability of the TensorRF model with respect to these three components.

4.2.3 Design of Augmentations. Employing a high-resolution and high-rank tensor \mathcal{G}_σ enables TensorRF to learn significantly different σ values for two nearby points in \mathbb{R}^3 leading to undesired depth discontinuities in smooth regions. We constrain the augmented TensorRF to learn only smooth and continuous depth surfaces by reducing the number of components to $R_\sigma^s < R_\sigma$ and also reducing the number of voxels of \mathcal{G}_σ from $N_{vox}^s < N_{vox}$. We note that modifying only one of these components is insufficient to achieve the desired smoothing. For example, only reducing the resolution of the grid allows TensorRF to learn sharp changes in σ at the voxel edges, leading to block artifacts. On the other hand, only reducing the number of components allows TensorRF to learn sharp changes in σ on account of the high-resolution grid.

Further, we find that reducing R_σ to R_σ^s leads to the augmented TensorRF learning cloudy volumes instead of hard object surfaces. We encourage the augmented TensorRF to learn hard surfaces by employing a mass concentration loss that minimizes the entropy of mass, grouped into N_{mc} intervals as

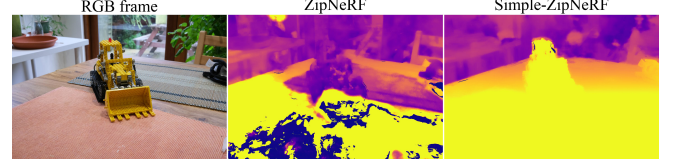
$$\mathcal{L}_{mc} = H \left(\left\{ \sum_{i=(j-1)(N/N_{mc})+1}^{j(N/N_{mc})} w_i \right\}_{j=1}^{N_{mc}} \right), \quad (15)$$

where $H(w_1, w_2, \dots, w_n) = -\sum_{i=1}^n w_i \log w_i$ is the entropy operator, N is the number of 3D points \mathbf{p}_i along a ray and w_i is the weight corresponding to \mathbf{p}_i as described in Eq. (2).

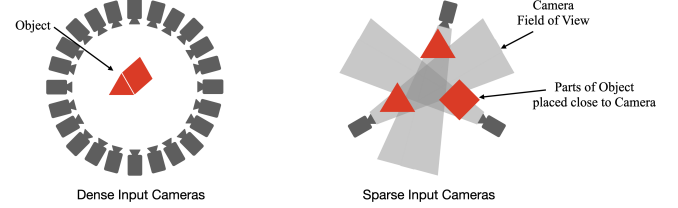
Objects that are incorrectly placed close to the camera due to a large bounding box are typically smooth in depth and hence are not mitigated by the above augmentation. We design a second augmentation to mitigate such distortions by reducing the size of the bounding box \mathbf{b} along the z-axis by increasing b_{z_1} to $b_{z_1}^s$. We note that replicating the same in the main TensorRF model could lead to distortions in objects that are truly close to the camera. In practice, we find that including both the augmentations in a single augmented TensorRF model works reasonably well, and hence we employ a single augmented model. We then use the reliable depth estimates from the augmented model to supervise the main model as in Eq. (9).

4.3 Simple-ZipNeRF

ZipNeRF [Barron et al. 2023] integrates the iNGP model [Müller et al. 2022], which achieves significant improvements in optimization and rendering times, with the anti-aliasing ability of MipNeRF [Barron et al. 2021] and the ability to handle unbounded 360° scenes of MipNeRF360 [Barron et al. 2022]. Our contributions to enable the training of ZipNeRF with sparse input views are mainly with respect to the components of the iNGP model, and hence, we believe that the augmentations designed for ZipNeRF are relevant to iNGP as well. We discuss the specific components of ZipNeRF that are relevant in our augmentations in Sec. 4.3.1, analyze the limitations of ZipNeRF



(a) *Floater artifacts*: We visualize the depth learned by the ZipNeRF model for an input frame from the MipNeRF360 kitchen scene.



(b) *Objects close to camera*: We illustrate ZipNeRF incorrectly placing objects close to the cameras using a toy example.

Fig. 8. We show two shortcomings of ZipNeRF when trained with few input views. In Fig (a), while the RGB frame for an input view is reconstructed perfectly, we observe floaters in the depth image, shown by the dark-blue regions. For reference, we also show the depth learned by Simple-ZipNeRF, which is free from floaters and better reconstructs the scene. In Fig (b), the image on the left depicts the true scene, which can be accurately learned by the ZipNeRF model provided with dense input views. The image on the right illustrates how the sparse-input ZipNeRF model can incorrectly place parts of the object close to the cameras, yet perfectly reconstruct the input views.

with sparse input views in Sec. 4.3.2 and then discuss the design of our augmentations in Sec. 4.3.3.

4.3.1 ZipNeRF Preliminaries. ZipNeRF employs a multi-resolution grid and a hash function that maps every vertex of the grid to an entry in a hash table. The hash table contains the latent features representing the volume density and the radiance. Concretely, given a point $\mathbf{p}_i \in \mathbb{R}^3$, the vertices of the voxel enclosing \mathbf{p}_i are mapped to an entry in a hash table of length T through the use of hash function \mathcal{H}_1 as,

$$\mathcal{H}_1(\mathbf{p}) = \left(\bigoplus_{j=1}^3 p_j \pi_j \right) \bmod T, \quad (16)$$

where \oplus denotes the bit-wise XOR operation, π_j is a prime number, and p_j is the j -th coordinate of \mathbf{p} . The feature vectors corresponding to the eight vertices of the voxel are trilinearly interpolated. The same procedure is repeated for every level of the multi-resolution grid, and the corresponding interpolated features are concatenated to form the latent feature $\mathcal{H}_1(\mathbf{p}_i)$. Two tiny MLPs are employed to decode $\mathcal{H}_1(\mathbf{p}_i)$ into the volume density and the radiance as

$$\sigma_i, \mathbf{h}_i = \mathcal{N}_1(\mathcal{H}_1(\mathbf{p}_i)); \quad \mathbf{c}_i = \mathcal{N}_2(\mathbf{h}_i, \gamma(\mathbf{v}, 0, l_o)), \quad (17)$$

where γ is the positional encoding as defined in Eq. (4). We note that \mathcal{F}_1 and \mathcal{F}_2 in Eq. (1) are thus represented as $\mathcal{F}_1 = \mathcal{H}_1 \circ \mathcal{N}_1$ and $\mathcal{F}_2 = \gamma \circ \mathcal{N}_2$. The color of the pixel is then obtained through volume rendering using σ_i and \mathbf{c}_i as in Sec. 3. Note that multiple vertices in the grid could map to the same entry in the hash table at every level.

iNGP and ZipNeRF rely on the MLP \mathcal{N}_1 to resolve such collisions based on multi-resolution features. Unbounded scenes are handled by employing a contraction function that maps the distance along the ray from $z \in [z_{\text{near}}, z_{\text{far}}]$ to a normalized distance $s \in [0, 1]$. The 3D points \mathbf{p}_i are then sampled in s -domain. For further details, we refer the readers to iNGP [Müller et al. 2022] and ZipNeRF [Barron et al. 2023].

4.3.2 Analysing Sparse Input ZipNeRF. We find that two components of ZipNeRF tend to cause overfitting when trained with sparse input views, leading to distortions in novel views. Firstly, employing a hash table with a large size T enables ZipNeRF to learn sharp depth edges, but introduces undesired depth discontinuities in smooth regions, causing floaters as shown in Fig. 8a. Secondly, since ZipNeRF handles unbounded 360° scenes, it learns the radiance field over the entire 3D space. Similar to TensorRF, ZipNeRF tends to incorrectly place multiple copies of objects very close to the camera without learning the correct geometry as shown in Fig. 8b. Thus, we design the augmentations to reduce the capability of the ZipNeRF model with respect to these two components.

4.3.3 Design of Augmentations. Employing a hash table of larger size T allows ZipNeRF to avoid collisions and not share features across multiple 3D points. This enables ZipNeRF to map two nearby points in \mathbb{R}^3 to two independent entries in the hash table, thus mapping them to two farther away points in the latent feature space. This allows the MLP \mathcal{N}_1 to learn discontinuities in the volume density, resulting in sharp depth edges. We encourage the augmented model to share features across more 3D points by reducing the size of the hash table to $T^s < T$.

To mitigate the objects being placed close to the camera incorrectly, we cannot reduce the size of the bounding box as in TensorRF, since ZipNeRF handles unbounded scenes. We achieve a similar effect by sampling the 3D points \mathbf{p}_i along the ray in s -domain in the range $s \in [s_{\text{near}}, 1]$ instead of $s \in [0, 1]$. This ensures that the objects are placed at least at a certain distance away from the camera in the augmented model. However, we note that employing the above modification in the main model is detrimental to learning or rendering any objects that are truly close to the camera. In practice, we find that including both the augmentations in a single augmented ZipNeRF model works reasonably well, and hence, we employ a single augmented model. We then use these depth estimates as in Eq. (9).

4.4 Overall Loss

Let \mathcal{L}_m denote the combination of the losses employed by the corresponding radiance fields. For example, while the NeRF employs only the photometric loss \mathcal{L}_{ph} , TensorRF employs a total variation regularization in addition to \mathcal{L}_{ph} . We refer the readers to the corresponding papers for the details of all the losses imposed. We impose all such losses on the augmented models as well and denote them by \mathcal{L}_a . In addition, we also include the sparse depth loss on both the main and augmented models as,

$$\mathcal{L}_{sd} = \|z_m - \hat{z}\|^2 + \|z_a - \hat{z}\|^2, \quad (18)$$

where z_m and z_a are the depths obtained from the main and augmented models respectively, and \hat{z} is the sparse depth given by the

Table 1. 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

SfM model [Deng et al. 2022]. Our final loss is a combination of all the losses as,

$$\mathcal{L} = \lambda_m \mathcal{L}_m + \lambda_a \mathcal{L}_a + \lambda_{sd} \mathcal{L}_{sd} + \lambda_{aug} \mathcal{L}_{aug} + \lambda_{cfc} \mathcal{L}_{cfc} + \lambda_{mc} \mathcal{L}_{mc}, \quad (19)$$

where \mathcal{L}_{cfc} and \mathcal{L}_{mc} are respectively imposed for the main NeRF and augmented TensorRF models only, and $\lambda_m, \lambda_a, \lambda_{sd}, \lambda_{aug}, \lambda_{cfc}$ and λ_{mc} are hyper-parameters.

5 EXPERIMENTAL SETUP

5.1 Datasets

We evaluate the performance of our models on four popular datasets, namely NeRF-LLFF [Mildenhall et al. 2019], RealEstate-10K [Zhou et al. 2018], MipNeRF360 [Barron et al. 2022] and NeRF-Synthetic [Mildenhall et al. 2020]. We assume the camera parameters are known for the input images, since in applications such as robotics or extended reality, external sensors or a pre-calibrated set of cameras may provide the camera poses.

NeRF-LLFF dataset contains eight real-world forward-facing scenes typically consisting of an object at the centre against a complex background. Each scene contains a varying number of images ranging from 20 to 60, each with a spatial resolution of 1008×756 . Following prior work [Niemeyer et al. 2022], we use every 8th view as the test view and uniformly sample 2, 3 or 4 input views from the remaining.

RealEstate-10K dataset contains a large number of real-world forward-facing scenes, from which we select 5 test scenes for our experiments. We include both indoor and unbounded outdoor scenes and select 50 temporally continuous frames from each scene. The frames have a spatial resolution of 1024×576 . Following prior work [Somraj and Soundararajan 2023], we reserve every 10th frame for training and choose 2, 3 or 4 input views among them. In the remaining 45 frames, we use those frames that are not very far from the input frames for testing. Specifically, we choose all the frames between the training views that correspond to interpolation and five frames on either side that correspond to extrapolation. Tab. 1 shows the train and test frame numbers we use for the three different settings.

MipNeRF360 dataset contains seven publicly available unbounded 360° real-world scenes including both indoor and outdoor scenes. Each scene contains 100 to 300 images. The four indoor scenes have a spatial resolution of approximately 1560×1040 , and the three outdoor scenes have an approximate spatial resolution of 1250×830 . Following prior work [Barron et al. 2023], we reserve every 8th view for testing and uniformly sample 12, 20 and 36 input views from the remaining. We use more input views on this dataset as compared to the other datasets owing to the significantly larger fields of view.

NeRF-Synthetic dataset contains eight bounded 360° synthetic scenes, each containing 100 train images and 200 test images. All the scenes have a spatial resolution of 800×800 . For training, we uniformly sample 4, 8 and 12 input views from the training set and test on all 200 test images.

