

Disentangled Generation and Aggregation for Robust Radiance Fields

Shihe Shen^{1*}, Huachen Gao^{1*}, Wangze Xu¹, Rui Peng^{1,2},
Luyang Tang^{1,2}, Kaiqiang Xiong^{1,2}, Jianbo Jiao³, and Ronggang Wang^{1,2,✉}

¹ School of Electronic and Computer Engineering, Peking University

² Peng Cheng Laboratory

³ School of Computer Science, University of Birmingham

{shshen0308, gaohuachen712}@gmail.com rgwang@pkusz.edu.cn

Abstract. The utilization of the triplane-based radiance fields has gained attention in recent years due to its ability to effectively disentangle 3D scenes with a high-quality representation and low computation cost. A key requirement of this method is the precise input of camera poses. However, due to the local update property of the triplane, a similar joint estimation as previous joint pose-NeRF optimization works easily results in local minima. To this end, we propose the Disentangled Triplane Generation module to introduce global feature context and smoothness into triplane learning, which mitigates errors caused by local updating. Then, we propose the Disentangled Plane Aggregation to mitigate the entanglement caused by the common triplane feature aggregation during camera pose updating. In addition, we introduce a two-stage warm-start training strategy to reduce the implicit constraints caused by the triplane generator. Quantitative and qualitative results demonstrate that our proposed method achieves state-of-the-art performance in novel view synthesis with noisy or unknown camera poses, as well as efficient convergence of optimization. Project page: <https://gaohchen.github.io/DiGARR/>.

Keywords: NeRF · Disentangle · Pose Estimation · Novel View Synthesis

1 Introduction

Recently, adopting neural networks has become increasingly popular for Novel View Synthesis (NVS), where the Neural Radiance Fields (NeRF) [36] has brought a great surge in high-fidelity synthesis quality. NeRF represents a 3D radiance field by multi-layer perceptrons (MLPs), and renders novel views by differentiable volume rendering [22]. A crucial prerequisite for achieving promising rendering results with NeRF is the precise annotated camera parameters. However, accurate camera poses are not easily attainable, which heavily relies on external Structure-from-Motion (SfM) algorithms like COLMAP [47].

* Equal contribution.

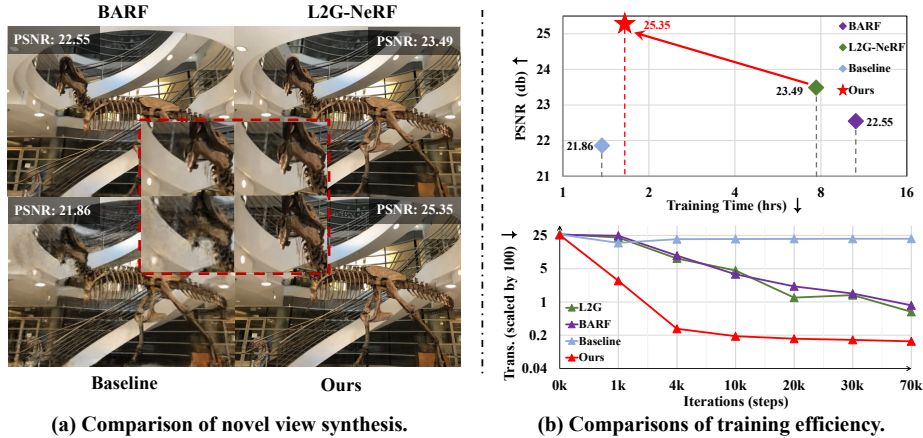


Fig. 1: Comparisons of novel view synthesis and training efficiency on *Trex*. (a) shows novel view synthesis results of different works, where the baseline is the direct joint pose-triplane optimization; (b) shows efficiency comparisons where the upper shows training time and PSNR, and the lower shows iterations and translation errors.

To this end, there has been a growing focus on methods for reducing the reliance on camera parameters by optimizing NeRF and camera parameters simultaneously [4, 9, 11, 28, 54]. However, based on the vanilla NeRF-MLP representation, most existing approaches demand hours even days for training with modern powerful GPUs, severely limiting their use in practical applications. Therefore, it is significant to improve training efficiency. An intuitive idea for acceleration is to explore alternative representations for MLP, such as voxels [14, 50], hash-grids [37], triplanes [5, 13] and Gaussian splats [23]. In this paper, we propose a novel triplane-based 3D representation to estimate camera poses and the 3D scene, with both high efficiency and quality. Triplane is a highly disentangled explicit representation of 3D scenes, which has been widely applied in many recent works [5, 7, 13, 19]. With a high data compression ratio, low computational cost and comparable performance, triplane is more concise and scalable than volumes, thus appropriate to accelerating pose-NeRF optimization. Meanwhile, as a structured 3D representation with fixed-size feature maps, triplane is more controllable and suitable for bundle adjusting tasks, while concurrent pose-3DGS joint optimizations [15, 26] are primarily restricted to video streams or ordered image collection.

However, introducing a triplane to the joint optimization is nontrivial. We observe that the naive direct combination of pose estimation and triplane radiance field is prone to get trapped into local-minima, especially in the early training stage, and could hardly get rectified, resulting in unsatisfactory view synthesis results (Fig. 1 (a)). The main reason lies in the following two aspects. (1) The triplane radiance field follows a local updating policy on each feature grid as shown in Fig. 2 (a), since features are derived from interpolation and planes are

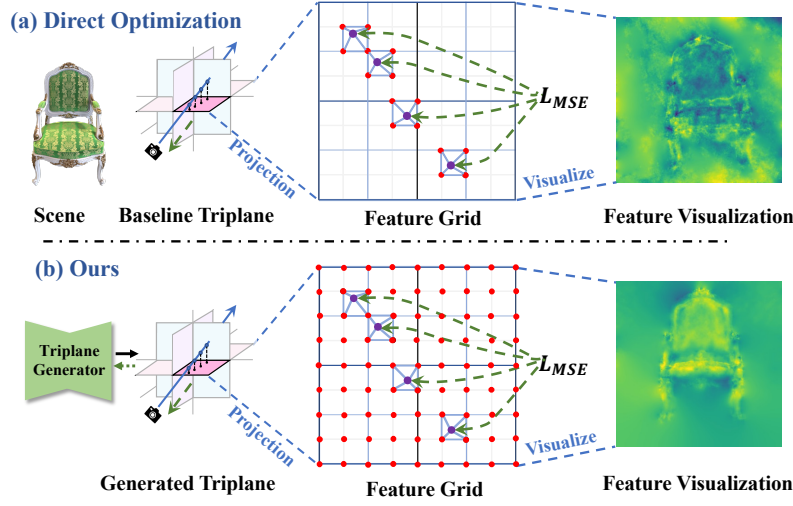


Fig. 2: Illustration of the local updating on feature grids. Feature points in **red** on the planes will be updated during one forward pass, lines in **green** represent gradient flow. (a) shows the updated feature points by direct optimization, and (b) shows the updated feature points after adopting the triplane generator.

