

SteerNeRF: Accelerating NeRF Rendering via Smooth Viewpoint Trajectory

Sicheng Li Hao Li Yue Wang Yiyi Liao* Lu Yu**

Zhejiang University

Abstract

Neural Radiance Fields (NeRF) have demonstrated superior novel view synthesis performance but are slow at rendering. To speed up the volume rendering process, many acceleration methods have been proposed at the cost of large memory consumption. To push the frontier of the efficiency-memory trade-off, we explore a new perspective to accelerate NeRF rendering, leveraging a key fact that the viewpoint change is usually smooth and continuous in interactive viewpoint control. This allows us to leverage the information of preceding viewpoints to reduce the number of rendered pixels as well as the number of sampled points along the ray of the remaining pixels. In our pipeline, a low-resolution feature map is rendered first by volume rendering, then a lightweight 2D neural renderer is applied to generate the output image at target resolution leveraging the features of preceding and current frames. We show that the proposed method can achieve competitive rendering quality while reducing the rendering time with little memory overhead, enabling 30FPS at 1080P image resolution with a low memory footprint.

1. Introduction

Novel View Synthesis (NVS) is a long-standing problem in computer vision and computer graphics with many applications in navigation [39], telepresence [58], and free-viewpoint video [49]. Given a set of posed images, the goal is to render the scene from unseen viewpoints to enable viewpoint control interactively.

With its ability to render high-fidelity images at novel viewpoints, Neural Radiance Fields (NeRF) have recently emerged as a popular representation for NVS. NeRF represents a scene as a continuous function, parameterized by a multilayer perceptron (MLP), that maps a continuous 3D position and a viewing direction to a density and view-dependent radiance [22]. A 2D image is then obtained via volume rendering, i.e., accumulating colors along each ray.

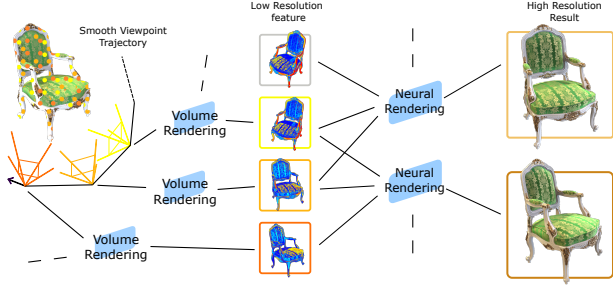


Figure 1. **Illustration.** We exploit smooth viewpoint trajectory to accelerate NeRF rendering, achieved by performing volume rendering at a low resolution and recovering the target image guided by multiple viewpoints. Our method enables fast rendering with a low memory footprint.

However, the rendering process of NeRF is relatively slow as the MLP needs to be queried at millions of samples to render a single image, preventing NeRF from interactive view synthesis. Many recent works have focused on improving the rendering speed of NeRF, yet there is a trade-off between rendering speed and memory cost. State-of-the-art acceleration approaches typically achieve fast rendering at the expense of large memory consumption [13] [54], e.g., by pre-caching the intermediate output of the MLP, leading to hundreds of megabytes to represent a single scene. While there are some attempts to accelerate NeRF rendering with a low memory footprint [16] [24], the performance has yet to reach cache-based methods. In practice, it is desired to achieve faster rendering at a lower memory cost.

To push the frontier of this trade-off, we propose to speed up NeRF rendering from a new perspective, leveraging the critical fact that the viewpoint trajectory is usually smooth and continuous in interactive control. Different from existing NeRF acceleration methods that reduce the rendering time of each viewpoint *individually*, we accelerate the rendering by exploiting the information overlap between *multiple* consecutive viewpoints.

Fig. 1 illustrates our SteerNeRF, a simple yet effective framework leveraging the Smooth ViEwpoint trajEctoRy to speed up NeRF rendering. Exploiting preceding view-

* Corresponding author. ** Co-corresponding author.

points, we can accelerate volume rendering by reducing the number of sample points and maintain the image fidelity using efficient 2D neural rendering.

More specifically, our method consists of a rendering buffer, neural feature fields, and a lightweight 2D neural renderer. At a given viewpoint, we first render a low-resolution feature map via volume rendering. The sampling range along each ray is reduced by fetching a depth map from the rendering buffer and projecting it to the current view. This effectively reduces the volume rendering computation as both the number of pixels and the number of samples for the remaining pixels are reduced. Next, we combine preceding and current feature maps to recover the image at the target resolution using a 2D neural renderer, i.e., a 2D convolutional neural network. The neural feature fields and the 2D neural renderer are trained jointly in an end-to-end manner.

The combination of the low-resolution volume rendering and high-resolution neural rendering leads to lower rendering time compared to directly performing volume rendering at a high resolution. It maintains high fidelity and temporal consistency at a low memory cost.

Our method is inspired by existing super-sampling methods for real-time rendering [50]. In contrast to these methods, we consider the more challenging NVS task without access to the perfect geometry. In this case, we demonstrate joint training of the neural feature fields and the neural renderer is crucial to achieve high image fidelity.

We summarize our contributions as follows.

