# Efficient Neural Radiance Fields with Learned Depth-Guided Sampling

Haotong Lin*    Sida Peng*    Zhen Xu    Hujun Bao    Xiaowei Zhou
State Key Lab of CAD&CG, Zhejiang University

## Abstract

*This paper aims to reduce the rendering time of generalizable radiance fields. Some recent works equip neural radiance fields with image encoders and are able to generalize across scenes, which avoids the per-scene optimization. However, their rendering process is generally very slow. A major factor is that they sample lots of points in empty space when inferring radiance fields. In this paper, we present a hybrid scene representation which combines the best of implicit radiance fields and explicit depth maps for efficient rendering. Specifically, we first build the cascade cost volume to efficiently predict the coarse geometry of the scene. The coarse geometry allows us to sample few points near the scene surface and significantly improves the rendering speed. This process is fully differentiable, enabling us to jointly learn the depth prediction and radiance field networks from only RGB images. Experiments show that the proposed approach exhibits state-of-the-art performance on the DTU, Real Forward-facing and NeRF Synthetic datasets, while being at least 50 times faster than previous generalizable radiance field methods. We also demonstrate the capability of our method to synthesize free-viewpoint videos of dynamic human performers in real-time. The code will be available at https://zju3dv.github.io/enerf/.*

## 1. Introduction

Photorealistic novel view synthesis has a variety of applications such as virtual tourism, telepresence, and sports broadcasting. Recently, NeRF [28] represents scenes as density and color fields, which works particularly well with volume rendering techniques. By optimizing the radiance fields from images, it achieves state-of-the-art performance on novel view synthesis. However, to synthesize novel views of a new scene, NeRF requires a time-consuming optimization process. Moreover, the rendering process of NeRF is extremely slow, which takes about 13 seconds to render a $512 \times 640$ image on an RTX 2080Ti GPU.

To avoid the per-scene optimization, recent works [42, 43, 49] introduce CNN-based encoders to extract features
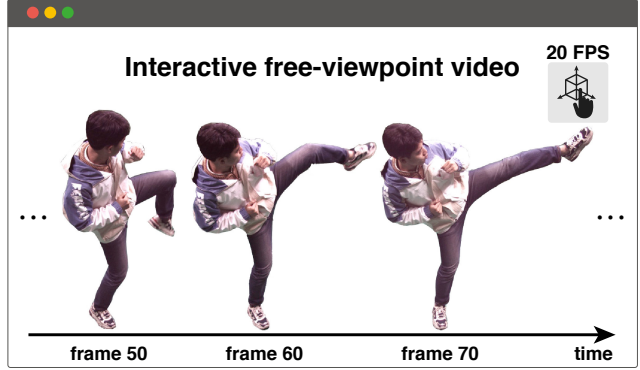
---

*Authors contributed equally



Figure 1. The proposed method achieves real-time photorealistic view synthesis without caching techniques. Please refer to https://zju3dv.github.io/enerf/ for the real-time demos that render free-viewpoint videos of dynamic scenes.

from input views and feed multi-view features into NeRF-like networks to reconstruct radiance fields, which thus can generalize across scenes. However, they still take a large amount of time to render an image at test time. To overcome this problem, some methods [9, 14, 48] propose to cache radiance fields into highly efficient data structures, which avoids the network forwarding during inference. Although this significantly improves the rendering speed, the caching process requires a lot of time, and the cached models are typically very large. For instance, for one frame of a ZJU-MoCap [32] sequence, it takes about 1.5 hours to convert a vanilla trained NeRF model to a PlenOctree [48] model, and the storage space of this PlenOctree model is about 1.45 GB. The expensive time and storage costs prevent caching-based methods from synthesizing free-viewpoint videos of dynamic scenes.

A major reason for the low rendering speed of implicit representations [28, 43, 49] is that they require hundreds of forwarding passes through a neural network for each pixel in the rendered image. In contrast, explicit 3D representations such as volumes and multi-plane images typically require only one pass to render the whole image, which makes the rendering fast. However, explicit representations are discrete and generally have limited resolution, while implicit representations are continuous and therefore have a high resolution to achieve photorealistic rendering [20].

In this paper, we present a hybrid scene representation that combines implicit neural radiance fields with explicit depth maps for fast view synthesis. Our core innovation is introducing explicit depth maps as coarse scene geometry to guide the rendering process. Specifically, for the target view, we construct the cascade cost volume, which is used to predict a depth probability distribution. The depth probability distribution gives an interval along each ray where the surface may locate in. With the depth interval, we only need to sample few 3D points along the ray and thus greatly improves the rendering speed of previous methods [2, 28]. Similar to MVSNeRF [2], we predict radiance fields using the cost volume feature, which makes our method generalize well to new scenes. Moreover, this whole process is fully differentiable, so the depth probability distribution can be jointly learned with NeRF from only RGB images.

We evaluate our approach on the DTU [15], Real Forward-facing [27, 28], NeRF Synthetic [28] and ZJU-MoCap [32] datasets, which are widely-used benchmark datasets for novel view synthesis. Across all datasets, our approach achieves state-of-the-art performance. Meanwhile, our approach significantly improves the rendering speed, which runs at least 50 times faster than previous generalizable radiance field methods [2, 43, 49]. Furthermore, experiments on the ZJU-MoCap [32] dataset prove the efficiency of our method in the storage and caching time compared to caching-based methods [9, 14, 48]. We also show that our approach is able to produce reasonable depth maps by supervising the networks with only images.

In summary, this work has the following contributions:

- We present a hybrid scene representation that combines the best of implicit neural radiance fields and explicit depth maps for fast high-quality view synthesis. It utilizes learned depth-guided sampling to improve the rendering efficiency of NeRF.

- We show that the depth-guided sampling can be jointly learned with NeRF from only RGB images.

- We demonstrate that our approach achieves state-of-the-art rendering quality while being significantly faster than previous methods on several view synthesis benchmarks without caching strategies.

- We demonstrate the capability of our method to synthesize novel views of dynamic scenes in real-time.

## 2. Related work

**Novel view synthesis.** Novel view synthesis is a long-standing problem in computer vision and computer graphics. Early works [7, 10, 18] achieve impressive rendering results based on light field interpolation techniques, which recover plenoptic functions from dense camera views. For

photorealistic rendering of extrapolated views, image-based rendering methods [1, 8, 16, 33, 35] have attempted to leverage explicit depth maps as proxy geometry. Specifically, these methods first reconstruct depth maps of input images using multi-view stereo techniques [37, 38]. Then, they warp input images to the target view and perform image blending. To utilize the power of deep learning, some methods [6, 13, 45] introduce learnable components to infer depth maps or blend images. Although image-based rendering methods have a big range of renderable viewpoints, they tend to be sensitive to the quality of precomputed depth maps. Recently, instead of predicting proxy geometry individually, some methods [19, 21, 22, 39, 41] propose to optimize 3D representations jointly with differentiable renderers from input images. An emerging direction is using MLPs as implicit representations, where the MLP networks map continuous spatial points to target values, including signed distance [30], occupancy [26], and radiance [28]. With differentiable renderers, recent methods [23, 29, 40] optimize implicit representations from images and achieve impressive performance on 3D reconstruction and novel view synthesis. NeRF [28] represents scenes as continuous color and density fields, which yields state-of-the-art rendering results.

**Improving efficiency of neural radiance fields.** In spite of impressive results of NeRF [28], it requires a long per-scene optimization process and its rendering speed is extremely slow. [2, 42, 43, 49] attempt to reduce the training time by generalizing NeRF across scenes. They design 2D CNNs to encode the content of input images and then decode the multi-view features to the target radiance fields through NeRF-like MLP networks. More recently, [24] proposes to consider the visibility of input views when predicting radiance fields, which improves the rendering quality. However, the rendering process of these methods [2, 24, 43, 49] is still slow. Another line of works aims to improve the rendering speed by using more efficient data structures to represent the radiance fields [20, 34, 48] or adopting the caching mechanism that pre-computes the color and density values [9, 14, 48]. However, the time and storage cost of the caching process is expensive, which makes them incapable of rendering dynamic scenes.