directly updated only in a few grids used for interpolation. When the camera poses and the triplane are ambiguous, the inaccurate poses lead to ray transmission bias. The local updating property of the triplane incorrectly updates the projected feature of sampled 3D points on biased rays. As the training proceeds, the error from local updating can not be rectified, where the feature parameters on the triplane fail to perceive the global contextual information. This causes additional ambiguity between camera poses and scene reconstructions, finally resulting in local minima of the joint optimization. (2) As triplane decoupling 3D scene into three orthogonal planes, anisotropic features need to be aggregated among independent planes. With different learning complexities of planes and the introduced entanglement with pose and planes, commonly used aggregations as sum [6] and production [13] impose conflicting signals on pose optimization, thus failing to achieve accurate pose estimation and expressive scene representation simultaneously.

To address the above issues of pose and triplane joint optimization, we propose a triplane-based representation with disentangled generation and aggregation. Firstly, we obtain triplane features from a disentangled generator, with frozen noise tokens of each plane as input, possessing both global smoothness in implicit representation and training efficiency in explicit representation. As shown in Fig. 2 (b), by parameterizing the triplane with triplane generator, network parameters are shared by all grids on the plane, and the information for feature updating will radiate around the plane, introducing global context to eliminate local error accumulating without additional constraints or post-processing steps

on triplane. Secondly, digging into how triplane feature aggregation impacts camera poses during joint optimization, we design a novel Disentangled Plane Aggregation (DPA) to relieve optimizing collision among three feature planes for a single 3D sampled point. The DPA disentangles pose with triplane feature by distributing updating signals to pose and triplane respectively, realizing robust and unambiguous joint optimization of pose and scene representation. Additionally, the smoothness introduced by the triplane generator has a side effect on high-frequency details of scene representation. To preserve our high-frequency feature and promote scene expressiveness, we present a two-stage warm-start training strategy. Inheriting the relatively accurate poses and a coarse triplane along with the MLP decoder from the first-stage training, we seamlessly transform to a direct feature optimization joint with pose refinement in the second stage.

In summary, the primary contributions are as follows:

- We present a hybrid disentangled scene representation based on triplane for joint estimation of camera poses and novel view synthesis.
- We leverage the triplane generator for joint estimation to address local errors in feature grids caused by the local updating in baseline pose-triplane optimization. To the best of our knowledge, our method is the first attempt to incorporate the deep network prior as a hybrid representation in the field of joint pose-NeRF optimization.
- We analyze the effect of different aggregation approaches on the pose estimation, and design a new disentangled plane aggregation method. We also introduce a two-stage warm-start training strategy to mitigate the smoothing on triplane grids caused by the triplane generator.
- Qualitative and quantitative results on real-world LLFF [35] and NeRF-synthetic [36] dataset show that our method achieves state-of-the-art performance on both novel view synthesis and pose estimation.

2 Related Work

Novel View Synthesis. NeRF [36] has become a popular representation in the field of novel view synthesis for its high-fidelity rendering results. Many follow-up works are proposed to improve NeRF’s performance. [1, 20] replace the ray casting with anti-aliased cone tracing, and [2] is further proposed to handle the unbounded scene with non-linear scene parameterization. [61] separates foreground and background using proposed sampling algorithms, [24, 38, 48] uses additional constraints, while [12, 46, 49, 55] adopts depth priors to improve NeRF’s geometry learning. NeRF has also gained many vision or graphic applications, such as surface reconstruction [39, 43, 52], dynamic scene reconstruction [5, 13, 27, 41, 42, 45], and single view reconstruction [8, 19, 30, 57]. Most recently, 3D Gaussian Splatting [23] also demonstrated strong capability in novel view synthesis.

To address the slow rendering speed of the vanilla MLP-based NeRF, many efforts tried to apply explicit grid-based representations, such as multi-resolution hash encoding [3, 37], voxel grids [14, 29, 50, 59] and triplanes [6, 7, 13, 20]. The

triplane representation leads to faster optimization and more compact modeling than other representations, which is widely adopted in recent works [19, 53, 62].

Joint Pose and NeRF Optimization. Pose-NeRF joint optimization has been widely studied recently. iNeRF [58] first shows the ability of pose estimation using reconstructed NeRF models. [54] first jointly estimates camera intrinsic, extrinsic and NeRF. [28] proposed a coarse-to-fine positional encoding method. [21] estimated camera distortion and proposed a geometric regularization. [33, 60] adopted GANs, while [11, 56] modified different activations for more suitable pose estimation. Recently, [4, 17] utilized mono-depth estimation [16, 44] for geometry guidance, which is designed for long-sequence images. [9] proposed a local-to-global registration to estimate poses from learned multiple local poses. [51] used pre-computed correspondences as priors for sparse view settings. [10, 25] optimized NeRF without any pose initialization. [31, 34] are proposed for the long sequence static and dynamic videos respectively. [18] accelerated previous pose-NeRF joint optimization using multi-resolution hash encoding. In this paper, we propose to use a fast and compact triplane for joint estimation.

3 Preliminaries

In this section, we first introduce the common pipeline of pose and NeRF joint learning, then we present the pipeline of our baseline directly optimization of triplane with noisy or unknown pose prior (denote as baseline).

Formulation of Pose-NeRF Optimization. Given a set of images $\{I_i\}_{i=1}^N$ with camera intrinsics and noisy or *unknown* extrinsics $\{\mathbf{T}_i\}_{i=1}^N$, our goal aims to reconstruct triplane-based radiance field and estimate the accurate camera poses \mathbf{T} . Denote a sampled 3D point $\mathbf{x} \in \mathbb{R}^3$ along a camera ray emitted from camera i , the color \mathbf{c} and density σ at \mathbf{x} can be derived from a MLP-based NeRF $f: \mathbb{R}^3 \rightarrow \mathbb{R}^4$ as $\mathbf{c}, \sigma = f(\mathbf{x}; \Theta)$. A synthesized image \hat{I} can be rendered along camera rays $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ between near and far plane t_n and t_f . The volume rendering function $\hat{\mathcal{I}}$ [22, 32] is formulated as:

$$\hat{\mathcal{I}}(\mathbf{r}) = \int_{t_n}^{t_f} W(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \quad (1)$$

where $W(t) = \exp(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds)$ is accumulated transmittance along the ray. The photometric loss is $\mathcal{L}_{Render} = \sum_{i=1}^N \|I_i - \hat{I}_i\|_2^2$. Taking the photometric loss as total loss \mathcal{L} , triplane NeRF parameters Θ and camera poses \mathbf{T} which are used to cast the ray can be optimized by minimizing the total loss:

$$\mathbf{T}^*, \Theta^* = \underset{\mathbf{T}, \Theta}{\operatorname{argmin}} \mathcal{L}(\hat{\mathbf{T}}, \hat{I}|I), \quad (2)$$

where $\hat{\mathbf{T}}$ denotes that learnable camera parameters are updated during optimization, and \mathbf{T}^*, Θ^* are the final optimized parameters.