- We provide a new perspective on NeRF rendering acceleration based on the assumption of smooth viewpoint trajectory. Our method is orthogonal to existing NeRF rendering acceleration methods and can be combined with existing work to achieve real-time rendering at a low memory footprint.
- To fully exploit information of preceding viewpoints, we propose a simple framework that combines low-resolution volume rendering and high-resolution 2D neural rendering. With end-to-end joint training, the proposed framework maintains high image fidelity.
- Our experiments on synthetic and real-world datasets show that our method achieved a rendering speed of nearly 100 FPS at an image resolution of 800×800 pixels and 30 FPS at 1920×1080 pixels. It is faster than other low-memory NeRF acceleration methods and narrows the speed gap between low-memory and cache-based methods.

2. Related Work

Advances in NeRF: Neural radiance fields [22] have received significant attention with photorealistic novel view

synthesis performance. Meanwhile, the vanilla NeRF has several limitations.

Many works have been conducted to address the limitations of NeRF, including unseen scene generalization [6] [55] [20] [45], dynamic scene representation [33] [17] [29] [30] [31] [28] [18], sparse view training [25] [8], surface reconstruction [52] [44] [51] [27], and training acceleration [53] [38] [5]. In addition to representing a single scene, NeRF is also demonstrated to have many applications in generative modeling [36] [26] [4] [12] [3] [9], and robotics [59] [37].

Another important question is accelerating the inference time of NeRF, which is critical for practical applications, e.g., interactive viewing control. Thus, many works focus on rendering acceleration.

NeRF Rendering Acceleration: Existing NeRF Rendering acceleration methods could be categorized into two groups: one reduces the computation at each sample point, and one reduces the number of sampling points.

In the first category, one line of works reduce the computation by pre-caching the intermediate output of the MLP [13] [11] [40] [56] [47] [54] [43] [14] [7] or completely omit the network by representing the scene using a voxel grid [38] [53]. During rendering, these methods retrieve the pre-stored information directly from the table instead of querying a deep network, thus accelerating the rendering. Another line of works reduce the computation at each sample by replacing the large MLPs with smaller ones [34] [10] [48], and maintains the image fidelity by using many small MLPs to represent a single scene. Despite achieving fast rendering, all these methods scarify memory over time, indicating the trade-off of rendering speed and memory cost.

The second category of methods speeds up rendering without increasing the memory cost. The core idea is to adaptively allocate different number of sampling points on each ray based on the content, thus reducing the number of sample points and accelerating rendering. Existing works in this area demonstrate that the number of sampling points can be effectively reduced while maintaining competitive rendering quality [24] [16] [32]. Neural light fields based methods [42] [1], which are equivalent to reducing the number of sampling points to one, also fall into this category. However, simply reducing the number of sampling points has not yet achieved the same level of speed as tabulation-based methods.

Our method is compatible with these two types of methods as we reduce the time from a new perspective, i.e., by combining low-resolution volume rendering and high-resolution neural rendering aided by preceding frames. The low-resolution volume rendering of our method can benefit from existing acceleration methods. The additional memory cost of our method is small as our neural ren-

derer is lightweight. More importantly, when combined with tabulation-based approaches for volume rendering, the resolution of voxel grid for pre-caching could be reduced sufficiently, since there is no need to render high-resolution content during the volume rendering stage.

Superresolution for Rendered Content: A few research works exist that focus on leveraging superresolution techniques for rendering acceleration, especially content rendered from game engine. One of the representative work is NVIDIA DLSS [2]. DLSS is a technology tailored to increase the rendering frame rate while maintaining high fidelity. Even though its technical details have not been fully disclosed, in general, DLSS collects the raw low-resolution texture input, accurate motion vectors, and depth buffers from the rendering engine and outputs high-quality full-resolution content. Xiao et al. [50] proposed a method similar to DLSS. They achieve a new state of the art in super-resolving rendered videos with extreme aliasing by using a new temporal super resolution design, which takes texture, depth, motion vector from the game engine as input as well. Different from above works, our method addresses the more challenging NVS task, and thus cannot obtain high precision motion vector for superresolution. Instead, we take a volume-rendered noisy depth map to warp the preceding frame to align with the current frame. Moreover, end-to-end joint training makes feature maps more expressive and enables 2D neural renderer to hallucinate high-fidelity images from current and warped feature. In the field of NeRF, NeRF-SR [41] is proposed to generate higher resolution images with low-resolution supervisions leveraging the sub-pixel information across different views. Note that NeRF-SR is not applicable for real-time rendering.

3. Method

In this work, we propose fully exploiting the smoothly changing viewpoints to accelerate the rendering process of NeRF. In general, we achieve rendering acceleration by reducing the total number of 3D points that need to be queried in volume rendering for each frame.

Fig. 2 gives an overview of our proposed pipeline consisting of a rendering buffer, neural feature fields, and a 2D neural renderer. Specifically, the rendering buffer saves low-resolution feature maps and depth maps of previous viewpoints. At the current viewpoint, a low-resolution feature map and depth map is rendered accelerated by the rendering buffer. Next, the lightweight neural renderer takes the preceding and the current feature maps as input to generate the output image at the target resolution.

In the following, we first introduce preliminaries of NeRF model in Section 3.1. Next, we present the accelerated volume rendering in Section 3.2, the buffer-guided neural rendering in Section 3.3, and the training procedure

in Section 3.4. Finally, we describe implementation details in Section 3.5.