**Multi-view stereo methods.** The cost volume has been widely used for depth estimation in multi-view stereo (MVS) methods. MVSNet [46] proposes to build the 3D cost volume from 2D image features and regularizes the cost volume with a 3D CNN. This design enables end-to-end training of the network with ground truth depth and achieves impressive results. However, memory consumption of MVSNet is huge. To overcome this problem, following works improve it with recurrent plane sweeping [47] or
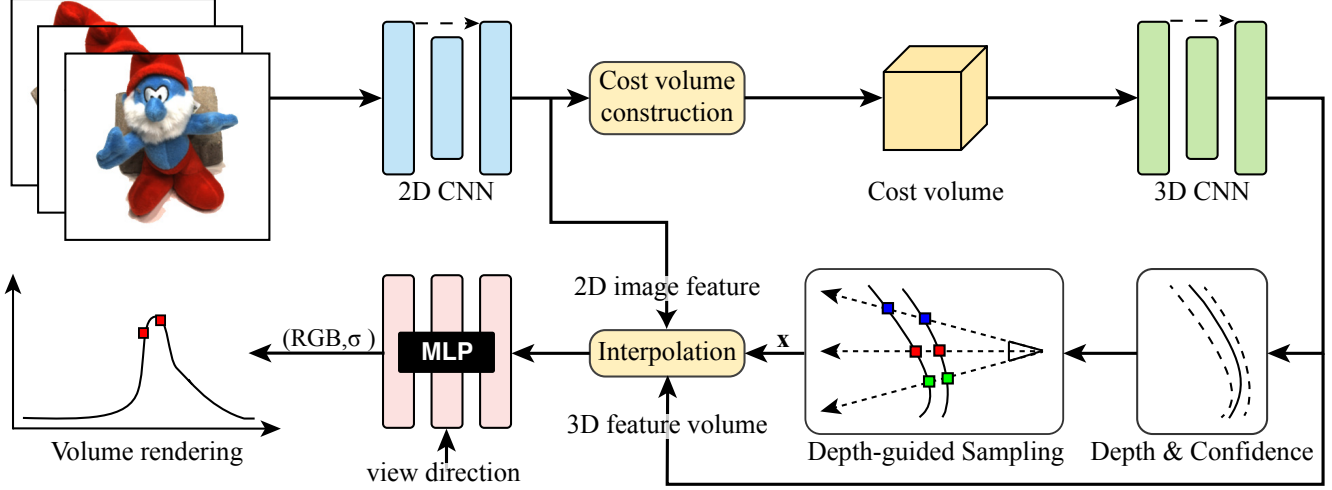
Figure 2. **Illustration of the proposed approach.** We first construct the cascade cost volume from multi-view images, which is processed to output the 3D feature volume and the coarse scene geometry (represented by depth and confidence maps). The estimated geometry guides us to sample around the surface, which significantly accelerates the rendering process. Also, the 3D feature volume provides rich geometry-aware information for radiance fields construction. All network components are trained end-to-end using only RGB images.

coarse-to-fine architectures [3, 4, 11, 50]. In addition, Luo et al. [25] propose a patch-wise matching confidence aggregation mechanism for the cost volume and enhance the performance of MVSNet. Multi-view stereo methods establish the cost volume to estimate the depth for accurate geometry reconstruction, while we take the depth as the coarse scene geometry to guide the volume rendering process for fast view synthesis. Thanks to the differentiable volume rendering technique, we jointly learn the depth prediction and the NeRF networks from only RGB images. [2,5] also attempt to combine multi-view stereo methods with neural radiance field methods. They [2,5] aim to facilitate the construction of the radiance fields leveraging stereo features, while we focus on accelerating the rendering process with explicit geometry from the cost volume.

## 3. Method

Given several captured images of a scene, our task is to generate images of novel views. Although previous methods [2, 43, 49] have shown impressive performance, they generally take a long time to render an image. To solve this problem, we propose a hybrid scene representation that combines implicit neural radiance fields with explicit depth maps for fast view synthesis.

To generate a novel view of the captured scene, we first select $N$ spatially nearby source views from captured images. The proposed approach takes as input $N$ source views and outputs the novel view image. As shown in Figure 2, we first extract multi-scale feature maps from input images, which are used to estimate coarse scene geometry and neural radiance fields (Sec. 3.1). Given extracted feature maps,

we construct the cascade cost volume to obtain the 3D feature volume and the coarse scene geometry (Sec. 3.2). The feature volume provides geometry-aware information for radiance fields construction, and the estimated coarse geometry allows us to sample around the surface, which significantly accelerates the rendering process (Sec. 3.3). All networks are trained end-to-end with the view synthesis loss using only RGB images (Sec. 3.4).

### 3.1. Multi-scale image feature extraction

To build the cascade cost volume, we first extract multi-scale image features from input views $\{I_i\}_{i=1}^N$ using a 2D UNet [36]. As shown in Figure 2, we first feed an input image $I_i \in \mathbb{R}^{H \times W \times 3}$ into the encoder to obtain low-resolution feature maps $F_{i,1} \in \mathbb{R}^{H/4 \times W/4 \times 32}$. Then we use two deconvolution layers to upsample the feature maps and obtain the other two-stage feature maps $F_{i,2} \in \mathbb{R}^{H/2 \times W/2 \times 16}$ and $F_{i,3} \in \mathbb{R}^{H \times W \times 8}$. $F_{i,1}$ and $F_{i,2}$ are used to construct cost volumes, and $F_{i,3}$ is used to reconstruct neural radiance fields.

### 3.2. Coarse-to-fine depth prediction

To efficiently obtain the coarse scene geometry (represented by a high-resolution depth map) from multi-view images, we borrow the design of MVS methods [4, 46] to construct the cascade cost volume. Specifically, we first construct a coarse-level low-resolution cost volume and recover a low-resolution depth map from this cost volume. Then we construct a fine-level high-resolution cost volume utilizing the depth map estimated in the last step. The fine-level cost volume is processed to produce a high-resolution depth map and a 3D feature volume.

**Coarse-level cost volume construction.** Given initial scene depth range, we first sample a set of depth planes $\{\mathbf{L}_i | i = 1, ..., D\}$. Following learning-based MVS methods [46], we construct the cost volume by warping image features $F_{i,1}$ into $D$ sweeping planes. Given the camera parameters $[\mathbf{K}_i, \mathbf{R}_i, \mathbf{t}_i]$ of input view $I_i$ and $[\mathbf{K}_t, \mathbf{R}_t, \mathbf{t}_t]$ of target view, the homography warping is defined as:

$$\mathbf{H}_i(z) = \mathbf{K}_i\mathbf{R}_i\Big(\mathbf{I} + \frac{(\mathbf{R}_i^{-1}\mathbf{t}_i - \mathbf{R}_t^{-1}\mathbf{t}_t)\mathbf{n}^T\mathbf{R}_t}{z}\Big)\mathbf{R}_t^{-1}\mathbf{K}_t^{-1},$$
(1)

where $\mathbf{n}$ denotes the principal axis of the target view camera, $\mathbf{I}$ is the identity matrix and $\mathbf{H}_i(z)$ warps a pixel $(u, v)$ in the target view at depth $z$ to the input view. The warped feature maps $F_i^w \in \mathbb{R}^{D \times H/8 \times W/8 \times 32}$ is defined as:

$$F_i^w(u, v, z) = F_{i,1}(\mathbf{H}_i(z)[u, v, 1]^T).$$
(2)

Based on the warped feature maps, we construct the cost volume by computing the variance of multi-view features $\{F_i^w(u, v, z) | i = 1, ..., N\}$ for each voxel.

**Depth probability distribution.** Given the constructed cost volume, we use a 3D CNN to process it into a depth probability volume $\mathbf{P} \in \mathbb{R}^{D \times H/8 \times W/8}$. Similar to [4], we compute a depth distribution based on the depth probability volume. For a pixel $(u, v)$ in the target view, we can obtain its probability at a certain depth plane $\mathbf{L}_i$ by linearly interpolating the depth probability volume, which is denoted as $\mathbf{P}_i(u, v)$. Then the depth value of pixel $(u, v)$ is defined as the mean $\hat{\mathbf{L}}$ of the depth probability distribution:

$$\hat{\mathbf{L}}(u, v) = \sum_{i=1}^{D} \mathbf{P}_i(u, v)\mathbf{L}_i(u, v),$$
(3)

and its confidence is defined as the standard deviation $\hat{\mathbf{S}}$:

$$\hat{\mathbf{S}}(u, v) = \sqrt{\sum_{i=1}^{D} \mathbf{P}_i(u, v)(\mathbf{L}_i(u, v) - \hat{\mathbf{L}}(u, v))^2}.$$
(4)

**Fine-level cost volume construction and processing.** The depth probability distribution with the mean $\hat{\mathbf{L}}(u, v)$ and the standard deviation $\hat{\mathbf{S}}(u, v)$ determines where the surface may be located in. Specifically, the surface should be located in the depth range defined as follows,

$$\hat{\mathbf{U}}(u, v) = [\hat{\mathbf{L}}(u, v) - \lambda\hat{\mathbf{S}}(u, v), \hat{\mathbf{L}}(u, v) + \lambda\hat{\mathbf{S}}(u, v)],$$
(5)

where $\lambda$ is a hyper-parameter that determines how large the depth range is. We simply set $\lambda$ to 1. To construct the fine-level cost volume, we first upsample the estimated depth range map $\hat{\mathbf{U}} \in \mathbb{R}^{H/8 \times W/8 \times 2}$ four times. Given the depth range map, we uniformly sample $D'$ depth planes within

it and construct the fine-level cost volume by applying the homography warping to the feature maps $F_{i,2}$, similar to the Equation (2). Then we use a 3D CNN to process this cost volume to obtain a depth probability volume and a 3D feature volume. Following the processing step the coarse-level depth probability volume, we get a finer depth range map $\hat{\mathbf{U}}' \in \mathbb{R}^{H/2 \times W/2 \times 2}$, which will guide the sampling for volume rendering.

