
HG³-NeRF: Hierarchical Geometric, Semantic, and Photometric Guided Neural Radiance Fields for Sparse View Inputs

Zelin Gao¹, Weichen Dai², Yu Zhang^{1*}

¹ College of Control Science and Engineering, Zhejiang University

² School of Computer Science, Hangzhou Dianzi University

Abstract

Neural Radiance Fields (NeRF) have garnered considerable attention as a paradigm for novel view synthesis by learning scene representations from discrete observations. Nevertheless, NeRF exhibit pronounced performance degradation when confronted with sparse view inputs, consequently curtailing its further applicability. In this work, we introduce **Hierarchical Geometric, Semantic, and Photometric Guided NeRF (HG³-NeRF)**, a novel methodology that can address the aforementioned limitation and enhance consistency of geometry, semantic content, and appearance across different views. We propose Hierarchical Geometric Guidance (HGG) to incorporate the attachment of Structure from Motion (SfM), namely sparse depth prior, into the scene representations. Different from direct depth supervision, HGG samples volume points from local-to-global geometric regions, mitigating the misalignment caused by inherent bias in the depth prior. Furthermore, we draw inspiration from notable variations in semantic consistency observed across images of different resolutions and propose Hierarchical Semantic Guidance (HSG) to learn the coarse-to-fine semantic content, which corresponds to the coarse-to-fine scene representations. Experimental results demonstrate that HG³-NeRF can outperform other state-of-the-art methods on different standard benchmarks and achieve high-fidelity synthesis results for sparse view inputs.

1 Introduction

Novel View Synthesis (NVS) is one of the crucial tasks in computer vision, aiming to generate images of unseen views through visual information, similar to how humans perceive and visualize their surroundings. Recent methods have opted to recover intermediate dense 3D-aware representation [57, 22, 37], multi-plane images [55, 46, 15], or volume density [3, 30, 29], followed by neural rendering theorem [38, 47, 11] to synthesize images. In particular, Neural Radiance Fields (NeRF) [29] have demonstrated remarkable potential with state-of-the-art performance in generating high-fidelity view synthesis results. However, NeRF suffer from significant challenges with sparse view inputs, primarily due to their reliance on dense scene coverage to mitigate the shape ambiguity problem [52, 60] arising from only photometric supervision.

Existing methods address this challenging issue with a variety of strategies, which can be classified into pre-training methods [25, 20, 4, 6, 58] and per-scene optimization methods [35, 21, 17]. Pre-training methods train the model on large-scale datasets and further fine-tune it for each scene at test time. However, the generalization ability heavily depends on the quality of datasets, and it is too expensive to obtain necessary datasets by capturing many different scenes. The alternative methods

*Corresponding authors.

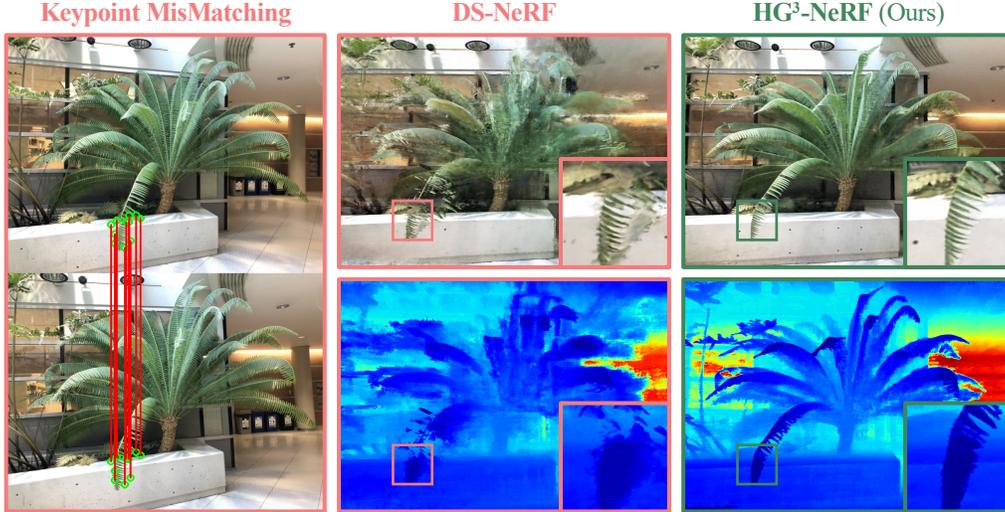


Figure 1: **Novel View Synthesis Results from 3 View Inputs.** Left: The bias in the sparse depth prior is caused by keypoint mismatching during the multi-view stereo process of SfM. Mid: Bias in the depth prior is introduced to NeRF through depth supervision and leads to geometric misalignment. Right: Our HG^3 -NeRF leverage additional guidance to learn the scene representations, showing great performance in both appearance and geometry.

are to train the network from scratch for each scene. Though most of these methods introduce additional regularization to prevent overfitting issues, they still suffer from geometric misalignment without real-world geometric supervision. Therefore, some methods [10] leverage the supervision function contributed by the rough depth prior estimated by Structure from Motion (SfM) [42, 1] to improve the geometric consistency by enforcing the distribution of density and color around the depth prior. However, as shown in the left and mid of Fig. 1, direct depth supervision can also introduce bias generated from the multi-view stereo process [32, 31, 16] into the scene representations, resulting in geometric misalignment and poor view synthesis quality.