5.2 Evaluation measures

We quantitatively evaluate the predicted frames from various models using the peak signal-to-noise ratio (PSNR), structural similarity (SSIM) [Wang et al. 2004] and LPIPS [Zhang et al. 2018] measures. For LPIPS, we use the v0.1 release with the AlexNet [Krizhevsky et al. 2012] backbone, as suggested by the authors. We also employ depth mean absolute error (MAE) to evaluate the models on their ability to predict absolute depth in novel views. In addition, we also evaluate the models with regard to their ability to predict better relative depth using the spearman rank order correlation coefficient (SROCC). Obtaining better relative depth might be more crucial in downstream applications such as 3D scene editing. Since the ground truth depth is not provided in the datasets, we train NeRF and ZipNeRF models with dense input views on forward-facing and 360° datasets respectively and use their depth predictions as pseudo ground truth. On the NeRF-LLFF and MipNeRF360 datasets, we normalize the predicted depths by the median ground truth depth, since the scenes have different depth ranges. With very few input views on forward-facing datasets, the test views could contain regions that are not visible in the input views, and hence, we also evaluate both the view synthesis and depth performance in visible regions only. 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. On the other hand, the input views cover most of the scene in the 360° datasets, and hence we evaluate the performance on full frames. We do not evaluate the rendered depth on the NeRF-Synthetic dataset since the depth estimated with the dense input ZipNeRF is unreliable, especially in the white background regions, and the ground truth depth is not provided in the dataset either.

6 EXPERIMENTAL RESULTS

We present the main results of our work with Simple-NeRF in Sec. 6.1 and then show the extension of our ideas to explicit models in Secs. 6.2 and 6.3.

6.1 Simple-NeRF

6.1.1 Comparisons. We evaluate the performance of our Simple-NeRF on the two forward-facing datasets only since the NeRF does not natively support unbounded 360 scenes. We evaluate the performance of our model against various sparse input NeRF models. We compare with DS-NeRF [Deng et al. 2022], DDP-NeRF [Roessle et al. 2022] and RegNeRF [Niemeyer et al. 2022] which regularize the depth estimated by the NeRF. We also evaluate DietNeRF [Jain et al. 2021] and InfoNeRF [Kim et al. 2022] that regularize the NeRF in hallucinated viewpoints. We also 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.

6.1.2 Implementation details. We develop our code in PyTorch and on top of DS-NeRF [Deng et al. 2022]. We employ the 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 similar orders. For the first 10k iterations of the training, we only impose \mathcal{L}_m , \mathcal{L}_a and \mathcal{L}_{sd} . \mathcal{L}_{aug} and \mathcal{L}_{cfc} are imposed after 10k iterations. We set the hyper-parameters as follows: $l_p = 10$, $l_v = 4$, $l_p^s = 3$, $k = 5$, $e_r = 0.1$, $\lambda_m = \lambda_a = 1$, $\lambda_{sd} = \lambda_{aug} = \lambda_{cfc} = 0.1$ and $\lambda_{mc} = 0$. The network architecture is exactly the same as DS-NeRF. For the augmented models, we only change the input dimension of the MLPs \mathcal{N}_1 and \mathcal{N}_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 2080 Ti GPU for 100k iterations.

6.1.3 Quantitative and Qualitative Results. Tabs. 2 and 3 show the view-synthesis performance of Simple-NeRF and other prior art on NeRF-LLFF and RealEstate-10K datasets respectively. We find that Simple-NeRF 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 NeRF-LLFF dataset. Fig. 9 shows predictions of various models on an example scene from the RealEstate-10K dataset, where we observe that Simple-NeRF is the best in reconstructing the novel view. Figs. 10 to 15 show more comparisons on both datasets with 2, 3, and 4 input views. Further, Simple-NeRF improves significantly in estimating the depth of the scene as seen in Tab. 4 and Fig. 16. We provide video comparisons in the supplementary.

We note that the quantitative results in Tabs. 2 and 3 differ from the values reported in ViP-NeRF [Somraj and Soundararajan 2023] on account of the following two differences. Firstly, the quality evaluation metrics are computed on full frames in ViP-NeRF. However, we exclude the regions not seen in the input views as explained in Sec. 5.2. Secondly, while we use the same train set as that of ViP-NeRF on the RealEstate-10K dataset, we modify the test set as shown in Tab. 1. We change the test set since the test views that are very far away from the train views may contain large unobserved regions. We provide more reasoning and details in the supplementary.

6.1.4 Ablations. We test the importance of each of the components of our model by disabling them one at a time. We disable the smoothing and Lambertian 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 efficiently 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}_{aug} and \mathcal{L}_{cfc} . In addition, we also analyze the effect of including residual positional encodings $\gamma(\mathbf{p}_i, l_p^s, l_p)$ while predicting the color in the smoothing augmentation model. Tab. 5 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 performance. Further, using all the depths for supervision instead of only the

Table 2. Quantitative results of NeRF based models on the NeRF-LLFF dataset.

Model	2 views			3 views			4 views		
	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
InfoNeRF	0.6024	0.2219	9.16	0.6732	0.1953	8.37	0.6985	0.2270	9.18
DietNeRF	0.5465	0.3283	11.94	0.6120	0.3405	11.76	0.6506	0.3496	11.86
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
Simple-NeRF	0.2688	0.6501	19.57	0.2559	0.6940	21.37	0.2633	0.7016	21.99

Table 3. Quantitative results of NeRF based models on the RealEstate-10K dataset.

Model	2 views			3 views			4 views		
	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
InfoNeRF	0.5924	0.4342	12.27	0.6561	0.3792	10.57	0.6651	0.3843	10.62
DietNeRF	0.4381	0.6534	18.06	0.4636	0.6456	18.01	0.4853	0.6503	18.01
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
Simple-NeRF	0.0635	0.8942	33.10	0.0726	0.8984	33.21	0.0847	0.8987	32.88



Fig. 9. **Qualitative examples of NeRF based models 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. Simple-NeRF does not suffer from these distortions and synthesizes a clean frame. For reference, we also show the input images.

reliable depths leads to a significant drop in performance. Finally, disabling \mathcal{L}_{fc} 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.

Since we design our regularizations on top of DS-NeRF [Deng et al. 2022] baseline, our framework can 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 radiance field 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. However, our augmented models have complementary abilities as compared to the main radiance field model. We now analyze if there is a need to design augmentations that learn “simpler” solutions by replacing our novel augmentations with identical

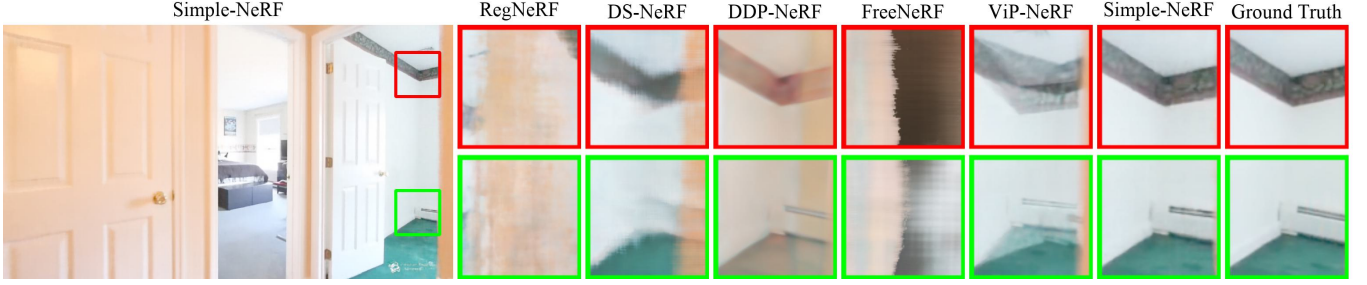


Fig. 10. **Qualitative examples of NeRF based models on RealEstate-10K dataset** with three input views. Simple-NeRF predictions are closest to the ground truth among all the models. In particular, DDP-NeRF predictions have a different shade of color and ViP-NeRF suffers from shape-radiance ambiguity,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

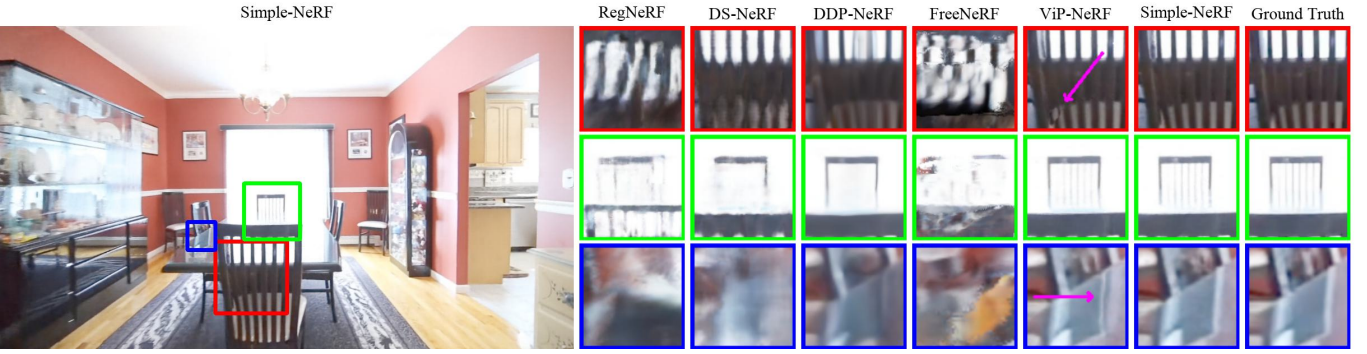


Fig. 11. **Qualitative examples of NeRF based models on the RealEstate-10K dataset with four input views.** We find that Simple-NeRF and ViP-NeRF perform the best among all the models. However, ViP-NeRF predictions contain minor distortions, as pointed out by the magenta arrow, which is rectified by Simple-NeRF.

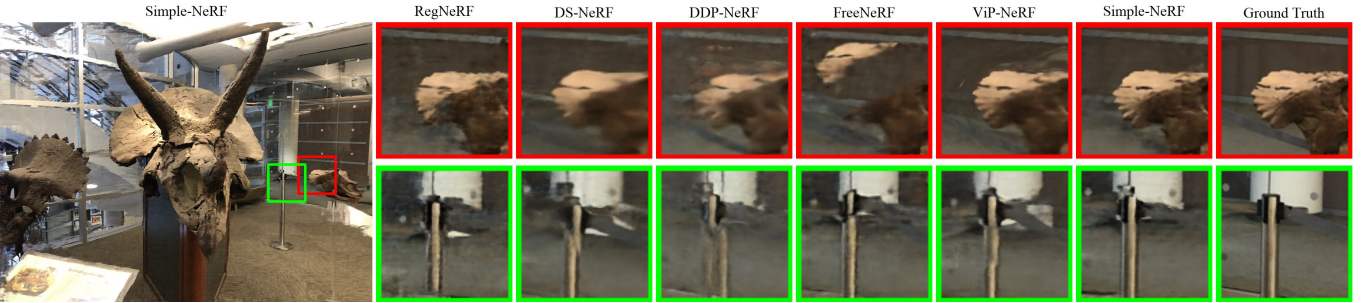


Fig. 12. **Qualitative examples of NeRF based models on the NeRF-LLFF dataset with two input views.** DDP-NeRF and ViP-NeRF synthesize frames with broken objects in the second row, and FreeNeRF breaks the object in the first row due to incorrect depth estimations. Simple-NeRF produces sharper frames devoid of such artifacts.

replicas of the NeRF as augmentations. The seventh row of Tab. 5 shows a performance drop when using identical augmentations.

Finally, we analyze the need for an augmentation that explicitly achieves depth smoothing. In other words, we ask if naively reducing the model capacity in the augmented model achieves a similar effect to that of our smoothing augmentation. We test this by replacing the smoothing augmentation with an augmented model that has a smaller MLP \mathcal{N}_1 . Specifically, we reduce the number of layers from

eight to four in the augmented model. From the results in the last row of Tab. 5, we conclude that reducing the positional encoding degree is more effective, perhaps because the MLP with fewer layers may still be capable of learning floaters on account of using all the positional encoding frequencies.

