

# Anisotropic Neural Representation Learning for High-Quality Neural Rendering

Yifan Wang\*, Jun Xu\*, Yi Gong, and Yuan Zeng

Southern University of Science and Technology

**Abstract.** Neural radiance fields (NeRFs) have achieved impressive view synthesis results by learning an implicit volumetric representation from multi-view images. To project the implicit representation into an image, NeRF employs volume rendering that approximates the continuous integrals of rays as an accumulation of the colors and densities of the sampled points. Although this approximation enables efficient rendering, it ignores the direction information in point intervals, resulting in ambiguous features and limited reconstruction quality. In this paper, we propose an anisotropic neural representation learning method that utilizes learnable view-dependent features to improve scene representation and reconstruction. We model the volumetric function as spherical harmonic (SH)-guided anisotropic features, parameterized by multilayer perceptrons, facilitating ambiguity elimination while preserving the rendering efficiency. To achieve robust scene reconstruction without anisotropy over-fitting, we regularize the energy of the anisotropic features during training. Our method is flexible and can be plugged into NeRF-based frameworks. Extensive experiments show that the proposed representation can boost the rendering quality of various NeRFs and achieve state-of-the-art rendering performance on both synthetic and real-world scenes.

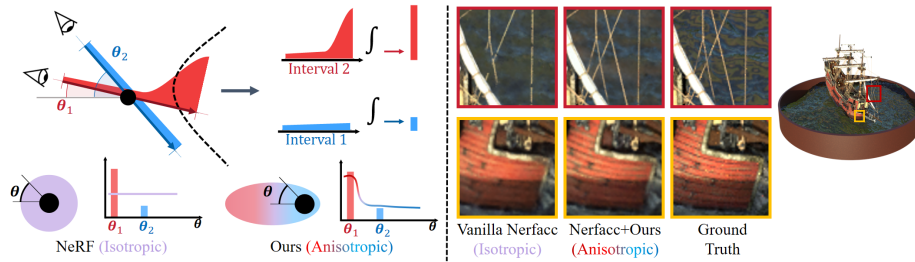
**Keywords:** Neural radiance fields · Anisotropic implicit representation · Neural rendering

## 1 Introduction

Neural radiance field (NeRF) [34] models the 3D scene geometry and view-dependent appearance by two cascaded multi-layer perceptrons (MLPs) and uses volume rendering to reconstruct photorealistic novel views. The advent of NeRF has sparked a flurry of work on neural rendering and has opened the way to many applications [2, 12, 15, 52, 58–60]. One of the key underpinnings of NeRF is differentiable volume rendering which facilitates learning the MLP-based radiance field solely from a 2D photometric loss on synthesized images in an end-to-end manner. However, NeRF still suffers from blurring and aliasing when reconstructing complex scenes. Specifically, since NeRF is done by sampling a set of points along a ray and approximating the piece-wise continuous integration as an accumulation of the estimated volumetric features of sampled

---

\* Authors contributed equally to this work.



**Fig. 1:** Left: Vanilla NeRF uses point sampling and view-independent functions to estimate  $\sigma$  and  $\mathbf{e}$ , resulting in directional ambiguity when representing intervals. To eliminate the ambiguity, we introduce anisotropic functions to model the integration along various directions. Right: Our anisotropic neural representation enables to assist Nerfacc [28] in capturing more geometric details and realistic textures.

intervals, MLP could only be queried at a fixed discrete set of positions and the same ambiguous point-sampled feature is used to represent the opacity at all points with the interval between two samples, leading to ambiguity in neural rendering. The features of samples can be very different across different rays cast from different viewing directions. Therefore, simultaneously supervising these rays to produce the isotropic volumetric features can result in artifacts like fuzzy surfaces and blurry texture, as shown in Fig. 1.

To overcome the limitations of the neural scene representation and enhance rendering quality of NeRF, several works [2–4, 23, 24] introduce shape-based sampling into the scene representation. By embedding novel spatial parameterization schemes, such as Gaussian ellipsoid, frustum, and sphere, into the encoding, these models can reduce representation ambiguity and improve rendering quality with less blurring and aliasing artifacts. However, the representation ambiguities still exist in their radiance fields, since the directional intervals used for rendering are parameterized by the non-directional features. A straightforward solution is to take the viewing direction as input of the first MLP, enabling to represent the scene geometry with view-dependent features. While this method can reduce directional ambiguity in representation, the radiance field is now a high-degree view-dependent function without regularization, and it is prone to over-fitting to the appearance of training images but fails to capture the correct geometry [61].

In this paper, we propose a novel radiance field representation to diverse a different model for the density and appearance in neural rendering, leading to better reconstruction of the scene’s geometry and appearance while producing photorealistic novel views. Instead of introducing a different spatial sampling strategy and parameterizing sampled-shapes, we model the density and latent features at a location as view-dependent functions using spherical harmonic (SH) basis. The spherical harmonics can be evaluated at arbitrary query viewing directions to capture anisotropy. Although this can be done by converting an existing NeRF into such anisotropic representations via projection onto the SH basis functions, we simply modify the first MLP in NeRF to

predict the geometry explicitly in terms of spherical harmonics, resulting in a compact and generalizable anisotropic neural representation. Specifically, we train the first MLP that produces coefficients for the SH functions instead of the density and latent features. The predicted values can later be directly used for appearance estimation and pixel rendering.

Additionally, training the anisotropic neural representation using only a color reconstruction loss will cause strong anisotropy and suffer from shape-radiance ambiguity in rendering [61]. Although existing regularization methods, such as Patch-based Consistency [14], warp-based loss [63], depth or multi-view stereo prior [6, 50, 53, 56] can be a solution, they rely on the geometric prior and structural output and significantly increase the computational cost. We aim to introduce an effective anisotropy regularization according to our SH function-based representation. To this end, we decompose the anisotropic elements from the predicted density and latent features and employ point-wise operation to penalize the anisotropy in the geometry representation. This improves the memory and computational efficiency of our method and allows us to render high-quality novel views. Moreover, our anisotropic neural representation can be used as a sub-model to replace the first MLP in various NeRFs, resulting in anisotropic geometric features and better rendering quality.

We extensively evaluate the effectiveness and generalizability of our anisotropic neural representation on benchmark datasets including both synthetic and real-world scenes. Both quantitative and qualitative comparison results demonstrate that plugging our anisotropic neural representation can further improve the rendering quality of various NeRFs. Our contributions are summarized as follows:

- A novel anisotropic neural representation is proposed to model the scene geometry by view-dependent density and latent features. It effectively reduces directional ambiguity in neural rendering and results in better geometric and appearance reconstruction and rendering quality.
- A modified NeRF network is trained to predict view-dependent geometry in terms of spherical basis functions, which is flexible and generalizable to various existing NeRFs.
- A point-wise anisotropy regularization loss enables highly efficient view-dependent penalties during model training to avoid over-fitting.

