

# NeuSample: Neural Sample Field for Efficient View Synthesis

Jiemin Fang<sup>1,2</sup>, Lingxi Xie<sup>3</sup>, Xinggang Wang<sup>2†</sup>, Xiaopeng Zhang<sup>3</sup>, Wenyu Liu<sup>2</sup>, Qi Tian<sup>3</sup>

<sup>1</sup>Institute of Artificial Intelligence, Huazhong University of Science & Technology

<sup>2</sup>School of EIC, Huazhong University of Science & Technology <sup>3</sup>Huawei Inc.

{jainfong, xgwang, liuwy}@hust.edu.cn

{198808xc, zxphistory}@gmail.com tian.qil@huawei.com

[jainfong.cn/neusample/](http://jainfong.cn/neusample/)

## Abstract

Neural radiance fields (NeRF) have shown great potentials in representing 3D scenes and synthesizing novel views, but the computational overhead of NeRF at the inference stage is still heavy. To alleviate the burden, we delve into the coarse-to-fine, hierarchical sampling procedure of NeRF and point out that the coarse stage can be replaced by a lightweight module which we name a neural sample field. The proposed sample field maps rays into sample distributions, which can be transformed into point coordinates and fed into radiance fields for volume rendering. The overall framework is named as NeuSample. We perform experiments on Realistic Synthetic 360° and Real Forward-Facing, two popular 3D scene sets, and show that NeuSample achieves better rendering quality than NeRF while enjoying a faster inference speed. NeuSample is further compressed with a proposed sample field extraction method towards a better trade-off between quality and speed.

## 1. Introduction

Novel view synthesis is a long-standing and important problem in computer vision [4, 36], aiming at reconstructing a 3D object/scene with sparsely sampled views and generating images from unseen views. It has a wide range of applications, such as rendering interactive objects in virtual reality and offering 3D preview of objects/scenes.

Recently, researchers propose the concept of neural radiance fields (NeRF) [17] that represent a scene with a continuous 5D function, which is often formulated using a deep neural network and hence can be optimized via gradient descent. The function takes point coordinates (3D) and the observer’s view direction (2D) as input and outputs the pixel radiance and opacity. To connect real 3D points with image pixels, a classical technique volume rendering [10] is used.

<sup>†</sup>Corresponding author.

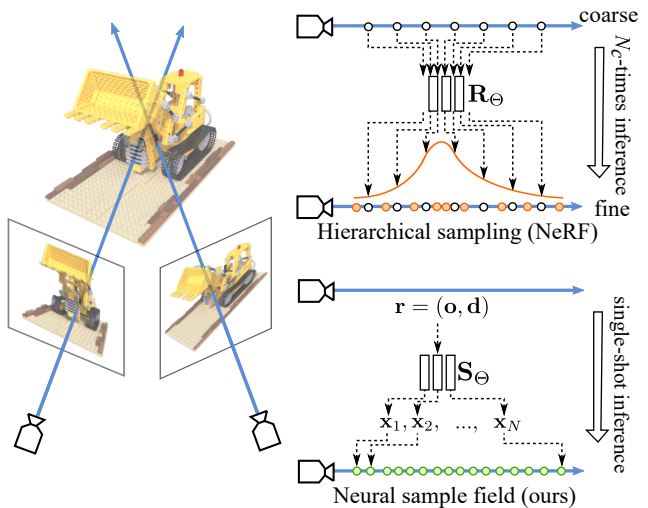


Figure 1. Comparisons between the conventional hierarchical sampling method and our NeuSample. Previous hierarchical sampling needs to first infer a set of coarse points. Based on opacity properties of coarse ones, fine samples are obtained for elaborate rendering. NeuSample directly maps a ray  $\mathbf{r} = (\mathbf{o}, \mathbf{d})$  into a sampling distribution with single inference by constructing a sample field. The obtained samples can directly render high-quality images.

It samples a number of points along the ray (starting from the observer, extending along the view direction), computes their density (*a.k.a.* opacity) values and color properties, and eventually accumulate them into the final output. Compared to conventional methods (*i.e.*, voxel-based or mesh-based representations), NeRF saves considerable storage by compressing each scene/object into a neural network.

To guarantee high rendering quality, NeRF and most subsequent methods [3, 7, 17, 29] adopted a hierarchical sampling method. As shown in Fig. 1, the algorithm first uniformly samples a set of  $N_c$  (usually taken as 64) coarse-level points along the ray, and feeds them into the radiance field to obtain the densities and color properties of these points, as well as an updated distribution to sample more

fine-level points (for another round of inference). Finally, the outputs of coarse-level and fine-level points are accumulated as the final output. Although the strategy leads to improved rendering performance, the two-stage inference incurs heavy computational overheads, making it difficult to apply NeRF to real-life or interactive scenarios.

To alleviate the burden, we propose a novel framework named **NeuSample** that only requires **single-shot** inference for sampling – in other words, we discard the costly coarse stage and instead use a **neural sample field**, which takes a ray representation as input and outputs a sampling distribution along the ray. As shown in Fig. 1, the sample field is also a 5D function and formulated using a neural network, which takes the observer’s coordinate  $(x_o, y_o, z_o)$  and view direction  $(\theta, \phi)$  as input and outputs  $N$  numbers, corresponding to the distances between the points-to-sample and the observer. These points are then fed into the neural radiance field to accomplish the volume rendering procedure.

We perform experiments on two commonly used benchmarks, namely, Realistic Synthetic 360° and Real Forward-Facing. Compared to NeRF, NeuSample demonstrates superior rendering quality and saves around a quarter of inference time. In addition, a sample field extraction method assists NeuSample to achieve competitive rendering quality with 25% of inference time. Diagnostic studies show that the improved quality-complexity trade-off owes to the neural sample field that learns an efficient way of sampling points. The contribution of this work is summarized as follows.

- We propose a sample field to map one ray into a sample distribution. This field is parameterized as a neural network with fully connected layers.
- The proposed sample field can be integrated with radiance fields to perform volume rendering, which not only saves computation cost from coarse networks used in the conventional hierarchical sampling strategy but also shows stronger rendering quality than NeRF.
- A sample field extraction method is proposed to reduce the number of sampled points. This method further accelerates rendering while maintaining competitive rendering quality.