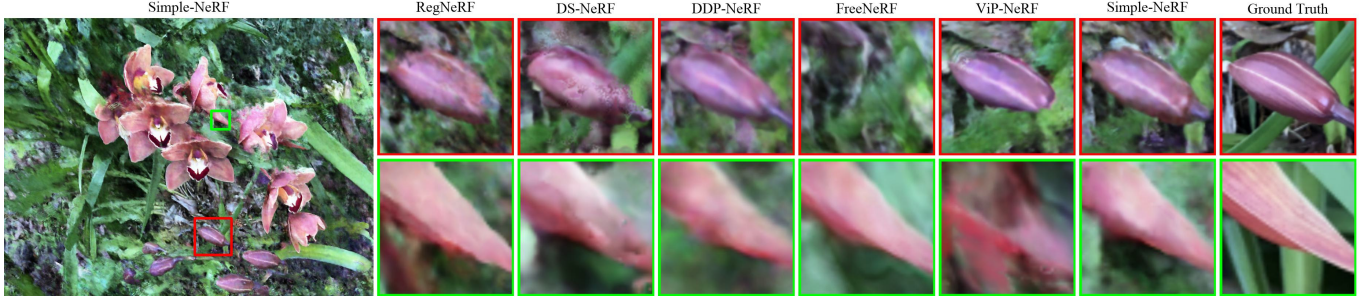


Fig. 13. **Qualitative examples of NeRF based models on the NeRF-LLFF dataset with three input views.** In the first row, the orchid is displaced out of the cropped box in the FreeNeRF prediction, due to incorrect depth estimation. ViP-NeRF and RegNeRF fail to predict the complete orchid accurately and contain distortions at either end. In the second row, ViP-NeRF prediction contains severe distortions. Simple-NeRF reconstructs the best among all the models in both examples.

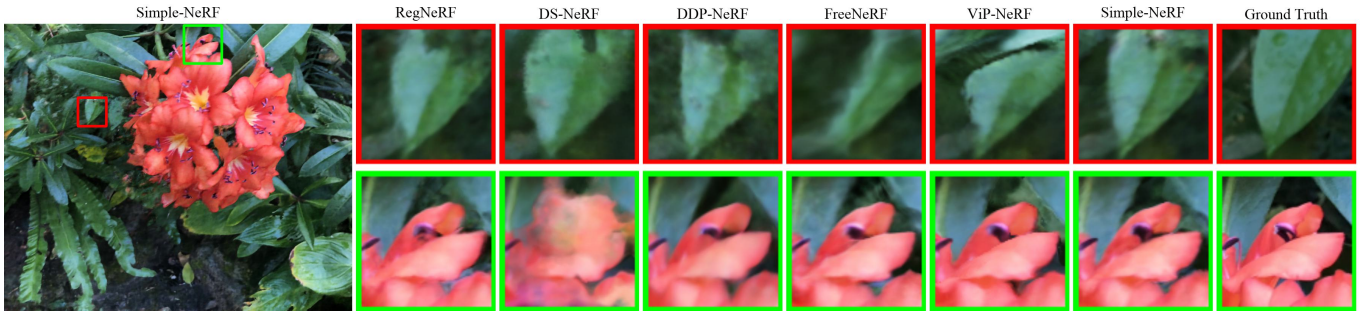


Fig. 14. **Qualitative examples of NeRF based models on the NeRF-LLFF dataset with four input views.** In the first row, we find that ViP-NeRF, FreeNeRF, and DDP-NeRF struggle to reconstruct the shape of the leaf accurately. In the second row, DS-NeRF introduces floaters. Simple-NeRF does not suffer from such artifacts and reconstructs the shapes better.

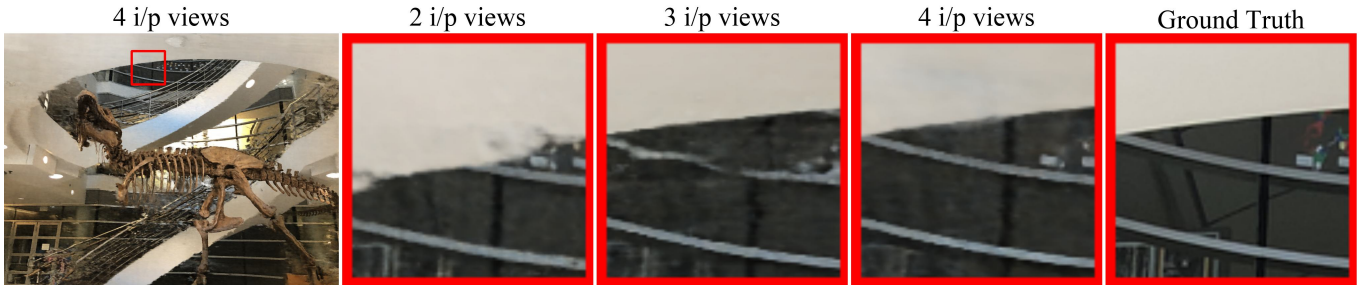


Fig. 15. **Qualitative examples of Simple-NeRF on the NeRF-LLFF dataset with two, three, and four input views.** We observe errors in depth estimation with two input views, causing a change in the position of the roof. While this is corrected with three input views, there are a few shape distortions in the metal rods. With four input views, even such distortions are corrected.

6.2 Simple-TensorRF

6.2.1 Implementation Details. Building on the original TensorRF code base, we employ Adam Optimizer with an initial learning rate of $2e - 2$ and $1e - 3$ for the tensor and MLP parameters respectively, which decay to $2e - 3$ and $1e - 4$. We employ the same hyper-parameters as the original implementation for the main model as follows: $R_\sigma = 24$, $R_c = 72$, $\mathbf{b} = \{(-1.5, 1.5), (-1.67, 1.67), (-1.0, 1.0)\}$, $N_{vox} = 640^3$, $D = 27$, and $l_v = 0$. We set $R_\sigma^s = 12$, $b_{z_1}^s = -0.5$ and $N_{vox}^s = 160^3$, $N_{mc} = 5$, $k = 5$, $e_\tau = 0.1$ for the augmented model and

the remaining hyper-parameters are the same as the main model. We weigh the losses as $\lambda_m = \lambda_a = 1$, $\lambda_{sd} = \lambda_{aug} = 0.1$, $\lambda_{mc} = 0.01$ and $\lambda_{cfc} = 0$. We train the models on a single NVIDIA RTX 2080 Ti 11GB GPU for 25k iterations and enable \mathcal{L}_{aug} after 5k iterations.

6.2.2 Quantitative and Qualitative Results. Tab. 6 shows the view-synthesis performance of Simple-TensorRF on the NeRF-LLFF and RealEstate-10K datasets. We compare the performance of our model against the vanilla TensorRF and a baseline we create by adding

Table 4. Evaluation of depth estimated by different NeRF based models with two input views. The reference depth is obtained using NeRF with dense input views. The depth MAE on the two datasets is of different orders on account of different depth ranges.

model	NeRF-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
Simple-NeRF	0.1420	0.8480	0.3269	0.9215

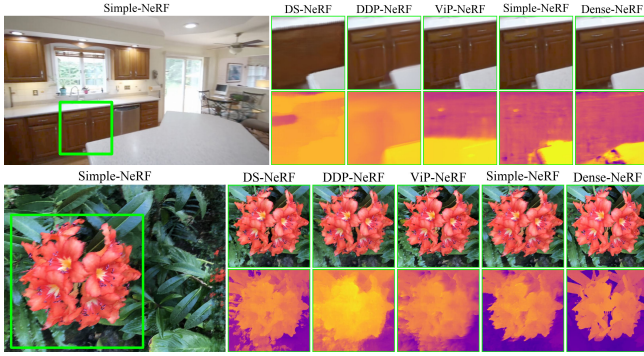


Fig. 16. **Estimated depth maps of NeRF based models** on RealEstate-10K and NeRF-LLFF datasets with two input views. In both examples, the two rows show the predicted images and the depths respectively. We find that Simple-NeRF is significantly better at estimating the scene depth. Also, DDP-NeRF synthesizes the left table edge at a different angle due to incorrect depth estimation.

Table 5. SimpleNeRF ablation experiments on RealEstate-10K and NeRF-LLFF datasets with two input views.

model	RealEstate-10K		NeRF-LLFF	
	LPIS ↓	MAE ↓	LPIS ↓	MAE ↓
Simple-NeRF	0.0635	0.33	0.2688	0.14
w/o smoothing augmentation	0.0752	0.38	0.2832	0.15
w/o Lambertian 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/o residual pos enc	0.0790	0.40	0.2837	0.16
w/ identical augmentations	0.0777	0.40	0.2849	0.15
w/ smaller n/w as smoothing aug	0.0740	0.38	0.2849	0.15

sparse depth loss on TensorRF, which we refer to as DS-TensorRF. We find that Simple-TensorRF significantly improves performance over TensorRF and DS-TensorRF on both datasets. Fig. 18 compares the three models visually, where we observe that Simple-TensorRF mitigates multiple distortions observed in the renders of TensorRF and DS-TensorRF. From Tab. 6 and Fig. 19, we observe that Simple-TensorRF is significantly better at estimating the scene depth than both TensorRF and DS-TensorRF. While we observe that Simple-NeRF performs marginally better than Simple-TensorRF in most cases,

Simple-TensorRF achieves a lower depth MAE on the RealEstate-10K dataset.

We test the need for the different components of our augmentation by disabling them one at a time and show the quantitative results in the second half of Tab. 6. Specifically, we disable the reduction in the number of tensor decomposition components and the number of voxels in the first two rows respectively. In the third row, we disable both the components, where the augmented model consists of the reduction in the bounding box size only. We find that disabling either or both of the components leads to a drop in performance. In particular, Fig. 20 shows that reducing only the tensor resolution and not reducing the number of tensor decomposition components leads to translucent blocky floaters. On the other hand, reducing only the number of components causes small and completely opaque floaters. These effects can be better observed in the supplementary videos. Further, we find that reducing the number of components R_σ is more crucial in obtaining simpler solutions on the RealEstate-10K dataset.

6.3 Simple-ZipNeRF

6.3.1 Implementation Details. We build our code in PyTorch on top of an unofficial ZipNeRF implementation¹. For the main model, we retain the hyper-parameters of the original ZipNeRF. For the augmented model, we reduce the size of the hash table from $T = 2^{21}$ to $T^s = 2^{11}$ and set $s_{\text{near}} = 0.3$. We impose \mathcal{L}_{aug} after 5k iterations and use $k = 5$, $e_\tau = 0.2$. The rest of the hyper-parameters for the augmented model are the same as the main model. We weigh the losses as $\lambda_m = \lambda_a = 1$, $\lambda_{\text{aug}} = 10$ and $\lambda_{sd} = \lambda_{cfc} = \lambda_{mc} = 0$. We do not impose the sparse depth loss \mathcal{L}_{sd} since we find that Colmap either fails in sparse reconstruction or provides noisy sparse depth for 360° scenes. We train the models on a single NVIDIA RTX 2080 Ti 11GB GPU for 25k iterations.

6.3.2 Quantitative and Qualitative Results. We compare the performance of ZipNeRF with and without our augmentations on the MipNeRF360 and NeRF-Synthetic datasets in Tabs. 7 and 8 respectively. We observe that including our augmentations improves performance significantly on both datasets in terms of all the evaluation measures. This observation is further supported by the qualitative examples in Figs. 21 to 23, where we observe a clear improvement in the quality of the rendered novel views and depth when employing our augmentations. In addition, Tab. 7 and Fig. 24 also shows the performance of the augmented model on the MipNeRF360 dataset. We observe a significant reduction in distortions in the renders of the augmented model; however, the same does not reflect in the quantitative evaluation due to the blur introduced by the augmented model.