Baseline Joint Optimization of Pose and Triplane. In our baseline, the 3D scene can be factorized by the triplane representation which contains three axis-aligned feature planes \mathcal{P}_{XY} , \mathcal{P}_{YZ} and \mathcal{P}_{XZ} . The features corresponding to sample \mathbf{x} on plane k can be queried by projection and interpolation:

$$F_k = \psi(\mathcal{P}_k, \pi_i(\mathcal{N}(\mathbf{x}))), k \in \{XY, XZ, YZ\}, \quad (3)$$

where π_k projects \mathbf{x} onto k -th plane, \mathcal{N} normalizes \mathbf{x} 's range to $[-1, 1]$, and ψ represents bilinear interpolations. Following [13], different features from triplane are combined by Hadamard product as $F = \prod_k F_k$ to produce a final feature F . The aggregated features will be decoded into color \mathbf{c} and density σ by an MLP decoder \mathcal{M} . Thus \mathcal{M} and \mathcal{P} represent Θ above in pose-NeRF optimization. Besides photometric loss, a commonly used total variation loss \mathcal{L}_{TV} [13, 14] is applied, and the whole loss term \mathcal{L} is formulated as:

$$\mathcal{L} = \mathcal{L}_{render} + \mathcal{L}_{TV}. \quad (4)$$

4 Method

As mentioned above, the baseline that naively combines pose estimation with triplane has drawbacks on local errors in feature plane updating and inappropriate plane aggregation. To address that, we provide a novel pipeline as illustrated in Fig. 3. Our method can be divided into two stages. In the first stage, we input random triplane noise to the proposed triplane generator to generate different feature grids for scene representation while optimizing the camera poses (Sec. 4.1). Features interpolated from different planes are aggregated through the proposed DPA (Sec. 4.2). Colors and densities are derived from an MLP decoder. In the second stage, we adopt a warm-start training strategy for both efficient training and better scene representation (Sec. 4.3).

4.1 Disentangled Triplane Generation Module

Triplane Generator. As we mentioned in Sec. 1, the baseline joint estimation is prone to fall into local minima. To mitigate the local error, we propose a disentangled triplane generator for robust learning.

As shown in Fig. 3, given a frozen triplane noise following a Gaussian distribution as input, we introduce a convolution-based neural network \mathcal{G} as our triplane generator for extracting the plane features \mathcal{P} . Our motivation is to parameterize spatial feature grids as deep neural networks. Since the feature planes are highly decoupled for different perspectives of 3D scenes, features among different planes are anisotropic. Therefore, we apply three individual neural networks with the same structure but not shared parameters as disentangled generators to produce three different planes respectively.

A convolution-based network is applied as our triplane generator. The generator is randomly initialized for each scene and *not* pre-trained, which still requires per-scene optimization. Detailed architecture of the generator can be

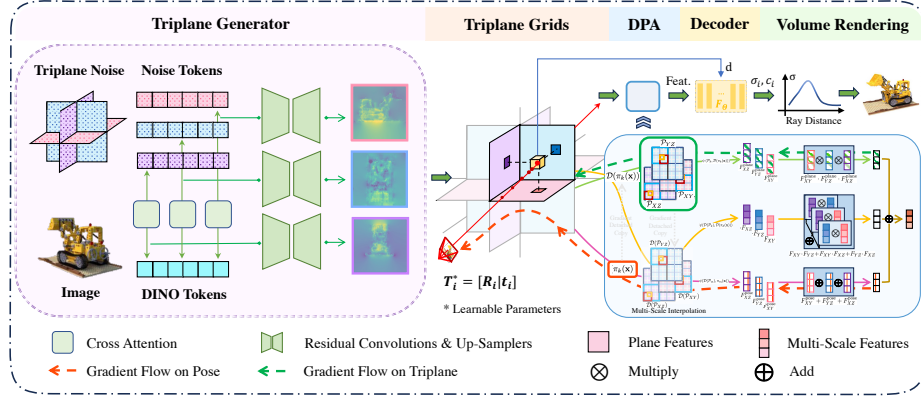


Fig. 3: Overview pipeline of our proposed method. In the 1st stage, we obtain triplane features from the triplane generator. In the 2nd stage, we discard the generator and transform it to direct updates on triplane feature grids.

found in the supplementary. As shown in Fig. 2, compared to baseline, our triplane generator reduces errors from local updates and exhibits much less noise in the visualization of planar feature maps. In addition, the deep network prior will excavate the scene’s geometric and texture features and embed them into triplane representation, instead of simply over-fitting the novel viewpoints.

It’s worth noting that the input of triplane generator \mathcal{G} is randomly initialized noise tokens t_{init} in order to introduce *spatial priors*. During generation, the core features are maintained in a triplane structure from input to output, suiting the disentangled scene representation.

Scene Texture Embedding Module. After the parameterization of the triplane with the generator, we present the Scene Texture Embedding module to enhance the triplane texture representation of the scene, thus mitigating the pose-NeRF ambiguous [9] problem. Inspired by [19, 53], we apply a DINOv2 [40] to encode the input image into patch-wise feature tokens. We found that without the need for several extracted feature maps for triplane generation, it is still capable of reducing the optimization ambiguity effectively. Therefore, we extract the feature h_0 of the first image I_0 for the scene texture prior.

Next, a cross-attention is applied to incorporate the 2D feature into the triplane. During optimization, 2D feature tokens and 3D triplane tokens will autonomously learn the alignment between modalities, which allows for better integration. Specifically, we use the frozen triplane tokens as query, and the extracted feature tokens as key and value to obtain the attended triplane tokens t_{in} . Finally, t_{in} is input to the triplane generator for the subsequent processing.