It should also be noted that NeRF is actually able to recover reasonable geometric results with dense view inputs [52, 4]. However, given the sparse viewpoints, it is too difficult for NeRF to learn the consistent color and density distribution across different views with only photometric supervision. Due to the fact that the sparse depth prior can indicate the approximate distribution of the color and density along the ray, the sparse depth prior should be used to guide the volume sampling to improve the geometric consistency. As shown in the right of Fig. 1, sampling volume points with the guidance of the sparse depth prior is a more precise strategy to improve geometric consistency rather than direct depth supervision.

In this paper, we exploit the geometric, semantic, and photometric guidance to represent the neural radiance fields from sparse view inputs. We propose hierarchical geometric guidance (HGG) to sample volume points with the depth prior, which is generated as a common attachment in NeRF pipelines by running SfM. Since the depth prior provides approximate locations of the density and color along the rays, we first sample points around the depth prior in a local region to initialize the neural radiance fields and then gradually extend the sampling region to full scene bounds for learning the global scene representations. Different from direct depth supervision, the HGG method utilizes a local-to-global sampling strategy to incorporate the depth prior into the representations and mitigate the potential geometric misalignment caused by bias in the depth prior. Furthermore, we propose hierarchical semantic guidance (HSG) to supervise semantic consistency of the complex real-world scenarios using CLIP [40]. We draw inspiration from the diverse variations in semantic consistency observed across images of different resolutions, where the semantic content of low-resolution images is hard to match with that of high-resolution images. Since the scene representations are learned from coarse to fine, the rendered images are blurred like low-resolution images at the start of training [23]. Therefore, the HSG method first leverages coarse feature vectors from the down-sampled images to supervise the semantic consistency. As the iteration increases, it gradually aggregates more content into the feature vectors by reducing the down sampling rate. Finally, we adopt the hierarchical photometric guidance proposed in NeRF to supervise the appearance consistency. Combined, we

call our method HG^3 -NeRF, which incorporate hierarchical geometric, semantic, and photometric guidance to represent the neural radiance fields from sparse view inputs.

In the experiments, we evaluate the effectiveness of the proposed HG^3 -NeRF in comparison to state-of-the-art baselines on various standard benchmarks. Furthermore, we conduct the model analysis including a comprehensive ablation study to investigate the contributions of HGG and HSG, respectively, as well as a view synthesis comparison study to demonstrate that our HGG method can achieve comparable performance to represent forward-facing scenarios in real-world space instead of using NDC (Normalized Device Coordinate) space.

To summarize, our contributions are as follows:

- We propose HG^3 -NeRF, a novel methodology that can exploit the hierarchical geometric, semantic, and photometric guidance to maintain consistency across different views to represent the neural radiance fields from sparse view inputs.
- We propose hierarchical geometric guidance (HGG) that incorporates sparse depth prior to the scene representations without introducing bias and hierarchical semantic guidance (HSG) that guides the coarse-to-fine semantic supervision of complex real-world scenarios.
- We conduct the experiments on various datasets and show that the proposed HG^3 -NeRF can achieve realistic synthesis results for sparse viewpoint inputs.

2 Related Works

Neural Scene Representations. In computer vision, coordinate-based neural representations [44, 5, 26, 27] have become one of the most popular representation methods for various 3D vision tasks, such as 3D reconstruction [59, 48, 50], 3D-aware generation [57, 22, 37] and novel view synthesis [3, 30, 35, 29]. As opposed to explicit representations such as point clouds [39, 14, 12], voxels [8, 53] or meshes [9, 33, 2], this paradigm signifies that color and 3D geometric information can be represented by the implicit neural network, leading to a more compact representation format. There are several works [24, 29, 36, 45] learning the scene representations from multi-view images using neural volume rendering. Among these methods, Neural Radiance Fields (NeRF) have demonstrated remarkable potential in generating high-fidelity images from novel viewpoints as well as its simplicity to represent the scene as a continuous implicit function. Therefore, various follow-up works have been proposed to improve the performance and generality of NeRF such as large-scale scene representations from city or satellite viewpoints [49, 54], real-time neural volume rendering for fast training [30, 13], and unbounded scene representations [3, 60]. Although these methods have extended NeRF to various domains with impressive performance, they typically require dense viewpoints to learn accurate scene representations for synthesizing realistic images, limiting their further application in real-world scenarios [35, 21]. In this work, we focus on solving the issue where only sparse view inputs are available for NeRF, which is much closer to real-world applications.

Novel View Synthesis from Sparse View Inputs. To tackle the problem of representing the neural radiance fields from sparse view inputs, several works are proposed for the further real-world application of NeRF, which can be classified into two main categories: pre-training methods and per-scene optimization methods. The pre-training methods [25, 20, 18, 41, 4, 6, 58] circumvent the requirement of dense view inputs by pre-training a conditional model to aggregate sufficient prior knowledge for reconstructing the neural radiance fields. They propose to train a generalizable model with the high-dimensional image features extracted from a CNN backbone network [6, 58] or the 3D cost volume obtained by image warping [4, 20, 25]. Though these methods achieve great results for sparse view inputs, large-scale datasets of many different scenarios are required for pre-training. Moreover, these methods require fine-tuning the network parameters and suffer from quality degradation on untrained domains. The per-scene optimization methods [43, 35, 21, 56, 17] propose to supervise NeRF with additional regularization loss function instead of using expensive pre-training models. Although various regularization functions over geometry [56, 35], appearance [35], semantics [56, 17], and density distribution [21] are used to supervise the consistency between seen and unseen viewpoints, the performance on geometric consistency is still hard to be improved due to the lack of real-world geometric supervision. Some methods [10] directly supervise the geometry of NeRF with the depth prior estimated by SfM, however, geometric misalignment still exists since the bias in the depth prior is introduced into the scene representations during the optimization.