### 3.3. Neural radiance fields construction

NeRF [28] represents the scene as color and volume density fields. To generalize across scenes, similar to [2,43,49], our method assigns features to arbitrary point in 3D space. Inspired by PixelNeRF [49], for any 3D point, we project it into input images and then extract corresponding pixel-aligned features from $\{F_{i,3} | i = 1, ..., N\}$, which are denoted as $\{f_i | i = 1, ..., N\}$. These features are then aggregated with a pooling operator $\psi$ to output the final feature $f_{\text{img}} = \psi(f_1, ..., f_N)$. The design of this pooling operator $\psi$ is borrowed from IBRNet [43] and its details will be described in the supplementary material. To leverage the geometry-aware information provided by the MVS framework, we also extract the voxel-aligned feature from the 3D feature volume by transforming the 3D point into the query view and trilinearly interpolating the voxel features, which is denoted as $f_{\text{voxel}}$. Our method passes $f_{\text{img}}$ and $f_{\text{voxel}}$ into an MLP network to obtain the point feature and density, which is defined as:

$$f_p, \sigma = \phi(f_{\text{img}}, f_{\text{voxel}}),$$
(6)

where $\phi$ denotes the MLP network. We estimate the color $\hat{\mathbf{c}}_p$ of this 3D point viewed in direction $\mathbf{d}$ by predicting blending weights for the image colors $\{\mathbf{c}_i\}_{i=1}^{N}$ in the source views. Specifically, we concatenate the point feature $f_p$ with the image feature $f_i$ and $\Delta\mathbf{d}_i$, and feed them into an MLP network to yield the blending weight $w_i$ defined as:

$$w_i = \varphi(f_p, f_i, \Delta\mathbf{d}_i),$$
(7)

where $\varphi$ denotes the MLP network and $\Delta\mathbf{d}_i$ is the difference between the target view ray direction $\mathbf{d}$ and the source view ray direction $\mathbf{d}_i$. The color $\hat{\mathbf{c}}_p$ is blended via a soft-argmax operator as the following,

$$\hat{\mathbf{c}}_p = \sum_{i=1}^{N} \frac{\exp(w_i)\mathbf{c}_i}{\sum_{j=1}^{N} \exp(w_j)}.$$
(8)

The implementation details of $\phi$ and $\varphi$ are described in the supplementary material.

**Depth-guided sampling for volume rendering.** Given a viewpoint, our method renders the radiance field into an image with the volume rendering technique. For each image

pixel, we sample a set of points along the camera ray and accumulate their volume densities and colors using an approximated integral equation [28]. Previous methods [28,43,49] mostly adopt the hierarchical sampling strategy to improve the rendering efficiency. However, the coarse-level NeRF still takes a lot of time. In contrast, our method leverages MVS methods to localize the scene surface. For pixel $(u, v)$, we have its depth range $\hat{\mathbf{U}}'(u, v)$ estimated from depth prediction module (Sec. 3.2). Then we uniformly sample $N_k$ points $\{\mathbf{x}_k | k = 1, ..., N_k\}$ within the depth range $\hat{\mathbf{U}}'(u, v)$ for volume rendering.

## 3.4. Training

During training, the gradients are back-propagated to the estimated depth probability distribution through sampled 3D points, so that the depth probability volume can be jointly learned with neural radiance fields from only images. Following [28], we optimize our model with the mean squared error that measures the difference between the rendered and ground-truth pixel colors. The corresponding loss is defined as:

$$\mathcal{L} = \frac{1}{N_r} \sum_{i=1}^{N_r} \|\hat{\mathbf{C}}_i - \mathbf{C}_i\|_2^2, \qquad (9)$$

where $N_r$ is the number of sampled rays at each optimization iteration and $\hat{\mathbf{C}}_i$ and $\mathbf{C}_i$ are the rendered and ground-truth color, respectively.

## 3.5. Implementation details

Our generalizable rendering model is implemented with PyTorch [31] and trained on 4 RTX 2080 Ti GPUs using Adam [17] optimizer with an initial learning rate of $5e^{-4}$. We randomly sample 10240 pixels for 2 novel view images as a batch. The model tends to converge after about 40k iterations and it takes about 10 hours. Given dense views of a new scene, we can finetune our pre-trained model on this scene. Fine-tuning on the scene generally takes about 20 minutes on an RTX 2080 Ti GPU. In practice, we sample 64 and 8 depth planes for the coarse-level and fine-level cost volumes, respectively. The proposed rendering model takes 3 input source views and 2 samples per ray in all experiments unless otherwise specified.

## 4. Experiments

## 4.1. Experiments setup

**Datasets.** We train our generalizable rendering model on the DTU [15] dataset. DTU is composed of 124 real scenes and each scene has 49 or 64 images. We use the train-test split suggested by PixelNeRF [49] and take the evaluation setting proposed in MVSNeRF [2], where they select 16 views as input and 4 views for testing for each test scene.

To further show the generalization ability of our method, we also test the model (trained on DTU) on the Real Forward-facing [27] and NeRF Synthetic [28] datasets. They both include 8 complex scenes and have different scene and view distributions from DTU. We also conduct experiments on a 1200-frame sequence of the ZJU-MoCap [32] dataset to demonstrate the capability of our method to handle dynamic scenes. This sequence has 21 synchronized cameras, and we uniformly select 11 cameras as input and use the remaining cameras for testing.

**Baselines.** As a generalizable radiance field method, we first make comparisons with PixelNeRF [49], IBRNet [43] and MVSNeRF [2]. They introduce image features to avoid the slow per-scene optimization process. Then we show that the proposed method with a short fine-tuning process can achieve comparable results with NeRF [28] and other per-scene optimization methods [2,43]. Compared to caching-based methods [9,14,48], our method has the potential of handling dynamic scenes without expensive caching costs (storage and conversion) as shown in the comparisons with PlenOctree [48].

## 4.2. Performance on image synthesis

**Generalization.** We report metrics of PSNR, SSIM [44] and LPIPS [51] on the NeRF Synthetic [28], DTU [15] and Real Forward-facing [27] datasets in Table 1. Qualitative comparison results can be found in Figure 3. As shown in both qualitative and quantitative results, our method exhibits state-of-the-art performance across different metrics and benchmarks. Furthermore, our method also achieves a significant gain in rendering speed compared to other generalizable radiance field methods. Note that since PixelNeRF uses a much wider MLP than other methods, it takes almost 10x as much time for rendering to other methods [2,43].

**Per-scene optimization.** We compare our method with NeRF [28], IBRNet [43] and MVSNeRF [2] under the per-scene optimization setting. The quantitative and qualitative comparison results are reported in Table 1 and Figure 4, respectively. As shown in the results, we achieve comparable rendering quality with these methods and achieve a significant gain in rendering speed. Note that our rendering quality is slightly worse than $\text{NeRF}_{10.2h}$ [28] on the Real Forward-facing and NeRF Synthetic [28] dataset as shown in quantitative results. This is because instead of encoding an entire scene into a single network with all images as NeRF [28], our method synthesizes the novel view by interpolating nearby source views like IBRNet [43]. On complex scenes with only sparse views, the novel view may include contents that nearby source views have not seen. Anyway, we achieve significant gains in rendering and training efficiency compared to NeRF [28].