## 2. Related Work

**Neural Implicit Representations** Using neural representations to modelling 3D structures or geometry [2, 6, 13–15, 19, 23, 34, 35] has shown great success. Neural representations model 3D scenes in a continuous space and can be optimized with a differentiable manner. The storage can be saved with network inference. Most of above methods require explicit supervision. NeRF [17] and subsequent works [3, 9, 22, 25, 31, 44, 45] use neural radiance fields to map 3D coordinates into color and opacity

values, which achieve strong performance on synthesizing photo-realistic images from novel views. Besides neural radiance fields, some works adopt other types of neural fields to represent or model 3D scenes, *e.g.* textured material [8, 20, 26, 27], indirect illumination values [30], surfacing reconstruction [21, 40] and light fields [33]. Our work is inspired by the neural field concept and proposes a sample field, which maps rays into sample distributions. This field is efficient and the predicted samples help radiance fields render high-quality images.

**Accelerating Rendering with Caching** One main stream of methods accelerate rendering by pre-computing and storing explicit data structures. For inference, properties of points can be looked up and latency of neural network inference can be drastically reduced. Neural Sparse Voxel Fields (NSVF) [12] allows for empty space skipping and early ray termination by constructing an octree structure. KiloNeRF [29] represents a scene with thousands of tiny MLPs, each of which is only responsible for a cell in the space grid. SNeRG [7], FastNeRF [5] and PlenOctrees [43] *et al.* adopt a similar methodology by pre-computing properties/features with sparse voxel grid-like structures using a learned NeRF. Directly querying these properties during inference can effectively save cost. NeX [42] represents scenes based on the multiplane image (MPI) with MPLs. The MPI grid can be cached for fast rendering. Most of the above methods achieve real-time rendering speed. This methodology can be taken as using storage to complement inference. Though storage may be compressed via techniques like quantization, for high resolution, storage is still a non-negligible element. From a different perspective, our method aims at improving rendering speed without caching and additional storage. NeuSample is a parallel work to the above ones and can be integrated with caching-aided frameworks.

**Accelerating Rendering from Inference Times** A series of methods promote the rendering efficiency by focusing on the neural representation itself and reducing the inference times. AutoInt [11] proposes to automatically integrate output colors for rendering with neural networks, which approximates the integral along rays in piecewise sections. DONeRF [18] uses the depth information to guide sampling by training a “depth oracle” network, which greatly reduces the sample number for rendering. However, depth information is usually hard to obtain in real scenarios. Light Filed Networks [33] directly map a ray into the color, only requiring one single inference for rendering rather than mapping hundreds of points and significantly promoting efficiency. This method discards explicit 3D point modelling so additional supervision may be needed to construct multi-view consistency towards complicated geometry. NeRF-ID [1] reduces the sample number by training an importance predictor, where a proposal network with Transformer/MLP-

Mixer architectures is used for sampling based on coarse network features. TerminiNeRF [24] predicts weights of ray segments and further sample points based on the estimated weights. Our work maintains advantage of volume rendering for modelling multi-view consistency with explicit 3D point coordinates, but saves computation cost from cumbersome coarse fields. Samples are directly obtained by the proposed sample field.

### 3. Method

In this section, we first review formulations in NeRF [17]. Second, we introduce the proposed neural sample field and the overall framework for rendering. We finally introduce a sample field extraction method to produce fewer samples for further acceleration.

#### 3.1. Review of NeRF

**Field Construction and Optimization** The neural radiance field is first proposed in [17], which is designed to represent a scene with emitted radiance for each position in the space. The field is constructed by mapping the coordinate of a point into its volume density and color with neural networks. Denoting a point  $\mathbf{x}$  with a 5D-coordinate  $(x, y, z, \theta, \phi)$ , we predict its volume density and viewing color via neural networks as

$$\sigma, \mathbf{c} = \mathbf{R}_\Theta(x, y, z, \theta, \phi), \quad (1)$$

where  $x, y, z$  indicate the point location,  $\theta, \phi$  are the view direction, and  $\mathbf{R}_\Theta$  is the radiance field instantiated as a neural network. To connect these points in the real space and images taken from the camera, one pixel in an image can be rendered by computing densities and colors of points along the according ray and performing the classic volume rendering [10]. Specifically, the expected color  $\hat{C}(\mathbf{r})$  of the pixel along camera ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  is computed with the quadrature rule as:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i(1 - \exp(-\sigma_i\delta_i))\mathbf{c}(\mathbf{x}_i, \mathbf{d}), \quad (2)$$

$$T_i = \exp(-\sum_{j=1}^{i-1} \sigma(\mathbf{x}_j)\delta_j),$$

where  $\mathbf{o}, \mathbf{d}$  denote the origin and direction of the ray respectively,  $t$  denotes the distance of the sample from the origin, and  $\delta_i$  is defined as the distance between two adjacent samples, *i.e.*,  $\delta_i = t_{i+1} - t_i$ . The field defined in Eq. 1 can be optimized with a differentiable manner by minimizing the loss between the rendered and ground truth color in the image:

$$\mathcal{L} = \|\hat{C}(\mathbf{r}) - C(\mathbf{r})\|_2^2. \quad (3)$$

**Volume Sampling in Radiance Fields** In real implementation, it is impossible to achieve the rendering procedure

defined in Eq. 2 by traversing all the points along the ray. Therefore, samples carrying useful information need to be obtained for rendering. As shown in Fig. 1, NeRF [17] and most of its variants [3, 7, 17, 29] adopt a hierarchical sampling strategy, which requires two field networks. In the first stage,  $N_c$  coarse points are sampled by randomly drawing one from each of  $N_c$  evenly-partitioned bins:

$$t_i \sim \mathcal{U}[t_n + \frac{i-1}{N_c}(t_f - t_n), t_n + \frac{i}{N_c}(t_f - t_n)], \quad (4)$$

where  $t_n$  and  $t_f$  denote the near and far bound respectively. The  $N_c$  coarse samples are mapped by the first field network into colors and densities. In the second stage,  $N_f$  fine samples are generated based on properties from  $N_c$  coarse samples in the first stage. In [17], weight  $w_i = T_i(1 - \exp(-\sigma_i\delta_i))$  of each coarse sample computed in Eq. 2 serves as the probability for sampling fine points. Then both coarse and fine samples are fed into the second field network for final rendering.