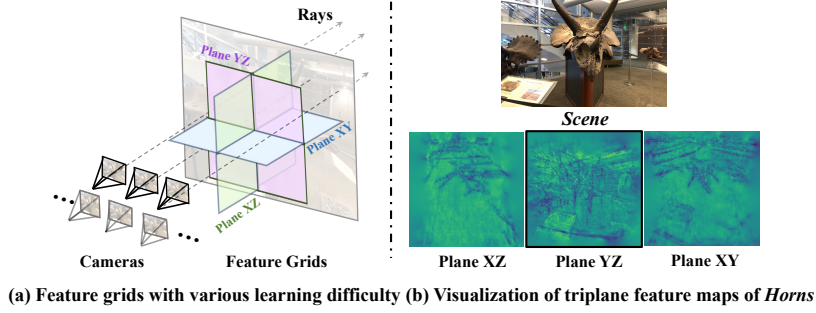


Fig. 4: Illustration of the training differences among feature grids. (a) With this camera distribution, Plane YZ obtains more scene information and converges with less difficulty than the other planes. (b) We visualize the feature maps on the triplane of the same training steps. Compared to the other two planes, Plane YZ shows realistic texture to the original scene, demonstrating sufficient training.

4.2 Disentangled Plane Aggregation

Features queried from different planes with a sampled 3D point \mathbf{x} will be aggregated before decoding and rendering. Here we analyze how the anisotropy in triplane influences pose in joint optimization, point out deficiencies in previous aggregation, and present our Disentangled Plane Aggregation as a solution.

Since the triplane representation is an interpretable explicit radiance field decomposing scenes into three orthogonal planes aligned with X, Y, and Z axes, each plane contains features of 3D scenes from orthogonal perspectives, which is influenced by different characteristics inherent to the scene. Adequate captured images from varied viewpoints furnish comprehensive scene information to triplane across all perspectives, while the images from limited or unevenly distributed viewpoints provide disparate information to each plane. For example, as depicted in Fig. 4, angles of images or videos are generally consistent, with all objects primarily facing toward cameras. As scenes emphasize the front view of objects, Plane YZ in Fig. 4 gains richer scene information and is easier to learn. During optimization, one plane that is easier to learn is more capable of fitting the scene and therefore provides better supervision for the pose. On the contrary, a plane that is more challenging to learn requires extra training iterations to accurately represent decoupled scene information, and therefore the probability of providing ambiguous supervision to the pose increases, resulting in pose optimizing collision during training.

As mentioned in Sec. 3 and Sec. 4.1, given an 2D sample $\mathbf{u} = (u, v)$ on image space of camera i and its homogeneous coordinates $\bar{\mathbf{u}} = [\mathbf{u}; 1]^T$, we can get $\mathbf{x} = \mathcal{W}(z\bar{\mathbf{u}}, \mathbf{T}_i)$ at depth z by space warping \mathcal{W} and obtain features F by interpolation and Hadamard product combination according to Eq. (3). Here, we can rewrite Eq. (1) as $\hat{I}(\mathbf{r}) = \hat{I}(\mathcal{M}(F))$, and derive the partial derivative with

respect to \mathbf{T}_i :

$$\frac{\partial \hat{I}_i}{\partial \mathbf{T}_i} = \frac{\partial \hat{I}_i}{\partial \mathcal{M}} \frac{\partial \mathcal{M}}{\partial F} \frac{\partial F}{\partial \mathbf{x}} \frac{\partial \mathcal{W}(z\bar{\mathbf{u}}, \mathbf{T}_i)}{\partial \mathbf{T}_i}, \quad (5)$$

which is formed via backpropagation and chain rules. Note that $\frac{\partial \hat{I}_i}{\partial \mathcal{M}}$ and $\frac{\partial \mathcal{M}}{\partial F}$ represent the rendering differentiation and weights of the MLP decoder relatively, $\frac{\partial \mathcal{W}(z\bar{\mathbf{u}}, \mathbf{T}_i)}{\partial \mathbf{T}_i}$ on behalf of warping, we expand the remaining part $\frac{\partial F}{\partial \mathbf{x}}$ as:

$$\frac{\partial F}{\partial \mathbf{x}} = \sum_m \left(\frac{\partial \psi(\mathcal{P}_m, \pi_m(\mathbf{x}))}{\partial \mathbf{x}} \cdot \prod_{k \neq m} \psi(\mathcal{P}_k, \pi_k(\mathbf{x})) \right), \quad (6)$$

where the subscript m and k are plane indexes, selected from set $\{XY, XZ, YZ\}$. When updating poses, especially in early steps, gradients are influenced by the fluctuation on all feature planes, thus planes with inadequate supervision lead to ambiguity. Another popular aggregation is sum, used in [5, 6]. However, as reported in Tab. 1, while sum relieves pose optimizing collision, its ability to learn triplane to express 3D scene is inferior to product, which is consistent with [13].

Aiming at both reducing collision to assist pose optimization and promoting the expressiveness of triplane, we design a disentangled aggregating strategy among planes. Dubbed as DPA, the expression of this aggregation is shown below:

Table 1: Ablation study on *Horns* over different aggregations.

| | Rot.(°) ↓ | Trans.(×100) ↓ | PSNR ↑ |
|-------|-----------|----------------|--------|
| Prod. | 1.002 | 0.348 | 21.64 |
| Sum | 0.794 | 0.172 | 21.29 |
| DPA | 0.296 | 0.102 | 23.66 |

$$\begin{aligned} \text{DPA}(\mathcal{P}, \mathbf{x}) = & \prod_k \psi(\mathcal{P}_k, \mathcal{D}(\pi_k(\mathbf{x}))) + \sum_m \psi(\mathcal{D}(\mathcal{P}_m), \pi_m(\mathbf{x})) \\ & + \sum_{\substack{k, m \\ k \neq m}} (\psi(\mathcal{D}(\mathcal{P}_k), \mathcal{D}(\pi_k(\mathbf{x}))) \cdot \psi(\mathcal{D}(\mathcal{P}_m), \mathcal{D}(\pi_m(\mathbf{x})))) + 1, \end{aligned} \quad (7)$$

where $\mathcal{D}(x)$ is the gradient-detached copy of x . Due to limited text space, we present detailed derivation of the above formula in supplemental materials. Noticing that $\mathcal{D}(x)$ keeps the value, the output of DPA is equal to $\prod_k (F_k + 1)$ in a Hadamard product form, and thus preserves the expressiveness in triplane. In practice, we change the addition item 1 to a hyperparameter λ and so as the coefficients of items in Eq. (7). Allowing the triplane gradient obtained from product while the pose gradient from sum, DPA distributes information with disentanglement to pose and triplane, leading to more robust joint optimization.