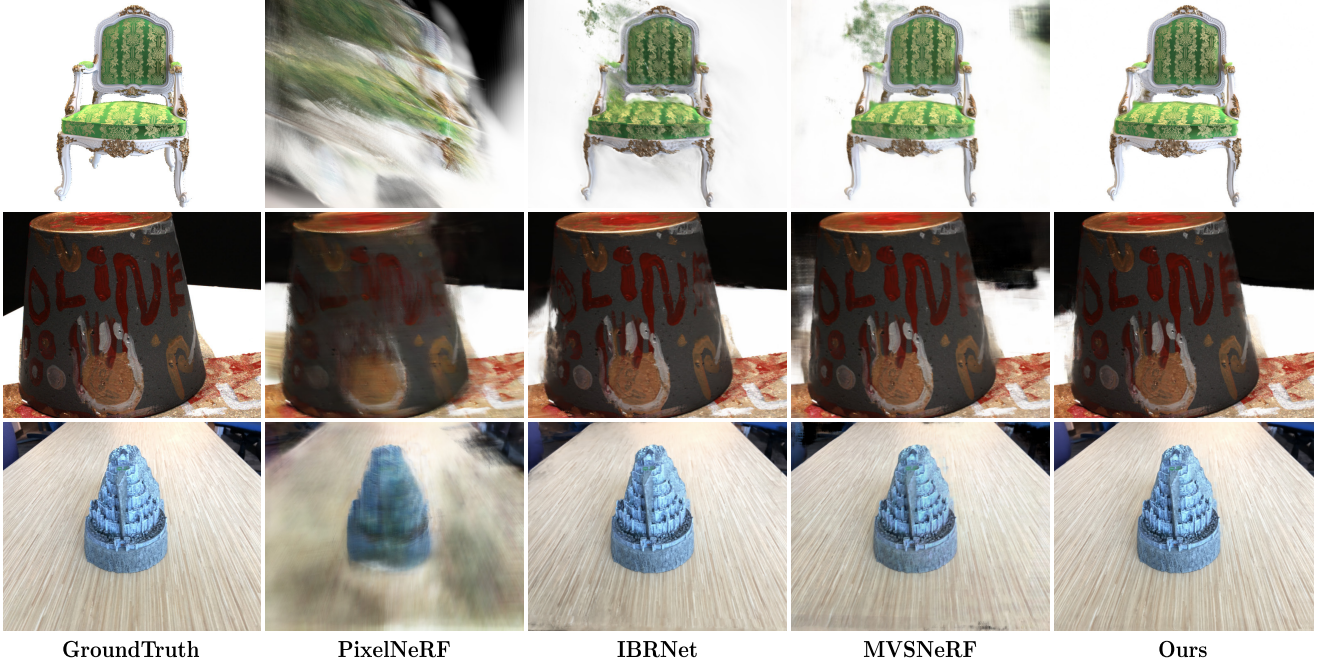| GroundTruth | PixelNeRF | IBRNet | MVSNeRF | Ours |

Figure 3. Qualitative comparison of image synthesis under **the generalization setting** on the NeRF Synthetic [28], DTU [15] and Real Forward-facing [28] datasets.

| Methods | Training settings | Rendering FPS↑ | NeRF Synthetic [28] | | | DTU [15] | | | Real Forward-facing [28] | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| PixelNeRF [49] | | 0.01 | 7.39 | 0.658 | 0.411 | 19.31 | 0.789 | 0.382 | 11.24 | 0.486 | 0.671 |
| IBRNet [43] | Generalization | 0.11 | 22.44 | 0.874 | 0.195 | 26.04 | 0.917 | 0.191 | 21.79 | 0.786 | 0.279 |
| MVSNeRF [2] | | 0.23 | 23.62 | 0.897 | 0.176 | 26.63 | 0.931 | 0.168 | 21.93 | 0.795 | 0.252 |
| Ours | | **12.31** | **25.47** | **0.936** | **0.097** | **27.46** | **0.951** | **0.101** | **22.89** | **0.807** | **0.229** |
| NeRF$_{10.2h}$ [28] | | 0.08 | **30.63** | **0.962** | 0.093 | 27.01 | 0.902 | 0.263 | **25.97** | 0.870 | 0.236 |
| IBRNet$_{ft-1.0h}$ [43] | Per-scene optimization | 0.11 | 25.62 | 0.939 | 0.111 | 31.35 | 0.956 | 0.131 | 24.88 | 0.861 | **0.189** |
| MVSNeRF$_{ft-15min}$ [2] | | 0.23 | 27.07 | 0.931 | 0.168 | 28.51 | 0.933 | 0.179 | 25.45 | **0.877** | 0.192 |
| NeRF$_{20min}$ [28] | | 0.08 | 21.42 | 0.861 | 0.312 | 24.45 | 0.842 | 0.372 | 21.36 | 0.692 | 0.513 |
| Ours$_{ft-20min}$ | | **12.31** | 26.55 | 0.942 | **0.088** | **31.43** | **0.961** | **0.103** | 24.91 | 0.856 | 0.201 |

Table 1. Quantitative comparison of image synthesis results on datasets of NeRF Synthetic [28], DTU [15] and Real Forward-facing [27]. The rendering quality results of [2, 28, 43, 49] other than "NeRF$_{20min}$" are taken from MVSNeRF [2], since our evaluation settings are the same as MVSNeRF. All methods are evaluated on an RTX 2080 Ti GPU to report the speed of rendering a $512 \times 640$ image.

**Dynamic scenes.** Compared to current real-time caching-based methods [9, 14, 48], the proposed method does not need expensive caching costs (caching time and storage). This indicates the potential of the proposed method to extend to dynamic scenarios, since it could produce novel view images in real-time without extra cost. We make quantitative and qualitative comparisons with PlenOctree [48] on the ZJU-MoCap [32] dataset in Table 2 and Figure 5, respectively. It takes about 1.5 hours to project a trained vanilla NeRF model to the PlenOctree model on one frame of a dynamic scene. The storage size of the cached model is

about 1.45 GB. Please refer to the supplementary material for experiment details of [48] on the ZJU-MoCap dataset. It is impractical to use caching-based methods for dynamic scenes because of the expensive cost. In contrast, as a generalizable radiance field method, our method simply trains one network on the whole sequence of the dynamic scene with only the images storage cost. Note that we experimentally find that the images storage cost could be compressed from 0.47GB to 30MB using mature video compression techniques. Furthermore, as our method trains one network on the whole sequence of the dynamic scene, our

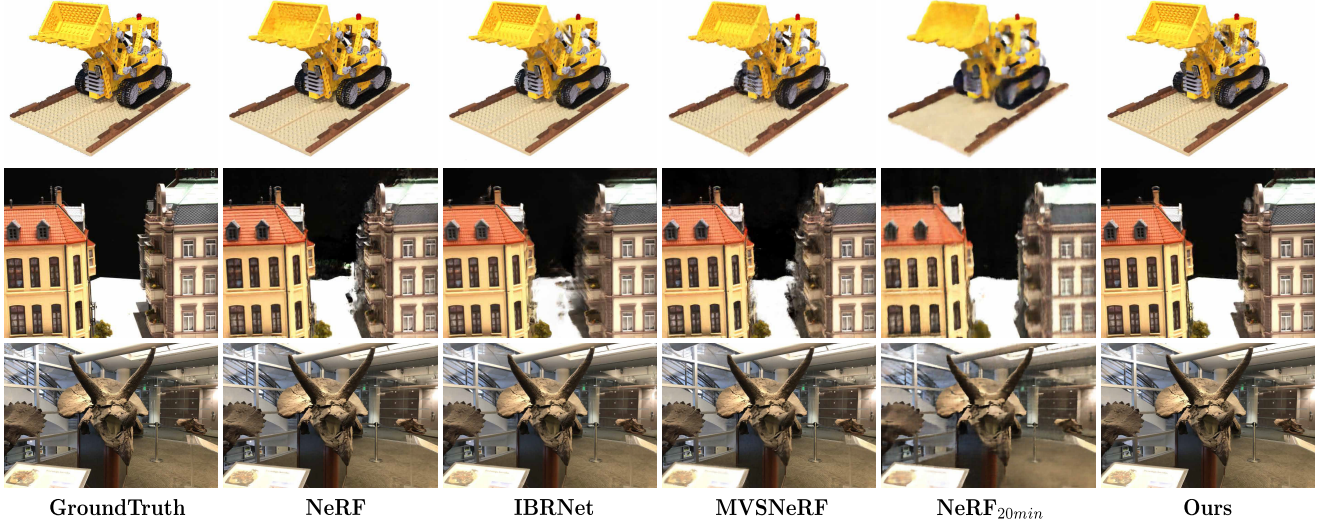|  | GroundTruth | NeRF | IBRNet | MVSNeRF | NeRF$_{20min}$ | Ours |