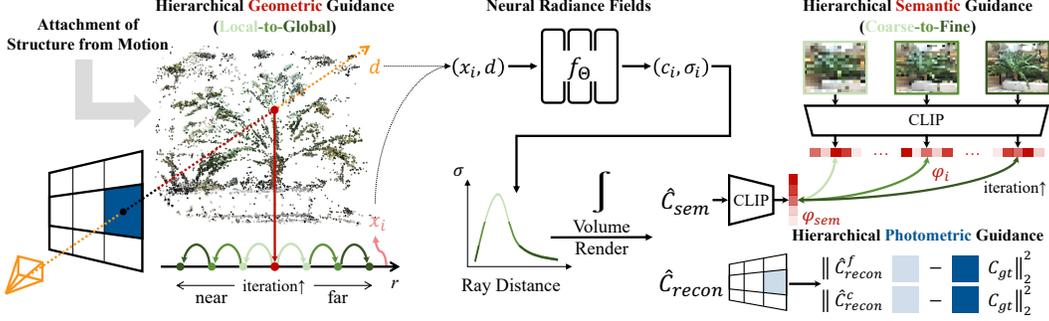


Figure 2: **Overview of HG³-NeRF.** The 3D volume location \mathbf{x}_i is first sampled within the local-to-global region setup by hierarchical geometric guidance and then fed into neural radiance fields along with the viewing direction \mathbf{d} to query color \mathbf{c}_i and density σ_i . Via the volume rendering theorem, the query results are integrated into $\hat{\mathbf{C}}_{recon}^c$, which contains $\hat{\mathbf{C}}_{recon}^c$ for the coarse model and $\hat{\mathbf{C}}_{recon}^f$ for the fine model. Moreover, we employ CLIP to encode the image $\hat{\mathbf{C}}_{sem}$ rendered from a randomly selected pose as a feature vector φ_{sem} . The scene representations are finally optimized by the coarse-to-fine cosine similarity between φ_{sem} and φ_i from hierarchical semantic guidance as well as the MSE between $\hat{\mathbf{C}}_{recon}$ and the observed color \mathbf{C}_{gt} .

3 Preliminaries

3.1 Neural Radiance Fields

In this work, we follow the framework proposed in NeRF to integrate the predicted color and density along the ray. Specifically, NeRF adopt multi-layer perceptrons (MLPs) to reconstruct a radiance field from dense input views, where the view-dependent appearance is modeled as a continuous function $f_{\Theta} : (\mathbf{x}_i, \mathbf{d}) \rightarrow (\mathbf{c}_i, \sigma_i)$, which maps a 3D location $\mathbf{x}_i = (x_i, y_i, z_i)$ with its 2D viewing direction $\mathbf{d} = (\theta, \phi)$ into a volume color $\mathbf{c}_i = (r_i, g_i, b_i)$ and a density σ_i .

Neural Volume Rendering. Given a ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ whose origin is at the camera’s center of projection \mathbf{o} , a volume location $t_i \in [t_n, t_f]$ is sampled within the near and far planes. By querying f_{Θ} , the final color $\hat{\mathbf{C}}_{recon}(\mathbf{r})$ is approximated via the volume rendering theorem as

$$\hat{\mathbf{C}}_{recon}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \text{ where } T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right) \quad (1)$$

Note that T_i denotes how much light is transmitted on ray \mathbf{r} up to sample i and $\delta_i = t_{i+1} - t_i$ denotes the interval between the i -th sample and its adjacent one.

Hierarchical Photometric Guidance. To improve the performance of NVS, NeRF propose to train a coarse model and a fine model, with parameters Θ^c and Θ^f , respectively. The hierarchical photometric guidance (HPG) can be formulated as:

$$\mathcal{L}_{hpg} = \sum_{\mathbf{r} \in \mathcal{R}} \left(\left\| \mathbf{C}_{gt}(\mathbf{r}) - \hat{\mathbf{C}}_{recon}^c(\mathbf{r}; \mathbf{t}^c) \right\|_2^2 + \left\| \mathbf{C}_{gt}(\mathbf{r}) - \hat{\mathbf{C}}_{recon}^f(\mathbf{r}; \mathbf{t}^c \cup \mathbf{t}^f) \right\|_2^2 \right) \quad (2)$$

where \mathcal{R} denotes the set of all rays across all images, $\hat{\mathbf{C}}_{recon}^c(\mathbf{r}; \mathbf{t}^c)$ denotes the predicted color of the coarse model with stratified volume samples \mathbf{t}^c , and $\hat{\mathbf{C}}_{recon}^f(\mathbf{r}; \mathbf{t}^c \cup \mathbf{t}^f)$ denotes the predicted color of the fine model with the union of \mathbf{t}^c and \mathbf{t}^f which is produced by inverse transform sampling.

3.2 Contrastive Language-Image Pre-Training

Contrastive Language-Image Pre-Training (CLIP) model is trained for learning joint representations between text and images [34]. CLIP consist of an image encoder and a text encoder, demonstrating the excellent few-shot transfer performance to image recognition tasks (for brevity, the following CLIP only refers to its image encoder f_{clip}). Though the semantic consistency implemented by CLIP can improve the performance of object-level scene representations [17, 56], it is still severely limited for complex real-world scenarios [35].

