

Improved Neural Radiance Fields Using Pseudo-depth and Fusion

Jingliang Li, Qiang Zhou, Chaohui Yu, Zhengda Lu, Jun Xiao, Zhibin Wang, Fan Wang

ABSTRACT

Since the advent of Neural Radiance Fields, novel view synthesis has received tremendous attention. The existing approach for the generalization of radiance field reconstruction primarily constructs an encoding volume from nearby source images as additional inputs. However, these approaches cannot efficiently encode the geometric information of real scenes with various scale objects/structures. In this work, we propose constructing multi-scale encoding volumes and providing multi-scale geometry information to NeRF models. To make the constructed volumes as close as possible to the surfaces of objects in the scene and the rendered depth more accurate, we propose to perform depth prediction and radiance field reconstruction simultaneously. The predicted depth map will be used to supervise the rendered depth, narrow the depth range, and guide points sampling. Finally, the geometric information contained in point volume features may be inaccurate due to occlusion, lighting, etc. To this end, we propose enhancing the point volume feature from depth-guided neighbor feature fusion. Experiments demonstrate the superior performance of our method in both novel view synthesis and dense geometry modeling without per-scene optimization.

CCS CONCEPTS

• Computing methodologies → Image-based rendering.

KEYWORDS

neural radiance fields, multi-scale, depth, feature fusion

ACM Reference Format:

Jingliang Li, Qiang Zhou, Chaohui Yu, Zhengda Lu, Jun Xiao, Zhibin Wang, Fan Wang. 2023. Improved Neural Radiance Fields Using Pseudo-depth and Fusion. In *Woodstock '18: ACM Symposium on Neural Gaze Detection*, Oct. 29–Nov. 3, 2023, Ontario, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Novel view synthesis is a fundamental problem for computer vision community, which aims to produce photo-realistic images of the same scene at novel viewpoints. This long-standing problem has recently received tremendous attention [12] due to the advent of neural rendering. Notably, the recent method of neural radiance fields (NeRF) [17] has shown impressive performance on novel view synthesis, represented using global MLPs. Although NeRF and its extensions [15, 35] can represent scenes faithfully and compactly,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference'23, October 2023, Ottawa, Ontario, Canada

© 2023 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

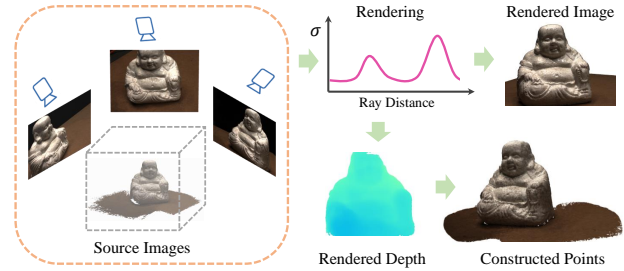


Figure 1: We present a more effective end-to-end method for both photo-realistic novel views rendering and high-quality dense geometry modeling. Our method utilizes multi-scale radiance field, auxiliary depth prediction head, and feature fusion for more realistic novel views, accurate depth maps, and dense points.

they typically require a very long per-scene optimization process to obtain high-quality radiance fields. Recently, some multi-view rendering literature [5, 10, 25, 29, 33] have been proposed to address these shortcomings, which can generalize well across scenes by taking advantage of the nearby input views. These methods construct encoding volumes from source images and extract volume features for each sample point via trilinear interpolation. The point volume feature contains the geometric information of the scene, which will be used as an additional input to the NeRF [17] model to improve the model's scene generalization ability.

Although these methods achieve promising results, they still have several limitations, especially when generalizing to scenes with diverse objects/structures and modeling delicate geometry. In this work, as depicted in Figure 1, we present an effective end-to-end method for both photo-realistic novel views rendering and high-quality geometry modeling.

The size of objects/structures in a realistic scene is diverse, and thus a single scale encoding volume is not adequate to provide geometric information for the entire scene. To this end, inspired by recent coarse-to-fine multi-view stereo (MVS) approaches, we propose constructing pyramid-structured encoding volumes to provide geometric information at all scales. The low-resolution encoding volume provides more geometric information for large objects/structures in the scene. In contrast, the high-resolution encoding volume focuses more on small-scale objects/structures. For each scale encoding volume, we sample points on the ray, trilinear interpolate point volume features and regress the volume density and view-dependent radiance in the scene. That is, we reconstruct the multi-scale radiance fields and combine the reconstructed radiance fields at all scales to perform the final volume rendering for view synthesis.

Building pyramid encoding volumes without depth prior guidance is a massive challenge for GPU memory usage and rendering rate, since we need to sample enough points across the entire depth range when building each scale volume. Although using RGB values