4.3 Two Stage Warm-Start Training

After adopting the triplane generator, we find that once the camera poses are well initialized, the baseline joint pose-triplane estimation outperforms the training of

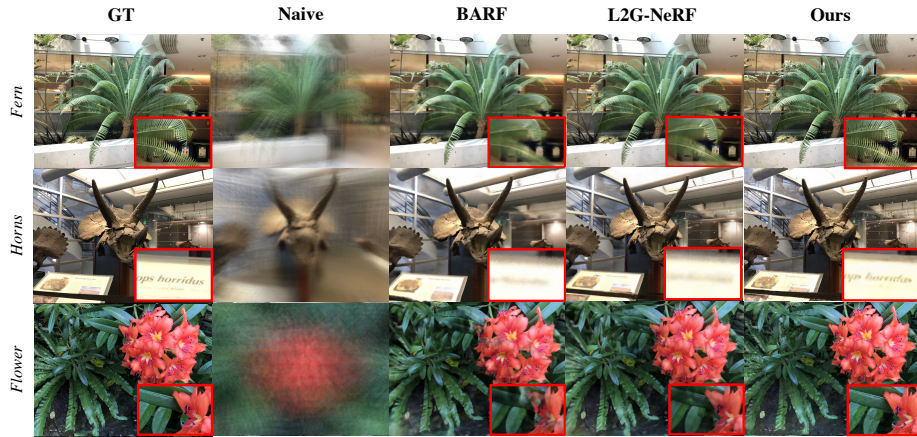


Fig. 5: Qualitative results of novel view synthesis on LLFF dataset. Naive represents reference K-Plane [13] trained with unknown camera poses.

adopting the generator from the beginning to the end, particularly in areas of high frequency. We argue that though the triplane generator mitigates local updating errors, it suffers from implicit constraints in different patches on the plane grids caused by the generator, thus introducing excessive smoothing on the feature grids. Therefore, we propose a warm-start two-stage training strategy to switch to direct plane optimization for better NVS quality. Instead of initializing the parameters from scratch in the second stage, we leverage the learned parameters in the first stage. The MLP decoder and camera poses are inherited directly as they work for the same triplane which contains no gaps between different stages. For the 2nd-stage triplane, we discard the generator and perform a forward inference to obtain the final 1st-stage triplane. Following [13, 37], we adopt multiple feature grids from different resolutions in a coarse-to-fine manner, thereby enhancing the structural and detailed representation. The 2nd-stage triplane with multiple scales will be obtained by bilinear interpolation from the 1st-stage triplane, where direct triplane optimization will be conducted in the second stage.

5 Experiments

In this section, we present the performance and analysis of our proposed method on two public datasets: LLFF [35] and NeRF-Synthetic dataset [36]. Next, we compare our results on both pose accuracy and novel view synthesis quality with other joint optimizing works containing BARF [28], GARF [11], L2G-NeRF [9] and HASH-BARF [18] in Sec. 5.1 and Sec. 5.2. Lastly, we conduct ablation studies to validate the effectiveness of our proposed modules in Sec. 5.3.

Table 2: Quantitative results of novel view synthesis on LLFF dataset.

| Scene | Novel View Synthesis Quality | | | | | | | | | |
|----------|------------------------------|-----------|---------|-----------|--------------|-----------------|-----------|---------|-----------|--------------|
| | PSNR \uparrow | | | | | SSIM \uparrow | | | | |
| | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours |
| Fern | 23.79 | 24.51 | 24.57 | 24.62 | 25.70 | 0.710 | 0.740 | 0.750 | 0.743 | 0.837 |
| Flower | 23.37 | 26.40 | 24.90 | 25.19 | 27.06 | 0.698 | 0.790 | 0.740 | 0.744 | 0.830 |
| Fortress | 29.08 | 29.09 | 29.27 | 30.14 | 30.79 | 0.823 | 0.820 | 0.840 | 0.901 | 0.905 |
| Horns | 22.78 | 22.54 | 23.12 | 22.97 | 23.66 | 0.727 | 0.690 | 0.740 | 0.736 | 0.828 |
| Leaves | 18.78 | 19.72 | 19.02 | 19.45 | 20.43 | 0.537 | 0.610 | 0.560 | 0.607 | 0.708 |
| Orchids | 19.45 | 19.37 | 19.71 | 20.02 | 20.24 | 0.574 | 0.570 | 0.610 | 0.610 | 0.660 |
| Room | 31.95 | 31.90 | 32.25 | 32.73 | 33.95 | 0.949 | 0.940 | 0.950 | 0.968 | 0.970 |
| T-Rex | 22.55 | 22.86 | 23.49 | 23.19 | 25.35 | 0.767 | 0.800 | 0.800 | 0.866 | 0.889 |
| Mean | 23.97 | 24.55 | 24.54 | 24.79 | 25.90 | 0.723 | 0.745 | 0.750 | 0.772 | 0.828 |

Table 3: Quantitative results of camera pose estimation on LLFF dataset.

| Scene | Camera Pose Estimation | | | | | | | | | |
|----------|--------------------------------------|--------------|---------|--------------|--------------|---|-----------|---------|--------------|--------------|
| | Rotation ($^{\circ}$) \downarrow | | | | | Translation ($\times 100$) \downarrow | | | | |
| | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours |
| Fern | 0.191 | 0.470 | 0.200 | 0.110 | 0.243 | 0.102 | 0.250 | 0.180 | 0.102 | 0.166 |
| Flower | 0.251 | 0.460 | 0.330 | 0.301 | 0.139 | 0.224 | 0.220 | 0.240 | 0.211 | 0.179 |
| Fortress | 0.479 | 0.030 | 0.250 | 0.211 | 0.612 | 0.364 | 0.270 | 0.250 | 0.241 | 0.385 |
| Horns | 0.304 | 0.030 | 0.220 | 0.049 | 0.296 | 0.222 | 0.210 | 0.270 | 0.209 | 0.102 |
| Leaves | 1.272 | 0.130 | 0.790 | 0.840 | 0.329 | 0.249 | 0.230 | 0.340 | 0.228 | 0.239 |
| Orchids | 0.627 | 0.430 | 0.670 | 0.399 | 0.226 | 0.404 | 0.410 | 0.410 | 0.386 | 0.235 |
| Room | 0.320 | 0.420 | 0.300 | 0.271 | 0.184 | 0.270 | 0.320 | 0.230 | 0.213 | 0.095 |
| T-Rex | 1.138 | 0.660 | 0.890 | 0.894 | 0.039 | 0.720 | 0.480 | 0.640 | 0.474 | 0.149 |
| Mean | 0.573 | 0.329 | 0.460 | 0.384 | 0.259 | 0.331 | 0.299 | 0.320 | 0.258 | 0.194 |

5.1 Results on LLFF dataset

Experimental Settings. For LLFF dataset [35], the resolution of images is set to 480×640 [18, 28] and the train/test splits are defined as [9, 11, 28]. We initialize the translation of each camera as a zero vector and the rotation matrix as an identity matrix. Our model is trained for 70k steps and switched to the second stage at step 4000. In the testing phase, our model performs the same test-time photometric pose optimization [9, 28] before evaluating view synthesis quality. The ground-truth poses are estimated from COLMAP [47]. We further compare our method to the reference baseline K-Planes [13] without camera poses. All methods are trained and evaluated under the same settings.