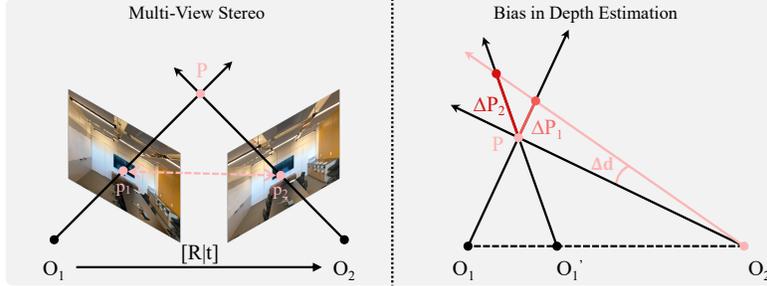


Figure 3: **Bias from Multi-View Stereo.** Classical stereo methods estimate depth through geometric constraints based on keypoint matching. The shift Δd in ray direction is caused by keypoint mismatching and thus introduces the bias into the depth estimation. Especially when the translation between two frames is small, the bias further increases.

4 Methodology

In this paper, we propose HG³-NeRF to represent the neural radiance fields for sparse view inputs. The overview of our method is illustrated in Fig. 2. Given the attachment of SfM, i.e., sparse depth prior, and a ray generated from the camera pose, we first sample volume points within the local-to-global region using hierarchical geometric guidance (Sec. 4.1) and then query the corresponding results to predict the final color via the volume rendering theorem. We also randomly select a camera pose to render an image and encode it into the corresponding feature vector using CLIP f_{clip} . Finally, we compute the coarse-to-fine semantic cosine similarity under hierarchical semantic guidance (Sec. 4.2) and the appearance MSE under hierarchical photometric guidance (Eq. 2) to optimize the neural radiance fields.

4.1 Hierarchical Geometric Guidance

We propose hierarchical geometric guidance (HGG), which incorporates the geometric consistency into the scene representations by using sparse depth prior from SfM. The HGG method guides NeRF to learn the approximate distribution of the density and color from a local-to-global sampling region determined by the depth prior.

As the common attachment of SfM, the sparse depth prior is estimated by the multi-view triangulation method, which can be formulated as:

$$s_1 p_1 = s_2 R p_2 + t \quad (3)$$

where p_1, p_2 denote the locations of two matched keypoints on the normalized coordinate, s_1, s_2 denote the depth values, and R, t denote rotation matrix and translation vector between these two points. As shown in Fig. 3, the accuracy of estimated 3D point P relies on the quality of keypoint matching. The keypoint mismatching can cause a shift Δd in the ray direction and thus introduces a bias ΔP_1 into the estimated depth. In addition, the small translation between two frames can further increase the bias to ΔP_2 .

The depth estimation of NeRF can be formulated as:

$$\hat{\mathbf{D}}_{recon}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) t_i \quad (4)$$

Although sparse depth prior from SfM can be used to supervise $\hat{\mathbf{D}}_{recon}$, the optimization function contributed by this sparse depth prior can introduce the aforementioned bias into the scene representations, resulting in geometric misalignment and poor view synthesis results.

To address this issue, the HGG method utilizes sparse depth prior to indicate the sampling region for learning the density distribution along the ray instead of direct supervision. We first initialize the scene representations by setting near planes t_n and far planes t_f close to the depth prior for sampling volume points from a local region. Then, t_n and t_f are gradually extended to the full scene bounds

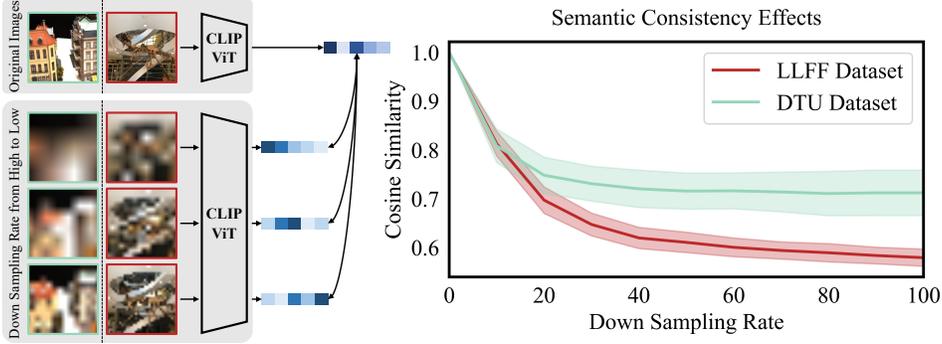


Figure 4: **Semantic Consistency Effects on Multi-Resolution Images.** By computing the cosine similarity of the feature vectors between the original image and its down sampling results, we can find that the cosine similarity decreases sharply with the increase of down sampling rate. Especially when the dataset is contributed by complex real-world scenes (e.g., LLFF dataset [28]), meaning that an image contains plenty of content, the effect on semantic consistency supervision is further limited.

for learning the global volume distribution along the ray, which can be formulated as:

$$\begin{aligned} t_n(i) &= t_{depth} + (t_n - t_{depth}) \gamma(i) \\ t_f(i) &= t_{depth} + (t_f - t_{depth}) \gamma(i) \end{aligned} \quad (5)$$

where $\gamma(i)$ denotes the region adjustment rate at i iteration, which can be formulated as:

$$\gamma(i) = \frac{1 - \cos((\min(\max(i/N_{hgg}, \epsilon_{hgg}), 1)) \pi)}{2} \quad (6)$$

where ϵ_{hgg} denotes the minimum adjusting rate, and N_{hgg} denotes how many iterations until the full bounds are reached. The HGG method can utilize the information provided by the depth prior, i.e., the approximate density distribution of the rays, and reconstruct the neural radiance fields with the local-to-global geometric guidance of the real world without introducing depth bias into the scene representations.