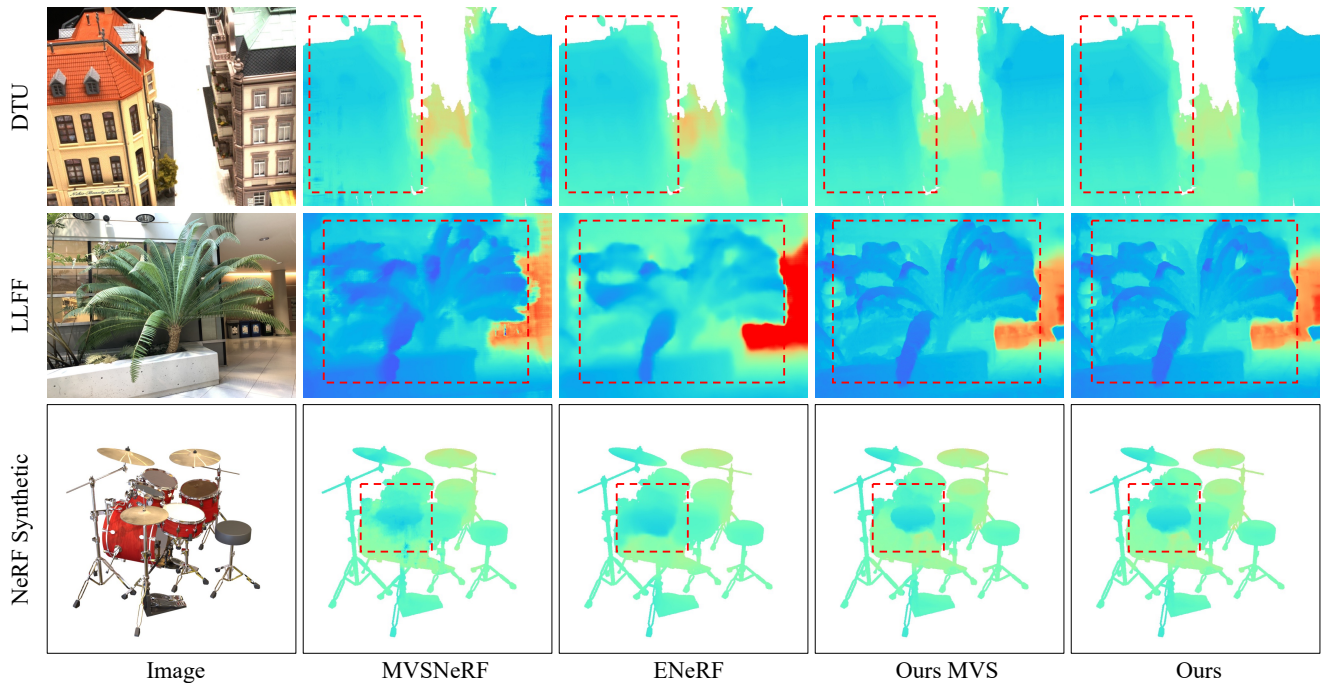


Figure 2: Visual illustrations of rendered source depths of different methods. Our method can learn more delicate geometry for neural radiance fields. “Ours MVS” means the output of the auxiliary depth prediction head.

to supervise the model helps to create more realistic novel views for certain scenes, there may still be challenges with ambiguous geometric information, particularly in texture-less scenes. Recently, some works [6, 18] have proposed depth-guided sampling for NeRF to reduce the number of sampling points, where depth information comes from third-party software or models, such as COLMAP [19] and depth completion models. In this work, we propose to train an additional depth prediction head along with the radiance field reconstruction, and use the predicted depth map of the depth head to supervise the rendered depth, reduce the depth range and guide point sampling. The depth head is trained in an unsupervised manner with the same loss function as [28]. As shown in Figure 2, our method using pseudo-depth guidance can learn more compact geometry compared to other counterparts.

Point volume features trilinearly interpolated from encoding volumes play an essential role in rendering precision across scenes. However, the geometric information contained in point volume features may be inaccurate due to occlusion, lighting, etc. Aggregating the volume features of nearby points from the same surface can effectively improve the correctness of the geometric information in point volume features. To this end, we propose a depth-guided adaptive neighbor feature fusion module. Utilizing the predicted depth map, we adaptively select adjacent points from the same surface and fuse these features using a shared multi-head cross-attention module. The projected features of these adaptive neighbor points on the source view are also input into NeRF as auxiliary information.

Our approach is fully differentiable and can therefore be trained end-to-end using multi-view images. Similar to CasMVSNet [7], we build multi-scale encoding volumes at the target view by warping

2D image features from source views onto sweeping planes in the target view. The multi-scale encoding volumes are used for depth maps prediction and radiance fields (including density σ and view-dependent radiance r) reconstruction at the same time. Finally, we combine the reconstructed multi-scale radiance fields and perform view synthesis via differentiable ray marching. In summary, our contributions are:

- We present a more efficient multi-scale framework for improving NeRF, which is capable of rendering photo-realistic novel views as well as high-quality dense geometry.
- We propose an auxiliary depth maps prediction head, used for supervising rendered depth, narrowing the depth range and guide points sampling.
- We propose enhancing the point volume feature, warped feature and color information from depth-guided neighbor feature fusion.
- Extensive experiments show that our framework achieves state-of-the-art results on various view synthesis datasets.

2 RELATED WORK

View synthesis. The long-standing problem of novel view synthesis is fundamental in computer vision. Recently, various neural scene representations have been presented [3, 13, 20, 36], due to the advent of neural rendering. NeRF [17] has shown impressive performance, which uses global MLPs to regress the volume density and view-dependent radiance at any arbitrary point in the space and applies volume rendering to synthesize images at novel viewpoints. Following works have highlighted different tasks such

as composing and editing [8, 26, 30, 34], large-scale scenes synthesizing [22, 23], relighting [2, 4, 21]. Like NeRF, most of these works construct radiance fields using global MLPs for the entire space, which requires a lengthy optimization process for each new scene before they can synthesize any novel views of that scene. To achieve cross-scene estimation for view synthesis, borrowing from the deep multi-view stereo, several approaches improve the cross-scene generalization ability of models by encoding scene geometry information.

Multi-view stereo. Multi-view stereo (MVS) is a classical computer vision problem, which aims to reconstruct the geometry from multiple viewpoints. Recently, learning-based multi-view stereo methods have been introduced and achieve impressive results. MVSNet [31] first leveraged the plane sweep-based cost volume formulation followed by 3D CNN for regularization to predict the depth maps. However, methods relying on 3D cost volumes are often limited to low-resolution input images, because 3D CNNs are generally time and GPU memory-consuming. Following works have extended this technique with recurrent plane sweeping [32], confidence-based aggregation [14] and multi-scale cost volumes [7, 24], improving the efficiency and reconstruction quality. In a coarse-to-fine fashion, early cost volumes are built on coarser-resolution features with sparsely sampled depth hypotheses, which result in relatively low volumetric resolution. Subsequently, estimated depth maps from earlier stages are used to narrow the depth range and construct thinner cost volumes on finer-resolution features. Drawing on these works, we propose to train depth prediction and radiance field reconstruction simultaneously. With the predicted depth, we construct efficient pyramid encoding volumes that provide multi-scale geometric information for the NeRF module.