Further, in Fig. 25, we show the performance of ZipNeRF and Simple-ZipNeRF as the number of input views increases. We observe that the performance of ZipNeRF is too low with very few input images, where our augmentation does not help improve the performance. As the number of input views increases, the performance of ZipNeRF improves, and our augmentation helps improve the performance significantly. Further, with a large number of input

¹ ZipNeRF implementation: <https://github.com/SuLvXiangXin/zipnerf-pytorch>

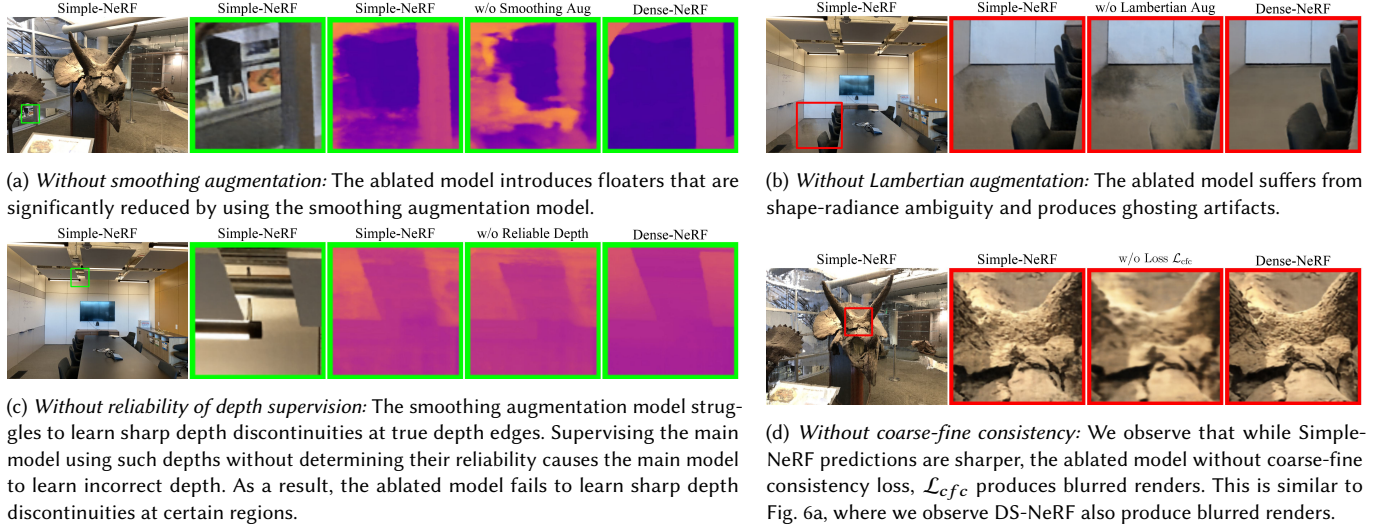


Fig. 17. **Qualitative examples for Simple-NeRF ablated models on the NeRF-LLFF dataset with two input views.** We also show the outputs of the dense-input NeRF for reference.

Table 6. Quantitative results of TensorRF based models with three input views.

Model	NeRF-LLFF					RealEstate-10K				
	LPIPS ↓	SSIM ↑	PSNR ↑	Depth MAE ↓	Depth SROCC ↑	LPIPS ↓	SSIM ↑	PSNR ↑	Depth MAE ↓	Depth SROCC ↑
TensorRF	0.5474	0.3163	12.29	0.67	0.03	0.0986	0.8532	29.62	0.44	0.63
DS-TensorRF	0.2897	0.6291	18.58	0.23	0.73	0.0739	0.8872	32.50	0.27	0.75
Simple-TensorRF	0.2461	0.6749	20.22	0.17	0.83	0.0706	0.8920	32.70	0.22	0.80
$R_\sigma^s = R_\sigma$	0.2536	0.6677	19.85	0.18	0.81	0.085	0.8821	30.94	0.27	0.77
$N_{vox}^s = N_{vox}$	0.2568	0.6579	19.95	0.19	0.79	0.0735	0.8896	32.22	0.22	0.82
$R_\sigma^s = R_\sigma; N_{vox}^s = N_{vox}$	0.2728	0.6424	19.50	0.22	0.74	0.0787	0.8871	31.73	0.23	0.79

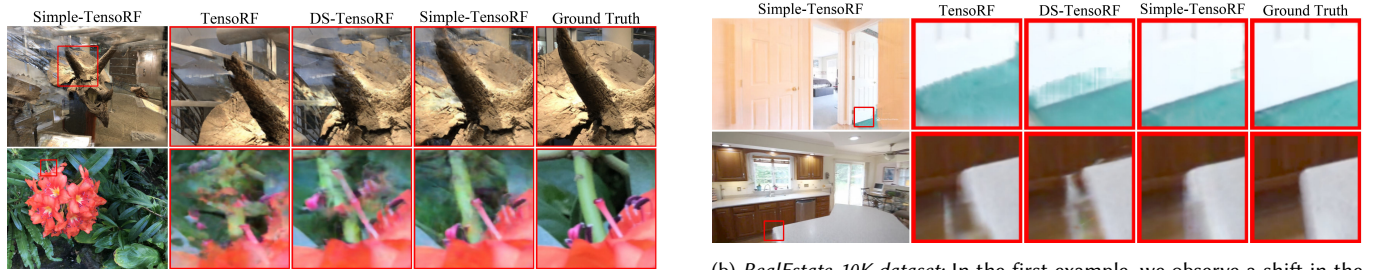


Fig. 18. **Qualitative examples of TensorRF based models with three input views.**

views, the performance of ZipNeRF saturates, and our augmentation does not help improve the performance. This shows that our augmentations are highly effective when the performance of the base model is moderately good.

Table 7. Quantitative results of ZipNeRF based models on the MipNeRF360 dataset.

Model	12 input views					20 input views			36 input views		
	LPIPS ↓	SSIM ↑	PSNR ↑	Depth MAE ↓	Depth SROCC ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
ZipNeRF	0.5614	0.4616	15.86	7.43	0.28	0.435	0.5911	18.89	0.3316	0.6737	21.78
Augmented ZipNeRF	0.6825	0.4462	16.27	96.42	0.49	0.619	0.5244	19.31	0.5917	0.5646	21.21
Simple-ZipNeRF	0.4876	0.5245	17.60	3.54	0.51	0.3421	0.6456	21.03	0.239	0.7458	24.19

Table 8. Quantitative results of ZipNeRF based models on the NeRF-Synthetic dataset.

Model	4 input views			8 input views			12 input views		
	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑	LPIPS ↓	SSIM ↑	PSNR ↑
ZipNeRF	0.4263	0.7548	11.04	0.2877	0.7973	15.01	0.1625	0.8528	20.12
Simple-ZipNeRF	0.3878	0.7715	11.50	0.2461	0.8063	15.88	0.1532	0.8531	20.51

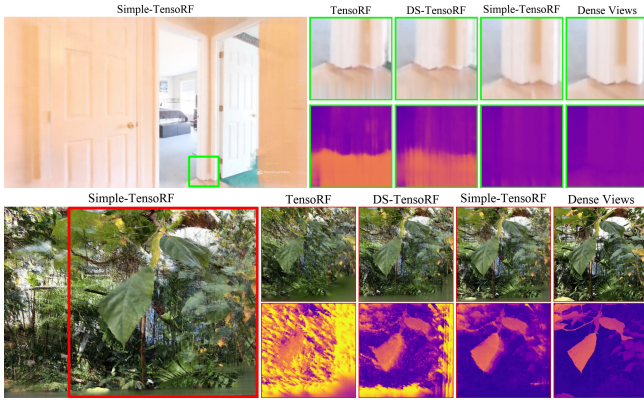


Fig. 19. **Estimated depth maps of TensorRF based models** on RealEstate-10K and NeRF-LLFF datasets with three input views. In both examples, the two rows show the predicted images and the depths respectively. In the first example, TensorRF and DS-TensorRF incorrectly estimate the depth of the floor as shown by the orange regions. In the second row, while TensorRF is unable to estimate the scene geometry, DS-TensorRF is unable to mitigate all the floaters in orange color. We find that Simple-TensorRF is significantly better at estimating the scene depth.

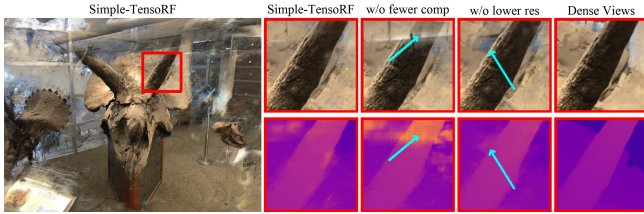


Fig. 20. **Qualitative examples of Simple-TensorRF ablations** on NeRF-LLFF dataset with three input views. Reducing the tensor resolution only leads to translucent floaters as shown by the arrows in the second column. On the other hand, only reducing the number of tensor decomposed components leads to small opaque floaters as shown by the arrows in the third column.

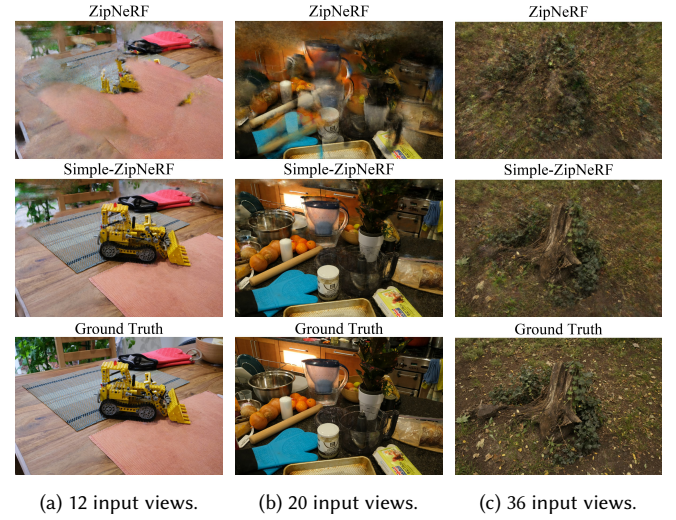


Fig. 21. **Qualitative examples of ZipNeRF and Simple-ZipNeRF on the MipNeRF360 dataset.** In the first column, we observe that ZipNeRF places large regions of the pink mat close to the camera, occluding the bulldozer. In the second example, we observe objects being broken or placed at incorrect positions due to incorrect depth estimation, as well as translucent floaters in ZipNeRF predictions. Finally, in the third column, we observe that ZipNeRF fails to reconstruct the tree stump. In all the cases, Simple-ZipNeRF produces very good reconstructions without any floaters.

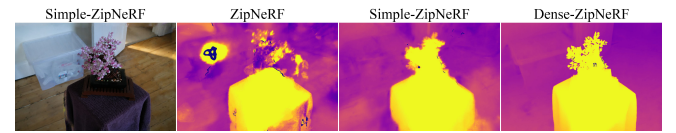


Fig. 22. **Simple-ZipNeRF estimated depth maps** on MipNeRF360 dataset with 20 input views. We observe that the depth map estimated by ZipNeRF contains floaters and that the depth estimates for the bonsai are incorrect. However, Simple-ZipNeRF does not suffer from such issues and the estimated depth is very close to that of ZipNeRF with dense input views.



Fig. 23. **Qualitative examples of ZipNeRF and Simple-ZipNeRF on the NeRF-Synthetic dataset.** While the renders of ZipNeRF contain multiple floaters, Simple-ZipNeRF outputs are cleaner and free from such artifacts.

Table 9. Training and inference (per frame) time and memory comparison of various models.

Model	Training		Inference	
	Time (hrs)	Mem (GB)	Time (sec)	Mem (GB)
NeRF	14	6.1	54	0.8
Simple-NeRF	21	8.8	54	0.8
TensorRF	2.1	6.8	21	4.0
Simple-TensorRF	3.7	7.2	21	4.0
ZipNeRF	2.0	6.7	13	4.8
Simple-ZipNeRF	4.2	8.6	13	4.8

7 DISCUSSION

7.1 Computational Complexity

We report the approximate GPU memory utilization and time taken for training and inference of our family of Simple-RF models in Tab. 9. We observe that Simple-NeRF with two augmentations takes only 1.5 times more time than NeRF for training on account of employing augmentations on the coarse NeRF only. While coarse NeRF queries the MLPs 64 times, the fine NeRF queries the MLPs 192 times, giving a combined 256 queries per pixel. Simple-NeRF queries the coarse MLPs 192 times and the fine MLPs 192 times, with a total of 384 queries per pixel. On the other hand, Simple-TensorRF and Simple-ZipNeRF take twice the time as TensorRF and ZipNeRF respectively on account of employing a single augmentation with

exactly the same number of queries as the main model. We note that it could be possible to further reduce the training time for ZipNeRF by employing the augmentation only on the proposal MLP. However, this requires the proposal MLPs to output color and to be trained with the photometric loss instead of the interval loss [Barron et al. 2023]. The effect of such a change is unclear and is left for future work. Interestingly, Simple-TensorRF requires only a little more memory than TensorRF during training, perhaps due to the low resolution tensor employed by the augmented model. Further, while the NeRF models require significantly less memory during inference, grid based models such as TensorRF and ZipNeRF require more memory due to the use of a voxel grid in place of MLPs. Finally, we note that at inference time, Simple-RF models take exactly the same time and memory as the baseline models since the augmentations are disabled during inference. All the above experiments are conducted on a single NVIDIA RTX 2080 Ti 11GB GPU.

7.2 Limitations and Future Work

Our approach of obtaining reliable depth supervision by learning simpler solutions is limited to the cases where the base model achieves a reasonable performance but suffers from various distortions due to overfitting. Our regularizations may not help significantly if the performance of the base model is very poor, as in the case of learning a highly complex scene with very few input views. In such cases, we can employ our regularizations on top of other sparse input radiance fields as we do for the NeRF-LLFF dataset.