4.2 Hierarchical Semantic Guidance

Inspired by the diverse variations in semantic consistency on multi-resolution images, we propose hierarchical semantic guidance to supervise the coarse-to-fine semantic content corresponding to the coarse-to-fine scene representations with the CLIP encoder.

As shown in Fig. 4, the effect of CLIP is limited for supervising the semantic consistency over multi-resolution images. NeRF first learn low-frequency information for coarse scene representations [23], resulting in images rendered at the beginning of training that resemble those low-resolution images containing little content. Therefore, it is difficult to match the feature vectors encoded from these rendered images with those encoded from high-resolution original images.

To solve this problem, the proposed HSG randomly selects a known camera pose and performs the semantic supervision between the rendered images and the original images of the same scale by sampling a set of pixels (rays) \mathcal{P} using a coarse-to-fine grid sampling strategy, which can be formulated as:

$$\mathcal{P}(s_i) = \{(u, v) \mid u \in [0, H, stride = s_i], v \in [0, W, stride = s_i]\} \quad (7)$$

where (u, v) denotes pixel location and s_i denotes sampling stride, which can be formulated as:

$$s_i = \max\left(\text{ceil}\left(s_{max} \cdot \frac{1 + \cos((i/N_{hsg})\pi)}{2}\right), 1\right) \quad (8)$$

Note that s_{max} denotes the maximum sampling stride, N_{hsg} denotes how many iterations until the full pixels are sampled, and ceil denotes rounding operation. Then, feature vectors of the sampled $\hat{\mathbf{C}}_{sem}$ and $\hat{\mathbf{C}}_i$ are encoded by CLIP f_{clip} , which can be formulated as:

$$\varphi_{sem} = f_{clip}(\hat{\mathbf{C}}_{sem}), \varphi_i = f_{clip}(\hat{\mathbf{C}}_i) \quad (9)$$

The coarse-to-fine semantic consistency supervision thus can be performed, since more scene content is aggregated into the feature vectors with the sampling stride s_i gradually decreasing to 1.

4.3 Reconstruction Loss Function

As we discuss in Sec. (4.1), HG³-NeRF represent the neural radiance fields without using explicit depth supervision for sparse view inputs. Therefore, our reconstruction loss is constructed with hierarchical photometric guidance and hierarchical semantic guidance.

HPG Loss Function. We follow the loss function \mathcal{L}_{hpg} (Eq. (2)) proposed in NeRF to supervise the reconstruction of the scene appearance.

HSG Loss Function. The HSG method allows HG³-NeRF to supervise the the coarse-to-fine semantic consistency by computing the cosine similarity \mathcal{L}_{hsg} between φ_{sem} and φ_i , which can be formulated as:

$$\mathcal{L}_{hsg} = \varphi_{sem}^T \cdot \varphi_i \quad (10)$$

Finally, we use these two loss functions to optimize the scene representations, which can be formulated as the following minimization problem:

$$\min_{\Theta_c, \Theta_f} (\mathcal{L}_{hpg} + \lambda \mathcal{L}_{hsg}) \quad (11)$$

where λ denotes the weighting factor for \mathcal{L}_{hsg} .

5 Experiments

Dataset. We validate the proposed HG³-NeRF using the standard datasets of different levels: DTU dataset [19] and LLFF dataset [28]. DTU dataset consists of object-level scenes where the image content contains objects on a white table with a black background. LLFF dataset consists of real-world-level scenes with sequentially captured images from hand-held cameras.

Baselines. We compare against state-of-the-art baselines including pre-training methods and per-scene optimization methods. Moreover, we adopt 'w *sdp*' to label the method using the sparse depth prior and 'ft' to label the method to fine-tune for each scene. Note that [†]DS-NeRF is our re-implementation, since there is no direct way to evaluate DS-NeRF in their released code.

Evaluation Metrics. We follow the official evaluation metrics and thus adopt the mean of PSNR, structural similarity index (SSIM) [51], and the LPIPS [61] perceptual metric.

5.1 Implement Details

We implement the framework of HG³-NeRF based on NeRF. The pre-training baselines are pre-trained on the large-scale DTU dataset. The regularization baselines are trained from scratch for each scene. We use the Adam optimizer with a learning rate of 1×10^{-3} exponentially decaying to 1×10^{-5} and train our method on two NVIDIA V100 GPUs. We sample 64 volume points for the coarse model and 128 volume points for the fine model. We set N_{hgg} and N_{hsg} as 10% and 50% of total training iterations. We set $s_{max} = 0.1 \cdot \min(H, W)$, $\epsilon_{hgg} = 0.2$, and $\lambda = 0.2$ where H, W denote the height and width of the image. Following experimental protocols of baselines [35, 58], we compare the performance of HG³-NeRF against these baselines for the scenarios of 3, 6, and 9 views.