## 2 Related Work

*Scene representations for Novel view synthesis.* Synthesizing views from a novel viewpoint is a long-standing problem in computer vision and computer graphics. Traditional methods typically synthesize novel views from a set of images [8–10, 17, 42] or light fields [27, 46, 55]. Although these methods work well on dense input images, they are limited by the sparse inputs and the quality of 3D reconstruction. To synthesize novel views from a sparse set of images, some methods leverage the geometry structure of the scene [20, 37]. Recent advances in deep learning have facilitated the use of neural networks for estimating scene geometry, such as voxel grids [22, 29, 39, 43], point

clouds [1, 26, 45, 51], and multi-plane images [16, 64]. Although these discrete representation methods can improve the rendering quality of novel views, the estimation of the scene geometry is not accurate enough for high-resolution scenes.

By modeling the scene geometry and appearance as a continuous volume, neural radiance fields (NeRFs) [34] have achieved state-of-the-art novel synthesis effects. Specifically, NeRF maps the input coordinate to density and radiance scene values and uses volume rendering [31] to synthesize the images. In addition, various schemes have been introduced to improve the robustness of NeRF to few-shot inputs [25, 36, 49, 60], anti-aliasing [2–4, 23], handle dynamics [13, 30, 33, 57] and speed up rendering [11, 19, 21, 35, 48], etc. Our method is more closely related to anti-aliasing NeRFs, which adopt interval-dependent features to assist the MLP in capturing more accurate geometry and reducing blurring artifacts in novel view synthesis. In contrast, we introduce a plug-and-play anisotropic neural representation that enables it to be plugged into various NeRFs to alleviate ambiguity and improve rendering quality.

*Anti-ambiguity in neural rendering.* Volume rendering is an important technique with a long history of research in the graphics community. Traditional graphics rendering methods include studying ray sampling efficiency and data structures for coarse-to-fine hierarchical sampling. Recent NeRF and its succeeding works have shown impressive results by replacing or augmenting the traditional graphics rendering with neural networks. The volume rendering used in these works is approximated to a discrete accumulation under the assumption of a piece-wise constant opacity and color, enabling to learn the NeRF-based implicit representations. However, the piece-wise constant assumption results in rendering results that are sensitive to the sampled points as well as the cumulative density function of the distribution of the sampled interval, introducing ambiguous features and aliasing artifacts in NeRF renderings.

To address these challenges, recent works have explored super-sampling or pre-filtering techniques. Super-sampling is done by casting multiple rays per pixel to approach the Nyquist frequency. Although this strategy works well for eliminating ambiguity, it is computationally expensive. Pre-filtering techniques are more computationally efficient, since the filtered versions of scene content can be pre-computed ahead of time. Recently, pre-filtering has been introduced into neural representation and rendering to reduce ambiguity and aliasing artifacts [7]. Mip-NeRF [2] samples the cone instead of rays to consider the shape and size of the volume viewed by each ray and optimizes a pre-filtered representation of the scene during training. With sampled shape-dependent inputs, the MLP of Mip-NeRF can capture various volumetric features to mitigate ambiguity, resulting in a high-quality multi-scale representation and anti-aliasing. Mip-NeRF 360 [3] extends Mip-NeRF with a novel distortion-based regularizer to tackle unbounded scenes. Zip-NeRF [4] adopts multi-sampling to approximate a cone with hash encoding. Tri-Mip [23] leverages multi-level 2D mipmaps to model the pre-filtered 3D feature space and projects parameterized spheres on three mipmaps to achieve anti-aliasing encoding. Although the shape and size of volumes at different scales can be fitted by introducing various sampling techniques in NeRF, shape-based features still exist in ambiguity since the non-directional shape intervals are used in accumulation for pixel rendering. Our work draws inspiration from the early work on anisotropic



volume rendering [41] and is the first to model scene geometry as an anisotropic neural representation based on spherical harmonics for volume rendering.

### 3 Preliminaries

*Neural Radiance Fields.* Given a 3D position  $\mathbf{x}$  and a 2D viewing direction  $\mathbf{d}$ , NeRF [34] first uses a multilayer perceptron (MLP) parameterized by weights  $\theta$  to predict the density  $\sigma$  and an intermediate vector  $\mathbf{e}$  from the input position  $\mathbf{x}$ :  $(\sigma, \mathbf{e}) = \mathcal{F}_\theta(\mathbf{x})$ . Then, a second MLP parameterized by weights  $\phi$  is employed to estimate the color  $\mathbf{c}$  from the direction  $\mathbf{d}$  and the vector  $\mathbf{e}$ :  $\mathbf{c} = \mathcal{F}_\phi(\mathbf{d}, \mathbf{e})$ .

*Volume Rendering.* The color of a pixel in NeRFs can be rendered by casting a ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  from the camera origin  $\mathbf{o}$  through the pixel along the direction  $\mathbf{d}$ , where  $t$  is the distance to the origin. The pixel’s color value can be computed by integrating colors and densities along a ray based on the volume rendering [31]:

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

where  $T(t) = \exp\left(-\int_0^t \sigma(\mathbf{r}(s)) ds\right)$  represents occlusions by integrating the differential density between 0 to  $t$ .

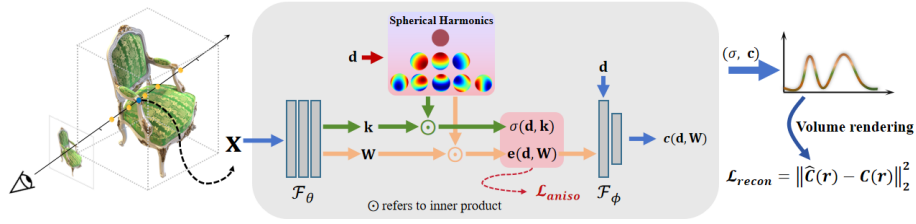
Since the volume density  $\sigma$  and radiance  $\mathbf{c}$  are the outputs of MLPs, NeRF rendering methods approximate this continuous integral using a sampling-based Riemann sum instead [40]. Within the near and far bounds,  $t_n$  and  $t_f$  of the cast ray, a subset of points on the ray is sampled in a near-to-far order. Let  $\mathbf{t}_N = t_1, \dots, t_N$  be  $N$  samples on the ray that define the intervals, e.g.,  $I_i = [t_i, t_{i+1}]$  is the  $i$ th interval, and  $I_0 = [0, t_1]$ ,  $I_N = [t_N, \infty]$ . The volume density for particles along the interval  $I_i$  is predicted under the assumption that opacity is constant along each interval, which indicates  $\sigma(\mathbf{r}(t)) = \sigma(\mathbf{r}(t_i)), \forall t \in [t_i, t_{i+1}]$  for particles of constant radius and material. We denote  $\sigma(\mathbf{r}(t_i)) = \sigma_i$  for notation convenience. Under this assumption, the rendered color can be written as an approximation of the  $N$  sampled points:

$$\hat{C}(\mathbf{r}) = \sum_{i=0}^N \mathbf{c}_i T_i (1 - \exp(-\sigma_i(t_{i+1} - t_i))), \quad (2)$$

where  $T_i = \exp\left(-\sum_{j=0}^{i-1} \sigma_j(t_{j+1} - t_j)\right)$  represents the transmittance accumulated along the ray until the  $i$ th sample. For more detailed derivation, please refer to [32]. The NeRF model is optimized by minimizing the  $L_2$  reconstruction loss between the ground truth and synthesized images, which can be expressed as follows:

$$\mathcal{L}_{recon} = \frac{1}{|\mathcal{R}|} \sum_{\mathbf{r} \in \mathcal{R}} \left\| \hat{C}(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2, \quad (3)$$

where  $\mathcal{R}$  is a set of rays sampled during training, and  $C$  is the ground truth pixel color value.



**Fig. 2:** Overview of our anisotropic neural representation. In contrast to output isotropic  $\sigma$  and  $\mathbf{e}$  directly as in vanilla NeRF, we composite anisotropic features by learning the SH coefficients predicted from MLP  $\mathcal{F}_\theta$ . During training, the model is optimized end-to-end by minimizing a joint loss of  $\mathcal{L}_{recon}$  and  $\mathcal{L}_{aniso}$ .

*Limitations.* While images can be efficiently rendered using NeRF renderings, it is non-trivial to represent directional intervals with the point-sampled features. In addition, leveraging the same predicted features to represent the distribution along the interval between two samples can hardly capture the correct geometry of different rays cast from different directions. For example, under piece-wise constant color, the contribution  $\hat{C}_i(\mathbf{r})$  of the  $i$ th interval to  $\hat{C}(\mathbf{r})$ , i.e., the volume rendering integral of the interval  $I_i$ , can be written as [31, 34]:

$$\begin{aligned}
 \hat{C}_i(\mathbf{r}) &= \int_{t_i}^{t_{i+1}} \mathbf{c}(t) \sigma(\mathbf{r}(t)) e^{-\int_0^t \sigma(\mathbf{r}(s)) ds} dt \\
 &= \mathbf{c}_i e^{-\int_0^{t_i} \sigma(\mathbf{r}(s)) ds} (1 - e^{-\int_{t_i}^{t_{i+1}} \sigma(\mathbf{r}(s)) ds}) \\
 &= \mathbf{c}_i T(t_i) (1 - e^{-\int_{t_i}^{t_{i+1}} \sigma(\mathbf{r}(s)) ds}).
 \end{aligned} \tag{4}$$

To achieve efficient rendering, NeRF uses point-sampled density  $\sigma_i$  to approximate the distribution along the directional interval  $(\int_{t_i}^{t_{i+1}} \sigma(\mathbf{o} + s\mathbf{d}) ds / (t_{i+1} - t_i))$  and simplifies the rendering computation in equation (4) as  $\mathbf{c}_i T_i (1 - \exp(-\sigma_i (t_{i+1} - t_i)))$  in equation (2). This indicates that the integral of  $\sigma(\mathbf{o} + s\mathbf{d})$  is related to the viewing direction, which has also been studied in [41] for volume rendering. Approximating the continuous integral of the interval as an accumulation of the isotropic features of sampled points has no guarantee for correct approximation. Fig. 1 shows that the vanilla NeRF using an isotropic representation has limited rendering quality and suffers from blurring artifacts. This insight motivates our anisotropic neural representation.

## 4 Method

Given a set of scenes with a collection of images and their camera parameters, we aim to learn an anisotropic neural representation for high-quality neural rendering. In this section, we first represent the scene geometry and appearance using spherical harmonic-based implicit representations to capture the anisotropy of surfaces (Sec. 4.1). Then we introduce an anisotropy regularization to encourage a sparse anisotropic neural representation and summarize our overall training procedure (Sec. 4.2).

#### 4.1 Anisotropic Neural Representation

Fig. 2 illustrates an overview of our method. While vanilla NeRF takes a 3D position  $\mathbf{x}$  as an input of a density MLP  $\mathcal{F}_\theta$  to estimate the scene geometry, including volume density  $\sigma$  and latent feature  $\mathbf{e}$ , we leverage both the 3D position and 2D viewing direction to introduce anisotropic features to represent the scene geometry. Although naively embedding the direction  $\mathbf{d}$  into the density MLP  $\mathcal{F}_\theta$  can realize anisotropic implicit geometry representation, this design will directly increase the dimension of input data and cause the model more easily fit the training data and estimate incorrect scene geometry [61]. To capture anisotropy in scene representation, our method utilizes spherical harmonics (SH), which have been used to model Lambertian surfaces [5, 38], or even glossy surfaces [44]. We query the SH functions  $Y_l^m: \mathbb{S}^2 \mapsto \mathbb{R}$  at a viewing direction  $\mathbf{d}$  and then fit the anisotropic neural representations by finding the corresponding coefficients. We use low-degree SH functions to compute ideal values of view-independent density and latent components, and high-degree SH functions for view-dependent components.

For any sampled point  $\mathbf{x}$  in the space, we adapt  $\mathcal{F}_\theta$  to estimate the spherical harmonic coefficients  $\mathbf{k}$  and  $\mathbf{W}$ , rather than the volume density and latent feature:

$$(\mathbf{k}, \mathbf{W}) = \mathcal{F}_\theta(\mathbf{x}), \quad (5)$$

where the spherical harmonic coefficients  $\mathbf{k} = (k_l^m)_{l:0 \leq l \leq L}^{m:-l \leq m \leq l}$  is related to calculate the view-dependent volume density and  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K]^\top$  consists of  $K$  sets of SH coefficients that used to determine the  $K$ -dimensional latent feature vector. For  $n \in \{1, \dots, K\}$ , we have  $\mathbf{w}_n = (w_{nl}^m)_{l:0 \leq l \leq L}^{m:-l \leq m \leq l}$ . For  $\mathbf{k}$  or  $\mathbf{w}_n$ , there are  $(L+1)^2$  spherical harmonics of degree at the most  $L$ . The view-dependent density  $\sigma$  and latent feature  $\mathbf{e}$  at position  $\mathbf{x}$  are then determined by querying the SH functions  $Y_l^m$  at the desired viewing direction  $\mathbf{d}$ :

$$\sigma(\mathbf{d}, \mathbf{k}) = \sum_{l=0}^L \sum_{m=-l}^l k_l^m Y_l^m(\mathbf{d}), \quad (6)$$

and

$$e_n(\mathbf{d}, \mathbf{w}_n) = \sum_{l=0}^L \sum_{m=-l}^l w_{nl}^m Y_l^m(\mathbf{d}). \quad (7)$$