Training augmented models adds compute and memory overhead during training. It would be interesting to explore deriving augmentations from the main model itself, without training an augmented model separately. For example, instead of training a separate augmented model with a lower resolution grid in Simple-TensorRF, one could try downsampling the grid from the main model. However, achieving the same with other augmentations is non-trivial and is left for future work.

While we perform elementary modifications of radiance fields to obtain simpler solutions, it would be interesting to explore more sophisticated augmentations to learn simpler solutions. For example, one could explore inducing smoothness through different hash functions in Simple-ZipNeRF. Such sophisticated augmentations could help obtain larger improvements in performance. However, such explorations are beyond the scope of this work.

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.

Our model requires accurate camera poses of the sparse input images. While the joint optimization of camera poses and scene geometry is explored when dense input views are available [Jeong et al. 2021; Park et al. 2023], it would be helpful to explore the same with sparse input views. In our experiments, we found that a trivial combination of a pose optimization radiance field such as NeRF- [Wang et al. 2021c] and a sparse input radiance field such as Simple-NeRF or DS-NeRF leads to very poor performances. While

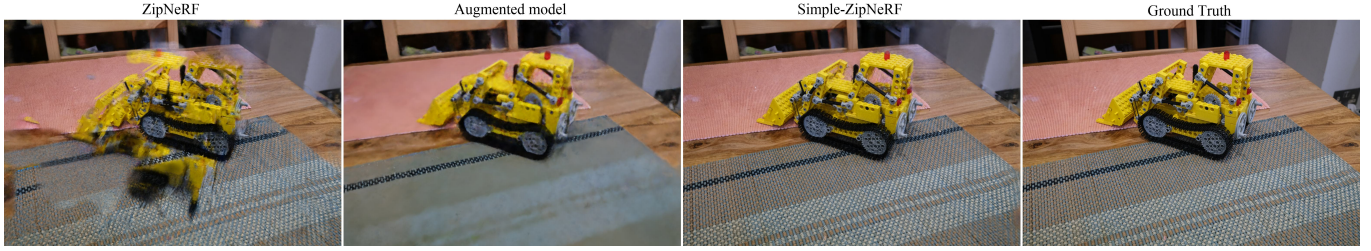


Fig. 24. **Qualitative examples to visualize the effect of our augmentation.** We observe that the ZipNeRF render contains severe distortions. The output of our augmented model is significantly better in reconstructing the scene, but the render contains severe blur on account of smoothing introduced by small hash table. Learning from the depth provided by the augmented model, Simple-ZipNeRF is able to reconstruct the scene better as well as retain sharpness by utilizing the larger hash table.

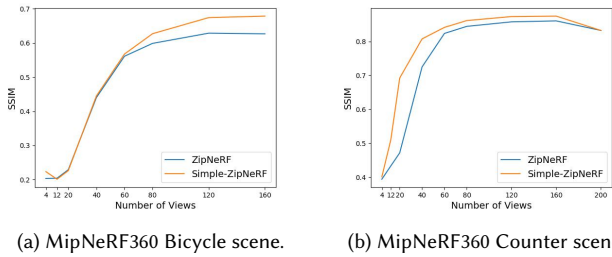


Fig. 25. **Performance of ZipNeRF and Simple-ZipNeRF with increasing number of input views.** We observe that our augmentation improves performance significantly over ZipNeRF, when the performance of the base model is moderate. When the performance of the base model is extremely poor or extremely good, the augmentation does not have a significant impact. However, our augmentation does not lead to significant degradation in performance in either case.

some approaches [Lin et al. 2021] for camera pose optimization are limited to NeRF and do not extend to explicit radiance fields, other approaches [Bian et al. 2023; Han et al. 2024; Truong et al. 2023] require pre-training on a large dataset. We believe this problem requires a deeper study and is a very important direction to be pursued in the future.

Finally, our approach of employing augmentations to obtain better supervision can be explored in sparse input novel view synthesis of highly specular objects [Verbin et al. 2022], refractive objects [Deng et al. 2024], low-light scenes [Mildenhall et al. 2022], blurred images [Ma et al. 2022] and high dynamic range images [Lu et al. 2024]. Further, it would be interesting to explore the use of augmentations in related inverse problems such as surface reconstruction [Long et al. 2022; Wang et al. 2021a], dynamic view synthesis [Fridovich-Keil et al. 2023; Pumarola et al. 2021; Somraj et al. 2024, 2022] and style transfer [Wang et al. 2024] when only sparse input viewpoints are available.

8 CONCLUSION

We address the problem of few-shot radiance fields by obtaining depth supervision from simpler solutions learned by lower capability augmented models that are trained in tandem with the main radiance field model. We show that augmentations can be designed

for both implicit models, such as NeRF, and explicit radiance fields, such as TensorRF and ZipNeRF. Since the shortcomings of various radiance fields are different, we design the augmentations appropriately for each model. We show that our augmentations improve performance significantly on all three models, and we achieve state-of-the-art performance on forward-facing scenes as well as 360° scenes. Notably, our models achieve a significant improvement in the depth estimation of the scene, which indicates a superior geometry estimation.

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.

REFERENCES

- Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. 2021. Mip-NeRF: A Multiscale Representation for Anti-Aliasing Neural Radiance Fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- 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)*.
- Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2023. Zip-NeRF: Anti-Aliased Grid-Based Neural Radiance Fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. 2020. Neural Reflectance Fields for Appearance Acquisition. *arXiv e-prints* (2020). arXiv:2008.03824
- Wenjing Bian, Zirui Wang, Kejia 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* (2022). arXiv:2210.04214
- Joel Carranza, Christian Theobalt, Marcus A. Magnor, and Hans-Peter Seidel. 2003. Free-Viewpoint Video of Human Actors. *ACM Transactions on Graphics (TOG)* 22, 3 (July 2003), 569–577. <https://doi.org/10.1145/882262.882309>
- Jin-Xiang Chai, Xin Tong, Shing-Chow Chan, and Heung-Yeung Shum. 2000. Plenoptic Sampling. In *Proceedings of the ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/344779.344932>
- Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. 2013. Depth Synthesis and Local Warps for Plausible Image-Based Navigation. *ACM Transactions on Graphics (TOG)* 32, 3 (July 2013). <https://doi.org/10.1145/2487228.2487238>
- 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* (March 2021). 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 ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/166117.166153>
- Yuedong Chen, Hao-fei Xu, Qianyi Wu, Chuanxia Zheng, Tat-Jen Cham, and Jianfei Cai. 2023. Explicit Correspondence Matching for Generalizable Neural Radiance Fields. *arXiv e-prints* (2023). arXiv:2304.12294
- Yuedong Chen, Hao-fei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. 2024. MVSPat: Efficient 3D Gaussian Splatting from Sparse Multi-View Images. *arXiv e-prints* (2024). arXiv:2403.14627
- 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.
- Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. 2015. High-Quality Streamable Free-Viewpoint Video. *ACM Transactions on Graphics (TOG)* 34, 4 (July 2015). <https://doi.org/10.1145/2766945>
- Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. 2022. Depth-Supervised NeRF: Fewer Views and Faster Training for Free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Weijian Deng, Dylan Campbell, Chunyi Sun, Shubham Kanitkar, Matthew Shaffer, and Stephen Gould. 2024. Ray Deformation Networks for Novel View Synthesis of Refractive Objects. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*.
- Christoph Fehn. 2004. Depth-Image-Based Rendering (DIBR), Compression and Transmission for a New Approach on 3D-TV. In *Proceedings of the Stereoscopic Displays and Virtual Reality Systems XI*. <https://doi.org/10.1117/12.524762>
- John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. 2016. DeepStereo: Learning to Predict New Views From the World’s Imagery. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbek Warburg, Benjamin Recht, and Angjoo Kanazawa. 2023. K-Planes: Explicit Radiance Fields in Space, Time, and Appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. 2022. Plenoxels: Radiance Fields Without Neural Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zelin Gao, Weichen Dai, and Yu Zhang. 2024. HG3-NeRF: Hierarchical Geometric, Semantic, and Photometric Guided Neural Radiance Fields for Sparse View Inputs. *arXiv e-prints* (2024). arXiv:2401.11711
- Steven J. Gortler, Radek Grzeszczek, Richard Szeliski, and Michael F. Cohen. 1996. The Lumigraph. In *Proceedings of the ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/237170.237200>
- Shuai Guo, Qiuwen Wang, Yijie Gao, Rong Xie, and Li Song. 2024. Depth-Guided Robust and Fast Point Cloud Fusion NeRF for Sparse Input Views. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 3 (March 2024), 1976–1984. <https://doi.org/10.1609/aaai.v38i3.27968>
- Xinyang Han, Zelin Gao, Angjoo Kanazawa, Shubham Goel, and Yossi Gandelsman. 2024. The More You See in 2D, the More You Perceive in 3D. *arXiv e-prints* (2024). arXiv:2404.03652
- 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)*.
- Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. 2021. Self-Calibrating Neural Radiance Fields. 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)*.
- Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. 2016. Learning-Based View Synthesis for Light Field Cameras. *ACM Transactions on Graphics (TOG)* 35, 6 (December 2016). <https://doi.org/10.1145/2980179.2980251>
- 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)*.
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics (TOG)* 42, 4 (2023).
- 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)*.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*.
- Minseop Kwak, Jiuhun Song, and Seungryong Kim. 2023. GeCoNeRF: Few-shot Neural Radiance Fields via Geometric Consistency. *arXiv e-prints* (2023). 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* (2023). arXiv:2303.11728
- Seoyoung Lee and Joonseok Lee. 2024. PoseDiff: Pose-Conditioned Multimodal Diffusion Model for Unbounded Scene Synthesis From Sparse Inputs. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*.
- Marc Levoy and Pat Hanrahan. 1996. Light Field Rendering. In *Proceedings of the ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/237170.237199>
- Jiahe Li, Jiawei Zhang, Xiao Bai, Jin Zheng, Xin Ning, Jun Zhou, and Lin Gu. 2024. DNGaussian: Optimizing Sparse-View 3D Gaussian Radiance Fields with Global-Local Depth Normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. 2021. BARF: Bundle-Adjusting Neural Radiance Fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- 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)*.
- Xinhang Liu, Shiu-hong Kao, Jiaben Chen, Yu-Wing Tai, and Chi-Keung Tang. 2023. Deceptive-NeRF: Enhancing NeRF Reconstruction using Pseudo-Observations from Diffusion Models. *arXiv e-prints* (2023). arXiv:2305.15171
- Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. 2022. Neural Rays for Occlusion-Aware Image-Based Rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Xiaoxiao Long, Cheng Lin, Peng Wang, Taku Komura, and Wenping Wang. 2022. SparseNeuS: Fast Generalizable Neural Surface Reconstruction from Sparse Views. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Zhan Lu, Qian Zheng, Boxin Shi, and Xudong Jiang. 2024. Pano-NeRF: Synthesizing High Dynamic Range Novel Views with Geometry from Sparse Low Dynamic Range Panoramic Images. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 4 (March 2024), 3927–3935. <https://doi.org/10.1609/aaai.v38i4.28185>
- 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)*.
- William R Mark. 1999. *Post-Rendering 3D Image Warping: Visibility, Reconstruction, and Performance for Depth-Image Warping*. The University of North Carolina at Chapel Hill.
- Leonard McMillan and Gary Bishop. 1995. Plenoptic Modeling: An Image-Based Rendering System. In *Proceedings of the ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/218380.218398>
- Leonard McMillan Jr. 1997. *An Image-Based Approach to Three-Dimensional Computer Graphics*. The University of North Carolina at Chapel Hill.
- 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.
- Zhangkai Ni, Peiqi Yang, Wenhan Yang, Hanli Wang, Lin Ma, and Sam Kwong. 2024. Col-NeRF: Collaboration for Generalizable Sparse Input Neural Radiance Field. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 5 (March 2024), 4325–4333. <https://doi.org/10.1609/aaai.v38i5.28229>