3.1. Background

NeRF represents a scene as a continuous function f_θ parameterized by learnable parameters θ that maps a 3D point $\mathbf{x} \in \mathbb{R}^3$ and a viewing direction $\mathbf{d} \in \mathbb{S}^2$ and to a volume density σ and a color value \mathbf{c} :

$$f_\theta : (\mathbf{x} \in \mathbb{R}^3, \mathbf{d} \in \mathbb{S}^2) \mapsto (\sigma \in \mathbb{R}^+, \mathbf{c} \in \mathbb{R}^3) \quad (1)$$

Given a target viewpoint, the color \mathbf{c}_r and depth d_r at a camera ray r is obtained via volume rendering integral approximated by the numerical quadrature [21]:

$$\mathbf{c}_r = \sum_{i=1}^N T_r^i \alpha_r^i \mathbf{c}_r^i \quad d_r = \sum_{i=1}^N T_r^i \alpha_r^i t_r^i \quad (2)$$

$$\alpha_r^i = 1 - \exp(-\sigma_r^i \delta_r^i) \quad T_r^i = \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (3)$$

where T_r^i and α_r^i denote transmittance and alpha value of a sample point \mathbf{x}_i .

Rendering Time: The rendering time of NeRF is proportional to the amount of computation required to render an image. Let $H \times W$ denote the target image resolution, N the number of samples on each ray and F the FLOPs of querying one sample’s color and density. We can roughly estimate the amount of computation as $H \times W \times N \times F$. Existing NeRF acceleration strategies mainly focus on how to decrease the FLOPs F of each query by pre-caching the output of f_θ or directly use a voxel grid [13] [11] [40] [54] [53], thus leading to large memory consumption. There are a few attempts to reduce the sampling points N along the ray via early ray termination, empty space skipping [47] or adaptive sampling [16], yet using these techniques alone has not yet reached the performance of pre-cache based methods.

Our solution can elevate rendering speed from a new perspective, that is, reduce the number of pixels $H \times W$ and the number of samples N for volume rendering via fully utilizing the smooth viewpoint trajectory. Besides, our solution can cooperate with the existing work to achieve high-speed rendering, leading to a higher rendering framerate while maintaining visual quality.

3.2. Accelerating Volume Rendering

We propose to learn a neural feature fields that renders a low-resolution feature map suited for the subsequent neural renderer. The acceleration of our framework comes from 1) rendering the feature map at a lower resolution and 2) reducing the sampling range guided by the rendering buffer.

Low-Resolution Feature Rendering: We render a feature map via volume rendering. Our neural feature fields maps

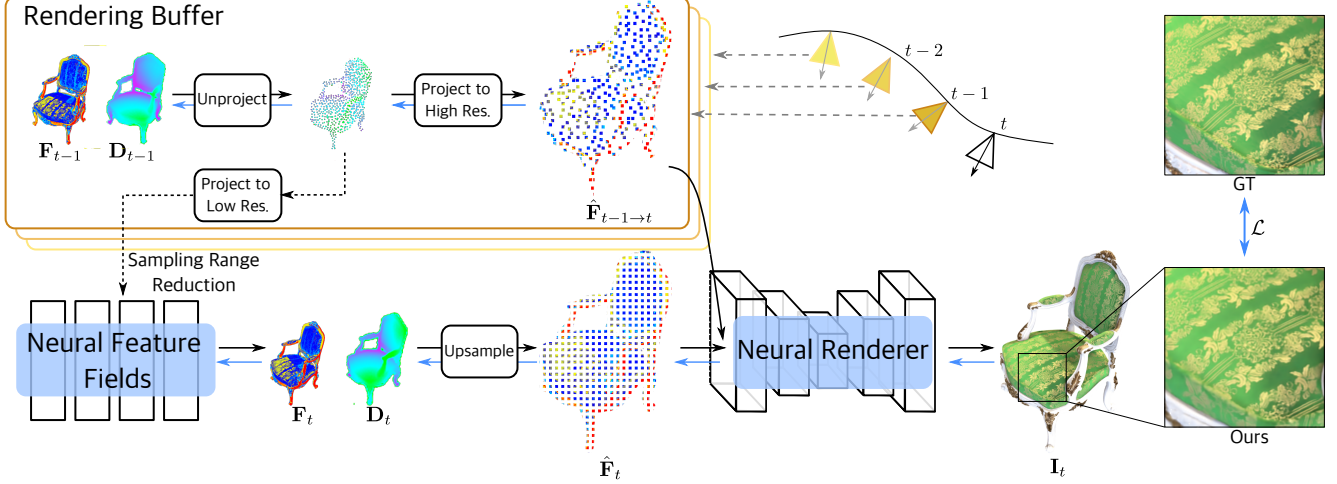


Figure 2. **SteerNeRF**. The rendering buffer saves low-resolution feature maps $\{\mathbf{F}_{t-L}, \dots, \mathbf{F}_{t-1}\}$ and depth maps $\{\mathbf{D}_{t-L}, \dots, \mathbf{D}_{t-1}\}$ of previous L viewpoints. At the current viewpoint t , a low-resolution feature map \mathbf{F}_t and a depth map \mathbf{D}_t are rendered accelerated by the rendering buffer. Next, the lightweight neural renderer takes as input the reprojected features maps at the high resolution $\{\hat{\mathbf{F}}_{t-L \rightarrow t}, \dots, \hat{\mathbf{F}}_{t-1 \rightarrow t}\}$ and the upsampled feature map $\hat{\mathbf{F}}_t$ to generate the output image \mathbf{I}_t . As illustrated by the blue arrows, during training, we apply the reconstruction loss \mathcal{L} to an image patch and jointly optimize the entire model in an end-to-end manner, including preceding frames in the rendering buffer.