5.2 Comparison with SOTA Baselines

Comparison on DTU Dataset. We first evaluate the performance of our method in the object-level DTU dataset. Fig. 5 presents the qualitative results on object-level DTU dataset. HG³-NeRF can synthesize images with more fine details and estimate depth maps with distinct edges in comparison to other SOTA baselines. As we discuss in Sec. 4.1, bias can be introduced into the scene representations through the loss supervision function contributed by the depth prior. Our HG³-NeRF incorporate geometric consistency using the HGG method, where the sparse depth prior is employed to guide a local-to-global volume sampling rather than a loss function. Additionally, our HSG method can further improve the semantic consistency of the scene representations by computing a coarse-to-fine cosine similarity. Therefore, as illustrated in Table. 1, HG³-NeRF can outperform most SOTA baselines with high consistency between seen and unseen views even with only 3 views available (PSNR=19.37). Moreover, as the semantic content is aggregated into the scene representation from coarse to fine, the synthesis quality can be further improved (PSNR=25.87, 9 views) in comparison to DietNeRF (23.83, 9 views), which directly compute the cosine similarity between rendered images and original images.

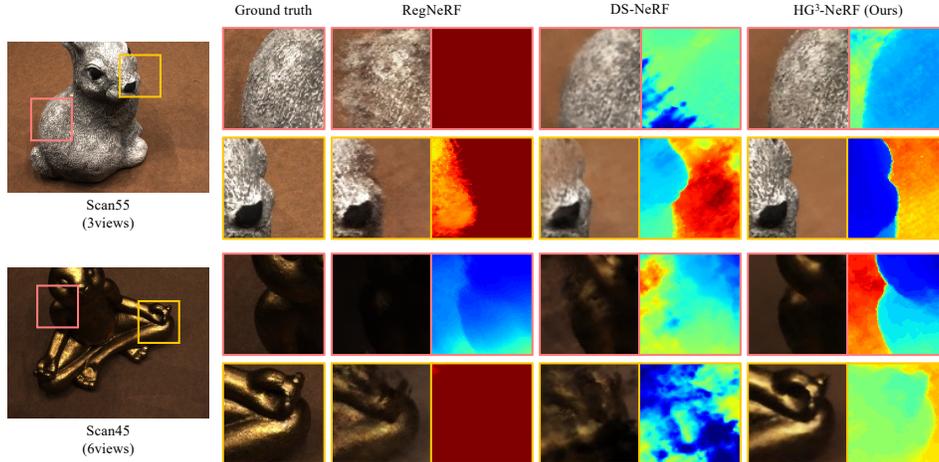


Figure 5: **Qualitative Results on DTU Dataset.** HG^3 -NeRF maintain more fine details in the synthetic images and more sharper edges in the estimated depth maps.

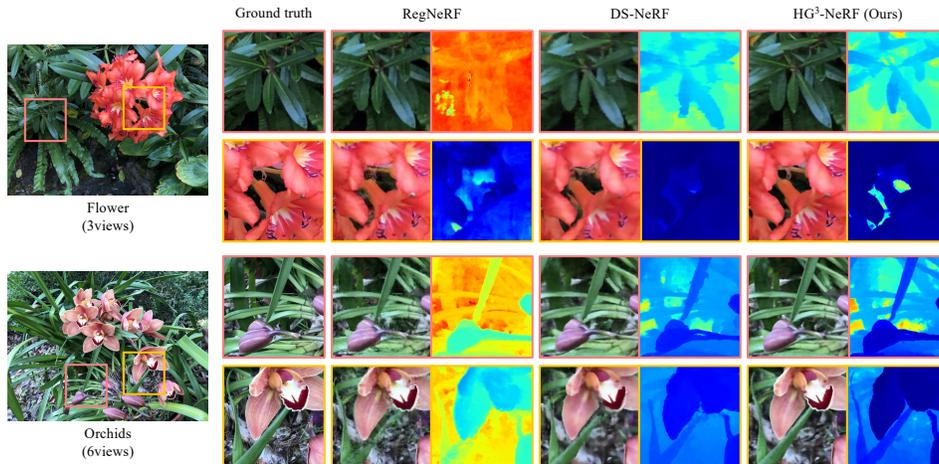


Figure 6: **Qualitative Results on LLFF Dataset.** Our method synthesizes high-fidelity images and estimates accurate depth maps in the forward-facing scene without NDC space.

Comparison on LLFF Dataset. We then evaluate HG^3 -NeRF on complex real-world LLFF dataset. As the qualitative results presented in Fig. 6, different from some baselines (e.g. RegNeRF) that require the NDC space to represent the forward-facing scene, our method can realize high-fidelity view synthesis under the real-world space. Moreover, as shown in Table. 2, the effect of semantic consistency for complex scenes is severely limited (PSNR=14.94, DietNeRF for 3 views). The HSG method can utilize the coarse-to-fine semantic consistency to gradually aggregate the semantic content into the representations. As a result, our HG^3 -NeRF is able to employ the hierarchical geometric, semantic, and photometric guidance to represent the complex scene with high consistency in appearance, geometry, and semantics even from sparse view inputs (PSNR=20.98, 3 views).

5.3 Model Analysis

Ablation Study. We conduct an ablation study to investigate the effect of HGG and HSG under 3,6,9 view inputs on LLFF dataset. As shown in Table. 3, since we represent the scene in the real-world space, the HGG method can maintain the most important geometric consistency for different viewpoints and thus plays a key role in our method, especially for only 3 view inputs. Moreover, the HSG method can supervise the semantic consistency for the coarse-to-fine representations to further improve the synthesis quality based on the HGG method. These two methods can be regarded as the convenient tools for improving the performance of novel view synthesis.