#### 3.2. Neural Sample Field

As aforementioned, color and density properties of  $N_c$  coarse samples need to be first inferred by a field network to generate fine samples for final rendering. This procedure takes a large amount of computation cost, *i.e.* 25% of the total cost for 64 coarse samples and 128 fine samples. Inspired by the neural field concept as Eq. 1, we propose a neural sample field  $\mathbf{S}_\Theta$  which maps a ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$  directly into a series of samples for volume rendering:

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \leftarrow \mathbf{S}_\Theta(\mathbf{o}, \mathbf{d}), \quad (5)$$

where  $\mathbf{x}_i$  denotes the coordinates of the  $i$ th sample,  $N$  denotes the number of desired samples. Specifically, we first obtain  $N$  scalars by feeding the ray origin coordinates and direction into  $\mathbf{S}_\Theta$ :

$$\hat{t}_1, \hat{t}_2, \dots, \hat{t}_N = \mathbf{S}_\Theta(x_o, y_o, z_o, \theta, \phi), \quad (6)$$

where  $\hat{t}_i \sim [0, 1]$  represents the relative sample position between the near and far bound along the ray.  $\hat{t}_i$  is mapped to a absolute position by computing  $t_i = (1 - \hat{t}_i)t_n + \hat{t}_i t_f$ . Then we compute the coordinates of the  $i$ th sample with

$$\mathbf{x}_i = \mathbf{o} + t_i\mathbf{d}. \quad (7)$$

**Architecture** We show the architecture of the sample field network  $\mathbf{S}_\Theta$  in Fig. 3. The ray vector with origin coordinates and direction are first mapped into a higher dimension with position encoding, which is widely used in previous works [3, 17, 32, 38]. Then the mapped input is passed through 8 fully connected (FC) layers with ReLU activations. A skip connection is included by concatenating input with the 4-th layer’s output. The hidden dimension

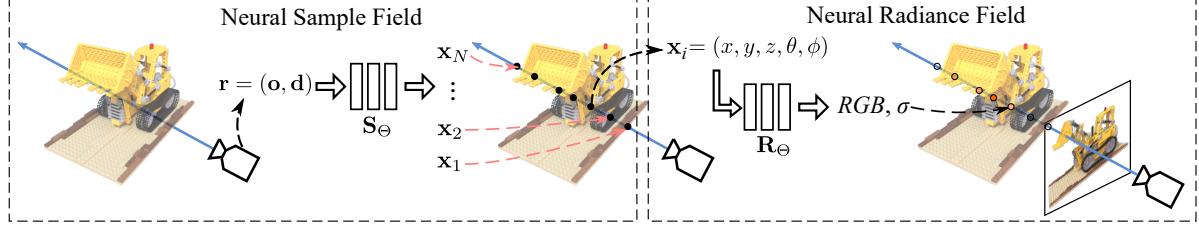


Figure 2. Overall framework of NeuSample. To render a pixel in the image, the ray passing through the pixel is first fed into a sample field network  $\mathbf{S}_\Theta$ , which is mapped to a distribution along the ray. The distribution is then transformed to 3D-point coordinates, which are fed into the radiance field  $\mathbf{R}_\Theta$  to obtain colors and densities. Finally, volume rendering is performed on these points.

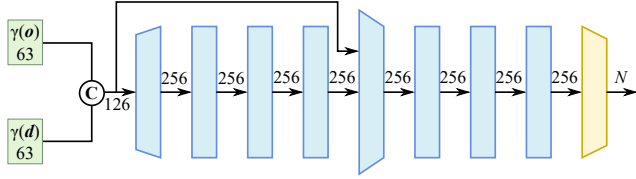


Figure 3. Architecture of the neural sample field network. “ $\gamma$ ” denotes position encoding, “C” denotes concatenation. The blue blocks indicate fully connected layers followed by ReLU activation and yellow blocks indicate layers with sigmoid activation.

of features between FC layers is set as 256. At the end of the network, an additional FC layer maps 256-dimension features into an  $N$ -dimension vector. The vector is finally fed into a sigmoid activation layer and becomes the relative sample positions  $(\hat{t}_1, \hat{t}_2, \dots, \hat{t}_N)$  defined in Eq. 6.

**Overall Framework** As shown in Fig. 2, we integrate the proposed neural sample fields with radiance fields as the overall framework. To render a pixel in the image, we first compute the camera pose and transform the pose to the ray origin and direction according to the pixel position. Then we feed the ray origin coordinates and normalized direction vector into the neural sample field network  $\mathbf{S}_\Theta$ . The sample field network outputs distributions within  $[0, 1]$  which are transformed into 3D coordinates. Samples are then fed into the second neural radiance field to obtain corresponding colors and densities. Finally, the pixel color is generated using volume rendering as Eq. 2. Noting that both the sample field and radiance field network are optimized by minimizing the final rendered color loss as Eq. 3. The whole framework can be trained end to end via gradient descent. The sample field is intuitively learning how to sample points along rays, which is intrinsically modelling geometry structures of the scene.

### 3.3. Sample Field Extraction

Besides saving computation cost from coarse fields, we propose to extract the learned sample field for further acceleration. As shown in Fig. 4, we first train a **regular** sample field network  $\mathbf{S}_\Theta$  which predicts  $N$  substantial samples (e.g.

$N = 192$  as in the fine network of NeRF [17]) for the radiance field  $\mathbf{R}_\Theta$  learning. Then we reduce the output number of  $\mathbf{S}_\Theta$  to  $N_e < N$  and obtain an **extracted** sample field  $\mathbf{S}_\Theta^e$  followed by radiance field  $\mathbf{R}_\Theta^e$ . Parameters of the regular sample field  $\mathbf{S}_\Theta$  and radiance field  $\mathbf{R}_\Theta$  are mapped to the extracted ones  $\mathbf{S}_\Theta^e$  and  $\mathbf{R}_\Theta^e$ . The two radiance field networks share the same architecture, so parameters from  $\mathbf{R}_\Theta$  can be directly copied to  $\mathbf{R}_\Theta^e$ . For the sample fields, only the final FC layers for sample prediction differ where parameters are evenly mapped from  $\mathbf{S}_\Theta$  to  $\mathbf{S}_\Theta^e$  on the output channel dimension. All the other parameters are directly copied. With parameters mapped, we fine-tune the extracted fields  $\mathbf{S}_\Theta^e$  and  $\mathbf{R}_\Theta^e$  only for a few iterations to fit the new distribution.

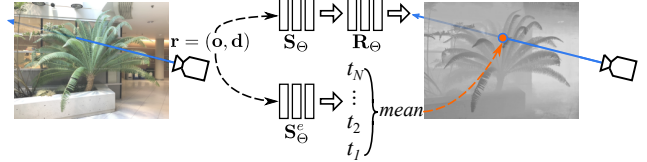


Figure 4. Depth boost for initializing an extracted sample field. The output mean value of the extracted field is forced to fit the depth predicted by a learned regular field.