The equations (6) and (7) can be seen as the factorization of the density and latent feature with the isotropic and anisotropic SH basis functions, respectively. This eliminates the input of view direction to the density MLP and enables efficient generation of view-dependent geometry features. Then, a color MLP takes the inputs of the estimated latent feature  $\mathbf{e}(\mathbf{d}, \mathbf{W})$  and the direction  $\mathbf{d}$  to predict the color value  $\mathbf{c}$ :

$$\mathbf{c}(\mathbf{d}, \mathbf{W}) = \mathcal{F}_\phi(\mathbf{d}, \mathbf{e}(\mathbf{d}, \mathbf{W})). \quad (8)$$

Given the estimated volume density  $\sigma(\mathbf{d}, \mathbf{k})$  and color value  $\mathbf{c}(\mathbf{d}, \mathbf{W})$ , a pixel color  $\hat{C}(\mathbf{r})$  in the radiance field along the ray  $\mathbf{r}$  can be predicted using the volume rendering in equation (2). Note that since the proposed scene geometry representation only

adapts the original density MLP to learn SH coefficients and efficiently convert the input position into anisotropic features, making it easy to be plugged into various existing NeRF-based scene representation models and assists the models in capturing more precise geometries for high-quality novel view synthesis.

## 4.2 Training

NeRF adopts pixel-wise RGB reconstruction loss  $\mathcal{L}_{recon}$  to optimize view-independent geometry density and view-dependent color for scene reconstruction. However, it is difficult to optimize a correct geometry purely from the input RGB images, especially when view-dependent components are used in the geometry representation, since a high degree of anisotropy will be learned to fit the training images, leading to shape-radiance ambiguity and the degradation of rendering quality. To facilitate robust mapping under our anisotropic neural representation, we propose a point-wise anisotropy constraint, which penalizes the anisotropy of the model to mitigate shape-radiance ambiguity efficiently.

*Anisotropy Regularization.* Without any regularization, the model is free to fit a set of training images by exploiting view-dependent anisotropic neural representation rather than recovering the correct geometry. The representation with strong anisotropy would generate shape-radiance ambiguity, resulting in blurring artifacts and incorrect geometries in rendering novel test views. We therefore introduce a new regularization method that suppresses the anisotropy in geometry representations.

To apply regularization techniques for penalizing anisotropy, we first define the anisotropic features according to the SH basis functions. Since the 0-degree SH function  $Y_0^0(\mathbf{d})$  is view-independent, we remove the view-independent component from the estimated density  $\sigma(\mathbf{d}, \mathbf{k})$  and latent feature  $\mathbf{e}(\mathbf{d}, \mathbf{W})$  to compute the view-dependent component:

$$\sigma^{aniso}(\mathbf{d}, \mathbf{k}) = \sum_{l=1}^L \sum_{m=-l}^l k_l^m Y_l^m(\mathbf{d}), \quad (9)$$

and

$$\mathbf{e}_n^{aniso}(\mathbf{d}, \mathbf{w}_n) = \sum_{l=1}^L \sum_{m=-l}^l w_{nl}^m Y_l^m(\mathbf{d}). \quad (10)$$

We formulate our anisotropy regularization loss as:

$$\mathcal{L}_{aniso} = \frac{1}{N} \sum_{i=1}^N \left( \|\sigma^{aniso}(\mathbf{d}_i, \mathbf{k}_i)\|_2^2 + \|\mathbf{e}^{aniso}(\mathbf{d}_i, \mathbf{W}_i)\|_2^2 \right), \quad (11)$$

where  $\mathbf{d}_i$ ,  $\mathbf{k}_i$  and  $\mathbf{W}_i$  represent the direction, density-related SH coefficients and latent feature-related SH coefficients at the  $i$ th point sampled on ray  $\mathbf{r}$ , respectively.

*Full Objective Loss.* To learn the high-fidelity scene reconstruction, we optimize the following total loss in each iteration:

$$\mathcal{L} = \mathcal{L}_{recon} + \lambda \mathcal{L}_{aniso}, \quad (12)$$

where  $\lambda$  is a hyperparameter to scale the anisotropy regularization loss. In addition, to plug our anisotropy neural representation in different existing NeRFs for novel view synthesis, we remain the original losses and add the anisotropy regularization loss  $\mathcal{L}_{aniso}$  in the full objective loss for model training.

## 5 Experiments

In this section, we evaluate the effectiveness and generalizability of our method on novel view synthesis. We plug our anisotropic neural representation into existing state-of-the-art NeRFs and present quantitative and qualitative comparisons between the baseline NeRFs and our models on both synthetic and real-world benchmark datasets in Sec. 5.1. A comprehensive ablation study that supports our design choices is also provided in Sec. 5.2. More details and per-scene results are provided in our supplementary materials.

*Datasets and metrics* We report our results on two datasets: Blender [34] and Mip-360 [3]. Blender consists of eight synthetic scenes (*lego, chair, hotdog, ficus, drums, materials, mic, and ship*), where each has 400 synthesized images. Mip-360 is an unbounded real-world dataset including three outdoor scenes (*garden, bicycle, stump*) and four indoor scenes (*room, kitchen, bonsai, counter*), and each scene contains a complex central object or region along with intricate background details. We follow the default split in [3] to produce training and testing views. In addition, we follow previous NeRF methods and report our quantitative results in terms of peak signal-to-noise ratio (PSNR), structural similarity index (SSIM) [54], learning perceptual image patch similarity (LPIPS) [62] and average error (Avg.) [2] which summarizes three above metrics.

*Baselines.* We adopt the following four recently proposed neural rendering methods as baselines: two point-sampling methods Nerfacc [28] and K-Planes [18], and two shape-sampling methods Tri-Mip [23] and Zip-NeRF [4]. We use the official implementation of Nerfacc, K-Planes, and Tri-Mip and retrain the three models on the Blender and the PyTorch implementation of Zip-NeRF [47] and retrain Zip-NeRF and Nerfacc on Mip-360.

*Implementation Details.* For our anisotropic neural representation, we set the maximal degree of spherical harmonic basis  $L = 3$  and the hyperparameter  $\lambda = 1e^{-4}$ . To keep the dimension of the predicted density  $\sigma$  and latent feature  $\mathbf{e}$  the same as the baselines, we modify the output dimension of the first MLP  $\mathcal{F}_\theta$  of the baselines to  $(L + 1)^2$  times the original output dimension. We keep other settings the same as our baselines. In addition, we use Nerfacc as our implementation backbone to verify our design choices in ablation studies.

**Table 1:** Quantitative findings derived from the Blender dataset reveal significant advancements in average metrics with the incorporation of our method into the baselines.