- 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)*.
- Keunhong Park, Philipp Henzler, Ben Mildenhall, Jonathan T. Barron, and Ricardo Martin-Brualla. 2023. CamP: Camera Preconditioning for Neural Radiance Fields. *ACM Transactions on Graphics (TOG)* 42, 6 (December 2023). <https://doi.org/10.1145/3618321>
- Eric Penner and Li Zhang. 2017. Soft 3D Reconstruction for View Synthesis. *ACM Transactions on Graphics (TOG)* 36, 6 (November 2017), 1–11. <https://doi.org/10.1145/3130800.3130855>
- Julien Philip and Valentin Deschaintre. 2023. Floaters No More: Radiance Field Gradient Scaling for Improved Near-Camera Training. In *Proceedings of the Eurographics Symposium on Rendering*. <https://doi.org/10.2312/sr.20231122>
- Malte Prinzler, Otmar Hilliges, and Justus Thies. 2023. DINER: Depth-Aware Image-Based NEural Radiance Fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. 2021. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Ravi Ramamoorthi. 2023. NeRFs: The Search for the Best 3D Representation. *arXiv e-prints* (2023). [arXiv:2308.02751](https://arxiv.org/abs/2308.02751)
- Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. 2021. KiloNeRF: Speeding Up Neural Radiance Fields With Thousands of Tiny MLPs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- 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, Yeonjin Chang, and Nojun Kwak. 2023a. FlipNeRF: Flipped Reflection Rays for Few-shot Novel View Synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Seunghyeon Seo, Donghoon Han, Yeonjin Chang, and Nojun Kwak. 2023b. 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)*.
- Jonathan Shade, Steven Gortler, Li-wei He, and Richard Szeliski. 1998. Layered Depth Images. In *Proceedings of the ACM Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/280814.280882>
- Ruoxi Shi, Xinyue Wei, Cheng Wang, and Hao Su. 2024. ZeroRF: Fast Sparse View 360° Reconstruction with Zero Pretraining. 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* (2022). [arXiv:2212.02280](https://arxiv.org/abs/2212.02280)
- Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 2020. 3D Photography Using Context-Aware Layered Depth Inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Aljoscha Smolic, Karsten Mueller, Philipp Merkle, Christoph Fehn, Peter Kauff, Peter Eisert, and Thomas Wiegand. 2006. 3D Video and Free Viewpoint Video - Technologies, Applications and MPEG Standards. In *Proceedings of the IEEE International Conference on Multimedia and Expo (ICME)*. <https://doi.org/10.1109/ICME.2006.262683>
- Nagabhushan Somraj, Kapil Choudhary, Sai Harsha Mupparaju, and Rajiv Soundararajan. 2024. Factorized Motion Fields for Fast Sparse Input Dynamic View Synthesis. In *Proceedings of the ACM Special Interest Group on Computer Graphics and Interactive Techniques (SIGGRAPH)*. <https://doi.org/10.1145/3641519.3657498>
- Nagabhushan Somraj, Pranali Sancheti, and Rajiv Soundararajan. 2022. Temporal View Synthesis of Dynamic Scenes through 3D Object Motion Estimation with Multi-Plane Images. In *Proceedings of the IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. <https://doi.org/10.1109/ISMAR55827.2022.00100>
- 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>
- Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. 2021. NeRV: Neural Reflectance and Visibility Fields for Relighting and View Synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Pratul P. Srinivasan, Richard Tucker, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. 2019. Pushing the Boundaries of View Extrapolation With Multiplane Images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct Voxel Grid Optimization: Super-Fast Convergence for Radiance Fields Reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wenxiu Sun, Lingfeng Xu, Oscar C Au, Sung Him Chui, and Chun Wing Kwok. 2010. An Overview of Free View-Point Depth-Image-Based Rendering (DIBR). In *Proceedings of the APSIPA Annual Summit and Conference*.
- 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)*.
- 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)*.
- Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. 2023. SPARF: Neural Radiance Fields From Sparse and Noisy Poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. 2018. Deep Image Prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Mikaela Angelina Uy, Ricardo Martin-Brualla, Leonidas Guibas, and Ke Li. 2023. SCADe: NeRFs from Space Carving With Ambiguity-Aware Depth Estimates. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- 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* (2023). [arXiv:2303.16196](https://arxiv.org/abs/2303.16196)
- Peihao Wang, Xuxi Chen, Tianlong Chen, Subhashini Venugopalan, Zhangyang Wang, et al. 2022. Is Attention All That NeRF Needs? *arXiv e-prints* (2022). [arXiv:2207.13298](https://arxiv.org/abs/2207.13298)
- Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. 2021a. NeuS: Learning Neural Implicit Surfaces by Volume Rendering for Multi-view Reconstruction. *arXiv e-prints* (2021). [arXiv:2106.10689](https://arxiv.org/abs/2106.10689)
- Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. 2021b. IBRNet: Learning Multi-View Image-Based Rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Yifan Wang, Ang Gao, Yi Gong, and Yuan Zeng. 2024. Stylizing Sparse-View 3D Scenes with Hierarchical Neural Representation. *arXiv e-prints* (2024). [arXiv:2404.05236](https://arxiv.org/abs/2404.05236)
- 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.
- Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. 2021c. NeRF-: Neural Radiance Fields without Known Camera Parameters. *arXiv e-prints* (2021). [arXiv:2102.07064](https://arxiv.org/abs/2102.07064)
- Rundi Wu, Ben Mildenhall, Philipp Henzler, Keunhong Park, Ruiqi Gao, Daniel Watson, Pratul P. Srinivasan, Dor Verbin, Jonathan T. Barron, Ben Poole, et al. 2024. ReconFusion: 3D Reconstruction with Diffusion Priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jamie Wynn and Daniyar Turmukhambetov. 2023. DiffusionNeRF: Regularizing Neural Radiance Fields with Denoising Diffusion Models. *arXiv e-prints* (2023). [arXiv:2302.12231](https://arxiv.org/abs/2302.12231)
- HaoLin Xiong, Sairisheek Muttukuru, Rishi Upadhyay, Pradyumna Chari, and Achuta Kadambi. 2023. SparseGS: Real-Time 360° Sparse View Synthesis using Gaussian Splatting. *arXiv e-prints* (2023). [arXiv:2312.00206](https://arxiv.org/abs/2312.00206)
- 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)*.
- Chen Yang, Sikuang Li, Jiemin Fang, Ruofan Liang, Lingxi Xie, Xiaopeng Zhang, Wei Shen, and Qi Tian. 2024. GaussianObject: Just Taking Four Images to Get A High-Quality 3D Object with Gaussian Splatting. *arXiv e-prints* (2024). [arXiv:2402.10259](https://arxiv.org/abs/2402.10259)
- Jiawei Yang, Marco Pavone, and Yue Wang. 2023. FreeNeRF: Improving Few-shot Neural Rendering with Free Frequency Regularization. (June 2023).
- Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. 2021a. PlenOctrees for Real-Time Rendering of Neural Radiance Fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021b. pixelNeRF: Neural Radiance Fields From One or Few Images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. 2021b. 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* (2020). [arXiv:2010.07492](https://arxiv.org/abs/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)*.

- Xiuming Zhang, Pratul P. Srinivasan, Boyang Deng, Paul Debevec, William T. Freeman, and Jonathan T. Barron. 2021a. NeRFactor: Neural Factorization of Shape and Reflectance Under an Unknown Illumination. *ACM Transactions on Graphics (TOG)* 40, 6 (December 2021). <https://doi.org/10.1145/3478513.3480496>
- 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. 2023b. VDN-NeRF: Resolving Shape-Radiance Ambiguity via View-Dependence Normalization. *arXiv e-prints* (2023). arXiv:2303.17968
- Hanxin Zhu, Tianyu He, Xin Li, Bingchen Li, and Zhibo Chen. 2024. Is Vanilla MLP in Neural Radiance Field Enough for Few-shot View Synthesis?. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zehao Zhu, Zhiwen Fan, Yifan Jiang, and Zhangyang Wang. 2023a. FSGS: Real-Time Few-shot View Synthesis using Gaussian Splatting. *arXiv e-prints* (2023). arXiv: 2312.00451

Supplement

The contents of this supplement include

- A. Details on evaluation measures.
- B. Video examples on LLFF, RealEstate-10K and MipNeRF360 datasets.
- C. Additional analysis - positional encoding frequency, and depth reliability masks.
- D. Extensive quantitative evaluation reports.

A DETAILS ON EVALUATION MEASURES

A.1 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 radiance fields 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, radiance fields are also ill-equipped to predict the depth of objects seen in only one of the input views, again by design. Thus, radiance fields 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 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. 10 to 17.

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. For the MipNeRF360 dataset, we render the videos along the elliptical trajectory that is commonly used in the literature.

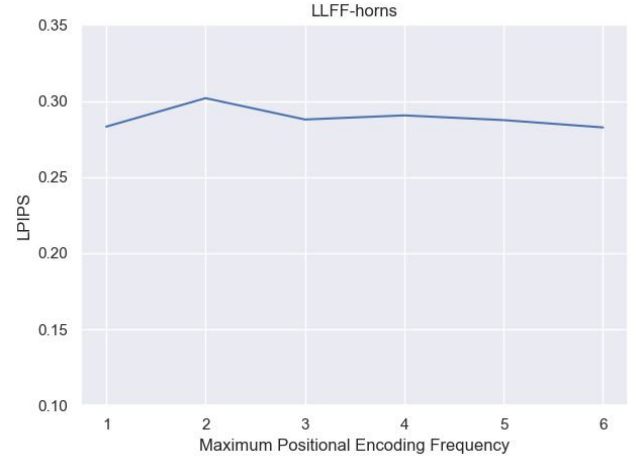


Fig. 26. LPIPS scores on the horns scene as l_p^{ap} of Simple-NeRF is varied.

We group the video comparisons based on Simple-NeRF, Simple-TensoRF and Simple-ZipNeRF models. In each group, we divide the video comparisons into two sets. In the first set, we show how our regularizations reduce various artifacts 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. Finally, we also include a few videos to support certain arguments made in the main paper. The videos are available on our project website <https://nagabhushansn95.github.io/publications/2024/Simple-RF.html>

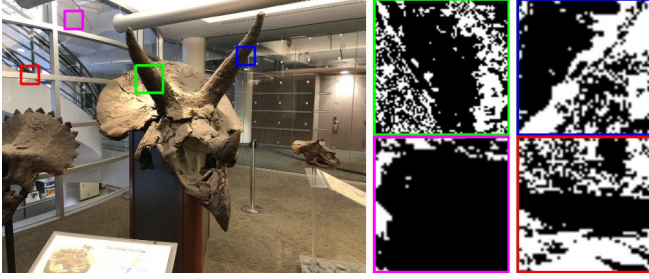
C ADDITIONAL ANALYSIS

C.1 Ablation on Positional Encoding Frequency in Simple-NeRF

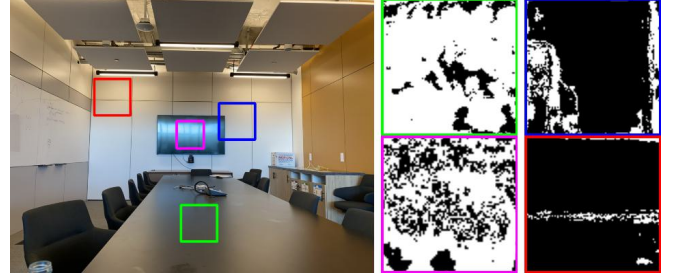
We analyze the variation in the performance of Simple-NeRF as l_p^{as} varies. We vary l_p^{as} from 1 to 6 and test the performance of Simple-NeRF on the horns scene of the NeRF-LLFF dataset. We show the quantitative performance in terms of LPIPS in Fig. 26. We observe only small variations in the performance as l_p^{ap} is varied, and thus, we conclude that our framework is robust to the choice of l_p^{ap} . Further, we note that using $l_p^{ap} = l_p = 10$ is equivalent to using an identical augmentation.

C.2 Visualization of Depth Reliability Masks

In Fig. 27, we present visualizations that motivate the design of our augmentations in Simple-NeRF, namely the smoothing and Lamberian augmentations. We train our model without augmentations and the individual augmentations separately with only \mathcal{L}_{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 4.4 of the main paper [Verify if this is correct](#). For two scenes from the



(a) **Smoothing 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) **Lambertian 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.

Fig. 27. Visualizations of depth reliability mask for the two augmentations of Simple-NeRF. 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.

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

We observe that the smoothing augmentation is determined to have estimated better depths in smooth regions. At edges, the depth estimated by the main model is more accurate. Similarly, the Lambertian augmentation estimates better depth in Lambertian regions, while 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.

D 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. 18 to 33.