For real-world scenes which usually have complicated geometry structures or depth distribution, we use the depth information predict by the regular fields to help initialize the extracted sampling filed network. We name this procedure as **depth boost**. Specifically, we sample some camera poses of the scene and feed the corresponding rays to the regular fields. For a ray  $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ , the depth  $d_{\mathbf{r}}$  is predicted as:

$$\begin{aligned} t_1, t_2, \dots, t_N &= \mathbf{S}_\Theta(\mathbf{o}, \mathbf{d}), \\ \sigma_1, \sigma_2, \dots, \sigma_N &= \mathbf{R}_\Theta(t_1, t_2, \dots, t_N), \\ d_{\mathbf{r}} &= \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) t_i. \end{aligned} \quad (8)$$

We make the mean value of the extracted sample field output fit the predicted depth value  $d_{\mathbf{r}}$  by minimizing the loss:

$$\begin{aligned} t_1^e, t_2^e, \dots, t_N^e &= \mathbf{S}_\Theta^e(\mathbf{o}, \mathbf{d}), \\ \mathcal{L}_d &= |(t_1^e + t_2^e + \dots + t_N^e)/N^e - d_{\mathbf{r}}|. \end{aligned} \quad (9)$$



Table 1. PSNR comparisons on the Realistic Synthetic 360° dataset.

Method	Inf. Cost	Average	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship
NeRF [17]	1.00	31.01	33.00	<b>25.01</b>	30.13	36.18	32.54	29.62	<b>32.91</b>	28.65
NeuSample	0.76	<b>31.15</b>	<b>33.02</b>	24.99	<b>30.72</b>	<b>36.29</b>	<b>33.17</b>	<b>29.66</b>	32.68	<b>28.65</b>
AutoInt [11] ( $N = 32$ )	0.31	26.83	25.82	22.02	25.51	31.84	27.26	<b>28.58</b>	28.42	25.18
NeuSample ( $N_e = 64$ )	0.25	<b>28.39</b>	<b>29.96</b>	<b>23.43</b>	<b>27.53</b>	<b>34.41</b>	<b>29.14</b>	27.76	<b>29.42</b>	<b>25.47</b>

\* “Inf. Cost” denotes the relative inference cost compared with NeRF [17], *i.e.* time for rendering one image which is measured on one V100 GPU.

\*  $N_e$  of NeuSample denotes the sample number of the extracted sample field.

\* “ $N = 32$ ” for AutoInt [11] denotes the number of piecewise sections.

Depth boost helps sample fields with fewer points quickly converge to positions with useful information along the ray. Noting that this procedure only requires gradients for the extracted sample field and is not used in subsequent fine-tuning, which is efficient and can be finished with negligible cost.

## 4. Experiments

In this section, we first describe the implementation and experimental details (Sec. 4.1). Then we show results on two widely used benchmarks, *i.e.* Realistic Synthetic 360° and Real Forward-Facing datasets, and compare with other methods (Sec. 4.2). Some ablation studies are performed and shown in Sec. 4.3.

### 4.1. Implementation Details

**Architecture Settings** Our framework is implemented based on PyTorch. We use the same network architecture for all the evaluated scenes. The sample field network is set as Fig. 3 depicts and the radiance field network is used as the same one in NeRF [17]. For the sample field input, we apply 10-frequency position encoding to both ray origin  $\mathbf{o}$  and direction  $\mathbf{d}$ , and the same position encoding as NeRF for radiance field input. The output number of the sample field is set as 192 without specified in the following part.

**Training Hyperparameters** We use the Adam optimizer with a batch size of 4,096 rays. The learning rate decays from  $5 \times 10^{-4}$  to  $5 \times 10^{-6}$  following a polynomial strategy with power 1. For the Realistic Synthetic 360° dataset, we train each scene for 400k iterations in total. For the Real Forward-Facing scenes, we train each one for 100 epochs, where each epoch is completed by randomly sampling 4,096 rays from the whole training set for each iteration. Random noise with 0 mean and unit variance is added to the radiance field’s output densities as NeRF.

**Extraction Hyperparameters** For Realistic Synthetic 360° scenes, the geometry structure is not complicated so we directly map parameters of regular fields to the extracted

ones without depth boost (which we find no additional gain in experiments). For Real Forward-Facing scenes, we sample 120 camera poses with a spiral path to perform depth boost. 8,192 rays are randomly sampled in each iteration. This procedure takes only one epoch with a  $5 \times 10^{-5}$  learning rate. The subsequent fine-tuning takes 40k iterations for synthetic scenes and 20 epochs for real scenes.

**Datasets** We use two datasets, Realistic Synthetic 360° and Real Forward-Facing, to evaluate our method, which are also used in NeRF [17] and most subsequent related methods [1, 3, 7, 11, 29, 42]. The Realistic Synthetic 360° dataset consists of 8 scenes and each one includes 100 views for training and 200 views for testing. We take all views with an  $800 \times 800$  resolution. The Real Forward-Facing dataset contains 8 complex real-world scenes from NeRF [17] and LLFF [16]. Each scene includes 20 - 62 images. Following [17], we hold out  $\frac{1}{8}$  images for testing and the rest are for training. All images are at  $1008 \times 756$  pixels for experiments if unspecified.

### 4.2. Results and Comparisons

**Realistic Synthetic 360°** We show main PSNR results in Tab. 1 and compare with NeRF [17] and AutoInt [11], which is a very related work focusing on reducing inference cost. Though  $\sim 25\%$  computation cost of course network inference is saved, our NeuSample still achieves similar or better PSNR performance compared with NeRF [17], *e.g.* +0.59 PSNR for the Ficus scene and +0.63 for Lego. When the sample field is extracted to output 64 points, NeuSample achieves significantly better performance than AutoInt, *i.e.* 1.56 better average PSNR. The qualitative visualization results are shown in Fig. 7. Noting that though NeRF has produced high-quality synthesis results, NeuSample shows advantages in modelling some details.

**Real Forward-Facing** The Real Forward-Facing dataset is more challenging as it contains complex scenes in real-world scenarios, which requires more elaborate geometry structure modelling than synthetic ones. As shown in Tab. 2,

Table 2. PSNR comparisons on the Real Forward-Facing dataset.