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg. $\downarrow$
Nerfacc	33.06	0.961	0.053	0.017
Nerfacc+Ours	<b>34.08</b>	<b>0.966</b>	<b>0.046</b>	<b>0.015</b>
K-Planes	32.34	0.962	0.052	0.018
K-Planes+Ours	<b>33.00</b>	<b>0.964</b>	<b>0.052</b>	<b>0.017</b>
Tri-Mip	33.78	0.963	0.051	0.016
Tri-Mip+Ours	<b>34.69</b>	<b>0.965</b>	<b>0.049</b>	<b>0.015</b>

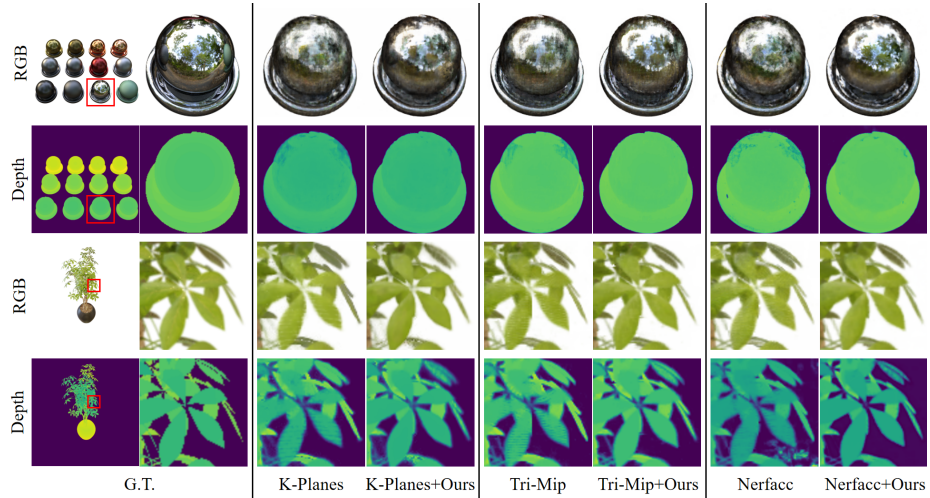
**Table 2:** Quantitative comparison results on Mip-360 dataset. Using our representation, the baseline models can achieve better rendering quality.

	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg. $\downarrow$
Nerfacc	27.47	0.771	0.294	0.063
Nerfacc+Ours	<b>28.32</b>	<b>0.782</b>	<b>0.281</b>	<b>0.058</b>
Zip-NeRF	28.67	0.838	0.270	0.053
Zip-NeRF+Ours	<b>29.14</b>	<b>0.846</b>	<b>0.261</b>	<b>0.050</b>

## 5.1 Comparisons

*Quantitative Comparison.* Quantitative results on Blender are reported in Table 1. We observe that our representation can significantly improve the numerical performance of the three baselines Nerfacc, K-Planes, and Tri-Mip. Quantitative results on Mip-360 are summarized in Table 2. We see that our representation is also effective for real-world scenes. The rendering performance of Zip-NeRF+Ours indicates that our method can be used in combination with multi-sampling techniques to achieve further better rendering details.

*Qualitative Comparison.* Fig. 3 shows the visual comparison and depth results of two synthetic scenes, including *materials* and *ficus* scenes in Blender. It shows that our method can assist the baseline models in recovering the finer appearance and geometric details with less blur and artifacts, such as reflections and concentrated surface of *materials*, tiny leaves, and semi-translucent edges of *ficus*. The rendering results of real-world scenes are illustrated in Fig. 4 and Fig. 5. As shown in Fig. 4, although the vanilla Nerfacc can reconstruct the overall scene geometry well, the tiny objects and local textures are blurry. Our method can help the model capture more geometric details like the leaves and the potted plants in the *garden*, and generate better appearance, such as cleaner leaves and plates in the *room*. Fig. 5 shows that Zip-NeRF renders photo-realistic images but fails to reconstruct some challenging small structures like banners



**Fig. 3:** Visual comparison of our method with baselines K-Planes [18], Tri-Mip [23] and Nerfacc [28] on synthetic scenes. Our representation works well on the three baselines, enabling to reduce blurring and artifacts and reconstruct better geometry and appearance on challenging details.

of *bicycle* and *lego of kitchen*. Our method can improve the visual performance of Zip-NeRF, producing sharper and more accurate renderings.

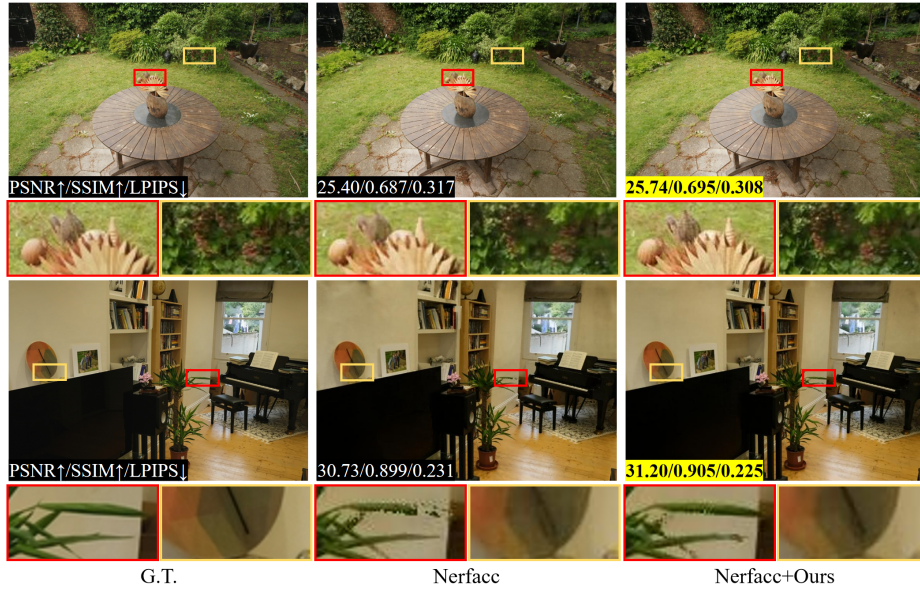
**Table 3:** Quantitative ablation study of the design choices of our anisotropic neural representation on Mip-360 and Blender.