Results. Quantitative results of novel view synthesis and pose estimation are shown in Tab. 2 and Tab. 3. The PSNR and SSIM are reported for novel view synthesis quality, and the rotation error (in degree) and translation error (scaled by 100) are reported for pose estimation accuracy. Qualitative results are shown in Fig. 5 for visual comparison. Without pose prior and optimization (Naive), it is impossible to reconstruct the scene with common triplane radiance fields.

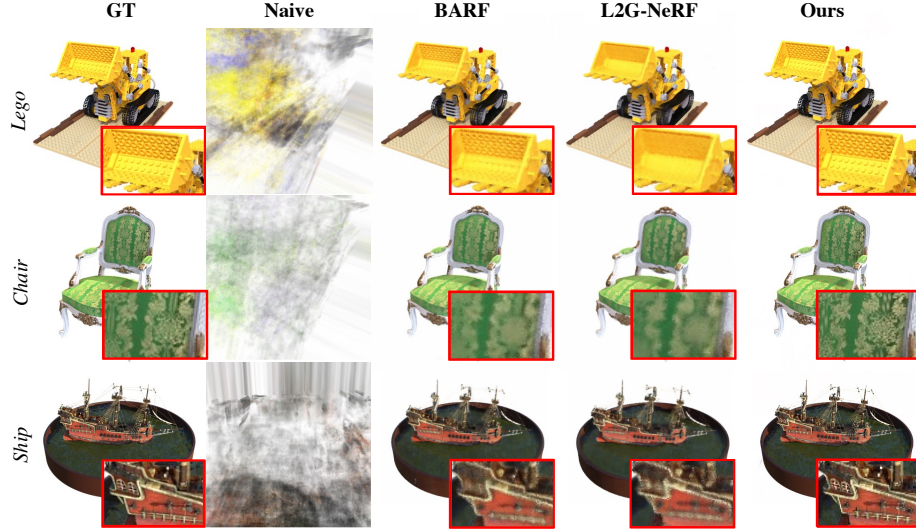


Fig. 6: Qualitative results of novel view synthesis on NeRF-synthetic dataset. Naive represents reference baseline K-Plane [13] is trained with noisy camera poses.

5.2 Results on NeRF-Synthetic dataset

Experimental Settings. In our implementation, the rendered images are resized to 400×400 resolution, and the train/test splits are the same as [9, 11, 18, 28]. We follow [28] to impose additive Gaussian noise $\xi \in \mathfrak{sc}(3)$ and $\xi \sim \mathcal{N}(\mathbf{0}, 0.15\mathbf{I})$ to the ground truth poses as initialization. Our model is trained for 60k steps and switched to the second stage at 2000 steps. The settings of test-time pose optimization are the same as the LLFF dataset. Note that L2G-NeRF [9] is implemented with a different way of the pose noise perturbation, we re-implement L2G-NeRF’s noise perturbation according to BARF to ensure the fairness of comparisons. We further show the results of the reference baseline K-Planes with the same noisy pose initialization.

Results. The quantitative results of NVS and pose estimation are presented in Tab. 4 and Tab. 5, where the same metrics are reported as in the LLFF dataset. The qualitative NVS results are shown in Fig. 6. The original triplane method K-Planes fails to perform reliable scene reconstruction without accurate camera poses, while our proposed method still performs robust novel view synthesis.

5.3 Ablation Study

In this section, we analyze the effectiveness of our proposed core components. The ablation studies are shown in Tab. 7.

Table 4: Quantitative results of novel view synthesis on NeRF-Synthetic.

| Scene | Novel View Synthesis | | | | | | | | | |
|-----------|----------------------|--------------|---------|-----------|--------------|-----------------|--------------|---------|--------------|--------------|
| | PSNR \uparrow | | | | | SSIM \uparrow | | | | |
| | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours |
| Chair | 31.16 | 31.32 | 34.28 | 31.95 | 37.78 | 0.954 | 0.959 | 0.982 | 0.962 | 0.992 |
| Drum | 23.91 | 24.15 | 25.39 | 24.16 | 26.10 | 0.900 | 0.909 | 0.901 | 0.912 | 0.926 |
| Ficus | 26.26 | 26.29 | 27.45 | 28.31 | 33.12 | 0.934 | 0.935 | 0.940 | 0.943 | 0.979 |
| Hotdog | 34.54 | 34.69 | 34.03 | 35.41 | 35.98 | 0.970 | 0.972 | 0.967 | 0.981 | 0.976 |
| Lego | 28.33 | 29.29 | 27.66 | 31.65 | 32.19 | 0.927 | 0.925 | 0.922 | 0.973 | 0.976 |
| Materials | 27.84 | 27.91 | 26.01 | 27.14 | 25.78 | 0.936 | 0.941 | 0.920 | 0.911 | 0.920 |
| Mic | 31.18 | 31.39 | 32.61 | 32.33 | 33.88 | 0.969 | 0.971 | 0.971 | 0.975 | 0.987 |
| Ship | 27.50 | 27.64 | 28.30 | 27.92 | 29.80 | 0.849 | 0.862 | 0.789 | 0.879 | 0.903 |
| Mean | 28.84 | 28.96 | 29.47 | 29.86 | 31.83 | 0.930 | 0.935 | 0.924 | 0.943 | 0.957 |

Table 5: Quantitative results of pose estimation on NeRF-Synthetic.

| Scene | Camera Pose Estimation | | | | | | | | | |
|-----------|--------------------------------------|--------------|---------|--------------|--------------|---|-----------|---------|--------------|--------------|
| | Rotation ($^{\circ}$) \downarrow | | | | | Translation ($\times 100$) \downarrow | | | | |
| | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours | BARF [28] | GARF [11] | L2G [9] | HASH [18] | Ours |
| Chair | 0.096 | 0.113 | 0.118 | 0.085 | 0.036 | 0.428 | 0.549 | 0.495 | 0.365 | 0.186 |
| Drum | 0.043 | 0.052 | 0.070 | 0.041 | 0.032 | 0.225 | 0.232 | 0.340 | 0.214 | 0.197 |
| Ficus | 0.085 | 0.081 | 0.168 | 0.079 | 0.058 | 0.474 | 0.461 | 1.037 | 0.479 | 0.299 |
| Hotdog | 0.248 | 0.235 | 0.661 | 0.229 | 0.226 | 1.308 | 1.123 | 4.283 | 1.123 | 0.976 |
| Lego | 0.082 | 0.101 | 0.088 | 0.071 | 0.033 | 0.291 | 0.299 | 0.399 | 0.272 | 0.141 |
| Materials | 0.844 | 0.842 | 0.805 | 0.852 | 0.486 | 2.692 | 2.688 | 2.510 | 2.743 | 1.537 |
| Mic | 0.071 | 0.070 | 0.080 | 0.068 | 0.124 | 0.301 | 0.293 | 0.331 | 0.287 | 1.554 |
| Ship | 0.075 | 0.073 | 0.163 | 0.079 | 0.255 | 0.326 | 0.310 | 0.585 | 0.287 | 0.340 |
| Mean | 0.193 | 0.195 | 0.269 | 0.189 | 0.156 | 0.756 | 0.744 | 1.248 | 0.722 | 0.654 |