Figure 4. Qualitative comparison of image synthesis under **the per scene optimization setting** on the NeRF Synthetic [28], DTU [15] and Real Forward-facing [28] datasets.

| Methods | Storage | Preprocess time | Rendering FPS ↑ | PSNR↑ |
|---|---|---|---|---|
| PlenOctree [48] | 1740 | 1800 | **150** | 34.85 |
| NeuralBody [32] | 0.018 | – | 5.49 | 33.90 |
| Ours | 0.47 | – | 19.04 | 33.07 |
| Ours$_{ft}$ | 0.47 | – | 19.04 | **36.22** |

Table 2. Quantitative comparison on the ZJU-MoCap [32] dataset. The units of "Storage" and "Preprocess time" are in GB and hours, respectively. "Ours" represents the proposed method trained on the DTU [15] dataset. Ours$_{ft}$ means we finetune the pre-trained model on this dynamic scene. The rendered image resolution is $512 \times 512$. To obtain the results of [48], we test it only on 4 frames and multiply the average cost by 1200 frames, because of its expensive cost.

| Methods | Abs err↓ | Acc (0.01)↑ | Acc (0.05)↑ |
|---|---|---|---|
| MVSNet [46] | **0.018** / – | 0.603/ – | **0.955**/ – |
| PixelNeRF [49] | 0.245/0.239 | 0.037/0.039 | 0.176/0.187 |
| IBRNet [43] | 1.69 / 1.62 | 0.000/0.000 | 0.000/0.001 |
| MVSNeRF [2] | 0.023/0.035 | 0.746/0.717 | 0.913/0.866 |
| Ours-MVS | 0.029/0.034 | 0.452/0.390 | 0.907/0.879 |
| Ours-NeRF | 0.023/**0.026** | **0.760/0.762** | 0.918/**0.899** |

Table 3. Quantitative comparison of depth results on the DTU [15] dataset. The two numbers of each cell refer to the depth at reference/novel views. The qualitative results of [2,43,46,49] are taken from [2].

method indeed integrate observations across video frames into the network. Thus we also compare our method with NeuralBody [32], which proposes to integrate information across frames into the structured latent codes. As shown in Table 2, our method exhibits state-of-the-art performance with real-time rendering speed. Please refer to the supplementary material for more discussion on integrating observations across video frames and more real-time novel view synthesis results on the DynamicCap [12] dataset.

### 4.3. Quality of reconstructed depth

We compare our depth results reconstructed from volume densities (denoted by "Ours-NeRF") with generalizable radiance field methods [2, 43, 49] and the classic deep MVS method MVSNet [46] on the DTU [15] testing set. In addition, we also report depth results from the cost volume

in our method (denoted by "Ours-MVS"), which are computed from the fine-level depth probability volume. Note that MVSNet is trained with depth supervision while other methods are trained with only image supervision. As shown in Table 3, "Ours-NeRF" outperforms baseline methods. Moreover, "Ours-MVS" produces reasonable depth results from the cost volume, which is a critical factor that leads to our high-quality rendering with only few samples. Please refer to the supplementary material for visual results.

### 4.4. Ablation studies and analysis

**Ablation of main proposed components.** We firstly execute experiments to show the power of depth-guided sampling. As shown in the first three rows of Table 4, when we reduce the number of samples from 128 to 2, the rendering quality of the method without depth-guided sampling drops a lot, while the method with depth-guided sampling can maintain almost the same performance. Secondly, we

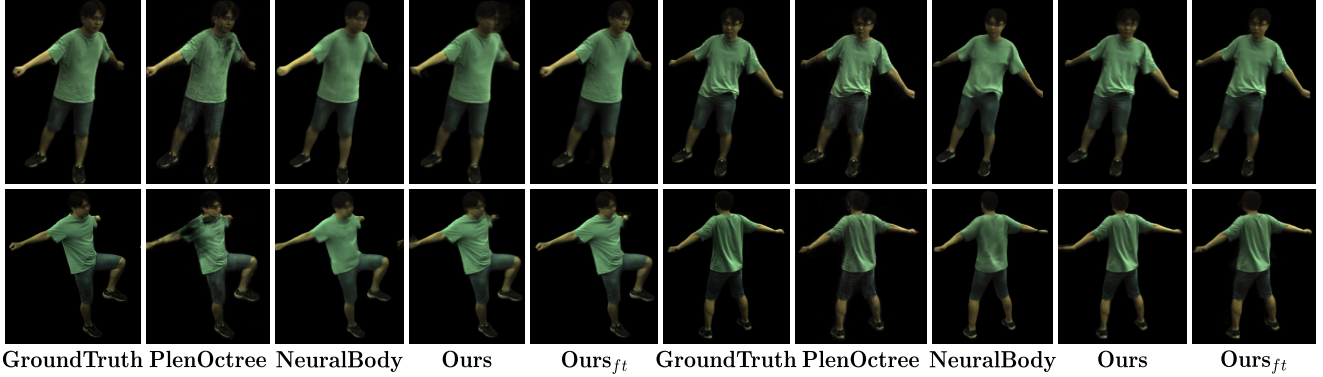| GroundTruth | PlenOctree | NeuralBody | Ours | $\text{Ours}_{ft}$ | GroundTruth | PlenOctree | NeuralBody | Ours | $\text{Ours}_{ft}$ |

Figure 5. Qualitative comparison of image synthesis results on the ZJU-MoCap [32] dataset.

| Samples | Depth-gui. | Cascade | Depth-sup. | PSNR↑ | FPS↑ |
|---------|-----------|---------|-----------|-------|------|
| 128 | | | | 27.39 | 0.32 |
| 2 | | | | 17.75 | 7.94 |
| 2 | ✓ | | | 26.97 | 7.21 |
| 2 | ✓ | ✓ | | **27.45** | **12.31** |
| 2 | ✓ | ✓ | ✓ | 27.11 | 12.31 |

Table 4. Quantitative ablation analysis of the design choices on the DTU [15] dataset. "Depth-gui." and "Depth-sup." means "Depth-guided" and "Depth-supervision", respectively.

| Samples | PSNR↑ | FPS↑ | Views | PSNR↑ | FPS↑ |
|---------|-------|------|-------|-------|------|
| 1 | 26.89 | 16.09 | 2 | 25.45 | 15.23 |
| 2 | 27.45 | 12.31 | 3 | 27.45 | 12.26 |
| 4 | 27.49 | 7.65 | 4 | 27.80 | 10.24 |
| 8 | 27.54 | 4.15 | 5 | 27.84 | 8.73 |

Table 5. Quantitative ablation analysis of the number of samples and input source views on the DTU [15] dataset.

conduct experiments to analyze the effect of the cascade cost volume. As shown in the fourth row of Table 4, the cascade design greatly improves the rendering speed and shows the same good performance. Finally, we additionally supervise the depth probability volume using ground truth depth. As shown in the last row of Table 4, using depth supervision does not improve the rendering performance. This demonstrates the effectiveness of end-to-end training with the view synthesis loss using RGB images.

**Sensitivity to the number of input views and samples.** We study how the number of input source views and samples affects the final rendering quality and speed in Table 5. Thanks to the effective depth-guided sampling strategy, the number of samples does not affect the view synthesis quality a lot. When the number of input source views is increased from 2 to 3, the rendering quality is greatly improved. We simply take 2 samples and 3 views for all experiments (2 views for the ZJU-MoCap scene and 4 views for a few difficult scenes under the per-scene optimization setting, e.g., T-Rex in Real Forward-facing.)

Please refer to the supplementary material for more qualitative ablation results.

### 4.5. Running time analysis

To render a $512 \times 640$ image, our method with 3 input views and 2 samples per ray runs at 12.31 FPS on a desktop with an RTX 2080 Ti GPU. Specifically, our method takes 8 ms to extract image features of 3 input images, 12 ms to construct the cost volumes, 16 ms to process the cost volume using 3D CNN, and 45 ms for inferring radiance fields and volume rendering.

### 5. Conclusion and discussion

This paper introduced a hybrid scene representation that combines the best of implicit radiance fields and explicit depth maps for fast view synthesis. The core innovation is utilizing explicit depth maps as coarse scene geometry to guide the rendering process of implicit radiance fields. We demonstrated state-of-the-art performance of our method and the capability to synthesize novel views of dynamic scenes in real-time.