Table 1: **Quantitative Results on DTU Dataset.** The best is in bold. Conducted with sparse view inputs, the proposed HG³-NeRF can achieve better view synthesis quality than most SOTA baselines.

Configures		PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
		3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
PixelNeRF	Pre-training	16.82	19.11	20.40	0.695	0.745	0.768	0.270	0.232	0.220
SRF		15.32	17.54	18.35	0.671	0.730	0.752	0.304	0.250	0.232
MVSNeRF		18.63	20.70	22.40	0.769	0.823	0.853	0.197	0.156	0.135
PixelNeRF <i>ft</i>	Pre-training and Fine-tune	18.95	20.56	21.83	0.710	0.753	0.781	0.269	0.223	0.203
SRF <i>ft</i>		15.68	18.87	20.75	0.698	0.757	0.785	0.281	0.225	0.205
MVSNeRF <i>ft</i>		18.54	20.49	22.22	0.769	0.822	0.853	0.197	0.155	0.135
DietNeRF	Per-scene Optimization	11.85	20.63	23.83	0.633	0.778	0.823	0.314	0.201	0.173
RegNeRF		18.89	22.20	24.93	0.745	0.841	0.884	0.190	0.117	0.089
[†] DS-NeRF <i>w sdp</i>		16.29	19.07	21.98	0.559	0.807	0.828	0.451	0.211	0.195
HG³-NeRF <i>w sdp</i>		19.37	23.35	25.87	0.759	0.855	0.891	0.177	0.094	0.061

Table 2: **Quantitative Results on LLFF Dataset.** The best is in bold. For the complex real-world scenes, HG³-NeRF can outperform other SOTA baselines without using NDC space.

Configures		PSNR \uparrow			SSIM \uparrow			LPIPS \downarrow		
		3-view	6-view	9-view	3-view	6-view	9-view	3-view	6-view	9-view
PixelNeRF	Pre-training	7.93	8.74	8.61	0.272	0.280	0.274	0.682	0.676	0.665
SRF		12.34	13.10	13.00	0.250	0.293	0.297	0.591	0.594	0.605
MVSNeRF		17.25	19.79	20.47	0.557	0.656	0.689	0.356	0.269	0.242
PixelNeRF <i>ft</i>	Pre-training and Fine-tune	16.17	17.03	18.92	0.438	0.473	0.535	0.512	0.477	0.430
SRF <i>ft</i>		17.07	16.75	17.39	0.436	0.438	0.465	0.529	0.521	0.503
MVSNeRF <i>ft</i>		17.88	19.99	20.47	0.584	0.660	0.695	0.327	0.264	0.244
DietNeRF	Per-scene Optimization	14.94	21.75	24.28	0.370	0.717	0.801	0.496	0.248	0.183
RegNeRF		19.08	23.10	24.86	0.587	0.760	0.820	0.336	0.206	0.161
[†] DS-NeRF <i>w sdp</i>		19.68	22.45	23.78	0.615	0.674	0.722	0.403	0.356	0.324
HG³-NeRF <i>w sdp</i>		20.98	24.60	25.51	0.682	0.823	0.844	0.198	0.107	0.075

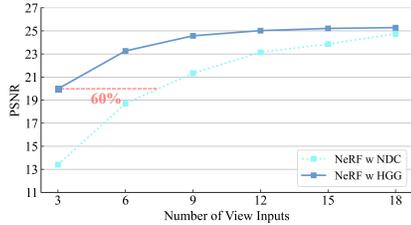


Figure 7: **Data and Space Efficiency.** In sparse settings, NeRF with HGG (in real-world space) require up to 60% fewer images than NeRF using NDC space to achieve comparable performance on LLFF dataset.

Table 3: **Ablation Study.** We investigate the effect of HGG and HSG on LLFF dataset. For sparse view inputs, HGG serves the crucial role of representations in the real-world space and avoiding overfitting problems. HSG can further improve the performance with the coarse-to-fine semantic guidance.

Configures		PSNR		
HGG	HSG	3-view	6-view	9-view
✓		19.97	23.46	24.57
	✓	15.82	22.71	23.29
✓	✓	20.98	24.06	25.51

NDC space vs Real-World Space. We evaluate the data and space efficiency by comparing the performance of NeRF with HGG (real-world space) and NeRF with NDC space. As shown in Fig. 7, for sparse view inputs, NeRF with HGG only requires up to 60% fewer view inputs to achieve the view synthesis quality that is comparable to NeRF with NDC space on the test image set. Therefore, in addition to improving the performance under sparse view inputs, the proposed HGG method can also replace the NDC space for representing the forward-facing scenes.

6 Conclusion

In this paper, the proposed HG³-NeRF exploit the hierarchical geometric, semantic, and photometric guidance to represent the scene from sparse views. We propose HGG to leverage the sparse depth prior to indicate a local-to-global region for volume sampling without introducing geometric bias. The HSG method maintains the semantic consistency by computing the coarse-to-fine cosine similarly. The experimental results demonstrate that the proposed HG³-NeRF can synthesize high-fidelity images and estimate accurate depth maps with distinct edges.