NeRF based on multi-view. Although MVS methods can reconstruct the geometry of the scene, they are sensitive to potential inaccuracies in point clouds from corrupted depth, especially when there are thin structures and textureless regions. In contrast, NeRF-based methods model scenes as neural volumetric radiance fields and can reproduce the faithful scene appearance, producing photo-realistic novel views. To improve the generalization of NeRF, some multi-view rendering methods are proposed. PixelNeRF [33] takes spatial image features aligned to each pixel as an input and learns a scene prior for reconstructing. IBRNet [25] generates a continuous scene radiance field on-the-fly from multiple source views for rendering novel views. MVSNeRF [5] and ENeRF [10] utilize a plane swept 3D cost volume for geometric-aware scene understanding, using only a few images as input. Point-NeRF [29] models a volumetric radiance field with a neural point cloud using multi-view images. These works use a single-scale encoding volume, which is insufficient to provide adequate geometric information for scenes with objects of various scales. In this work, we propose constructing pyramid-structured encoding volumes to provide geometric information at all scales.

3 APPROACH

Our framework has the same input settings as MVSNeRF [5], i.e., additional N source images and their camera parameters are provided when performing view synthesis. As shown in Fig. 3, our

framework constructs pyramid neural encoding volumes to provide multi-scale geometric information for the NeRF module (Sec. 3.1). The multi-scale encoding volumes will be utilized to regress the depth maps and reconstruct the multi-scale radiance fields. The key to efficient pyramid volume construction is effectively reducing the depth range. To this end, we propose to train depth prediction and radiance field reconstruction simultaneously (Sec. 3.2). The additional input f^l to the NeRF module in Eq. 6 is trilinear interpolated from the l -th level encoding volume V^l . To make the geometric information in the point features more robust, we further explore enhancing the point volume feature f^l by fusing multi-scale volume features, and adaptive neighbor volume features, respectively (Sec. 3.3).

3.1 Multi-scale neural field reconstruction

In this section, we introduce our multi-scale neural field reconstruction framework. Our framework learns an encoding volume and reconstructs the density field and radiance field separately. During training, the pixel is rendered at each level to facilitate the acquisition of multi-scale geometric information. Conversely, during testing, the pixel is solely rendered at the final level.

Feature pyramid. Similar to MVSNet [31], a shared eight-layer 2D CNN is applied to extract deep features for all source images $\{I_i\}_{i=1}^N$, $I_i \in \mathbb{R}^{H \times W}$. Afterward, we use a feature pyramid network [11] to extract hierarchical feature maps F^1, \dots, F^L with different spatial resolutions, where L denotes the number of levels. These hierarchical feature maps will be used to construct multi-scale encoding volumes.

Encoding volume. Taking the l -th level as an example. For each pixel p in the reference view, we sample $|\mathcal{S}_A^l|$ points, where \mathcal{S}_A^l denotes the l -th level sampling points for encoding volume, within the depth range \mathcal{R}^l (note that the depth range can be different for different pixels) and build encoding features for each sampled point. Following learning-based methods [31], we construct encoding volume by warping the source features into reference views. Specifically, for a sample point with pixel coordinate p and depth value d , we first warp it to the i -th source view via inverse warping:

$$p_i = K_i T_i (d \cdot K^{-1} p), \quad (1)$$

where K and K_i are the intrinsics for the reference and the i -th source image, respectively. T_i is the relative transformation from reference to i -th source image. p_i is the warped pixel location in the i -th source image. Based on the inverse warping, we construct the encoding volume $\bar{V}^l \in \mathbb{R}^{|\mathcal{S}_A^l| \times \frac{H}{2^{L-l}} \times \frac{W}{2^{L-l}} \times C}$ by computing the variance of multi-view source features for each point.

$$\bar{V}^l = \text{Var}(F_1^l[p_1], \dots, F_N^l[p_N]), \quad (2)$$

where N denotes the number of source images, $F_N^l[p_N]$ represents the bilinearly interpolated feature at the l -th level feature map of the N -th source image, using the warped pixel coordinate p_N .

To smooth the noise contaminated in the raw encoding volume, we apply a deep 3D U-Net to regularize the encoding volume and build a new encoding volume, and represent the similarity of the source features on different depth values. This process is expressed