Method	Inf. Cost	Average	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	T-Rex
<b>1008 × 756 Resolution</b>										
NeRF [17]	1.00	26.50	<b>25.17</b>	27.40	31.16	27.45	20.92	20.36	32.70	26.80
NeRF-ID <sup>†</sup>	>1.00	26.76	25.01	27.85	<b>31.51</b>	27.88	21.09	<b>20.38</b>	32.93	27.45
NeuSample	0.76	<b>26.83</b>	24.99	<b>28.14</b>	31.26	<b>28.32</b>	<b>21.10</b>	20.08	<b>33.26</b>	<b>27.46</b>
NeuSample ( $N_e = 64$ )	0.25	26.50	24.77	28.03	31.09	27.45	21.06	20.03	32.71	26.82
<b>504 × 378 Resolution</b>										
NeRF [17]	1.00	27.93	<b>26.92</b>	<b>28.57</b>	<b>32.94</b>	29.26	22.50	<b>21.37</b>	33.60	28.26
NeuSample	0.76	<b>28.14</b>	26.84	28.36	32.76	<b>30.20</b>	<b>22.50</b>	20.99	<b>34.22</b>	<b>29.23</b>
AutoInt [11] ( $N = 32$ )	0.31	25.53	23.51	28.11	28.95	27.64	20.84	17.30	30.72	27.18
NeuSample ( $N_e = 64$ )	0.25	<b>27.80</b>	<b>26.74</b>	<b>28.24</b>	<b>32.33</b>	<b>29.37</b>	<b>22.45</b>	<b>20.89</b>	<b>33.61</b>	<b>28.79</b>
NeuSample ( $N_e = 32$ )	0.13	26.94	26.24	27.95	31.87	27.48	22.33	20.50	31.56	27.55

\* <sup>†</sup> denotes training with a larger batch size, *i.e.* 66k in NeRF-ID [1].

\* Inference cost for NeRF-ID includes the additional proposal network.

Table 3. Comparison with state-of-the-art methods on Real Forward-Facing scenes.

	SRN [34]	LLFF [16]	NeRF [17]	DeRF [28]	IBRNet [41]	GRF [39]	SNeRG [7]	NeRF-ID [1]	NeuSample
<b>PSNR</b> ↑	22.84	24.13	26.50	24.81	26.73	26.64	25.63	26.76	<b>26.83</b>
<b>SSIM</b> ↑	0.668	0.798	0.811	0.767	<b>0.851</b>	0.837	0.818	0.822	0.823
<b>LPIPS</b> ↓	0.378	0.212	0.250	0.274	<b>0.175</b>	0.178	0.183	\	0.231

we perform experiments on two resolutions, *i.e.*  $1008 \times 756$  and  $504 \times 378$ , for better comparisons. Under the large resolution, our method achieves a high average PSNR compared with both NeRF and NeRF-ID<sup>†</sup>. When we reduce the sample number to 64, NeuSample still achieves 26.50 PSNR as high as NeRF while the computation cost has been decreased to 25.4% of NeRF. For the half resolution setting, NeuSample shows 0.21 higher PSNR than NeRF. With fewer samples, NeuSample consistently outperforms AutoInt for all scenes by a large margin, +2.27 PSNR for 64 points +1.41 for 32 points. We comprehensively compare with other state-of-the-art methods on the Real Forward-Facing dataset in Tab. 3. NeuSample still achieves competitive rendering quality. The qualitative results are provided in Fig. 7 and NeuSample performs better in detail modelling as well.

**Evaluation with Diverse Sample Numbers** To comprehensively compare with the state-of-the-art method NeRF-ID [1], we evaluate NeuSample with diverse output numbers of the sample field, *i.e.* 192, 128, 64 and 32, which reduce the computation cost of NeRF to 75.4%, 50.4%, 25.4% and 12.9% respectively. As shown in Fig. 5, NeuSample achieves evidently better trade-off between rendering qual-

<sup>1</sup>NeRF-ID uses a very large batch size, *i.e.* 66k, on 16 Cloud TPUs. This training strategy can significantly promote the baseline PSNR by 0.82 on synthetic datasets and 0.16 on real datasets, which we find it hard to reproduce. Therefore we only compare with NeRF-ID on Real Forward-Facing scenes for fairness.

ity and computation cost than both NeRF and NeRF-ID for almost all scenes.

The above experiments demonstrate that the proposed sample field can not only save cost from coarse network inference, but also capture more elaborate samples for learning scenes and rendering images. Even though the sample field is extracted, the radiance field can still render high-quality images with samples which carry useful information.

### 4.3. Ablation Study

We perform ablation studies on two scenes, *i.e.*, Realistic Synthetic 360° lego at  $800 \times 800$  and Real Forward-Facing fern at  $1008 \times 756$ .

Table 4. Depth boost effectiveness study on the “fern” scene of the Real Forward-Facing dataset.

#Sample	Depth Boost	PSNR ↑	SSIM ↑	LPIPS ↓
128	-	24.99	0.798	0.272
64	✓	24.77	0.785	0.289
	✗	24.65	0.782	0.290
32	✓	24.33	0.765	0.307
	✗	24.05	0.753	0.318

**Depth Boost for Field Extraction** We evaluate effectiveness of depth boost proposed in Sec. 3.3. As shown in Tab. 4, we perform this ablation study on two sample num-

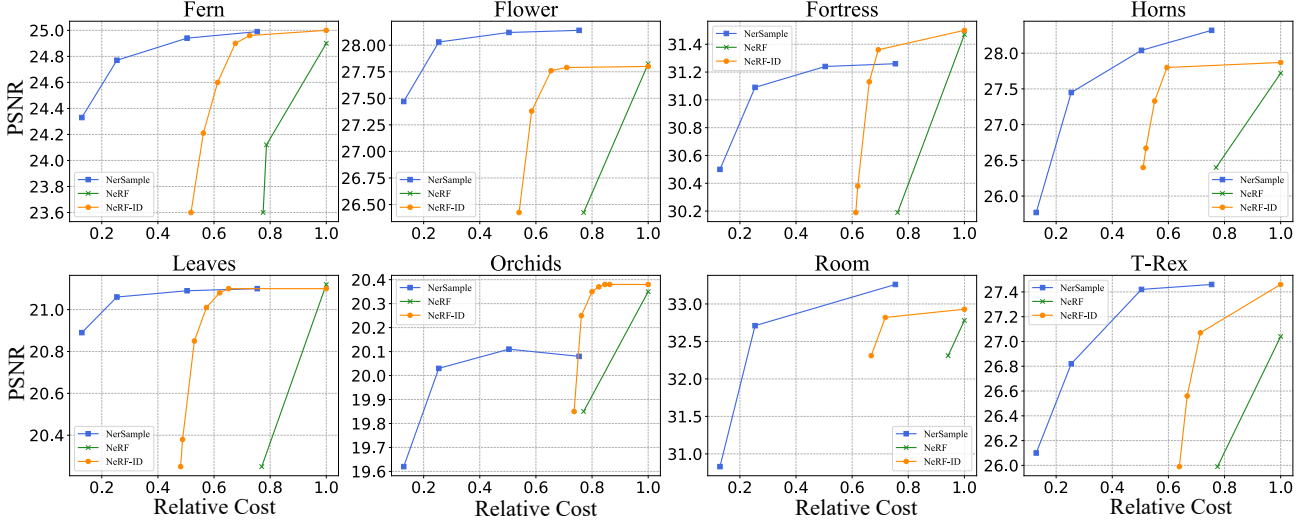


Figure 5. Speedup comparisons on the Real Forward-Facing dataset with NeRF [17] and NeRF-ID [1].