Effectiveness of Triplane Generator. The rows (e) and (c) show the results of the baseline and that with the triplane generator. The baseline method suffers from severe local minima due to incorrect local updating. After introducing the triplane generator, the quality of joint estimation substantially improved.

Effectiveness of Disentangled Plane Aggregation. After introducing DPA to the baseline, as shown in (e) and (d), indicating that the pose estimation is subject to the conflict between different planes. Therefore the proposed DPA can improve the accuracy of pose estimation and thus bring better NVS quality. Meanwhile, after removing DPA from (b) to (c), there is a certain decrease in the joint optimization results.

Effectiveness of Two-Stage Warm-Start Training. Without two-stage training, over-smoothness caused by the triplane generator degrades the image rendering quality. In addition, as shown in Tab. 8, we conduct runtime comparisons to show the trade-off between efficiency and performance. Our method incurs a larger time consumption when using only the first stage of training while the quality of viewpoint rendering is degraded. Our full model improves the image rendering quality with comparable time consumption to baseline direct optimization.

Table 6: Comparison of the reference K-Planes [13] (Ref.), naive joint pose-triplane optimization (Base), and our proposed method (Ours). # represents the reference K-Planes is trained with ground-truth camera poses.

| Scenes | PSNR \uparrow | | | SSIM \uparrow | | | Rot.($^\circ$) \downarrow | | Trans.($\times 100$) \downarrow | |
|----------|-----------------|-------|-------|-----------------|-------|-------|-------------------------------|--------------|-------------------------------------|--------------|
| | Ref.# | Base | Ours | Ref.# | Base | Ours | Base | Ours | Base | Ours |
| Fern | 24.44 | 23.94 | 25.70 | 0.828 | 0.782 | 0.837 | 1.852 | 0.243 | 0.428 | 0.166 |
| Flower | 27.42 | 24.88 | 27.06 | 0.872 | 0.764 | 0.830 | 8.624 | 0.139 | 0.733 | 0.179 |
| Fortress | 29.36 | 28.09 | 30.79 | 0.804 | 0.836 | 0.905 | 2.247 | 0.612 | 1.660 | 0.385 |
| Horns | 28.64 | 20.82 | 23.66 | 0.892 | 0.543 | 0.828 | 63.92 | 0.296 | 22.93 | 0.102 |
| Leaves | 20.22 | 16.46 | 20.43 | 0.746 | 0.427 | 0.708 | 173.6 | 0.329 | 6.799 | 0.239 |
| Orchids | 19.58 | 14.03 | 20.22 | 0.676 | 0.241 | 0.660 | 17.48 | 0.226 | 5.283 | 0.235 |
| Room | 34.07 | 21.68 | 33.95 | 0.957 | 0.716 | 0.970 | 174.2 | 0.184 | 6.139 | 0.095 |
| T-Rex | 24.14 | 21.86 | 25.35 | 0.915 | 0.675 | 0.889 | 84.62 | 0.039 | 21.26 | 0.149 |
| Mean | 25.98 | 21.47 | 25.90 | 0.847 | 0.623 | 0.828 | 65.82 | 0.259 | 8.154 | 0.194 |

Table 7: Ablation study of the proposed components on LLFF dataset [35].

| | TriPlane Generator | Disentangled Plane Agg. | Two-Stage | Rot. \downarrow | Trans. \downarrow | PSNR \uparrow |
|-----|--------------------|-------------------------|-----------|-------------------|---------------------|-----------------|
| (a) | ✓ | ✓ | ✓ | 0.259 | 0.194 | 25.90 |
| (b) | ✓ | ✓ | | 0.634 | 0.328 | 25.21 |
| (c) | ✓ | | | 0.660 | 0.356 | 24.93 |
| (d) | | ✓ | | 6.180 | 2.708 | 23.41 |
| (e) | | | | 65.82 | 8.154 | 21.47 |

Table 8: Runtime comparisons on *Drums*. Training Times are in the format of hh: mm.

| Methods | Iters. | Training Time | PSNR \uparrow |
|-------------------------|------------|---------------|-----------------|
| BARF [28] | 200k | 11:49 | 23.91 |
| HASH-BARF [18] | 200k | 00:36 | 24.16 |
| Ours-w/o 1^{st} stage | 60k | 00:59 | 19.93 |
| Ours-w/o 2^{nd} stage | 60k | 04:00 | 24.98 |
| Ours | 60k | 01:21 | 26.10 |

Comparisons with Baseline Joint Estimation. As shown in Tab. 6, we conduct experiments to compare our full model, reference K-Planes [13] with COLMAP poses, and direct joint pose-triplane optimization baseline to illustrate the effectiveness. Note that the train/test split of the original K-Planes implementation is different from ours, we re-implement the evaluation as [9, 11, 28] for fair comparison. We parameterize the camera poses and integrate them into K-Planes as our baseline implementation. Our approach exhibits significant improvements in pose estimation and novel view rendering and even surpasses K-Planes with COLMAP poses, which further demonstrates the effectiveness of our method.

6 Conclusion

In this paper, we propose a novel algorithm that allows joint estimation of camera pose and disentangled scene representation to reduce the computational cost of neural renderings. We first propose the Disentangled Triplane Generation Module to parameterize the triplane with a convolution-based generator to mitigate local updating errors in the direct joint optimization baseline. Then we analyze the effect of triplane feature aggregation on camera poses and propose Disentangled Plane Aggregation to reduce the entanglement between feature planes and camera poses. Finally, we adopt a two-stage warm-start strategy to tackle the oversmoothing on the feature planes caused by the generator. Comprehensive evaluations demonstrate that our proposed method achieves state-of-the-art performance and rapid convergence with noisy or unknown poses.