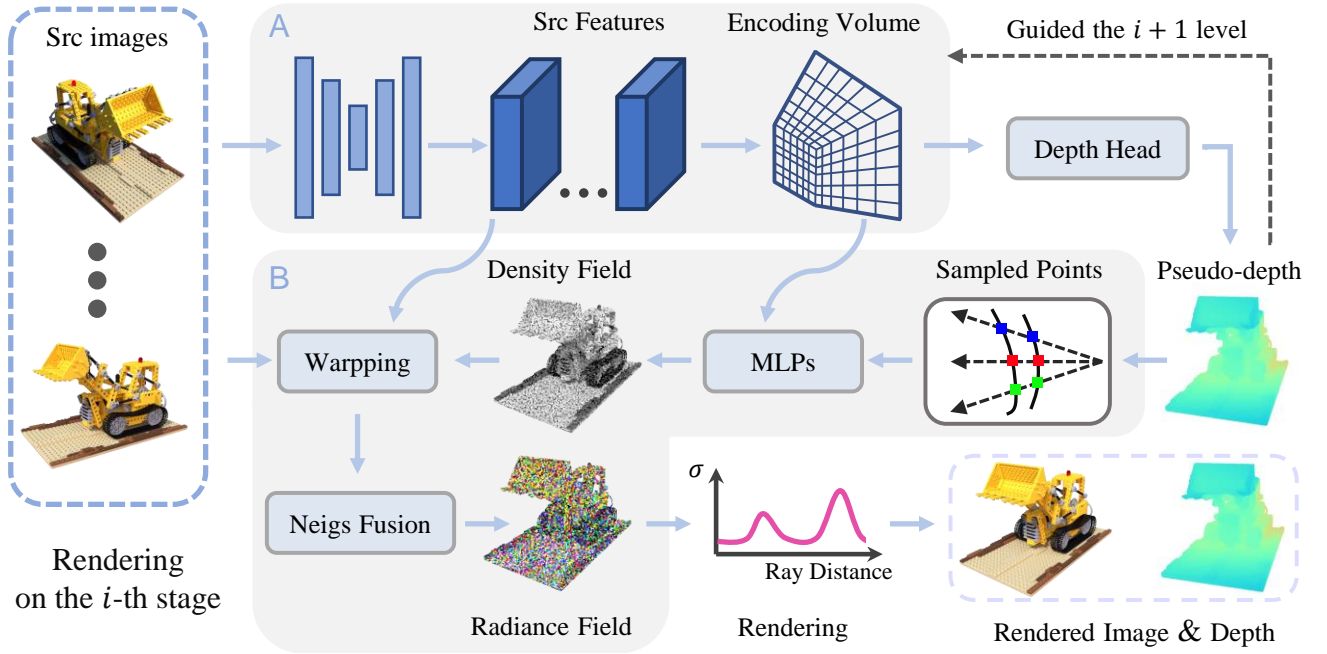


Figure 3: Overview of our pipeline, *A* represents the encoding volume reconstruction, while *B* represents the neural field reconstruction. Our framework first identify a set of neighboring source views and extract hierarchical feature maps. Then, we apply inverse warping, 3D CNN, and auxiliary depth heads to reconstruct encoding volumes and estimate depth maps at each stage. Afterward, the multi-scale neural fields (including density field and radiance field) are reconstructed by trilinear interpolation and MLPs. Finally, we apply depth-guided adaptive feature fusion to enhance the radiance field obtained by fusing neighbors features. Differentiable ray marching uses multi-scale neural fields for final rendering.

by:

$$V^l = \text{UNet}_{3D}(\bar{V}^l). \quad (3)$$

Density field reconstruction. The higher similarity, the greater the probability that the point is located on the surface of the object, which means that the encoding volume contains the density information of the scene. Existing methods construct the density field using the 3D position and viewing direction of the point, which lacks guidance from global and local information and leads to discontinuities in the rendered depth map. Different from these methods, we reconstruct the density by sharing the encoding volume, with a tiny MLP network.

Specifically, along a ray passing through the pixel p of the l -th level, we sampled a set of points S_B^l for neural field reconstruction. For each point $s_b \in S_B^l$, we regress the point's density σ using a tiny MLP network MLP_1 , which only contains two layers MLP. The first layer is used for fusing the feature volume and the direction information, and the second layer aims to regress the density:

$$\sigma, f_h = \text{MLP}_1(f_v, \Delta d), \quad (4)$$

where f_h is the hidden features, and f_v is the volume feature, obtained by trilinear interpolation in the l -th level encoding volume V^l at position x , 3d position of the point s_b . Δd is defined as the concatenation of the direction of $d_i - d$ from all source views, d and d_i are the ray directions of the point under the reference view

and corresponding source view.

$$\Delta d = \text{Concat}(d_1 - d, \dots, d_N - d). \quad (5)$$

Radiance field reconstruction. Same with density field reconstruction, we also perform volume attribute regression at each level. Following with IBRNet [25], we regress the point's s_b radiance, viewed in direction d , by predicting blending weights for the image colors $\{c_i\}_{i=1}^N$ in the source views. The weight of the i -th source view is predicted by an MLP network MLP_2 :

$$w_i = \text{MLP}_2(\Delta d, f_h, c_i, f_i), \quad (6)$$

where c_i, f_i are the RGB value and the image features of the pixel at the projected position of location x in the i -th source view. The color of the point s_r location in the direction d is blended via a soft-argmax operator, as the following:

$$r = \sum_{i=1}^N \frac{\exp(w_i)c_i}{\sum_{j=1}^N \exp(w_j)}. \quad (7)$$

Volume rendering. To render the pixel color and density, we perform the differentiable ray marching [17] based on the reconstructed radiance fields $\{\sigma^l, r^l\}_{l=1}^L$. During testing, we only perform the rendering at the final level. In contrast, we conduct rendering at each level during the training. By supervising the color and density at each level, the model can more effectively learn multi-scale information, thereby improving its ability to generalize.

Specifically, for the l -th level, the color \tilde{c}_l and depth \tilde{d}_l of the pixel p is calculated by accumulating the radiance of sampling points at each level, and is given by:

$$\tilde{c}_l = \sum_{k=1}^{|\mathcal{S}_B^l|} \tau_k (1 - \exp(-\sigma_k)) r_k, \quad \tilde{d}_l = \sum_{k=1}^{|\mathcal{S}_B^l|} \sigma_k z_k \quad (8)$$

$$\text{where } \tau_k = \exp\left(-\sum_{j=1}^{k-1} \sigma_j\right), \quad (9)$$

and z_k is the depth value of the point $s_b \in \mathcal{S}_B^l$ in the reference view. And τ represents the volume transmittance.

3.2 Auxiliary depth prediction head

Due to the lack of depth guidance, MVSNeRF sample points uniformly across the entire depth range for each pixel. The number of sample points is a trade-off between rendering precision and rate. Inspired by CasMVSNet [7], we append a depth prediction head F_{depth} (a convolutional layer) after each scale encoding volume, and predict the depth map D^l with confidence U^l in the target view:

$$D^l, U^l = F_{\text{depth}}(V^l). \quad (10)$$

Based on the encoding volume, the core of the depth prediction head compresses the similarity of source features. As a result, the predicted depth can have higher accuracy compared to the regression approach [31], particularly in regions with abundant texture information.