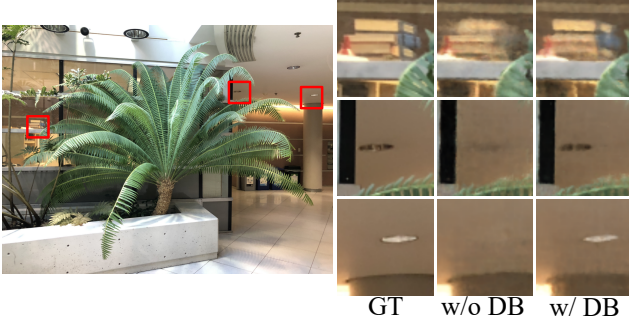


Figure 6. Rendering quality comparisons of sample field extraction *w/o* and *w/* depth boost on the  $1008 \times 756$  fern scene. “GT” denotes ground truth and “DB” denotes depth boost. The extracted field without depth boost fails to render some small objects in the scene.

bers, *i.e.* 64 and 32. For the 64-sample setting, extracting the sample field without depth boost leads to 0.12 PSNR decay. When the samples become fewer to 32, PSNR degrades more by 0.28. We show the rendering results in Fig. 6 and find without depth boost, some small objects in the scene are omitted while the field with depth boost models accurate outline of these objects. These defects though cause small changes to PSNR values but are evident in the final rendered image. This experiment reveals that a few iterations of depth boost to initialize the sample field effectively helps to locate points with high importance. Even with few samples, the real pixel color can be rendered accurately.

**Layer Numbers of Sample Field Network** We study the layer number design for constructing the sample field network. As shown Tab. 5, three layer settings are evaluated, *i.e.* 8 (the default setting), 4 and 2, on the two scenes of lego and fern. We observe that the layer number decrease causes

Table 5. Layer number study of the sample field network on the synthetic lego and real fern scenes.

Scene	#Layers	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
lego	8	33.17	0.965	0.048
	4	29.29	0.934	0.106
	2	28.81	0.937	0.090
fern	8	24.99	0.798	0.272
	4	24.70	0.778	0.303
	2	24.79	0.790	0.283

slight impact for the fern scene rendering within  $\sim 0.29$  PSNR decay. However, smaller layer numbers lead to drastically rendering quality degradation for the lego scene. Setting to 4 layers drops PSNR by 3.88; and setting to 2 causes 4.36 decay. We analyze this phenomenon and deduce the reason as follows. The synthetic lego scene contains more views in a larger range than fern, so it requires a deeper neural network with more parameters to fit diverse view-dependent sampling distributions.

Table 6. Frequency study of position encoding on the synthetic lego and real fern scenes.

Scene	#Frequencies	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
lego	5	33.09	0.966	0.047
	10	33.17	0.965	0.048
	15	33.02	0.965	0.047
fern	5	24.80	0.790	0.283
	10	24.99	0.798	0.272
	15	25.12	0.800	0.270

**Frequencies of Position Encoding** We study frequencies of position encoding in the sample field network. As shown



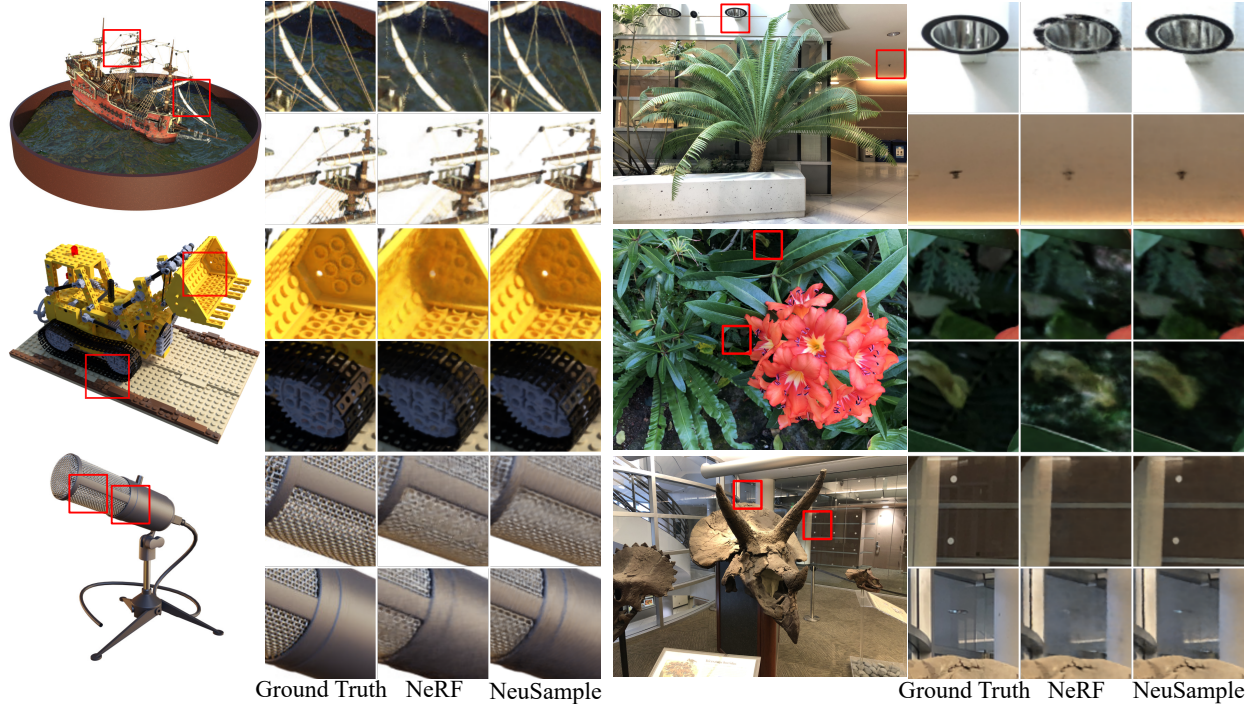


Figure 7. Visualization comparisons on the Realistic Synthetic 360° (left) and Real Forward-Facing (right) dataset. Our method shows higher quality on rendering detailed of thin structures in the scene.