Table 10. Quantitative Results of NeRF based models 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)
Simple-NeRF	0.2688(0.3899)	0.6501(0.5529)	19.57(17.57)	0.1420 (0.2777)	0.8480(0.7531)
Simple-NeRF w/o smoothing aug	0.2832(0.4100)	0.6402(0.5448)	19.33(17.38)	0.1505(0.2805)	0.8334(0.7465)
Simple-NeRF w/o Lambertian aug	0.2834(0.4115)	0.6396(0.5438)	19.27(17.37)	0.1529(0.2760)	0.8306(0.7481)
Simple-NeRF 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)
Simple-NeRF w/o reliable depth	0.3020(0.4212)	0.6012(0.5115)	18.41(16.68)	0.2186(0.3644)	0.7564(0.6724)
Simple-NeRF 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)
Simple-NeRF w/ identical augs	0.2849(0.4110)	0.6379(0.5433)	19.27(17.38)	0.1509(0.2803)	0.8354(0.7466)

Table 11. Quantitative Results of NeRF based models 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)
Simple-NeRF	0.2559(0.3259)	0.6940(0.6222)	21.37(19.47)	0.1199(0.2201)	0.8935(0.8153)

Table 12. Quantitative Results of NeRF based models 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)
Simple-NeRF	0.2633(0.3083)	0.7016(0.6521)	21.99(20.44)	0.1110(0.1741)	0.9355(0.8952)

Table 13. Quantitative Results of NeRF based models 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)
Simple-NeRF	0.0635(0.0745)	0.8942(0.8783)	33.10(30.30)	0.3269(0.4584)	0.9215(0.8781)
Simple-NeRF w/o smoothing aug	0.0752(0.0889)	0.8886(0.8722)	32.54(29.85)	0.3795(0.5109)	0.8973(0.8487)
Simple-NeRF w/o Lambertian aug	0.0790(0.0925)	0.8884(0.8714)	32.13(29.62)	0.3870(0.5110)	0.8837(0.8428)
Simple-NeRF 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)
Simple-NeRF w/o reliable depth	0.0687(0.0802)	0.8913(0.8754)	32.63(30.33)	0.4485(0.5778)	0.8729(0.8404)
Simple-NeRF 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)
Simple-NeRF w/ identical augs	0.0777(0.0916)	0.8875(0.8713)	32.31(29.85)	0.4037(0.5269)	0.8949(0.8453)

Table 14. Quantitative Results of NeRF based models 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)
Simple-NeRF	0.0726(0.0829)	0.8984(0.8879)	33.21(31.40)	0.2770(0.3885)	0.9266(0.8931)

Table 15. Quantitative Results of NeRF based models 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)
Simple-NeRF	0.0847(0.0891)	0.8987(0.8917)	32.88(31.73)	0.2692(0.3565)	0.9209(0.9035)

Table 16. Quantitative Results of TensorRF based models on LLFF dataset with three input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
TensorRF	0.5474(0.6177)	0.3163(0.3002)	12.29(12.14)	0.6682(0.7824)	0.0252(0.0324)
DS-TensorRF	0.2897(0.3549)	0.6291(0.5702)	18.58(17.41)	0.2276(0.3483)	0.7279(0.6478)
Simple-TensorRF	0.2461(0.3037)	0.6749(0.6099)	20.22(18.71)	0.1671(0.2897)	0.8272(0.7451)
Simple-TensorRF w/ $R_\sigma^s = R_\sigma$	0.2536(0.3136)	0.6677(0.6020)	19.85(18.44)	0.1777(0.2986)	0.8112(0.7298)
Simple-TensorRF w/ $N_{vox}^s = N_{vox}$	0.2568(0.3182)	0.6579(0.5946)	19.95(18.54)	0.1902(0.3106)	0.7945(0.7128)
Simple-TensorRF w/ $R_\sigma^s = R_\sigma; N_{vox}^s = N_{vox}$	0.2728(0.3356)	0.6424(0.5814)	19.50(18.25)	0.2190(0.3425)	0.7446(0.6625)

Table 17. Quantitative Results of TensorRF based models on RealEstate-10K dataset with three input views. The values within parenthesis show unmasked scores.

Model	LPIPS	SSIM	PSNR	Depth MAE	Depth SROCC
TensorRF	0.0986(0.1050)	0.8532(0.8427)	29.62(28.00)	0.4394(0.4841)	0.6314(0.6077)
DS-TensorRF	0.0739(0.0827)	0.8872(0.8748)	32.50(30.20)	0.2720(0.3321)	0.7527(0.7227)
Simple-TensorRF	0.0706(0.0780)	0.8920(0.8809)	32.70(30.79)	0.2229(0.2882)	0.7983(0.7621)
Simple-TensorRF w/ $R_\sigma^s = R_\sigma$	0.0850(0.0919)	0.8821(0.8709)	30.94(29.51)	0.2742(0.3317)	0.7650(0.7238)
Simple-TensorRF w/ $N_{vox}^s = N_{vox}$	0.0735(0.0809)	0.8896(0.8773)	32.22(30.29)	0.2200(0.2888)	0.8223(0.7751)
Simple-TensorRF w/ $R_\sigma^s = R_\sigma; N_{vox}^s = N_{vox}$	0.0787(0.0867)	0.8871(0.8743)	31.73(29.88)	0.2346(0.3083)	0.7939(0.7458)

Table 18. Per-scene performance of various NeRF based 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)
Simple-NeRF	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 19. Per-scene performance of various NeRF based 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)
Simple-NeRF	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 20. Per-scene performance of various NeRF based 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)
Simple-NeRF	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 21. Per-scene performance of Simple-NeRF 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
Simple-NeRF	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)
Simple-NeRF w/o Smoothing Aug	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)
Simple-NeRF w/o Lambertian Aug	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)
Simple-NeRF 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)
Simple-NeRF 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)
Simple-NeRF 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)
Simple-NeRF 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 22. Per-scene performance of various NeRF based 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)
Simple-NeRF	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 23. Per-scene performance of various NeRF based 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)
Simple-NeRF	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 24. Per-scene performance of various NeRF based 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)
Simple-NeRF	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 25. Per-scene performance of Simple-NeRF 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
Simple-NeRF	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)
Simple-NeRF w/o Smoothing 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)
Simple-NeRF w/o Lambertian 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)
Simple-NeRF 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)
Simple-NeRF 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)
Simple-NeRF 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)
Simple-NeRF 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)

Table 26. Per-scene performance of various TensorRF based 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
TensorRF	0.61(0.66)	0.47(0.52)	0.46(0.51)	0.55(0.64)	0.43(0.57)	0.56(0.60)	0.66(0.70)	0.62(0.70)	0.55(0.62)
	0.28(0.26)	0.37(0.33)	0.37(0.38)	0.25(0.25)	0.20(0.14)	0.23(0.17)	0.50(0.51)	0.28(0.27)	0.32(0.30)
	12.4(12.4)	14.4(14.2)	14.4(13.8)	11.7(11.6)	11.8(11.3)	11.8(11.8)	12.2(12.2)	10.3(10.5)	12.3(12.1)
	0.47(0.56)	0.59(0.70)	0.67(0.88)	0.84(1.02)	0.74(0.79)	0.54(0.75)	0.79(0.81)	0.54(0.58)	0.67(0.78)
	0.13(0.12)	-0.18(-0.10)	0.05(0.05)	0.02(0.03)	-0.04(-0.02)	-0.08(-0.09)	-0.00(-0.00)	0.23(0.21)	0.03(0.03)
DS-TensorRF	0.40(0.46)	0.24(0.30)	0.13(0.17)	0.29(0.38)	0.31(0.46)	0.47(0.52)	0.31(0.33)	0.28(0.34)	0.29(0.35)
	0.47(0.44)	0.68(0.59)	0.80(0.76)	0.62(0.55)	0.37(0.28)	0.35(0.26)	0.79(0.77)	0.70(0.65)	0.63(0.57)
	17.8(17.5)	20.8(18.8)	23.3(21.2)	17.2(15.8)	14.6(13.3)	14.8(14.2)	19.8(19.3)	18.4(17.5)	18.6(17.4)
	0.19(0.27)	0.19(0.32)	0.12(0.33)	0.31(0.49)	0.47(0.55)	0.30(0.47)	0.20(0.23)	0.12(0.17)	0.23(0.35)
	0.55(0.49)	0.80(0.59)	0.92(0.85)	0.62(0.55)	0.54(0.41)	0.57(0.49)	0.80(0.83)	0.85(0.77)	0.73(0.65)
Simple-TensorRF	0.38(0.43)	0.19(0.25)	0.11(0.15)	0.23(0.32)	0.24(0.36)	0.40(0.43)	0.25(0.27)	0.28(0.32)	0.25(0.30)
	0.50(0.46)	0.72(0.63)	0.81(0.77)	0.69(0.61)	0.48(0.37)	0.42(0.32)	0.83(0.81)	0.71(0.66)	0.67(0.61)
	18.3(17.8)	22.4(20.2)	26.0(23.7)	19.9(17.6)	16.5(14.9)	16.0(15.1)	21.7(20.9)	18.1(17.4)	20.2(18.7)
	0.20(0.30)	0.15(0.29)	0.10(0.32)	0.21(0.37)	0.27(0.35)	0.25(0.41)	0.14(0.19)	0.10(0.16)	0.17(0.29)
	0.51(0.42)	0.83(0.68)	0.97(0.95)	0.80(0.71)	0.77(0.59)	0.77(0.69)	0.88(0.88)	0.89(0.80)	0.83(0.75)
Simple-TensorRF w/ $R_{\sigma}^s = R_{\sigma}$	0.40(0.46)	0.21(0.27)	0.13(0.17)	0.24(0.34)	0.25(0.38)	0.40(0.43)	0.24(0.27)	0.27(0.32)	0.25(0.31)
	0.48(0.44)	0.71(0.62)	0.80(0.77)	0.67(0.59)	0.46(0.36)	0.42(0.32)	0.83(0.81)	0.71(0.66)	0.67(0.60)
	18.1(17.8)	21.9(19.6)	24.3(22.5)	18.8(17.0)	16.1(14.6)	16.1(15.1)	22.0(21.0)	19.0(18.0)	19.9(18.4)
	0.21(0.29)	0.17(0.32)	0.11(0.31)	0.25(0.43)	0.28(0.38)	0.23(0.39)	0.13(0.16)	0.10(0.15)	0.18(0.30)
	0.52(0.45)	0.80(0.65)	0.96(0.91)	0.74(0.66)	0.74(0.59)	0.78(0.70)	0.87(0.88)	0.90(0.82)	0.81(0.73)
Simple-TensorRF w/ $N_{vox}^s = N_{vox}$	0.38(0.44)	0.21(0.26)	0.12(0.16)	0.23(0.33)	0.29(0.42)	0.43(0.47)	0.25(0.28)	0.27(0.32)	0.26(0.32)
	0.50(0.46)	0.70(0.61)	0.81(0.77)	0.68(0.60)	0.38(0.29)	0.39(0.29)	0.82(0.80)	0.71(0.66)	0.66(0.59)
	18.4(18.0)	21.7(19.5)	25.3(23.1)	19.7(17.7)	15.2(14.0)	15.4(14.6)	22.0(21.1)	18.5(17.8)	19.9(18.5)
	0.20(0.29)	0.18(0.32)	0.10(0.31)	0.22(0.38)	0.39(0.46)	0.30(0.48)	0.14(0.18)	0.10(0.16)	0.19(0.31)
	0.51(0.41)	0.76(0.57)	0.96(0.92)	0.82(0.75)	0.59(0.47)	0.66(0.59)	0.87(0.87)	0.89(0.80)	0.79(0.71)
Simple-TensorRF w/ $R_{\sigma}^s = R_{\sigma};$ $N_{vox}^s = N_{vox}$	0.38(0.44)	0.21(0.27)	0.12(0.16)	0.31(0.40)	0.29(0.43)	0.44(0.49)	0.24(0.26)	0.28(0.33)	0.27(0.34)
	0.50(0.46)	0.70(0.61)	0.81(0.77)	0.59(0.52)	0.41(0.32)	0.37(0.28)	0.83(0.81)	0.70(0.65)	0.64(0.58)
	18.3(17.8)	21.5(19.6)	25.1(23.0)	17.5(16.2)	15.3(14.0)	15.1(14.5)	22.4(21.4)	18.5(17.7)	19.5(18.2)
	0.20(0.29)	0.20(0.37)	0.11(0.31)	0.36(0.53)	0.39(0.45)	0.31(0.49)	0.12(0.16)	0.11(0.17)	0.22(0.34)
	0.51(0.40)	0.72(0.55)	0.97(0.93)	0.59(0.51)	0.59(0.45)	0.63(0.54)	0.90(0.90)	0.88(0.78)	0.74(0.66)