Utilizing the predicted depth map, We first supervise the pixel depth calculated by volume rendering, as described in Sec 3.4. Besides, we can reduce the number of sampling points for radiance field reconstruction, and narrow down the depth range of the next level, thus resulting in better rendering precision with fewer sample points:

$$\begin{aligned} |\mathcal{S}_B^l| &= |\mathcal{S}_A^l| * \alpha_l, \\ \mathcal{R}^{l+1} &= \mathcal{R}^l * \beta_l, \end{aligned} \quad (11)$$

As described in Sec. 3.1, \mathcal{S}_A^l and \mathcal{S}_B^l denote the l -th level sampling points for encoding volume and radiance field reconstruction, respectively. Points \mathcal{S}_A^l , with $N_A^l = |\mathcal{S}_A^l|$ sampled points, are uniformly sampled in the depth range \mathcal{R}^l for the l -th level, and the reducing factor $\alpha_l (< 1)$ is used to decrease the number of points for radiance field reconstruction. Specifically, for \mathcal{S}_B^l , we select $N_B^l = |\mathcal{S}_A^l| * \alpha_l$ points closest to the predicted depth from \mathcal{S}_A^l . In addition, we narrow the depth range of the next level \mathcal{R}^{l+1} to $\mathcal{R}^l * \beta_l$ centered on the predicted depth value, in which $\beta_l < 1$ is the reducing factor. By default, we use three levels and set $[N_A^l], [\alpha_l], [\beta_l]$ to $[48, 32, 8], [1/6, 1/4, 1/2]$ and $[1/6, 1/16]$ respectively.

3.3 Depth-guided adaptive feature fusion

The geometric information contained in the point volume feature f^l may be inaccurate due to occlusion, lighting, etc. This section introduces our feature fusion strategies to enhance the point volume features.

The MVS methods [24, 27] adaptively aggregate cost volume features to improve matching robustness, using reference view

images as input to predict adjacent pixel offsets. However, reference view images are not available in the view synthesis task. To this end, we propose to use the predicted depth map $D^l \in \mathbb{R}^{H \times W \times 1}$ to adaptively enhance the hidden volume features f^h , projected colors c_i and projected features f_i in Eq. 6.

For the adaptive fusion of point volume features f_v in Eq. 6, we employ the depth map and inverse warped source image features as inputs to predict the pixel offsets. Specifically, we first obtain the warped features in the reference view by inverse warping, taking the predicted depth map and source image features as input:

$$\hat{F}^l = \text{Concat}(IW(F_1^l, D^l), \dots, IW(F_N^l, D^l)), \quad (12)$$

where F_i denotes the features of the i -th source image, N denotes the number of source images, D^l is the predicted depth map in the reference view, $IW(\cdot)$ denotes the inverse warping in Eq. 1, and \hat{F}^l is the warped image features in the reference view. Then, we concatenate the depth map and warped features and apply a 2D CNN to predict the pixel offsets $\{\Delta p_k\}_{k=1}^K$ of the pixel p . The final neighbor volume feature fusion is obtained using a single-layer multi-head cross-attention module F_{MHCA} , taking the volume feature at location p as the query and the volume features at locations $\{p + \Delta p_k\}_{k=1}^K$ as keys and values:

$$\tilde{f}_v = F_{\text{MHCA}}(V(p), \{V(p + \Delta p_k)\}_{k=1}^K). \quad (13)$$

For projected colors and features c_i, f_i in Eq. 6, we employ the same adaptive fusion approach, taking adaptive colors c_i fusion as an example. Specifically, for each predicted neighborhood pixel $\hat{p}_k = p + \Delta p_k$, we project it to the corresponding source images and obtain the project color value $c_{k,i}$:

$$c_{k,i} = I_i[p_{k,i}], \quad (14)$$

where $p_{k,i}$ is the projected pixel location in the i -th source image I_i . The final colors input to the NeRF module is the attention warped colors of predicted neighbors, which can be formulated as:

$$\tilde{c}_i = F_{\text{MHCA}}(c_i, \{c_{k,i}\}_{k=1}^K). \quad (15)$$

3.4 End-to-end training

In particular, we train our full pipeline on the DTU dataset [1], which can learn a powerful generalizable function, reconstructing radiance fields across scenes from only three input nearby source images.

For the auxiliary depth prediction, we adopt the same loss function as in unsupervised MVS tasks [9] to improve the accuracy of predicted depth. This, in turn, helps more effectively supervise the depth of the rendering branch, guide the sampling points, and adjust the depth range. The loss is defined as follows:

$$L_{mvs} = L_{PC} + L_{SSIM} + L_{Smooth} \quad (16)$$

where, L_{PC} , L_{SSIM} , and L_{Smooth} are the photometric consistency loss, structured similarity loss, and depth smoothness loss. Please refer to [9] for more details.

We optimize the outputs of our rendering approach, including the color \tilde{c} and depth \tilde{d} , with the ground truth target view I and the pseudo depth D . The auxiliary depth prediction head predicts the pseudo-depth based on the similarity of source features. Although it may have higher accuracy in texture-rich regions, it may be

Table 1: Quantitative results of novel view synthesis. We show averaged results of PSNRs, SSIMs, and LPIPs on three different datasets. For the Realistic Synthetic NeRF dataset, the two numbers in each item refer to the evaluation of the central/entire regions of the novel images.

Method	NeRF Synthetic			DTU			Real Forward-Facing		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
PixelNeRF	-/7.39	-/0.658	-/0.411	19.31	0.789	0.382	11.24	0.486	0.671
IBRNet	21.91/22.44	0.857/0.874	0.203/0.195	26.04	0.971	0.191	21.79	0.786	0.279
MVSNeRF	23.62/24.63	0.897/0.929	0.176/0.155	26.63	0.931	0.168	21.93	0.795	0.252
ENeRF	23.25/26.65	0.893/0.947	0.152/0.072	27.61	0.956	0.091	22.78	0.808	0.209
Ours	24.88/26.76	0.902/0.950	0.110/0.067	28.52	0.952	0.110	22.78	0.781	0.243

unreliable in texture-less regions. Therefore, we only consider pixels with a confidence score greater than δ to ensure the reliability of the predicted depth. The loss is defined as follows. For simplicity, we show the calculation for one pixel p .

$$L_{render} = \sum_{p \in I} \lambda_1 \|\tilde{c}(p) - I(p)\|_2^2 + \lambda_2 \mathbb{I}(U(p) > \delta) \|\tilde{d}(p) - D(p)\|, \quad (17)$$

where \mathbb{I} is the indicator function, λ_i sets the influence of each loss. The final objective from each level can be constructed as follows:

$$L = \sum_{l=1}^3 L_{render}^l + \lambda_3 L_{mos}^l. \quad (18)$$

4 EXPERIMENTS

In this section, we first describe the details including the datasets and implementation details, then we show the qualitative and quantitative comparisons of our network with state-of-the-art methods. Finally, we perform ablation studies to validate the effectiveness of our proposed method.