the input \mathbf{x} and the viewing direction \mathbf{d} to a density value and a feature vector \mathbf{f} :

$$f_\theta : (\mathbf{x} \in \mathbb{R}^3, \mathbf{d} \in \mathbb{S}^2) \mapsto (\sigma \in \mathbb{R}^+, \mathbf{f} \in \mathbb{R}^K) \quad (4)$$

where K is the number of channels of our feature vector. Despite providing more information, rendering extra channels leads to a very little overhead in time as we only expand the last layer of the MLP to predict more channels.

We can obtain a feature vector \mathbf{f}_r at each ray r via volume rendering.

$$\mathbf{f}_r = \sum_{i=1}^N T_r^i \alpha_r^i \mathbf{f}_r^i \quad (5)$$

We render the feature vector at a subset of the rays, yielding a feature map $\mathbf{F} \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4} \times K} = \{\mathbf{f}_r\}$. The corresponding low-resolution depth value $\mathbf{D} \in \mathbb{R}^{\frac{H}{4} \times \frac{W}{4}}$ is also rendered. Both the feature map \mathbf{F} and \mathbf{D} are stored into our rendering buffer for subsequent frames.

Buffer-Guided Sampling Range Reduction: The depth information of adjacent frames rendered previously provides a coarse scene geometry. Thus, warped depth from the preceding frame could be used as guidance to determine sampling positions and thus accelerate rendering. Specifically, given a low-resolution depth map $\mathbf{D}_{t', t' < t}$ of the previous frame, it is first unprojected to a 3D point cloud and then projected to the current viewpoint. More formally, the following unprojection and projection functions are applied

to each pixel (u, v) of $\mathbf{D}_{t'}$ with depth $d_{t'}$:

$$\mathbf{p} = \xi_{t'}^{-1} \mathbf{K}_l^{-1} d_{t'} [u, v, 1]^T \quad (6)$$

$$d_{t' \rightarrow t} [u_{t' \rightarrow t}, v_{t' \rightarrow t}, 1]^T = \mathbf{K}_l \xi_t \mathbf{p} \quad (7)$$

where \mathbf{p} denotes a 3D point and \mathbf{K}_l is the intrinsic matrix of the low-resolution image. This yields the reprojected depth map $\mathbf{D}_{t' \rightarrow t}$ where

$$\mathbf{D}_{t' \rightarrow t} (\lfloor u_{t' \rightarrow t} \rfloor, \lfloor v_{t' \rightarrow t} \rfloor) = d_{t' \rightarrow t}. \quad (8)$$

Note that here we simply round $(u_{t' \rightarrow t}, v_{t' \rightarrow t})$ and observe negligible impact on the performance. Given $\mathbf{D}_{t' \rightarrow t}$, the sampling range at frame t can be limited to the depth interval $[\mathbf{D}_{t' \rightarrow t} - \epsilon, \mathbf{D}_{t' \rightarrow t} + \epsilon]$ for rendering \mathbf{F}_t and \mathbf{D}_t at the camera viewpoint ξ_t . This simple strategy further decreases the number of 3D sample points and accelerates the rendering of the low-resolution feature map.

3.3. Buffer-Guided Neural Rendering

Given the rendering buffer consisting of L preceding feature maps $\{\mathbf{F}_{t-L}, \dots, \mathbf{F}_{t-1}\}$ and depth maps $\{\mathbf{D}_{t-L}, \dots, \mathbf{D}_{t-1}\}$, we combine them with the feature map at the current viewpoint t to recover the output image \mathbf{I}_t using a 2D neural renderer. We first project frames in the rendering buffer to the current viewpoint. Next, we use a 2D neural renderer to recover the target image.

Preceding Frames Reprojection: Inspired by natural video superresolution approaches, our method warps previous frames to align with the current frame to ease the task

of the subsequent neural renderer. Instead of reprojecting the depth map to the low-resolution image as in Eq. 7, we directly project the 3D point cloud to the target resolution to achieve higher precision, i.e., maintain sub-pixel precision in terms of the low-resolution image:

$$d_{t' \rightarrow t} [\hat{u}_{t' \rightarrow t}, \hat{v}_{t' \rightarrow t}, 1]^T = \mathbf{K}_h \boldsymbol{\xi}_t \mathbf{p} \quad (9)$$

where \mathbf{K}_h denotes the intrinsic matrix of the high-resolution image. This allows us to obtain the reprojected high-resolution feature map $\hat{\mathbf{F}}_{t' \rightarrow t} \in \mathbb{R}^{H \times W \times K}$:

$$\hat{\mathbf{F}}_{t' \rightarrow t}([\hat{u}_{t' \rightarrow t}], [\hat{v}_{t' \rightarrow t}]) = \mathbf{F}_{t'}(u, v). \quad (10)$$

Neural Renderer: We use a lightweight 2D convolutional network for fast inference. The reprojected high-resolution feature maps $\{\hat{\mathbf{F}}_{t' \rightarrow t}\}$ are concatenated with the upsampled feature map $\hat{\mathbf{F}}_t$ and mapped to the output target image:

$$g_\theta : (\{\hat{\mathbf{F}}_{t' \rightarrow t}\}, \hat{\mathbf{F}}_t) \mapsto \mathbf{I}_t \in \mathbb{R}^{H \times W \times 3} \quad (11)$$

In practice, we choose a simple modified U-Net as our neural renderer. Compared to traditional U-Net, we reduce the number of convolution layers for high-resolution features and increase the depth of convolution layers for low-resolution features. The simple adjustment allows us to greatly reduce the inference time and keep visual quality when the number of parameters is almost the same. We use an off-the-shelf inference acceleration toolbox, NVIDIA TensorRT, to optimize neural renderer to reduce the inference time.