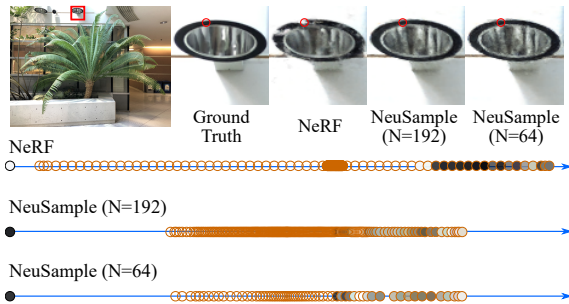


Figure 8. Sample visualization comparisons along the ray in the real scene “fern”. Samples are colored with “RGBA” values, which are set as predicted colors and densities. The edge of each sample is of the orange color, while the orange mass in the NeRF ray represents that lots of points are gathered. The start point is with the finally rendered color.

in Tab. 6, we evaluate three frequency numbers 5, 10 (the default setting) and 15 on two scenes of lego and fern. We find changing frequencies in this range does not affect the sample field much,  $< 0.2$  PSNR for both scenes.

**Sample Visualization** As shown in Fig. 8, we visualize the obtained samples of three models, *i.e.* NeRF [17], our NeuSample with 192 points and an extrated NeuSample with 64 points. We observe NeRF tends to sample points in a wide range from the near bound to far bound but locate too many samples at a similar position (the orange mass). On the contrary, NeuSample locates points in a narrower

range but the points are more even. With fewer points, NeuSample can still obtain the key samples which render accurate colors. We deduce that an extremely dense distribution around a local point is not beneficial for the field training, *e.g.* NeRF predicts wrong opacity values for key points in this example.

## 5. Conclusion

In this paper, we propose a neural sample field which maps a ray to sampling distributions. The proposed sample field can be integrated with neural radiance fields for volume rendering, which we name NeuSample. Our method obtains samples for rendering with a lightweight network module, which saves computation cost from widely used course networks and shows stronger ability of synthesizing novel views for 3D scenes. With a learned NeuSample, the sample field can be extracted for further acceleration with maintaining high rendering quality.

**Limitations** The proposed NeuSample approach has similar limitations as radiance fields in previous methods [3, 17, 45] which needs to be optimized towards every independent scene, as the sample field is directly related to the geometry structure of each specific scene. The methods for generalizing the learned field or speeding up the optimization procedure [9, 33, 37, 39, 44] could be integrated with our method for further reducing the training cost.

Besides, though NeuSample accelerates rendering, it is



unlikely that purely using NeuSample can achieve real-time rendering like methods with caching-based techniques [5, 7, 29, 42, 43]. In comparison, NeuSample and [1, 11, 18] accelerate rendering from another path that focuses on the field itself and requires no additional storage. Overall, this is a problem of achieving better trade-off between storage and speed, and we believe that integrating the above methodologies is a promising future research direction.

## Acknowledgements

We sincerely thank Liangchen Song, Yingqing Rao and Yuzhu Sun for their generous assistance and discussion.

## References

- [1] Relja Arandjelović and Andrew Zisserman. Nerf in detail: Learning to sample for view synthesis. *arXiv:2106.05264*, 2021. 2, 5, 6, 7, 9, 10, 11
- [2] Matan Atzmon, Niv Haim, Lior Yariv, Ofer Israelov, Haggai Maron, and Yaron Lipman. Controlling neural level sets. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 1, 2, 3, 5, 8
- [4] Joel Carranza, Christian Theobalt, Marcus A Magnor, and Hans-Peter Seidel. Free-viewpoint video of human actors. *ACM transactions on graphics (TOG)*, 2003. 1
- [5] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. 2, 9
- [6] Kyle Genova, Forrester Cole, Daniel Vlasic, Aaron Sarna, William T Freeman, and Thomas Funkhouser. Learning shape templates with structured implicit functions. In *ICCV*, 2019. 2
- [7] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021. 1, 2, 3, 5, 6, 9
- [8] Philipp Henzler, Niloy J Mitra, and Tobias Ritschel. Learning a neural 3d texture space from 2d exemplars. In *CVPR*, 2020. 2
- [9] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, 2021. 2, 8
- [10] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 1984. 1, 3
- [11] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *CVPR*, 2021. 2, 5, 6, 9, 10, 11
- [12] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 2
- [13] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020. 2
- [14] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, 2019. 2
- [15] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *ICCV*, 2019. 2
- [16] 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 (Proceedings of SIGGRAPH)*, 2019. 5, 6
- [17] 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 *ECCV*, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 10, 11
- [18] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton S. Kaplanyan, and Markus Steinberger. DONeRF: Towards Real-Time Rendering of Compact Neural Radiance Fields using Depth Oracle Networks. *Computer Graphics Forum*, 2021. 2, 9
- [19] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d reconstruction by learning particle dynamics. In *ICCV*, 2019. 2
- [20] Michael Oechsle, Lars Mescheder, Michael Niemeyer, Thilo Strauss, and Andreas Geiger. Texture fields: Learning texture representations in function space. In *ICCV*, 2019. 2
- [21] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *International Conference on Computer Vision (ICCV)*, 2021. 2
- [22] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021. 2
- [23] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *ECCV*, 2020. 2
- [24] Martin Piala and Ronald Clark. Terminerf: Ray termination prediction for efficient neural rendering. In *3DV*, 2021. 3
- [25] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2
- [26] Gilles Rainer, Abhijeet Ghosh, Wenzel Jakob, and Tim Weyrich. Unified neural encoding of BTFs. *Computer Graphics Forum (Proc. Eurographics)*, 2020. 2
- [27] Gilles Rainer, Wenzel Jakob, Abhijeet Ghosh, and Tim Weyrich. Neural btf compression and interpolation. In *Computer Graphics Forum*, 2019. 2
- [28] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. In *CVPR*, 2021. 6
- [29] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, 2021. 1, 2, 3, 5, 9

- [30] Peiran Ren, Jiaping Wang, Minmin Gong, Stephen Lin, Xin Tong, and Baining Guo. Global illumination with radiance regression functions. *ACM Transactions on Graphics (TOG)*, 2013. 2
- [31] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *NeurIPS*, 2020. 2
- [32] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *NeurIPS*, 2020. 3
- [33] Vincent Sitzmann, Semon Rezchikov, William T. Freeman, Joshua B. Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. In *NeurIPS*, 2021. 2, 8
- [34] Vincent Sitzmann, Michael Zollhoefer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. *NeurIPS*, 2019. 2, 6
- [35] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *CVPR*, 2021. 2
- [36] R Szeliski. Computer vision: Algorithms and applications. *Instructor*, 2019. 1
- [37] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *CVPR*, 2021. 8
- [38] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. In *NeurIPS*, 2020. 3
- [39] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *ICCV*, 2021. 6, 8
- [40] 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. *NeurIPS*, 2021. 2
- [41] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, 2021. 6
- [42] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2, 5, 9
- [43] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 2, 9
- [44] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. 2, 8
- [45] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv:2010.07492*, 2020. 2, 8