4.1 Dataset and network details

Dataset. We train our framework end-to-end on the DTU [1] dataset to learn a generalizable network, using the same training and testing splits as MVSNeRF [5]. The DTU dataset is divided into 88 training scenes and 16 testing scenes with an image resolution of 512×640 . We evaluate the generalization ability of our method on the Realistic Synthetic NeRF [17] data and the Forward-Facing data [16] by using the model merely trained on the DTU dataset. They both include 8 complex scenes that have a different distribution from DTU.

Network details. We set $N = 3$ for nearby source images, and implement the image feature extraction network using a FPN [11] like architecture. The number of depth hypotheses $[N_r^l]$ for encoding volume and sampling points $[N_d^l]$ for radiance field reconstruction at three stages is set to $[48, 32, 8]$ and $[8, 8, 4]$, respectively. And the number of adaptive neighbors K is set to 8. For the weight of each loss λ_i , λ_1 and λ_3 are set same value of 1.0, and λ_2 the weight of the depth loss grows linearly based on the training steps from $1e^{-4}$ to $1e^{-2}$. The MLP decoder in rendering is similar to the MVSNeRF [5].

We train our network using four V100 GPUs with a batch size of 1. During training, 2048 pixels are randomly sampled from one

Table 2: Quantitative results of reconstructed depth. We evaluate our depths and points reconstruction on the DTU and compare it with other methods. Our method significantly outperforms other neural rendering methods.

Novel-view depth metric			
Method	Abs err \downarrow	Acc(2mm) \uparrow	Acc(10mm) \uparrow
PixelNeRF	47.8	0.039	0.187
IBRNet	324	0.000	0.866
MVSNeRF	7.00	0.717	0.866
ENeRF	4.60	0.792	0.917
Ours	4.29	0.867	0.937
Reconstructed points metric			
Method	Overall \downarrow	Acc \downarrow	Comp \downarrow
MVSNeRF	0.588	0.641	0.534
ENeRF	2.055	2.354	1.755
Ours	0.368	0.355	0.380

novel viewpoint. The model is trained with the Adam optimizer for 8 epochs with a base learning of $5e^{-4}$. The learning rate is halved iteratively at the 2-th, 4-th, and 8-th epoch.

4.2 Evaluation results

To evaluate our model, we compare it against current top-performing techniques for view synthesis, detailed below. We adopt the PSNR, SSIM, and LPIPS as the quantitative results for the rendered images. We also evaluate the geometry reconstruction quality by comparing depth and points reconstruction results on the DTU dataset.

Quality of rendered image. Tab. 1 shows the comparison results with recent concurrent works with PSNR, SSIM, and LPIPS; For Realistic Synthetic NeRF [17], objects are primarily located in the central area, with the surrounding regions consisting of a white background. Existing methods typically utilize cropping to evaluate only the central portion of the image, while some methods evaluate the entire image region. In this paper, we explore two different approaches. Our results on the DTU are significantly better than the comparison methods, where also 3 views are provided for the source images. To show that our method has good generalizability, we also evaluate our model on the NeRF Synthetic and Real