Table 27. Per-scene performance of various TensorRF based 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
TensorRF	0.0385(0.0407)	0.0498(0.0563)	0.3072(0.3196)	0.0590(0.0628)	0.0386(0.0458)	0.0986(0.1050)
	0.9580(0.9428)	0.9709(0.9642)	0.4348(0.4228)	0.9378(0.9279)	0.9644(0.9559)	0.8532(0.8427)
	32.66(29.91)	32.19(30.83)	16.16(15.91)	31.18(29.52)	35.90(33.81)	29.62(28.00)
	0.2546(0.2708)	0.3033(0.3398)	0.6816(0.7829)	0.0969(0.1215)	0.8604(0.9056)	0.4394(0.4841)
	0.6713(0.6488)	0.8540(0.8215)	0.3840(0.3944)	0.8204(0.7474)	0.4272(0.4263)	0.6314(0.6077)
DS-TensorRF	0.0333(0.0390)	0.0303(0.0391)	0.2288(0.2480)	0.0346(0.0381)	0.0423(0.0489)	0.0739(0.0827)
	0.9616(0.9476)	0.9801(0.9731)	0.5777(0.5531)	0.9533(0.9464)	0.9633(0.9538)	0.8872(0.8748)
	34.15(31.09)	38.01(34.85)	18.79(18.29)	36.35(33.48)	35.21(33.32)	32.50(30.20)
	0.1580(0.1822)	0.2152(0.2471)	0.4200(0.5451)	0.0363(0.0516)	0.5304(0.6343)	0.2720(0.3321)
	0.7528(0.7036)	0.9446(0.9348)	0.7748(0.7419)	0.9258(0.8734)	0.3654(0.3596)	0.7527(0.7227)
Simple-TensorRF	0.0413(0.0454)	0.0324(0.0395)	0.1957(0.2165)	0.0367(0.0374)	0.0466(0.0514)	0.0706(0.0780)
	0.9535(0.9423)	0.9793(0.9732)	0.6197(0.5939)	0.9514(0.9480)	0.9562(0.9470)	0.8920(0.8809)
	33.57(30.98)	38.17(34.72)	19.65(19.17)	36.88(35.51)	35.25(33.58)	32.70(30.79)
	0.1979(0.2127)	0.1998(0.2438)	0.3293(0.4552)	0.0508(0.0572)	0.3366(0.4721)	0.2229(0.2882)
	0.7026(0.6628)	0.9102(0.8768)	0.8654(0.8077)	0.8866(0.8736)	0.6267(0.5896)	0.7983(0.7621)
Simple-TensorRF w/ $R_{\sigma}^s = R_{\sigma}$	0.0556(0.0585)	0.0784(0.0833)	0.1941(0.2149)	0.0472(0.0493)	0.0494(0.0535)	0.0850(0.0919)
	0.9429(0.9303)	0.9426(0.9406)	0.6248(0.5974)	0.9460(0.9415)	0.9542(0.9448)	0.8821(0.8709)
	31.96(29.66)	32.86(31.77)	19.59(19.01)	35.39(33.83)	34.93(33.26)	30.94(29.51)
	0.2579(0.2734)	0.3988(0.4114)	0.3190(0.4456)	0.0662(0.0714)	0.3290(0.4567)	0.2742(0.3317)
	0.5966(0.5550)	0.8298(0.7942)	0.8539(0.7936)	0.8752(0.8647)	0.6694(0.6118)	0.7650(0.7238)
Simple-TensorRF w/ $N_{vox}^s = N_{vox}$	0.0503(0.0545)	0.0360(0.0425)	0.1990(0.2183)	0.0329(0.0344)	0.0494(0.0548)	0.0735(0.0809)
	0.9483(0.9356)	0.9777(0.9714)	0.6132(0.5851)	0.9539(0.9495)	0.9547(0.9447)	0.8896(0.8773)
	32.49(30.12)	36.51(33.82)	19.46(18.91)	37.50(35.20)	35.15(33.41)	32.22(30.29)
	0.2154(0.2294)	0.1612(0.2106)	0.3467(0.4895)	0.0403(0.0456)	0.3362(0.4688)	0.2200(0.2888)
	0.7102(0.6776)	0.9554(0.9091)	0.8388(0.7518)	0.9083(0.8993)	0.6990(0.6379)	0.8223(0.7751)
Simple-TensorRF w/ $R_{\sigma}^s = R_{\sigma}; N_{vox}^s = N_{vox}$	0.0485(0.0523)	0.0522(0.0583)	0.2012(0.2236)	0.0405(0.0428)	0.0512(0.0564)	0.0787(0.0867)
	0.9495(0.9358)	0.9739(0.9681)	0.6079(0.5793)	0.9504(0.9454)	0.9537(0.9431)	0.8871(0.8743)
	32.55(30.02)	35.56(33.27)	19.33(18.80)	36.22(34.12)	34.97(33.18)	31.73(29.88)
	0.2143(0.2284)	0.2039(0.2483)	0.3522(0.5015)	0.0555(0.0626)	0.3474(0.5006)	0.2346(0.3083)
	0.7181(0.6866)	0.8937(0.8607)	0.7969(0.7146)	0.8851(0.8694)	0.6756(0.5976)	0.7939(0.7458)

Table 28. Per-scene performance of various ZipNeRF based models with twelve input views on the MipNeRF360 dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively.

Model \ Scene Name	Bicycle	Bonsai	Counter	Garden	Kitchen	Room	Stump	Average
ZipNeRF	0.6617	0.5722	0.7258	0.4130	0.5927	0.3978	0.6249	0.5614
	0.2044	0.5486	0.4330	0.4385	0.4072	0.7121	0.2591	0.4616
	14.77	15.79	12.41	18.65	14.05	18.65	17.19	15.86
	8.4860	2.9808	15.1470	1.5437	15.9720	0.4371	8.7365	7.4260
	-0.3807	0.4405	-0.0017	0.8949	-0.0473	0.7352	0.0952	0.2755
Augmented ZipNeRF	0.7995	0.6486	0.7180	0.8011	0.6239	0.5824	0.7060	0.6825
	0.2433	0.5089	0.4790	0.2863	0.4916	0.6213	0.2705	0.4462
	14.30	15.73	13.78	16.16	18.92	17.65	16.31	16.27
	428.2730	26.6151	50.3386	23.5793	7.0898	3.7203	356.3442	96.4194
	-0.0204	0.5007	0.2334	0.7987	0.9650	0.5605	0.0301	0.4869
Simple-ZipNeRF	0.6588	0.5123	0.6105	0.4698	0.2888	0.3918	0.6278	0.4876
	0.2015	0.5851	0.5094	0.3995	0.7017	0.7115	0.2616	0.5245
	14.77	17.46	14.79	17.80	20.50	19.09	17.31	17.60
	9.0387	3.7980	3.3017	0.5187	1.1201	2.9600	6.0815	3.5434
	-0.1401	0.5518	0.2898	0.8854	0.9621	0.6553	-0.0605	0.5103

Table 29. Per-scene performance of various ZipNeRF based models with twenty input views on the MipNeRF360 dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively.

Model \ Scene Name	Bicycle	Bonsai	Counter	Garden	Kitchen	Room	Stump	Average
ZipNeRF	0.6554	0.3793	0.7201	0.2472	0.2799	0.3114	0.6076	0.4350
	0.2293	0.6912	0.4716	0.6358	0.7406	0.7889	0.2732	0.5911
	14.97	18.43	13.45	22.49	20.23	23.15	17.58	18.89
	5.5415	2.0900	2.2907	0.3442	1.0369	0.3183	5.8841	2.1151
	-0.1092	0.6476	0.0093	0.9838	0.7730	0.8285	-0.0517	0.5032
Augmented ZipNeRF	0.7825	0.5268	0.5903	0.7503	0.5722	0.5380	0.7340	0.6190
	0.2756	0.6568	0.6164	0.3362	0.5333	0.6867	0.3021	0.5244
	15.44	21.23	18.90	19.04	20.83	20.00	17.14	19.31
	71.2764	6.4862	6.0847	6.7096	2.2221	1.2190	246.0264	31.2000
	0.3799	0.7704	0.7391	0.9282	0.9753	0.7951	-0.0616	0.7117
Simple-ZipNeRF	0.6415	0.2783	0.3624	0.2311	0.1777	0.3096	0.5890	0.3421
	0.2266	0.7630	0.6922	0.6359	0.8155	0.7796	0.2582	0.6456
	15.11	22.39	20.23	21.85	23.80	22.68	17.30	21.03
	7.2913	1.9279	1.0639	0.8546	6.4959	0.6978	7.3037	3.2887
	0.2979	0.8119	0.7890	0.9569	0.9466	0.8305	-0.0959	0.7190

Table 30. Per-scene performance of various ZipNeRF based models with thirty-six input views on the MipNeRF360 dataset. The five rows show LPIPS, SSIM, PSNR, Depth MAE and Depth SROCC scores, respectively.

Model \ Scene Name	Bicycle	Bonsai	Counter	Garden	Kitchen	Room	Stump	Average
ZipNeRF	0.5717	0.1689	0.7272	0.1485	0.1416	0.2159	0.5632	0.3316
	0.2810	0.8687	0.4546	0.7571	0.8636	0.8515	0.2733	0.6737
	16.09	25.54	13.01	24.52	25.15	25.56	17.75	21.78
	12.5203	0.9944	2.2470	0.2896	0.5732	0.5515	5.0126	2.6502
	0.0187	0.9025	0.0558	0.9905	0.9455	0.8953	0.2225	0.6353
Augmented ZipNeRF	0.7788	0.4916	0.5547	0.7130	0.5420	0.5022	0.7448	0.5917
	0.3422	0.6942	0.6623	0.3682	0.5589	0.7275	0.3397	0.5646
	18.81	22.93	20.98	20.40	22.11	22.07	18.63	21.21
	14.5803	4.3221	4.4459	5.4332	3.1258	0.7287	94.4808	11.8336
	0.9211	0.8651	0.8508	0.9660	0.9870	0.8210	0.1884	0.8414
Simple-ZipNeRF	0.4057	0.1688	0.2473	0.1663	0.1284	0.2318	0.4938	0.2390
	0.4448	0.8638	0.7872	0.7417	0.8702	0.8398	0.3699	0.7458
	20.49	26.04	23.23	24.06	26.71	25.16	19.73	24.19
	4.9626	2.2896	3.0641	0.4073	12.9376	4.1386	10.2534	5.2852
	0.8906	0.8646	0.8717	0.9800	0.9185	0.7992	0.6079	0.8591

Table 31. Per-scene performance of various ZipNeRF based models with four input views on NeRF-Synthetic dataset. The five rows show LPIPS, SSIM, and PSNR, respectively.

Model \ Scene Name	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship	Average
ZipNeRF	0.3172	0.4762	0.2500	0.4280	0.5128	0.4689	0.3186	0.6389	0.4263
	0.8372	0.7501	0.8430	0.7869	0.6811	0.7340	0.8562	0.5496	0.7548
	14.12	10.75	14.47	11.33	9.39	8.63	13.20	6.43	11.04
Simple-ZipNeRF	0.2848	0.4635	0.2114	0.4083	0.5021	0.4085	0.2903	0.5331	0.3878
	0.8454	0.7446	0.8500	0.8087	0.6867	0.7592	0.8609	0.6162	0.7715
	14.54	11.28	15.17	11.38	9.40	9.05	14.63	6.54	11.50

Table 32. Per-scene performance of various ZipNeRF based models with eight input views on NeRF-Synthetic dataset. The five rows show LPIPS, SSIM, and PSNR, respectively.

Model \ Scene Name	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship	Average
ZipNeRF	0.2227	0.3508	0.0813	0.2773	0.4562	0.4464	0.0453	0.4216	0.2877
	0.8471	0.7706	0.8980	0.8391	0.6991	0.7388	0.9571	0.6282	0.7973
	16.86	11.90	21.91	14.21	9.72	8.67	26.99	9.82	15.01
Simple-ZipNeRF	0.1552	0.2572	0.0879	0.2318	0.4900	0.2418	0.0405	0.4641	0.2461
	0.8687	0.8084	0.8964	0.8481	0.6871	0.7801	0.9541	0.6075	0.8063
	18.53	14.17	21.68	14.56	9.59	12.30	26.87	9.31	15.88

Table 33. Per-scene performance of various ZipNeRF based models with twelve input views on NeRF-Synthetic dataset. The five rows show LPIPS, SSIM, and PSNR, respectively.

Model \ Scene Name	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship	Average
ZipNeRF	0.1012	0.2181	0.0474	0.1278	0.1149	0.1659	0.0402	0.4844	0.1625
	0.9062	0.8143	0.9336	0.9055	0.8718	0.8403	0.9613	0.5898	0.8528
	23.40	14.72	25.18	21.09	20.24	18.52	27.39	10.39	20.12
Simple-ZipNeRF	0.1010	0.1458	0.0569	0.1227	0.1116	0.1123	0.0342	0.5411	0.1532
	0.9013	0.8475	0.9236	0.9055	0.8602	0.8617	0.9614	0.5634	0.8531
	22.27	18.27	24.01	21.24	20.12	20.59	28.32	9.27	20.51