Although our method can efficiently render high-quality images, it still has several limitations. 1) Our approach predicts depth maps based on cost volumes, which requires a lot of memory. It would be interesting to explore more efficient ways to estimate the scene geometry. 2) The proposed model generates novel views based on nearby images. Once some target regions under the novel view are invisible in input views, the rendering quality may degrade. It could be solved by considering the visibility of input views.

# References

[1] Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM TOG*, 2013. 2

[2] 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 *ICCV*, 2021. 2, 3, 4, 5, 6, 7, 11, 12

[3] Rui Chen, Songfang Han, Jing Xu, and Hao Su. Point-based multi-view stereo network. In *ICCV*, 2019. 3

[4] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *CVPR*, 2020. 3, 4

[5] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (srf): Learning view synthesis for sparse views of novel scenes. In *CVPR*, 2021. 3

[6] Inchang Choi, Orazio Gallo, Alejandro J. Troccoli, Min H. Kim, and Jan Kautz. Extreme view synthesis. In *ICCV*, 2019. 2

[7] Abe Davis, Marc Levoy, and Fredo Durand. Unstructured light fields. In *Eurographics*, 2012. 2

[8] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world's imagery. In *CVPR*, June 2016. 2

[9] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *ICCV*, 2021. 1, 2, 5, 6

[10] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *SIGGRAPH*, 1996. 2

[11] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *CVPR*, pages 2492–2501, 2020. 3

[12] Marc Habermann, Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Real-time deep dynamic characters. *ACM TOG*, 2021. 7

[13] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM TOG*, 2018. 2

[14] 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, 5, 6

[15] Rasmus Ramsbøl Jensen, Anders Lindbjerg Dahl, George Vogiatzis, Engil Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *CVPR*, 2014. 2, 5, 6, 7, 8, 11, 12

[16] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM TOG*, 2016. 2

[17] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5

[18] Marc Levoy and Pat Hanrahan. Light field rendering. In *SIGGRAPH*, 1996. 2

[19] Zhengqi Li, Wenqi Xian, Abe Davis, and Noah Snavely. Crowdsampling the plenoptic function. In *ECCV*, 2020. 2

[20] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 1, 2

[21] Lingjie Liu, Weipeng Xu, Marc Habermann, Michael Zollhöfer, Florian Bernard, Hyeongwoo Kim, Wenping Wang, and Christian Theobalt. Neural human video rendering by learning dynamic textures and rendering-to-video translation. *TVCG*, 2020. 2

[22] Lingjie Liu, Weipeng Xu, Michael Zollhoefer, Hyeongwoo Kim, Florian Bernard, Marc Habermann, Wenping Wang, and Christian Theobalt. Neural rendering and reenactment of human actor videos. *ACM TOG*, 2019. 2

[23] Shaohui Liu, Yinda Zhang, Songyou Peng, Boxin Shi, Marc Pollefeys, and Zhaopeng Cui. DIST: rendering deep implicit signed distance function with differentiable sphere tracing. In *CVPR*, 2020. 2

[24] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. *arXiv*, 2021. 2

[25] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *ICCV*, 2019. 3

[26] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *CVPR*, 2019. 2

[27] 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 TOG*, 2019. 2, 5, 6

[28] 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, 4, 5, 6, 7, 11

[29] Michael Niemeyer, Lars M. Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *CVPR*, 2020. 2

[30] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *CVPR*, 2019. 2

[31] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 5

[32] 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 *CVPR*, 2021. 1, 2, 5, 6, 7, 8, 11

[33] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM TOG*, 2017. 2

[34] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, pages 14335–14345, 2021. 2

[35] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *ECCV*, 2020. 2

[36] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 3

[37] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 2

[38] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *ECCV*, 2016. 2

[39] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deep-voxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019. 2

[40] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wet-zstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 2

[41] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM TOG*, 2019. 2

[42] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *ICCV*, 2021. 1, 2

[43] 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 *CVPR*, 2021. 1, 2, 3, 4, 5, 6, 7, 11

[44] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 2004. 5

[45] Zexiang Xu, Sai Bi, Kalyan Sunkavalli, Sunil Hadap, Hao Su, and Ravi Ramamoorthi. Deep view synthesis from sparse photometric images. *ACM TOG*, 2019. 2

[46] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018. 2, 3, 4, 7

[47] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *CVPR*, 2019. 2

[48] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 1, 2, 5, 6, 7, 11

[49] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 1, 2, 3, 4, 5, 6, 7

[50] Zehao Yu and Shenghua Gao. Fast-mvsnet: Sparse-to-dense multi-view stereo with learned propagation and gauss-newton refinement. In *CVPR*, 2020. 3

[51] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 5

## Supplementary material

In the supplementary material, we provide network architectures, details of experimental setup, and more experimental results.

## 1. Network architectures

**Pooling operator.** Given the multi-view point features $\{f_i\}_{i=1}^N$, the pooling operator $\psi$ aims to aggregate these features to obtain the feature $f_{\text{img}}$, which is used to infer the radiance field. Instead of simply concatenating these features like MVSNeRF [2], we use a weighted pooling operator proposed in IBRNet [43], which allows us to input any number of source views. Specifically, we first compute a per-element mean $\boldsymbol{\mu}$ and variance $\mathbf{v}$ of $\{f_i\}_{i=1}^N$ to capture global information. Then we concatenate each feature $f_i$ with $\boldsymbol{\mu}$ and $\mathbf{v}$, and feed the concatenated feature into a small shared MLP to obtain a weight $w_i$. The feature $f_{\text{img}}$ is blended via a soft-argmax operator using weights $\{w_i\}_{i=1}^N$ and multi-view features $\{f_i\}_{i=1}^N$.

**Architectures of MLPs.** The MLP $\phi$ is used to infer the density $\sigma$ from the image feature $f_{\text{img}}$ and the voxel feature $f_{\text{voxel}}$. To predict the color of the point, we use the MLP $\varphi$ to yield the blending weights for image colors in the source views. We illustrate the architectures of $\phi$ and $\varphi$ in Table 6.

## 2. Details of the experimental setup

**Evaluation details.** Our evaluation setup is taken from MVSNeRF [2] and is described as the following. To report the results on the DTU [15] dataset, we compute the metric score of foreground part in images. For metrics of SSIM and LPIPS, we set the background to black and calculate the metric score of the whole image. The segmentation mask is defined by whether there is ground-truth depth available at each pixel. Since marginal regions of images are usually invisible to input images on the Real Forward-facing [28] dataset, we only evaluate 80% area in the center of images. The image resolutions are set to $512 \times 640$, $640 \times 960$ and $800 \times 800$ on the DTU, Real forward-facing and NeRF Synthetic [28] datasets, respectively.

**Experimental details of PlenOctree.** We convert vanilla trained NeRF models to PlenOctree models following the suggestion by [48]. Training a NeRF model on one frame of the ZJU-MoCap [32] dataset takes about 2.5 hours. It takes about 1.45 hours to convert the trained NeRF model to the PlenOctree model. After converting, we optimize the PlenOctree model using the view synthesis loss with SGD optimizer as suggested by [48]. The optimization process takes about 0.1 hours.

| MLP | Layer | Chns. | Input | Output |
|---|---|---|---|---|
| | $LR_0$ | 8 + 16 / 128 | $f_{\text{img}}, f_{\text{voxel}}$ | hidden feature |
| $\phi$ | $LR_i$ | 128 / 128 | hidden feature | hidden feature |
| | $LR_3$ | 128 / 64 + 1 | hidden feature | $f_p, \sigma$ |
| | $LR_0$ | 64 + 16 + 4 / 128 | $f_p, f_i, \Delta\mathbf{d}_i$ | hidden feature |
| $\varphi$ | $LR_1$ | 128 / 64 | hidden feature | hidden feature |
| | $LR_2$ | 64 / 1 | hidden feature | $w_i$ |

Table 6. **The architectures of MLPs $\phi$ and $\varphi$.** We denote LR to be LinearRelu layer. "Chns." shows the number of input and output channels for each layer.

| Test frame id | 1 | 2 | 3 | 4 | Mean |
|---|---|---|---|---|---|
| Per-frame training | 36.73 | 33.07 | 34.46 | 33.24 | 34.38 |
| Per-sequence training | 37.68 | 35.55 | 36.33 | 35.30 | 36.22 |