	Blender				Mip-360			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg. $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg. $\downarrow$
Nerfacc	33.06	0.961	0.053	0.017	27.47	0.771	0.294	0.063
Nerfacc+aniso- $\sigma$	33.41	0.961	0.050	0.017	28.07	0.774	0.290	0.060
Nerfacc+aniso-e	33.68	0.963	0.049	0.016	28.21	0.779	0.286	0.059
Nerfacc+Ours	<b>34.08</b>	<b>0.966</b>	<b>0.046</b>	<b>0.015</b>	<b>28.32</b>	<b>0.782</b>	<b>0.281</b>	<b>0.058</b>

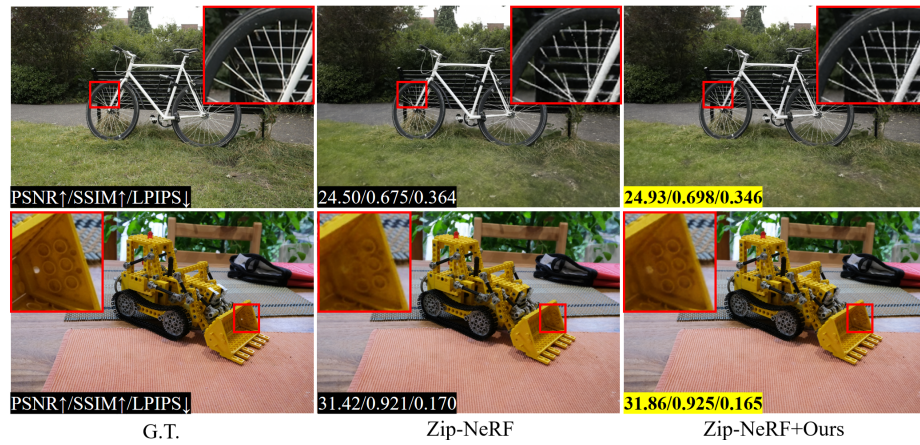
## 5.2 Ablation Study

*Anisotropic neural representation.* We first verify the effectiveness of our anisotropic neural representation. In Table 3, we compare our proposed representation to two variations and a baseline (*Nerfacc*). The first one (*Nerfacc+aniso- $\sigma$* ) is to replace the anisotropic latent feature  $\mathbf{e}(\mathbf{d}, \mathbf{W})$  with isotropic one  $\mathbf{e}$  and uses only anisotropic density  $\sigma(\mathbf{d}, \mathbf{k})$ .

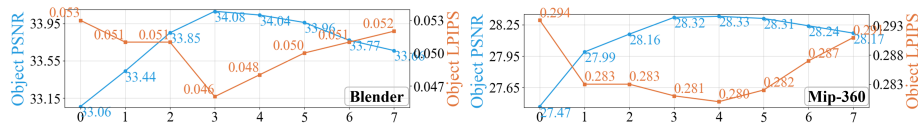




**Fig. 4:** Qualitative comparison of our method with Nerfacc on Mip-360. Our method enhances the rendering quality of vanilla Nerfacc, capturing more geometric details and producing more correct structure and texture.



**Fig. 5:** Qualitative comparison of our method with Zip-NeRF on Mip-360. Our method enables them to estimate opacity more precisely and reconstruct finer details.



**Fig. 6:** The impact of the maximal degree of SH functions on rendering quality in terms of PSNR and LPIPS. The  $x$ -axis refers to the maximal SH degree  $L$ .

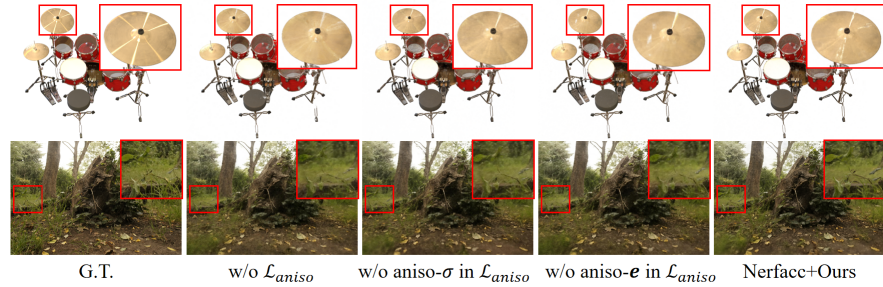
This change with anisotropic density can improve the rendering quality of the vanilla Nerfacc but it leads to large performance drops on all metrics compared to our design, showing the necessity of having the anisotropic latent feature for rendering. The second one (*Nerfacc+aniso-e*) is to replace the anisotropic density  $\sigma(\mathbf{d}, \mathbf{k})$  with the isotropic density  $\sigma$  and uses only the anisotropic latent feature  $\mathbf{e}(\mathbf{d}, \mathbf{W})$ . Although this also enhances the rendering quality of the vanilla Nerfacc, it produces inferior performance than our overall design, indicating that the anisotropic density can indeed help to learn geometry better. Additionally, The second variation (*Nerfacc+aniso-e*) yields slightly better numerical performance than the first variation (*Nerfacc+aniso- $\sigma$* ). This can be explained by the fact that the dimension of the anisotropic latent feature  $\mathbf{e}(\mathbf{d}, \mathbf{W})$  is higher than the density  $\sigma(\mathbf{d}, \mathbf{k})$ , which brings more performance gains.

*Maximal Degree of SH basis.* In Fig. 6, we show the impact of the maximal degree of SH basis on rendering quality in terms of PSNR and LPIPS. Our model is the same as vanilla Nerfacc when setting  $L = 0$ . Our model gets the best rendering quality when setting  $L = 3$  on Blender and  $L = 4$  on Mip-360, since higher maximal degree more anisotropy can be captured by SH functions but continuously increasing the maximal degree will lead to anisotropy over-fitting and can hardly be interpolated accurately. Meanwhile, the computation complexity of SH functions ( $O(L^2)$ ) increases exponentially as  $L$  increases. We set  $L = 3$  in our experiments.

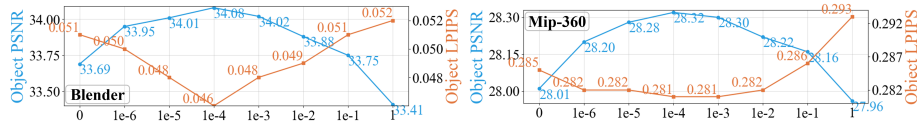
**Table 4:** Quantitative ablation study of anisotropy regularization loss on Mip-360 and Blender.

	Blender				Mip-360			
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg. $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	Avg. $\downarrow$
w/o $\mathcal{L}_{aniso}$ .	33.69	0.962	0.051	0.016	28.01	0.775	0.285	0.060
w/o aniso- $\sigma$ in $\mathcal{L}_{aniso}$	33.88	0.964	0.049	0.016	28.13	0.775	0.285	0.059
w/o aniso- $\mathbf{e}$ in $\mathcal{L}_{aniso}$	33.81	0.963	0.050	0.016	28.25	0.778	0.284	0.059
Nerfacc+Ours	<b>34.08</b>	<b>0.966</b>	<b>0.046</b>	<b>0.015</b>	<b>28.32</b>	<b>0.782</b>	<b>0.281</b>	<b>0.058</b>

*Anisotropy regularization.* We first verify the effectiveness of different anisotropy regularization loss terms for the mapping process from Sec. 4.2. The experimental results



**Fig. 7:** Qualitative ablation study of our anisotropy regularization loss. Using our full loss with  $\mathcal{L}_{recon}$  and  $\mathcal{L}_{aniso}$  renders novel views with sharper texture details and more accurate geometry.