3.4. Training

The training strategy is crucial to achieving high-quality novel view synthesis. In practice, we first pre-train our neural feature fields and then train the full model jointly in an end-to-end training fashion.

Pre-training: We pre-train our neural feature fields on the target resolution $H \times W$. Here, we render a high-resolution feature map $\hat{\mathbf{F}} \in \mathbb{R}^{H \times W \times K}$ and apply an $L2$ reconstruction loss on the first three channels supervised by the high-resolution ground truth image. We leave other output channels without constraint of supervision as pretraining on the first three channels is sufficient to learn reasonable volume density.

End-to-end Joint Training: With the pre-trained neural feature fields, we train our full model in an end-to-end manner using an $L2$ loss \mathcal{L} on the final output image $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$. Note that we do not apply reconstruction loss to the rendered feature map during end-to-end training and let the neural feature fields learn features suited for the 2D neural renderer. During training, the loss \mathcal{L} is applied to image patches. As the training viewpoints are scattered in

the space without a smooth trajectory, we generate a short sequence of preceding camera poses for each training image to train the buffer-based neural renderer. Note that this process does not introduce additional supervision as we only apply the loss \mathcal{L} to the training viewpoints despite taking preceding feature maps as input.

3.5. Implementation Details

Network Architecture: Our method is compatible with different NeRF approaches for learning the neural feature fields. In this work, we implement our neural feature fields based on Instant-NGP [23] using a third-party PyTorch implementation¹. This allows more efficient feature map rendering than the vanilla NeRF. We follow the original architecture of Instant-NGP that uses multi-resolution hash tables where the table length at each resolution is fixed to 2^{19} . Following Instant-NGP, empty space skipping and early ray termination are applied when rendering the low-resolution feature map. Regarding the 2D neural renderer, we adopt a shallow U-Net [35] with the detailed architecture described in the supplementary.

Distillation: When the number of training views is relatively small, the 2D neural renderer tends to overfit the training views, yielding degenerated performance on the test poses. In this case, we leverage a pre-trained NeRF model to synthesize more viewpoints as our pseudo ground truth by randomly sampling viewpoints within the available viewing zone. Adding the randomly sampled pseudo ground truth alleviates the overfitting problem.

Inference optimization: Optimizing trained neural renderer for real-time inference and lower memory footprint is necessary. Thus, we leverage NVIDIA TensorRT to optimize 2D neural renderer. Prior to testing, we optimize 2D neural renderer into two versions in FP16 and INT8 precision separately.

4. Experiments

In this section, we evaluate the performance of our method. We first quantitatively compare our work with prior work. Next, we report the runtime breakdown of our method in two representative scenes. Finally, we validate our design decisions with extensive ablation studies.

Datasets: We take NeRF-Synthetic dataset, and a subset of the Tanks & Temples dataset [15] for performance evaluation. NeRF-Synthetic dataset contains 8 synthetic scenes rendered by Blender at a resolution of 800×800 and full of high frequency texture. Tanks & Temples dataset is a real-world dataset at a resolution of 1920×1080 . We follow the