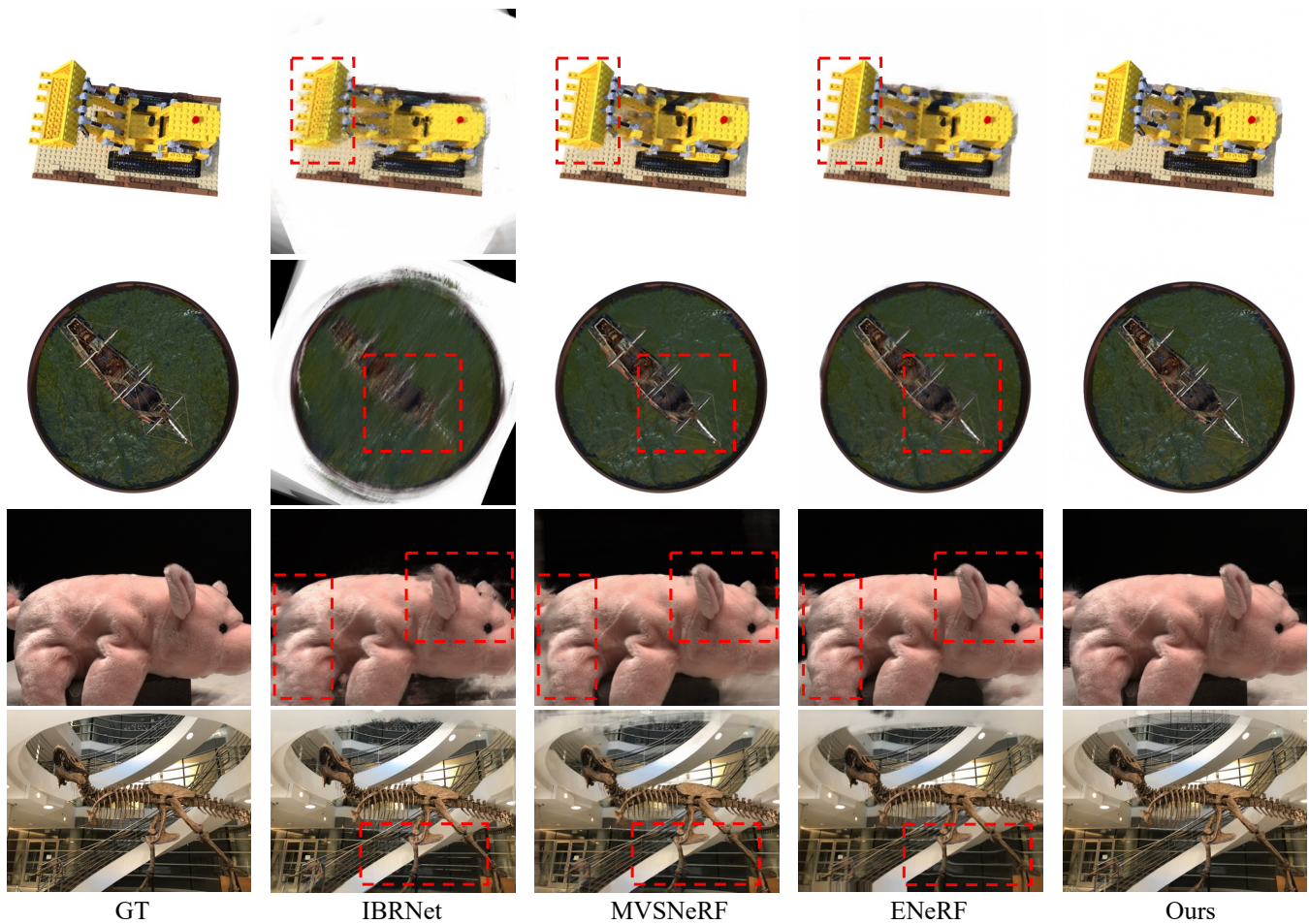


Figure 4: Qualitative comparison of rendering quality on diverse scenes between our method and state-of-the-art counterparts. We show the results of all three datasets. Our method renders better textures than other methods. This is particularly evident in the scenes highlighted in the red dashed boxes. (Best viewed by zooming in).

Forward-Facing datasets, following the experiment setting in MVSNeRF. For the NeRF Synthetic dataset, the quality of our results is further boosted significantly, leading to the best PSNR, SSIM, and LPIPS in all compared methods. And also generates results that are comparable to ENeRF on the Real Forward-Facing dataset. As shown in Fig. 4, our results on these scenes are of very high visual quality compared with other methods.

Quality of reconstructed depth. To improve the quality of geometry, our approach reconstructs a density field by sharing the encoding volume. We evaluate our geometry reconstruction quality by comparing depth reconstruction and points reconstruction results. To reconstruct the final point, we follow [31] to fuse the depth from multiple views. We compare our approach with the recent multi-view based methods on the DTU testing set.

As shown in Tab. 2, benefit to our shared encoding volume, our approach achieves significantly more accurate depth rendering than others. Especially for 2mm accuracy, the metric is improved by 9.3% compared with ENeRF. For the quantitative evaluation of point

cloud, we calculate the accuracy and completeness by the MATLAB code provided by the DTU dataset. We can see that our method outperforms other methods in both completeness and accuracy quality and rank the first place.

The qualitative results of depths and points are shown in Fig. 2 and Fig. 5. For the qualitative depths results, our method renders more accurate depth maps and includes more texture. Thanks to the accurate depths, ours generates more complete point clouds with finer details.

4.3 Ablations and analysis

To further demonstrate the contribution of our method, we introduce some groups of ablation studies on the DTU dataset in this section.

Contribution of each proposed module. In this section, we conduct experiments to verify the effectiveness of our proposed neural radiance field reconstruction. As shown in Tab. 3, when reconstructing the density field using the 3D position and viewing

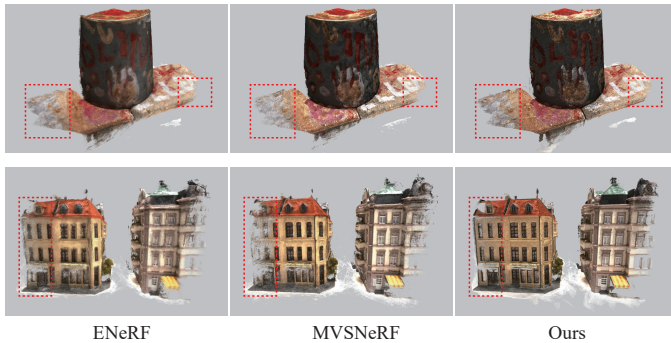


Figure 5: Qualitative comparisons between other rendering methods. we follow [31] to fuse the depth from multiple views and reconstruct points. Our method achieves increasingly dense reconstruction.

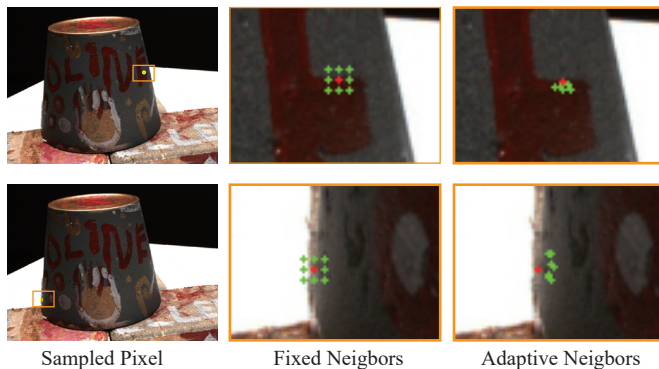


Figure 6: Visualization of adaptively sampled neighbors. The left column shows the center pixels, the middle column shows the eight fixed neighbors and the right column shows the learned adaptive locations.

direction of the point, not shared volume, the rendering quality drops a lot. Besides, we remove the feature fusion module for comparison, this module shows the same good performance. Finally, we evaluate the influence of the pseudo-depth on photo-realistic quantity. As shown in the 3-th row, using pseudo-depth supervision also improves the rendering performance.

Contribution of depth-guided adaptive feature fusion. To verify the effectiveness of the proposed depth-guided adaptive feature fusion, we conduct experiments comparing with no-neighbor and different numbers of neighbors setting. As shown in Tab. 4, our adaptive method performs the best in various settings of the number of neighbors. As shown in Fig. 6, the sampled neighbors tend to the locations which have salient features and have the same texture as the center pixel.