## A. Appendix

### A.1. Detailed Quantitative Results

As shown Tab. 7 and Tab. 8, we provide additional quantitative results for both Realistic Synthetic 360° scenes at  $800 \times 800$  and Real Forward-Facing scenes at  $1008 \times 756$  with three evaluation metrics of PSNR, SSIM and LPIPS. Our NeuSample consistently outperforms NeRF [17] with only 76% computation cost. With fewer samples for rendering, NeuSample still shows promising rendering quality with a significantly better trade-off between computation budget and quality than compared methods [1, 11].

Table 7. Per-scene quantitative comparisons on the Realistic Synthetic 360° dataset.

Method	Inf. Cost	Average	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship
<b>PSNR ↑</b>										
NeRF [17]	1.00	31.01	33.00	<b>25.01</b>	30.13	36.18	32.54	29.62	<b>32.91</b>	28.65
NeuSample	<b>0.76</b>	<b>31.15</b>	<b>33.02</b>	24.99	<b>30.72</b>	<b>36.29</b>	<b>33.17</b>	<b>29.66</b>	32.68	<b>28.65</b>
AutoInt [11] ( $N = 32$ )	0.31	26.83	25.82	22.02	25.51	31.84	27.26	<b>28.58</b>	28.42	25.18
NeuSample ( $N_e = 64$ )	<b>0.25</b>	<b>28.39</b>	<b>29.96</b>	<b>23.43</b>	<b>27.53</b>	<b>34.41</b>	<b>29.14</b>	27.76	<b>29.42</b>	<b>25.47</b>
<b>SSIM ↑</b>										
NeRF [17]	1.00	0.947	0.967	<b>0.925</b>	0.964	0.974	0.961	0.949	0.980	0.856
NeuSample	<b>0.76</b>	<b>0.949</b>	<b>0.968</b>	0.924	<b>0.968</b>	<b>0.977</b>	<b>0.965</b>	0.949	0.980	<b>0.863</b>
AutoInt [11] ( $N = 32$ )	0.31	<b>0.927</b>	0.926	0.885	0.926	<b>0.973</b>	0.929	<b>0.953</b>	0.951	<b>0.869</b>
NeuSample ( $N_e = 64$ )	<b>0.25</b>	0.923	<b>0.941</b>	<b>0.897</b>	<b>0.942</b>	0.966	<b>0.931</b>	0.927	<b>0.964</b>	0.819
<b>LPIPS ↓</b>										
NeRF [17]	1.00	0.081	0.046	0.091	0.044	0.121	0.050	<b>0.063</b>	0.028	0.206
NeuSample	<b>0.76</b>	<b>0.068</b>	<b>0.045</b>	0.091	<b>0.036</b>	<b>0.043</b>	<b>0.048</b>	0.073	<b>0.027</b>	<b>0.183</b>
AutoInt [11] ( $N = 32$ )	0.31	0.152	0.149	0.209	0.109	0.088	0.135	0.100	0.127	0.295
NeuSample ( $N_e = 64$ )	<b>0.25</b>	<b>0.105</b>	<b>0.077</b>	<b>0.133</b>	<b>0.073</b>	<b>0.067</b>	<b>0.097</b>	<b>0.095</b>	<b>0.057</b>	<b>0.238</b>

\* “Inf. Cost” denotes the relative inference cost compared with NeRF [17], *i.e.* time for rendering one image which is measured on one V100 GPU.

\*  $N_e$  of NeuSample denotes the sample number of the extracted sample field.

\* “ $N = 32$ ” for AutoInt [11] denotes the number of piecewise sections.

Table 8. Per-scene quantitative comparisons on the Real Forward-Facing dataset.

Method	Inf. Cost	Average	Fern	Flower	Fortress	Horns	Leaves	Orchids	Room	T-Rex
<b>PSNR ↑</b>										
NeRF [17]	1.00	26.50	<b>25.17</b>	27.40	31.16	27.45	20.92	20.36	32.70	26.80
NeRF-ID <sup>†</sup> [1]	>1.00 <sup>‡</sup>	26.76	25.01	27.85	<b>31.51</b>	27.88	21.09	<b>20.38</b>	32.93	27.45
NeuSample	0.76	<b>26.83</b>	24.99	<b>28.14</b>	31.26	<b>28.32</b>	<b>21.10</b>	20.08	<b>33.26</b>	<b>27.46</b>
NeuSample ( $N_e = 64$ )	<b>0.25</b>	26.50	24.77	28.03	31.09	27.45	21.06	20.03	32.71	26.82
<b>SSIM ↑</b>										
NeRF [17]	1.00	0.811	0.792	0.827	0.881	0.828	0.690	<b>0.641</b>	0.948	0.880
NeRF-ID <sup>†</sup> [1]	>1.00 <sup>‡</sup>	0.822	<b>0.800</b>	0.840	<b>0.890</b>	0.840	<b>0.710</b>	0.640	0.950	<b>0.900</b>
NeuSample	0.76	<b>0.823</b>	0.798	<b>0.845</b>	0.888	<b>0.859</b>	0.708	0.630	<b>0.958</b>	0.899
NeuSample ( $N_e = 64$ )	<b>0.25</b>	0.811	0.785	0.840	0.882	0.826	0.702	0.623	0.950	0.882
<b>LPIPS ↓</b>										
NeRF [17]	1.00	0.250	0.280	0.219	0.171	0.268	0.316	<b>0.321</b>	0.178	0.249
NeuSample	0.76	<b>0.231</b>	<b>0.272</b>	<b>0.192</b>	<b>0.158</b>	<b>0.218</b>	<b>0.296</b>	0.331	<b>0.154</b>	<b>0.224</b>
NeuSample ( $N_e = 64$ )	<b>0.25</b>	0.251	0.289	0.201	0.167	0.261	0.306	0.361	0.173	0.248

<sup>†</sup> denotes training with a larger batch size, *i.e.* 66k in NeRF-ID [1].

<sup>‡</sup> Inference cost for NeRF-ID includes the additional proposal network.