¹<https://github.com/kwea123/ngp-pl>

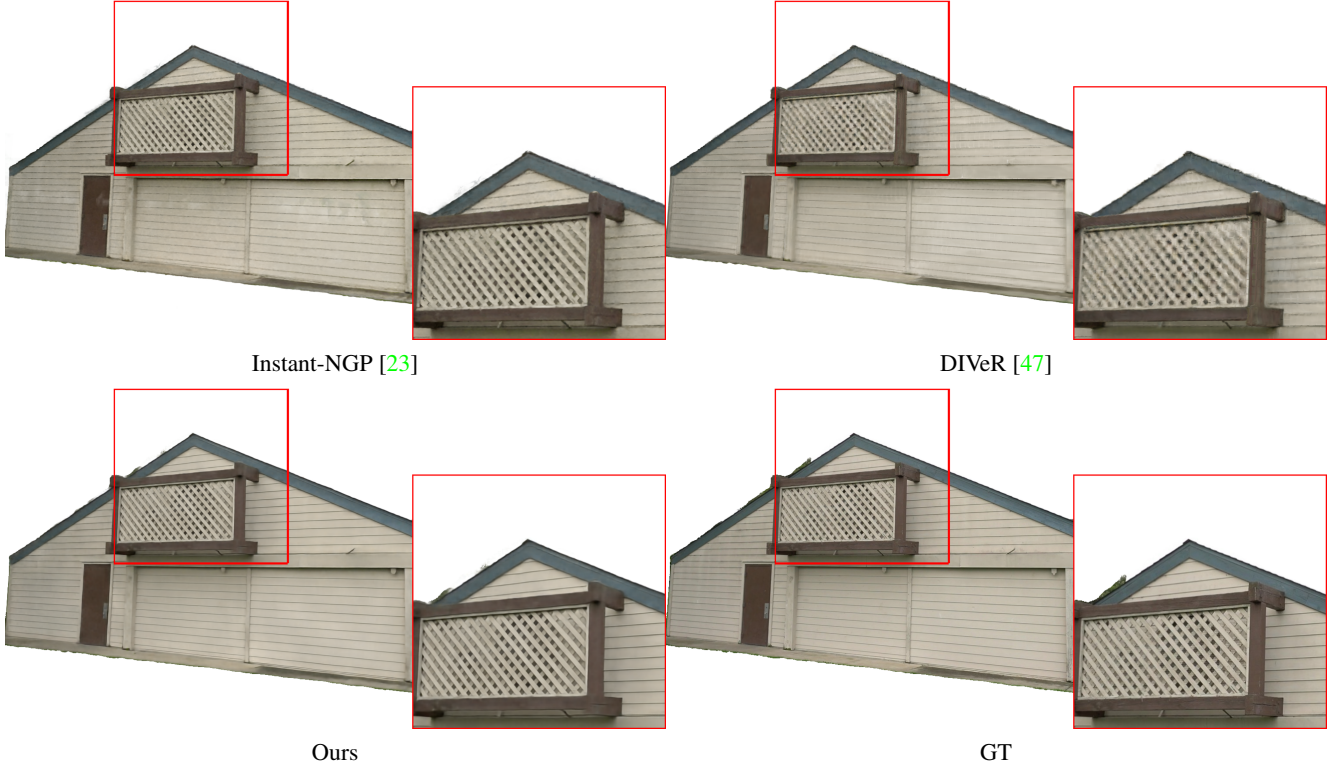


Figure 3. **Qualitative Comparison on Tanks & Temples.**

Method	Tanks & Temples				NeRF-Synthetic				Mem.(MB)↓
	PSNR(dB)↑	SSIM↑	LPIPS↓	FPS↑	PSNR(dB)↑	SSIM↑	LPIPS↓	FPS↑	
NeRF [22]	28.32	0.890	0.198	0.005	31.01	0.947	0.081	0.02	5
NSVF [19]	28.40	0.900	0.153	0.06	31.74	0.953	0.047	0.23	-
KiloNeRF [34]	28.41	0.900	0.092	10.95	31.00	0.950	0.030	38.50	161
PlenOctree [54]	-	-	-	-	31.71	0.958	0.053	167.7	1930
DIVEr [47]	28.18	0.912	0.116	-	<u>32.12</u>	0.958	0.033	74.00	68
Instant-NGP [23]	28.77	0.918	0.136	5.00	32.79	<u>0.957</u>	0.055	60.00	<u>25.2</u>
Ours (FP16)	<u>28.65</u>	0.924	0.121	<u>27.24</u>	31.60	0.954	0.058	75.19	29.4
Ours (INT8)	28.44	<u>0.919</u>	0.129	30.90	30.97	0.948	0.065	<u>86.97</u>	27.3

Table 1. **Quantitative results on Tanks & Temples and NeRF-Synthetic** show that our method could achieve high framerate rendering while keeping relatively low memory footprint. (**Best**, Second Best)

subset selection and crop the backgrounds of images as in NSVF [19].

Baselines: The baseline methods could be categorized into two groups: one is classical but unable to achieve real-time rendering, and another is capable of real-time rendering. As for classical methods, we compare with the original NeRF and NSVF. As for real-time rendering methods, we compare with PlenOctree, DIVEr, and Instant-NGP. Note that we take DIVEr32 (RT) for a fair comparison, a real-time version of DIVEr with better speed-quality trade-off.

Metrics: We evaluate our method from efficiency, quality, and memory usage. The efficiency is measured by the number of frames per second. The quality is measured by PSNR, SSIM [46] and LPIPS [57]. The memory usage is measured by megabytes.

4.1. Comparisons to Baseline

We first report the results on Tanks & Temples dataset in Tab. 1 (left). All baseline and our methods show similar visual quality. However, our method shows the best FPS, almost three times higher than the best baseline method, Kilo-



Figure 4. **Qualitative Comparison on NeRF-Synthetic.**

NeRF. In the meantime, the memory usage of our method is less than 20% of KiloNeRF. Instant-NGP, which has similar memory usage with ours, can only render views at 5 FPS. PlenOctree is not compared here as we fail to re-train the model to achieve reasonable results. NeRF and NSVF, both employing a large MLP, perform much slower in rendering efficiency. With such a trade-off between memory usage and rendering efficiency, our method satisfies 3D interaction along smooth viewport trajectories at 1080P resolution with 30FPS. We show corresponding qualitative results on Tanks & Temples in Fig. 3.

On NeRF-Synthetic dataset in Tab. 1 (right), all methods still show similar visual quality scores. For FPS, PlenOctree becomes the best with one to two orders of magnitude higher memory usage, which is mainly caused by its network-free nature. Among the others, all the methods keep the same memory usage, and our method still shows a superior FPS. There are two differences from the previ-

ous dataset in FPS. First, KiloNeRF is slower than Instant-NGP. Second, the margin of our method against the others is smaller. We explain both differences by the bandwidth of GPU, which may not be fully maximized due to the lower image resolution i.e. fewer times of ray marching or neural network inference. In summary, we consider that the parallelism of our method is able to push the frontier of the efficiency-memory trade-off, especially for high resolution rendering. We present corresponding qualitative results on NeRF-Synthetic in Fig. 4.