Limitations. In our method, the camera poses are still required to be estimated from SfM. Note that sparse view inputs can also affect the accuracy of pose estimation, and some works prove that noisy camera poses can degrade the view synthesis quality [23, 7]. Therefore, there is still scope to improve the performance of NeRF from sparse view inputs. In the future work, we will follow the methods proposed in this paper and extend them on solving the joint optimization problem about camera poses and NeRF from sparse view inputs.

References

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building rome in a day. *Communications of the ACM*, 54(10):105–112, 2011.
- [2] Timur Bagautdinov, Chenglei Wu, Jason Saragih, Pascal Fua, and Yaser Sheikh. Modeling facial geometry using compositional vaes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3877–3886, 2018.
- [3] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [4] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14124–14133, 2021.
- [5] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019.
- [6] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. 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*, pages 7911–7920, 2021.
- [7] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIII*, pages 264–280. Springer, 2022.
- [8] Christopher B Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VIII 14*, pages 628–644. Springer, 2016.
- [9] Boyang Deng, Kyle Genova, Soroosh Yazdani, Sofien Bouaziz, Geoffrey Hinton, and Andrea Tagliasacchi. Cvxnet: Learnable convex decomposition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 31–44, 2020.
- [10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12882–12891, 2022.
- [11] SM Ali Eslami, Danilo Jimenez Rezende, Frederic Besse, Fabio Viola, Ari S Morcos, Marta Garnelo, Avraham Ruderman, Andrei A Rusu, Ivo Danihelka, Karol Gregor, et al. Neural scene representation and rendering. *Science*, 360(6394):1204–1210, 2018.
- [12] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017.
- [13] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022.
- [14] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–224, 2018.
- [15] Yuxuan Han, Ruicheng Wang, and Jiaolong Yang. Single-view view synthesis in the wild with learned adaptive multiplane images. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–8, 2022.

- [16] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [17] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5885–5894, 2021.
- [18] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12949–12958, 2021.
- [19] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014.
- [20] Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. Geonerf: Generalizing nerf with geometry priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18365–18375, 2022.
- [21] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022.
- [22] Yeonkyeong Lee, Taeho Choi, Hyunsung Go, Hyunjoon Lee, Sunghyun Cho, and Junho Kim. Exp-gan: 3d-aware facial image generation with expression control. In *Proceedings of the Asian Conference on Computer Vision*, pages 3812–3827, 2022.
- [23] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021.
- [24] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3d supervision. *Advances in Neural Information Processing Systems*, 32, 2019.
- [25] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7824–7833, 2022.
- [26] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019.
- [27] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4743–4752, 2019.
- [28] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019.
- [29] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020.
- [30] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [31] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. Orb-slam: a versatile and accurate monocular slam system. *IEEE transactions on robotics*, 31(5):1147–1163, 2015.

- [32] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras. *IEEE transactions on robotics*, 33(5):1255–1262, 2017.
- [33] Charlie Nash, Yaroslav Ganin, SM Ali Eslami, and Peter Battaglia. Polygen: An autoregressive generative model of 3d meshes. In *International conference on machine learning*, pages 7220–7229. PMLR, 2020.
- [34] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- [35] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022.
- [36] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3504–3515, 2020.
- [37] Roy Or-El, Xuan Luo, Mengyi Shan, Eli Shechtman, Jeong Joon Park, and Ira Kemelmacher-Shlizerman. Stylesdf: High-resolution 3d-consistent image and geometry generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13503–13513, 2022.
- [38] Juewen Peng, Zhiguo Cao, Xianrui Luo, Hao Lu, Ke Xian, and Jianming Zhang. Bokehme: When neural rendering meets classical rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16283–16292, 2022.
- [39] Songyou Peng, Chiyu Jiang, Yiyi Liao, Michael Niemeyer, Marc Pollefeys, and Andreas Geiger. Shape as points: A differentiable poisson solver. *Advances in Neural Information Processing Systems*, 34:13032–13044, 2021.
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [41] Konstantinos Rematas, Ricardo Martin-Brualla, and Vittorio Ferrari. Sharf: Shape-conditioned radiance fields from a single view. In *International Conference on Machine Learning*, pages 8948–8958. PMLR, 2021.
- [42] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- [43] Seunghyeon Seo, Donghoon Han, Yeonjin Chang, and Nojun Kwak. Mixnerf: Modeling a ray with mixture density for novel view synthesis from sparse inputs. *arXiv preprint arXiv:2302.08788*, 2023.
- [44] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020.
- [45] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2437–2446, 2019.
- [46] Pavel Solovev, Taras Khakhulin, and Denis Korzhenkov. Self-improving multiplane-to-layer images for novel view synthesis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4309–4318, 2023.

- [47] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8269–8279, 2022.
- [48] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022.
- [49] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022.
- [50] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 34:27171–27183, 2021.
- [51] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [52] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5610–5619, 2021.
- [53] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.
- [54] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, pages 106–122. Springer, 2022.
- [55] Wenpeng Xing and Jie Chen. Temporal-mpi: Enabling multi-plane images for dynamic scene modelling via temporal basis learning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 323–338. Springer, 2022.
- [56] DeJia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, pages 736–753. Springer, 2022.
- [57] Yinghao Xu, Sida Peng, Ceyuan Yang, Yujun Shen, and Bolei Zhou. 3d-aware image synthesis via learning structural and textural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18430–18439, 2022.
- [58] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4578–4587, 2021.
- [59] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. In *Conference on Neural Information Processing Systems (NeurIPS 2022)*, 2022.
- [60] Kai Zhang, Gernot Riegler, Noah Snaveley, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.
- [61] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.