Table 7. **Image synthesis results on 4 frames of the ZJU-MoCap [32] dataset in terms of PSNR metric.** "Per-frame training" means we simply fine-tune the pre-trained network on the test frame. "Per-sequence training" means we fine-tune the pre-trained network on the whole sequence (1200 frames).

| Scan | #1 | #8 | #21 | #103 | #114 |
|---|---|---|---|---|---|
| | | | PSNR↑ | | |
| PixelNeRF | 21.64 | 23.70 | 16.04 | 16.76 | 18.40 |
| IBRNet | 25.97 | 27.45 | 20.94 | 27.91 | 27.91 |
| MVSNeRF | 26.96 | 27.43 | 21.55 | 29.25 | 27.99 |
| Ours | **28.86** | **28.98** | **22.69** | **30.64** | **29.00** |
| $NeRF_{10.2h}$ | 26.62 | 28.33 | 23.24 | 30.40 | 26.47 |
| $IBRNet_{ft-1h}$ | 31.00 | 32.46 | **27.88** | 34.40 | **31.00** |
| $MVSNeRF_{ft-15min}$ | 28.05 | 28.88 | 24.87 | 32.23 | 28.47 |
| $Ours_{ft-20min}$ | **31.93** | **32.69** | 27.21 | **34.66** | 30.66 |
| | | | SSIM↑ | | |
| PixelNeRF | 0.827 | 0.829 | 0.691 | 0.836 | 0.763 |
| IBRNet | 0.918 | 0.903 | 0.873 | 0.950 | 0.943 |
| MVSNeRF | **0.937** | **0.922** | **0.890** | **0.962** | **0.949** |
| Ours | 0.916 | 0.895 | 0.880 | 0.924 | 0.935 |
| $NeRF_{10.2h}$ | 0.902 | 0.876 | 0.874 | 0.944 | 0.913 |
| $IBRNet_{ft-1h}$ | 0.955 | 0.945 | 0.947 | 0.968 | 0.964 |
| $MVSNeRF_{ft-15min}$ | 0.934 | 0.900 | 0.922 | 0.964 | 0.945 |
| $Ours_{ft-20min}$ | **0.966** | **0.956** | **0.948** | **0.971** | **0.965** |
| | | | LPIPS ↓ | | |
| PixelNeRF | 0.373 | 0.384 | 0.407 | 0.376 | 0.372 |
| IBRNet | 0.190 | 0.252 | 0.179 | 0.195 | 0.136 |
| MVSNeRF | 0.155 | 0.220 | 0.166 | 0.165 | 0.135 |
| Ours | **0.105** | **0.149** | **0.121** | **0.128** | **0.094** |
| $NeRF_{10.2h}$ | 0.265 | 0.321 | 0.246 | 0.256 | 0.225 |
| $IBRNet_{ft-1h}$ | 0.129 | 0.170 | 0.104 | 0.156 | 0.099 |
| $MVSNeRF_{ft-15min}$ | 0.171 | 0.261 | 0.142 | 0.170 | 0.153 |
| $Ours_{ft-20min}$ | **0.088** | **0.133** | **0.092** | **0.119** | **0.086** |

Table 8. **Quantitative comparison on the DTU dataset.**

Ground-truth image      Ours image      Ours NeRF depth      Ours MVS depth      Ground-truth depth
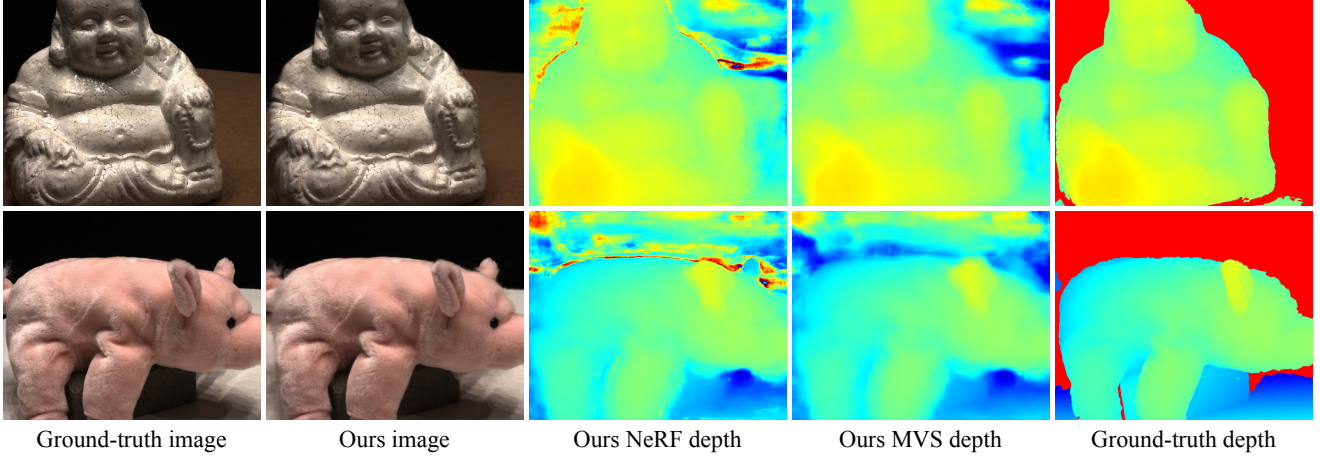
Figure 6. **Visual depth results on the DTU [15] dataset.** "Ours NeRF depth" represents the depth results recovered from volume densities. "Ours MVS depth" denotes the depth results from the fine-level cost volume.
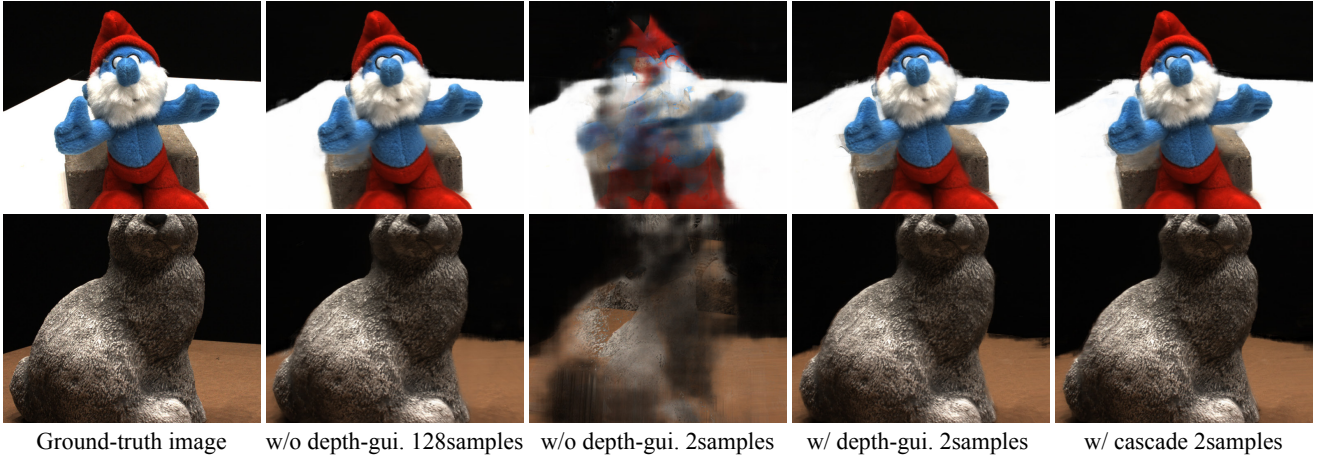


Ground-truth image    w/o depth-gui. 128samples    w/o depth-gui. 2samples    w/ depth-gui. 2samples    w/ cascade 2samples

Figure 7. **Visual ablation results on the DTU [15] dataset.** "w/o depth-gui." is similar to MVSNeRF [2].

## 3. Visual results

**Depth results.** As shown in Figure 6, the proposed method produces reasonable depth results by supervising networks with only images. The cost volume recovers high-quality depth, which allows us to place few samples around surfaces to achieve photorealistic rendering.

**Ablation results.** We provide visual ablation results in Figure 7. The results show that when we reduce the number of samples from 128 to 2, our method with depth-guided sampling almost maintains the same rendering quality. With the depth-guided sampling, the construction of a high-resolution cost volume becomes a bottleneck in the rendering speed. The cascade cost volume further speeds up the construction of the cost volume without loss of rendering quality as shown in Figure 7.