4.2. Runtime Breakdown

We report the average runtime of our method in Tab. 2 for Chair and Barn scenes, including the runtime of volume rendering, neural rendering and depth reprojection. For volume rendering, $\{\mathbf{D}_{t' \rightarrow t}\}_{t'=t-1}$ refers to the depth reprojection to enables buffer-guided sampling range reduction. When this function is turned on, it reduces the volume

Module	Time(ms)			
	Barn		Chair	
$\{\mathbf{D}_{t' \rightarrow t}\}_{t'=t-1}$	0	0.11	0	0.11
f_θ	20.28	17.15	4.81	3.76
$\{\hat{\mathbf{F}}_{t' \rightarrow t}\}_{t'=\{t-2, t-1\}}$	3.21		2.40	
g_θ	12		3.71	
Total	35.49	32.47	10.92	9.98

Table 2. **Runtime breakdown** for Chair and Barn scenes

L : # Previous frames	0	1	2	3
PSNR(dB)	32.94	33.05	33.14	33.10
SSIM	0.959	0.964	0.968	0.967
LPIPS	0.035	0.034	0.034	0.032
Runtime(ms)	8.46	9.8	11.14	12.48

Table 3. **Comparison of Number of Preceding Frames** on Chair.

rendering time (f_θ) to 85-90% of the original. As for the neural rendering part, the reprojection $\{\hat{\mathbf{F}}_{t' \rightarrow t}\}_{t'=\{t-2, t-1\}}$ takes longer as two frames are reprojected to a higher resolution. The reprojection can be further accelerated in the future work by using custom CUDA kernels.

4.3. Ablation Study

We conduct ablation studies on SteerNeRF using Chair scene from NeRF-Synthetic dataset.

Number of Preceding Frames: We first evaluate the visual quality and render time by taking a varying number of preceding frames L as input to neural renderer in Tab. 3. As more preceding frames are used, the reconstruction quality increases with a growing rendering time. The additional runtime mainly comes from the warping operation of preceding frames. Furthermore, we also observe that as the number of preceding frames increases, the quality gain brought by each additional frame decreases. Therefore, in practical applications, we can adjust this parameter more flexibly according to the scene content to achieve trade-off between quality and rendering efficiency.

Number of Feature Channels: We also verify the effect of the number of channels of feature images. Experiments show that as the number of channels increases, the gain of visual quality decreases gradually in Tab. 4. Therefore, we ended up choosing six channels in total for our implementation. Note that employing feature rendering almost causes no rendering time overhead.

Joint Training: The necessity of joint training is validated in Tab. 5. The method without joint training means the parameters of neural feature fields are frozen when training the neural renderer. Joint training helps neural feature fields

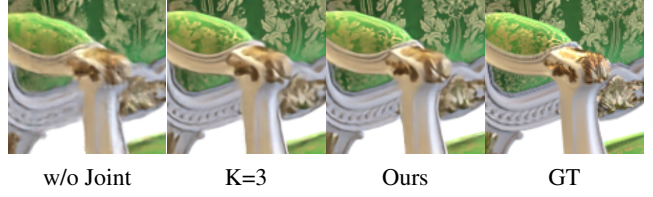


Figure 5. **Ablation study.** Impact of different training configurations has been shown in above visual examples.

K : # Feature channels	3	6	9
PSNR(dB)	32.53	33.14	33.21
SSIM	0.959	0.965	0.966
LPIPS	0.042	0.035	0.034

Table 4. **Comparison of Number of Feature Channels** on Chair.

	PSNR	SSIM	LPIPS
Ours w/ joint training	33.05	0.965	0.038
Ours w/o joint training	31.98	0.957	0.053

Table 5. **Joint training** on Chair.

Length of Hash Table	PSNR(dB)	SSIM	LPIPS	Mem.(MB)
2^{19}	33.05	0.965	0.038	29.2
2^{16}	33.03	0.961	0.039	8.1
2^{14}	32.70	0.953	0.039	5.0

Table 6. **Memory Usage of Neural Feature Fields** on Chair.

generate more expressive features compatible with following neural renderer to synthesize higher-quality textures. In Fig. 5, we show visual examples of different training configurations.

Memory Usage of Neural Feature Fields: In Tab. 6, we compare the visual quality and memory usage when neural feature fields is configured with different length of hash table, i.e., learnable parameters. We find that even when the length of the hash table is greatly reduced, the final visual quality is only slightly affected while the memory usage is significantly reduced.

5. Conclusion

We propose a new perspective for NeRF rendering acceleration by considering the smooth viewpoint trajectory during interaction. The main idea is to supersample the image rendered at the current viewpoint by taking preceding low-resolution features and depths. The experiments show that our method achieves real-time even when rendering 1080P images. As a limitation, training the 2D neural renderer is time-consuming. Moreover, we need to train a specialized neural renderer from scratch for each scene individually. In

the future, we plan to train a neural renderer that generalizes well on different scenes only after short-time fine-tuning or even without fine-tuning.