Ablation on pseudo-depth supervised. The auxiliary depth prediction head predicts a more accurate depth map compared to the regression approach, described in Sec. 3.2. The pseudo-depth supervising can help to reconstruct a more effective density field,

Table 3: Ablation studies for our proposed neural radiance field reconstruction of the photo-realistic quantity on DTU evaluation set.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
w/o shared volume	28.16	0.949	0.113
w/o feature fusion	28.26	0.950	0.110
w/o pseudo-depth supervised	28.08	0.948	0.115
Ours	28.52	0.952	0.110

Table 4: Ablation experiments for depth-guided adaptive feature fusion.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
w/o neigs	28.26	0.950	0.110
4-neigs	28.38	0.951	0.109
8-neigs (Ours)	28.53	0.952	0.108

Table 5: Ablation experiments for pseudo-depth supervision of the geometry quantity on the DTU evaluation set. $\delta = 0.0$ means all pixels are used to supervise the rendered depths, $\delta = 1.0$ means not using the pseudo-depth supervised, and value between $[0, 1]$ means using the pixels with a confidence score greater than δ .

Setting	Abs err \downarrow	Acc(2mm) \uparrow	Acc(10mm) \uparrow
$\delta = 0.0$	4.42	0.852	0.925
$\delta = 0.3$	4.32	0.867	0.935
$\delta = 0.5$ (Ours)	4.29	0.867	0.937
$\delta = 0.7$	4.33	0.858	0.937
$\delta = 1.0$	5.36	0.786	0.920

and the photo-realistic quantity is shown in Tab. 3. To evaluate the geometry quantity, we set five experiments in Tab 5. The results show that more accurate pixels can benefit to improve the geometry quantity.

5 CONCLUSION

This work proposes an end-to-end framework for generalizable radiance field reconstruction. In this framework, an auxiliary depth prediction head is additionally learned to provide depth priors when performing point sampling and pyramid encoding volumes construction. The constructed pyramid volumes provide multi-scale geometric information of the scene, thereby improving rendering performance across scenes. To mitigate the effects of factors such as occlusion and illumination on the constructed encoding volumes, we enhance the volume features by depth-guided adaptive feature fusion. Our method generalizes well across diverse testing datasets and can significantly outperform concurrent works on photo-realistic and geometry of rendered novel views.

REFERENCES

- [1] Henrik Aanaes, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. 2016. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision* 120, 2 (2016), 153–168.
- [2] Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. 2020. Neural reflectance fields for appearance acquisition. *arXiv preprint arXiv:2008.03824* (2020).
- [3] Sai Bi, Zexiang Xu, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. 2020. Deep reflectance volumes: Rightable reconstructions from multi-view photometric images. In *European Conference on Computer Vision*. Springer, 294–311.
- [4] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. 2021. Nerf: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 12684–12694.
- [5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. 2021. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14124–14133.
- [6] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. 2022. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12882–12891.
- [7] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuoqiuo Dai, Feitong Tan, and Ping Tan. 2020. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2495–2504.
- [8] Kacper Kania, Kwang Moo Yi, Marek Kowalski, Tomasz Trzciński, and Andrea Tagliasacchi. 2022. Conerf: Controllable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18623–18632.
- [9] Tejas Khot, Shubham Agrawal, Shubham Tulsiani, Christoph Mertz, Simon Lucey, and Martial Hebert. 2019. Learning unsupervised multi-view stereopsis via robust photometric consistency. *arXiv preprint arXiv:1905.02706* (2019).
- [10] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiao-wei Zhou. 2022. Efficient Neural Radiance Fields for Interactive Free-viewpoint Video. In *SIGGRAPH Asia 2022 Conference Papers*. 1–9.
- [11] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.
- [12] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. 2020. Neural sparse voxel fields. *Advances in Neural Information Processing Systems* 33 (2020), 15651–15663.
- [13] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. 2019. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751* (2019).
- [14] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. 2019. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10452–10461.
- [15] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. 2021. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7210–7219.
- [16] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. 2019. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* 38, 4 (2019), 1–14.
- [17] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.
- [18] Barbara Roessle, Jonathan T Barron, Ben Mildenhall, Pratul P Srinivasan, and Matthias Nießner. 2022. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12892–12901.
- [19] Johannes Lutz Schönberger and Jan-Michael Frahm. 2016. Structure-from-Motion Revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [20] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. 2019. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2437–2446.
- [21] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. 2021. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7495–7504.
- [22] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretschmar. 2022. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8248–8258.
- [23] Haitthem Turki, Deva Ramanan, and Mahadev Satyanarayanan. 2022. Mega-NeRF: Scalable Construction of Large-Scale NeRFs for Virtual Fly-Throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12922–12931.
- [24] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. 2021. Patchmatchnet: Learned multi-view patchmatch stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14194–14203.
- [25] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. 2021. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4690–4699.
- [26] Qianyi Wu, Xian Liu, Yuedong Chen, Kejie Li, Chuanxia Zheng, Jianfei Cai, and Jianmin Zheng. 2022. Object-compositional neural implicit surfaces. In *European Conference on Computer Vision*. Springer, 197–213.
- [27] Haofei Xu and Juyong Zhang. 2020. Aa-net: Adaptive aggregation network for efficient stereo matching. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1959–1968.
- [28] Hongbin Xu, Zhipeng Zhou, Yu Qiao, Wenxiong Kang, and Qiuxia Wu. 2021. Self-supervised Multi-view Stereo via Effective Co-Segmentation and Data-Augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [29] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. 2022. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5438–5448.
- [30] Tianhan Xu and Tatsuya Harada. 2022. Deforming Radiance Fields with Cages. In *European Conference on Computer Vision*. Springer, 159–175.
- [31] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. 2018. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*. 767–783.
- [32] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. 2019. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 5525–5534.
- [33] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. 2021. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4578–4587.
- [34] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yüewen Ma, Rongfei Jia, and Lin Gao. 2022. NeRF-editing: geometry editing of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 18353–18364.
- [35] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. 2020. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492* (2020).
- [36] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyfe, and Noah Snavely. 2018. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817* (2018).