## 4. Integrating information across video frames

We observe that training on a sequence produces higher rendering quality on the test frame compared to training on one frame of the sequence as shown in Table 7. This indicates that our method is able to integrate observations across video frames to produce higher quality images.

## 5. Per-scene breakdown

Tables 8, 9 and 10 present the per-scene comparisons. These results are consistent with the averaged results shown in the paper and show that our method achieves comparable performance to baselines.

| | Chair | Drums | Ficus | Hotdog | Lego | Materials | Mic | Ship |
|---|---|---|---|---|---|---|---|---|
| | | | | PSNR↑ | | | | |
| PixelNeRF | 7.18 | 8.15 | 6.61 | 6.80 | 7.74 | 7.61 | 7.71 | 7.30 |
| IBRNet | 24.20 | 18.63 | 21.59 | 27.70 | 22.01 | 20.91 | 22.10 | 22.36 |
| MVSNeRF | 23.35 | 20.71 | 21.98 | 28.44 | 23.18 | 20.05 | 22.62 | 23.35 |
| Ours | **27.01** | **22.67** | **23.22** | **31.66** | **24.24** | **23.05** | **25.00** | **25.20** |
| NeRF | **31.07** | **25.46** | **29.73** | **34.63** | **32.66** | **30.22** | **31.81** | **29.49** |
| IBRNet$_{ft-1h}$ | 28.18 | 21.93 | 25.01 | 31.48 | 25.34 | 24.27 | 27.29 | 21.48 |
| MVSNeRF$_{ft-15min}$ | 26.80 | 22.48 | 26.24 | 32.65 | 26.62 | 25.28 | 29.78 | 26.73 |
| Ours$_{ft-20min}$ | 27.81 | 24.01 | 24.25 | 33.15 | 25.16 | 24.79 | 27.38 | 25.81 |
| | | | | SSIM↑ | | | | |
| PixelNeRF | 0.624 | 0.670 | 0.669 | 0.669 | 0.671 | 0.644 | 0.729 | 0.584 |
| IBRNet | 0.888 | 0.836 | 0.881 | 0.923 | 0.874 | 0.872 | 0.927 | 0.794 |
| MVSNeRF | 0.876 | 0.886 | 0.898 | 0.962 | 0.902 | 0.893 | 0.923 | **0.886** |
| Ours | **0.942** | **0.912** | **0.899** | **0.963** | **0.908** | **0.901** | **0.950** | 0.791 |
| NeRF | **0.971** | **0.943** | **0.969** | 0.980 | **0.975** | **0.968** | **0.981** | **0.908** |
| IBRNet$_{ft-1h}$ | 0.955 | 0.913 | 0.940 | 0.978 | 0.940 | 0.937 | 0.974 | 0.877 |
| MVSNeRF$_{ft-15min}$ | 0.934 | 0.898 | 0.944 | 0.971 | 0.924 | 0.927 | 0.970 | 0.879 |
| Ours$_{ft-20min}$ | 0.961 | 0.932 | 0.927 | **0.981** | 0.933 | 0.941 | 0.975 | 0.889 |
| | | | | LPIPS ↓ | | | | |
| PixelNeRF | 0.386 | 0.421 | 0.335 | 0.433 | 0.427 | 0.432 | 0.329 | 0.526 |
| IBRNet | 0.144 | 0.241 | 0.159 | 0.175 | 0.202 | 0.164 | 0.103 | 0.369 |
| MVSNeRF | 0.282 | 0.187 | 0.211 | 0.173 | 0.204 | 0.216 | 0.177 | 0.244 |
| Ours | **0.060** | **0.097** | **0.101** | **0.066** | **0.108** | **0.102** | **0.058** | **0.205** |
| NeRF | **0.055** | 0.101 | **0.047** | 0.089 | **0.054** | 0.105 | **0.033** | 0.263 |
| IBRNet$_{ft-1h}$ | 0.079 | 0.133 | 0.082 | 0.093 | 0.105 | 0.093 | 0.040 | 0.257 |
| MVSNeRF$_{ft-15min}$ | 0.129 | 0.197 | 0.171 | 0.094 | 0.176 | 0.167 | 0.117 | 0.294 |
| Ours$_{ft-20min}$ | 0.056 | **0.092** | 0.085 | **0.048** | 0.095 | **0.081** | 0.038 | **0.207** |

Table 9. **Quantitative comparison on the NeRF Synthetic dataset.**

|  | Fern | Flower | Fortress | Horns | Leaves | Orchids | Room | Trex |
|---|---|---|---|---|---|---|---|---|
| | | | | PSNR↑ | | | | |
| PixelNeRF | 12.40 | 10.00 | 14.07 | 11.07 | 9.85 | 9.62 | 11.75 | 10.55 |
| IBRNet | 20.83 | 22.38 | 27.67 | 22.06 | **18.75** | 15.29 | 27.26 | 20.06 |
| MVSNeRF | **21.15** | 24.74 | 26.03 | 23.57 | 17.51 | **17.85** | 26.95 | **23.20** |
| Ours | 20.84 | **24.84** | **28.81** | **23.58** | 18.20 | 17.50 | **28.63** | 20.70 |
| $\text{NeRF}_{10.2h}$ | **23.87** | 26.84 | **31.37** | 25.96 | 21.21 | 19.81 | **33.54** | **25.19** |
| $\text{IBRNet}_{ft-1h}$ | 22.64 | 26.55 | 30.34 | 25.01 | **22.07** | 19.01 | 31.05 | 22.34 |
| $\text{MVSNeRF}_{ft-15min}$ | 23.10 | 27.23 | 30.43 | **26.35** | 21.54 | **20.51** | 30.12 | 24.32 |
| $\text{Ours}_{ft-20min}$ | 21.92 | **27.42** | 29.88 | 25.49 | 21.28 | 19.01 | 30.82 | 23.42 |
| | | | | SSIM↑ | | | | |
| PixelNeRF | 0.531 | 0.433 | 0.674 | 0.516 | 0.268 | 0.317 | 0.691 | 0.458 |
| IBRNet | **0.710** | 0.854 | **0.894** | 0.840 | **0.705** | 0.571 | 0.950 | 0.768 |
| MVSNeRF | 0.638 | **0.888** | 0.872 | **0.868** | 0.667 | **0.657** | 0.951 | 0.868 |
| Ours | 0.628 | 0.830 | 0.864 | 0.795 | 0.633 | 0.541 | 0.921 | 0.701 |
| $\text{NeRF}_{10.2h}$ | **0.828** | 0.897 | **0.945** | 0.900 | 0.792 | 0.721 | **0.978** | **0.899** |
| $\text{IBRNet}_{ft-1h}$ | 0.774 | 0.909 | 0.937 | 0.904 | **0.843** | 0.705 | 0.972 | 0.842 |
| $\text{MVSNeRF}_{ft-15min}$ | 0.795 | **0.912** | 0.943 | **0.917** | 0.826 | **0.732** | 0.966 | 0.895 |
| $\text{Ours}_{ft-20min}$ | 0.751 | 0.911 | 0.933 | 0.902 | 0.818 | 0.706 | 0.966 | 0.861 |
| | | | | LPIPS ↓ | | | | |
|  | Fern | Flower | Fortress | Horns | Leaves | Orchids | Room | Trex |
| PixelNeRF | 0.650 | 0.708 | 0.608 | 0.705 | 0.695 | 0.721 | 0.611 | 0.667 |
| IBRNet | 0.349 | 0.224 | 0.196 | 0.285 | 0.292 | 0.413 | **0.161** | 0.314 |
| MVSNeRF | **0.238** | 0.196 | 0.208 | 0.237 | 0.313 | **0.274** | 0.172 | **0.184** |
| Ours | 0.257 | **0.181** | **0.137** | **0.218** | **0.256** | 0.325 | 0.179 | 0.285 |
| $\text{NeRF}_{10.2h}$ | 0.291 | 0.176 | 0.147 | 0.247 | 0.301 | 0.321 | 0.157 | 0.245 |
| $\text{IBRNet}_{ft-1h}$ | 0.266 | 0.146 | 0.133 | 0.190 | **0.180** | 0.286 | **0.089** | 0.222 |
| $\text{MVSNeRF}_{ft-15min}$ | **0.253** | **0.143** | 0.134 | **0.188** | 0.222 | **0.258** | 0.149 | **0.187** |
| $\text{Ours}_{ft-20min}$ | 0.256 | 0.166 | **0.126** | 0.193 | 0.196 | 0.286 | 0.162 | 0.225 |

Table 10. **Quantitative comparison on the Real Forward-facing dataset.**