References

- [1] Benjamin Attal, Jia-Bin Huang, Michael Zollhöfer, Johannes Kopf, and Changil Kim. Learning neural light fields with ray-space embedding. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [2] Andrew Burnes. Nvidia dlss 2.0: A big leap in ai rendering. <https://www.nvidia.com/en-us/geforce/news/nvidia-dlss-2-0-a-big-leap-in-ai-rendering/>, 2020. 3
- [3] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3d generative adversarial networks. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [4] Eric R. Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. Pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [5] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2022. 2
- [6] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2
- [7] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *arXiv.org*, 2022. 2
- [8] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised nerf: Fewer views and faster training for free. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [9] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. GRAM: generative radiance manifolds for 3d-aware image generation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [10] Stefano Esposito, Daniele Baieri, Stefan Zellmann, André Hinkenjann, and Emanuele Rodolà. Kiloneus: Implicit neural representations with real-time global illumination. *arXiv.org*, 2022. 2
- [11] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien P. C. Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2, 3
- [12] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. Stylenerf: A style-based 3d-aware generator for high-resolution image synthesis. *Proc. of the International Conf. on Learning Representations (ICLR)*, 2022. 2
- [13] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul E. Debevec. Baking neural radiance fields for real-time view synthesis. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 1, 2, 3
- [14] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. Efficientnerf efficient neural radiance fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [15] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 5
- [16] Andreas Kurz, Thomas Neff, Zhaoyang Lv, Michael Zollhöfer, and Markus Steinberger. Adanerf: Adaptive sampling for real-time rendering of neural radiance fields. 2022. 1, 2, 3
- [17] Tianye Li, Mira Slavcheva, Michael Zollhöfer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard A. Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [18] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [19] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 6
- [20] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Theobalt Christian, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [21] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):99–108, 1995. 3
- [22] 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 *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 1, 2, 6
- [23] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. on Graphics*, 2022. 5, 6, 7
- [24] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of compact neural radiance fields using depth oracle networks. *Comput. Graph. Forum*, 40(4):45–59, 2021. 1, 2
- [25] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Reg-

- nerf: Regularizing neural radiance fields for view synthesis from sparse inputs. *arXiv.org*, 2021. 2
- [26] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [27] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2
- [28] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2
- [29] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: a higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6):238:1–238:12, 2021. 2
- [30] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2
- [31] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [32] Martin Píala and Ronald Clark. Terminerf: Ray termination prediction for efficient neural rendering. In *Proc. of the International Conf. on 3D Vision (3DV)*. IEEE, 2021. 2
- [33] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural radiance fields for dynamic scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [34] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2, 6
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015. 5
- [36] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: generative radiance fields for 3d-aware image synthesis. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2
- [37] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J. Davison. imap: Implicit mapping and positioning in real-time. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 2
- [38] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [39] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1
- [40] Krishna Wadhvani and Tamaki Kojima. Squeezenerf: Further factorized fastnerf for memory-efficient inference. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [41] Chen Wang, Xian Wu, Yuanchen Guo, Song-Hai Zhang, Yu-Wing Tai, and Shi-Min Hu. Nerf-sr: High quality neural radiance fields using supersampling. In *Communications of the ACM*, pages 6445–6454. ACM, 2022. 3
- [42] Huan Wang, Jian Ren, Zeng Huang, Kyle Olszewski, Menglei Chai, Yun Fu, and Sergey Tulyakov. R2l: Distilling neural radiance field to neural light field for efficient novel view synthesis. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2022. 2
- [43] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenotrees for dynamic radiance field rendering in real-time. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [44] 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. In Marc Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2
- [45] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [46] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. on Image Processing (TIP)*, 13(4):600–612, 2004. 6
- [47] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yu-Xiong Wang, and David A. Forsyth. Diver: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3, 6, 7
- [48] Xiuchao Wu, Jiamin Xu, Zihan Zhu, Hujun Bao, Qixing Huang, James Tompkin, and Weiwei Xu. Scalable neural indoor scene rendering. *ACM Transactions on Graphics (TOG)*, 41(4):1–16, 2022. 2
- [49] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1
- [50] Lei Xiao, Salah Nouri, Matthew Chapman, Alexander Fix, Douglas Lanman, and Anton Kaplanyan. Neural super-

- sampling for real-time rendering. *ACM Trans. Graph.*, 39(4):142, 2020. 2, 3
- [51] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2
 - [52] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. 2020. 2
 - [53] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
 - [54] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2021. 1, 2, 3, 6
 - [55] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
 - [56] Jian Zhang, Jinchi Huang, Bowen Cai, Huan Fu, Mingming Gong, Chaohui Wang, Jiaming Wang, Hongchen Luo, Rongfei Jia, Binqiang Zhao, et al. Digging into radiance grid for real-time view synthesis with detail preservation. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2022. 2
 - [57] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 6
 - [58] Yizhong Zhang, Jiaolong Yang, Zhen Liu, Ruicheng Wang, Guojun Chen, Xin Tong, and Baining Guo. Virtualcube: An immersive 3d video communication system. *IEEE Transactions on Visualization and Computer Graphics*, 28(5):2146–2156, 2022. 1
 - [59] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R. Oswald, and Marc Pollefeys. Nice-slam: Neural implicit scalable encoding for slam. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, June 2022. 2