**Fig. 8:** The effect of the strength of anisotropy regularization  $\lambda$  on rendering quality. The  $x$ -axis refers to  $\lambda$ .

in Table 4 and Fig. 7 show that using all losses together leads to the best overall performance. Without anisotropy regularization loss ( $w/o \mathcal{L}_{aniso}$ ) or without anisotropy density ( $w/o aniso-\sigma$ ) or latent feature ( $w/o aniso-e$ ) regularization in  $\mathcal{L}_{aniso}$ , the reconstruction performance drops significantly, indicating that these anisotropy regularization losses are important for the disambiguation of the optimization process. We also investigate the impact of anisotropy regularization strength in Fig. 8. Our method benefits more from a larger  $\lambda$  in terms of PSNR and LPIPS across two datasets, with  $\lambda = 1e^{-4}$  being the best. Note that  $\mathcal{L}_{aniso}$  is disabled when  $\lambda = 0$ . When  $\lambda \geq 1e^{-3}$ , the rendering quality decreases. Because the overly strong anisotropy penalization leads to limited anisotropy in representation. In our experiments, we set  $\lambda = 1e^{-4}$  to balance anisotropy capturing and over-fitting.

## 6 Conclusion

We introduced a novel anisotropic neural representation for NeRFs using spherical harmonic (SH) functions, which enables accurate representation and novel view rendering in complex scenes. We used SH functions and the corresponding coefficients that were estimated by a modified NeRF MLP to determine the anisotropic features. During training, an anisotropy regularization is introduced to alleviate the anisotropy over-fitting problem. We showed that our design results in a simple and flexible representation module that is easy to generalize to various NeRFs. Experiments on both synthetic and real-world datasets demonstrated the effectiveness of our method in improving the rendering quality for NeRFs.

## References

1. Achlioptas, P., Diamanti, O., Mitliagkas, I., Guibas, L.: Learning representations and generative models for 3d point clouds. In: International conference on machine learning. pp. 40–49. PMLR (2018) [4](#)
2. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. ICCV (2021) [1](#), [2](#), [4](#), [9](#)
3. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. CVPR (2022) [2](#), [4](#), [9](#)
4. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Zip-NeRF: Anti-aliased grid-based neural radiance fields. arXiv preprint arXiv:2304.06706 (2023) [2](#), [4](#), [9](#)
5. Basri, R., Jacobs, D.W.: Lambertian reflectance and linear subspaces. IEEE transactions on pattern analysis and machine intelligence **25**(2), 218–233 (2003) [7](#)
6. Bian, W., Wang, Z., Li, K., Bian, J., Prisacariu, V.A.: NoPe-NeRF: Optimising neural radiance field with no pose prior (2023) [3](#)
7. Blackman, R.B., Tukey, J.W.: The measurement of power spectra from the point of view of communications engineering—part i. Bell System Technical Journal **37**(1), 185–282 (1958) [4](#)
8. Buehler, C., Bosse, M., McMillan, L., Gortler, S., Cohen, M.: Unstructured lumigraph rendering. In: Seminal Graphics Papers: Pushing the Boundaries, Volume 2, pp. 497–504 (2023) [3](#)
9. Chaurasia, G., Duchene, S., Sorkine-Hornung, O., Drettakis, G.: Depth synthesis and local warps for plausible image-based navigation. ACM Transactions on Graphics (TOG) **32**(3), 1–12 (2013) [3](#)
10. Chaurasia, G., Sorkine, O., Drettakis, G.: Silhouette-aware warping for image-based rendering. In: Computer Graphics Forum. vol. 30, pp. 1223–1232. Wiley Online Library (2011) [3](#)
11. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: TensorRF: Tensorial radiance fields. In: European Conference on Computer Vision (ECCV) (2022) [4](#)
12. Chen, A., Xu, Z., Zhao, F., Zhang, X., Xiang, F., Yu, J., Su, H.: MVSNeRF: Fast generalizable radiance field reconstruction from multi-view stereo. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 14124–14133 (2021) [1](#)
13. Chen, X., Zhang, Q., Li, X., Chen, Y., Feng, Y., Wang, X., Wang, J.: Hallucinated neural radiance fields in the wild. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12943–12952 (2022) [4](#)
14. Chen, Z., Wang, C., Guo, Y.C., Zhang, S.H.: StructNeRF: Neural radiance fields for indoor scenes with structural hints. IEEE Transactions on Pattern Analysis and Machine Intelligence (2023) [3](#)
15. Chen, Z., Funkhouser, T., Hedman, P., Tagliasacchi, A.: MobileNeRF: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In: The Conference on Computer Vision and Pattern Recognition (CVPR) (2023) [1](#)
16. Choi, I., Gallo, O., Troccoli, A., Kim, M.H., Kautz, J.: Extreme view synthesis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7781–7790 (2019) [4](#)
17. Debevec, P.E., Taylor, C.J., Malik, J.: Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. In: Seminal Graphics Papers: Pushing the Boundaries, Volume 2, pp. 465–474 (2023) [3](#)
18. Fridovich-Keil, S., Meanti, G., Warburg, F.R., Recht, B., Kanazawa, A.: K-Planes: Explicit radiance fields in space, time, and appearance. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12479–12488 (2023) [9](#), [11](#)

19. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5501–5510 (2022) [4](#)
20. Goesele, M., Ackermann, J., Fuhrmann, S., Haubold, C., Klowsky, R., Steedly, D., Szeliski, R.: Ambient point clouds for view interpolation. In: ACM SIGGRAPH 2010 papers, pp. 1–6 (2010) [3](#)
21. Han, K., Xiang, W., Yu, L.: Volume feature rendering for fast neural radiance field reconstruction. *Advances in Neural Information Processing Systems* **36** (2024) [4](#)
22. He, T., Collomosse, J., Jin, H., Soatto, S.: DeepVoxels++: Enhancing the fidelity of novel view synthesis from 3d voxel embeddings. In: Proceedings of the Asian Conference on Computer Vision (2020) [3](#)
23. Hu, W., Wang, Y., Ma, L., Yang, B., Gao, L., Liu, X., Ma, Y.: Tri-MipRF: Tri-mip representation for efficient anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19774–19783 (2023) [2, 4, 9, 11](#)
24. Isaac-Medina, B.K., Willcocks, C.G., Breckon, T.P.: Exact-NeRF: An exploration of a precise volumetric parameterization for neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 66–75 (2023) [2](#)
25. Kim, M., Seo, S., Han, B.: InfoNeRF: Ray entropy minimization for few-shot neural volume rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12912–12921 (2022) [4](#)
26. Le, H.A., Mensink, T., Das, P., Gevers, T.: Novel view synthesis from single images via point cloud transformation. *arXiv preprint arXiv:2009.08321* (2020) [4](#)
27. Levoy, M., Hanrahan, P.: Light field rendering. In: *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pp. 441–452 (2023) [3](#)
28. Li, R., Gao, H., Tancik, M., Kanazawa, A.: NerfAcc: Efficient sampling accelerates nerfs. *arXiv preprint arXiv:2305.04966* (2023) [2, 9, 11](#)
29. Lombardi, S., Simon, T., Saragih, J., Schwartz, G., Lehrmann, A., Sheikh, Y.: Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751* (2019) [3](#)
30. Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duckworth, D.: NeRF in the wild: Neural radiance fields for unconstrained photo collections. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7210–7219 (2021) [4](#)
31. Max, N.: Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics* **1**(2), 99–108 (1995) [4, 5, 6](#)
32. Max, N., Chen, M.: Local and global illumination in the volume rendering integral. *Tech. rep., Lawrence Livermore National Lab.(LLNL), Livermore, CA (United States)* (2005) [5](#)
33. Mildenhall, B., Hedman, P., Martin-Brualla, R., Srinivasan, P.P., Barron, J.T.: 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. pp. 16190–16199 (2022) [4](#)
34. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: NeRF: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* **65**(1), 99–106 (2021) [1, 4, 5, 6, 9](#)
35. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* **41**(4), 1–15 (2022) [4](#)
36. Niemeyer, M., Barron, J.T., Mildenhall, B., Sajjadi, M.S., Geiger, A., Radwan, N.: RegNeRF: Regularizing neural radiance fields for view synthesis from sparse inputs. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5480–5490 (2022) [4](#)



37. Penner, E., Zhang, L.: Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)* **36**(6), 1–11 (2017) [3](#)
38. Ramamoorthi, R., Hanrahan, P.: On the relationship between radiance and irradiance: determining the illumination from images of a convex lambertian object. *JOSA A* **18**(10), 2448–2459 (2001) [7](#)
39. Rematas, K., Ferrari, V.: Neural voxel renderer: Learning an accurate and controllable rendering tool. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5417–5427 (2020) [3](#)
40. Riemann, B.: On the hypotheses which lie at the foundation of geometry, translated by wk clifford. *Nature* **8**, 1873 (1854) [5](#)
41. Schussman, G., Ma, K.L.: Anisotropic volume rendering for extremely dense, thin line data. In: *IEEE Visualization 2004*. pp. 107–114. IEEE (2004) [5, 6](#)
42. Sinha, S., Steedly, D., Szeliski, R.: Piecewise planar stereo for image-based rendering. In: *2009 International Conference on Computer Vision*. pp. 1881–1888 (2009) [3](#)
43. Sitzmann, V., Thies, J., Heide, F., Nießner, M., Wetzstein, G., Zollhofer, M.: DeepVoxels: Learning persistent 3d feature embeddings. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 2437–2446 (2019) [3](#)
44. Sloan, P.P., Kautz, J., Snyder, J.: Precomputed radiance transfer for real-time rendering in dynamic, low-frequency lighting environments. In: *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pp. 339–348 (2023) [7](#)
45. Song, Z., Chen, W., Campbell, D., Li, H.: Deep novel view synthesis from colored 3d point clouds. In: *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIV* 16. pp. 1–17. Springer (2020) [4](#)
46. Srinivasan, P.P., Wang, T., Sreelal, A., Ramamoorthi, R., Ng, R.: Learning to synthesize a 4d rgb-d light field from a single image. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 2243–2251 (2017) [3](#)
47. SuLvXiangXin: ZipNeRF-torch: a pytorch implementation of zipnerf. <https://github.com/SuLvXiangXin/zipnerf-pytorch> (2023) [9](#)
48. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5459–5469 (2022) [4](#)
49. Wang, G., Chen, Z., Loy, C.C., Liu, Z.: SparseNeRF: Distilling depth ranking for few-shot novel view synthesis. *arXiv preprint arXiv:2303.16196* (2023) [4](#)
50. Wang, J., Wang, P., Long, X., Theobalt, C., Komura, T., Liu, L., Wang, W.: NeuRIS: Neural reconstruction of indoor scenes using normal priors. In: *European Conference on Computer Vision*. pp. 139–155. Springer (2022) [3](#)
51. Wang, J., Sun, B., Lu, Y.: MVPNet: Multi-view point regression networks for 3d object reconstruction from a single image. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 33, pp. 8949–8956 (2019) [4](#)
52. Wang, Y., Gong, Y., Zeng, Y.: Hyb-NeRF: A multiresolution hybrid encoding for neural radiance fields. *arXiv preprint arXiv:2311.12490* (2023) [1](#)
53. Wang, Y., Li, Z., Jiang, Y., Zhou, K., Cao, T., Fu, Y., Xiao, C.: NeuralRoom: Geometry-constrained neural implicit surfaces for indoor scene reconstruction. *arXiv preprint arXiv:2210.06853* (2022) [3](#)
54. Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E.: Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* **13**(4), 600–612 (2004). <https://doi.org/10.1109/TIP.2003.819861> [9](#)
55. Wood, D.N., Azuma, D.I., Aldinger, K., Curless, B., Duchamp, T., Salesin, D.H., Stuetzle, W.: Surface light fields for 3d photography. In: *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pp. 487–496 (2023) [3](#)

56. Wu, H., Graikos, A., Samaras, D.: S-VolSDF: Sparse multi-view stereo regularization of neural implicit surfaces. arXiv preprint arXiv:2303.17712 (2023) [3](#)
57. Xin, H., Qi, Z., Ying, F., Hongdong, L., Xuan, W., Qing, W.: HDR-NeRF: High dynamic range neural radiance fields. arXiv preprint arXiv:2111.14451 [2](#) (2021) [4](#)
58. Xing, W., Chen, J.: MVSPlenOctree: Fast and generic reconstruction of radiance fields in plenotree from multi-view stereo. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 5114–5122 (2022) [1](#)
59. Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., Neumann, U.: Point-NeRF: Point-based neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5438–5448 (2022) [1](#)
60. Yang, J., Pavone, M., Wang, Y.: FreeNeRF: Improving few-shot neural rendering with free frequency regularization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8254–8263 (2023) [1](#), [4](#)
61. Zhang, K., Riegler, G., Snavely, N., Koltun, V.: NeRF++: Analyzing and improving neural radiance fields. arXiv preprint arXiv:2010.07492 (2020) [2](#), [3](#), [7](#)
62. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 586–595 (2018) [9](#)
63. Zhang, X., Zheng, Z., Gao, D., Zhang, B., Pan, P., Yang, Y.: Multi-view consistent generative adversarial networks for 3d-aware image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18450–18459 (2022) [3](#)
64. Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: Learning view synthesis using multiplane images. arXiv preprint arXiv:1805.09817 (2018) [4](#)