Acknowledgments

This work is financially supported by Outstanding Talents Training Fund in Shenzhen, Shenzhen Science and Technology Program-Shenzhen Cultivation of Excellent Scientific and Technological Innovation Talents project (Grant No. RCJC20200714114435057), Shenzhen Science and Technology Program-Shenzhen Hong Kong joint funding project (Grant No. SGDX20211123144400001), Guangdong Provincial Key Laboratory of Ultra High Definition Immersive Media Technology, National Natural Science Foundation of China U21B2012, R24115SG MIGU-PKU META VISION TECHNOLOGY INNOVATION LAB. Jianbo Jiao is supported by the Royal Society Short Industry Fellowship (SIF\R1\231009).

References

1. Barron, J.T., Mildenhall, B., Tancik, M., Hedman, P., Martin-Brualla, R., Srinivasan, P.P.: Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5855–5864 (2021)
2. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 5470–5479 (June 2022)
3. Barron, J.T., Mildenhall, B., Verbin, D., Srinivasan, P.P., Hedman, P.: Zip-nerf: Anti-aliased grid-based neural radiance fields. arXiv preprint arXiv:2304.06706 (2023)
4. Bian, W., Wang, Z., Li, K., Bian, J.W., Prisacariu, V.A.: Nope-nerf: Optimising neural radiance field with no pose prior. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4160–4169 (2023)
5. Cao, A., Johnson, J.: Hexplane: A fast representation for dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 130–141 (2023)
6. Chan, E.R., Lin, C.Z., Chan, M.A., Nagano, K., Pan, B., De Mello, S., Gallo, O., Guibas, L.J., Tremblay, J., Khamis, S., et al.: Efficient geometry-aware 3d generative adversarial networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16123–16133 (2022)
7. Chen, A., Xu, Z., Geiger, A., Yu, J., Su, H.: Tensorf: Tensorial radiance fields. In: European Conference on Computer Vision. pp. 333–350. Springer (2022)
8. Chen, H., Gu, J., Chen, A., Tian, W., Tu, Z., Liu, L., Su, H.: Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction. arXiv preprint arXiv:2304.06714 (2023)
9. Chen, Y., Chen, X., Wang, X., Zhang, Q., Guo, Y., Shan, Y., Wang, F.: Local-to-global registration for bundle-adjusting neural radiance fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8264–8273 (2023)
10. Cheng, Z., Esteves, C., Jampani, V., Kar, A., Maji, S., Makadia, A.: Lu-nerf: Scene and pose estimation by synchronizing local unposed nerfs. arXiv preprint arXiv:2306.05410 (2023)

11. Chng, S.F., Ramasinghe, S., Sherrah, J., Lucey, S.: Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In: European Conference on Computer Vision. pp. 264–280. Springer (2022)
12. Deng, K., Liu, A., Zhu, J.Y., Ramanan, D.: Depth-supervised nerf: Fewer views and faster training for free. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12882–12891 (2022)
13. Fridovich-Keil, S., Meanti, G., Warburg, F.R., Recht, B., Kanazawa, A.: K-planes: Explicit radiance fields in space, time, and appearance. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 12479–12488 (June 2023)
14. Fridovich-Keil, S., Yu, A., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5501–5510 (2022)
15. Fu, Y., Liu, S., Kulkarni, A., Kautz, J., Efros, A.A., Wang, X.: Colmap-free 3d gaussian splatting. arXiv preprint arXiv:2312.07504 (2023). <https://doi.org/10.48550/arXiv.2312.07504>
16. Gao, H., Liu, X., Qu, M., Huang, S.: Pdanet: Self-supervised monocular depth estimation using perceptual and data augmentation consistency. *Applied Sciences* **11**(12), 5383 (2021)
17. Gao, H., Shen, S., Zhang, Z., Xiong, K., Peng, R., Gao, Z., Wang, Q., Xie, Y., Wang, R.: Fdc-nerf: Learning pose-free neural radiance fields with flow-depth consistency. In: ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 3615–3619 (2024)
18. Heo, H., Kim, T., Lee, J., Lee, J., Kim, S., Kim, H.J., Kim, J.H.: Robust camera pose refinement for multi-resolution hash encoding. In: Proceedings of the 40th International Conference on Machine Learning (2023)
19. Hong, Y., Zhang, K., Gu, J., Bi, S., Zhou, Y., Liu, D., Liu, F., Sunkavalli, K., Bui, T., Tan, H.: Lrm: Large reconstruction model for single image to 3d. arXiv preprint arXiv:2311.04400 (2023)
20. Hu, W., Wang, Y., Ma, L., Yang, B., Gao, L., Liu, X., Ma, Y.: Tri-miprf: Tri-mip representation for efficient anti-aliasing neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19774–19783 (2023)
21. Jeong, Y., Ahn, S., Choy, C., Anandkumar, A., Cho, M., Park, J.: Self-calibrating neural radiance fields. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5846–5854 (2021)
22. Kajiya, J.T., Von Herzen, B.P.: Ray tracing volume densities. *ACM SIGGRAPH computer graphics* **18**(3), 165–174 (1984)
23. Kerbl, B., Kopanas, G., Leimkühler, T., Drettakis, G.: 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics* **42**(4) (2023)
24. Kim, M., Seo, S., Han, B.: Infonerf: Ray entropy minimization for few-shot neural volume rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12912–12921 (2022)
25. Levy, A., Matthews, M., Sela, M., Wetzstein, G., Lagun, D.: Melon: Nerf with unposed images using equivalence class estimation. arXiv preprint arXiv:2303.08096 (2023)
26. Li, H., Gao, Y., Zhang, D., Wu, C., Dai, Y., Zhao, C., Feng, H., Ding, E., Wang, J., Han, J.: Ggrt: Towards generalizable 3d gaussians without pose priors in real-time. arXiv preprint arXiv:2403.10147 (2024)

27. Li, T., Slavcheva, M., Zollhoefer, M., Green, S., Lassner, C., Kim, C., Schmidt, T., Lovegrove, S., Goesele, M., Newcombe, R., et al.: Neural 3d video synthesis from multi-view video. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5521–5531 (2022)
28. Lin, C.H., Ma, W.C., Torralba, A., Lucey, S.: Barf: Bundle-adjusting neural radiance fields. In: *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. pp. 5721–5731. IEEE Computer Society (2021)
29. Liu, L., Gu, J., Zaw Lin, K., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. *Advances in Neural Information Processing Systems* **33**, 15651–15663 (2020)
30. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 9298–9309 (2023)
31. Liu, Y.L., Gao, C., Meuleman, A., Tseng, H.Y., Saraf, A., Kim, C., Chuang, Y.Y., Kopf, J., Huang, J.B.: Robust dynamic radiance fields. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 13–23 (2023)
32. Max, N.: Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics* **1**(2), 99–108 (1995)
33. Meng, Q., Chen, A., Luo, H., Wu, M., Su, H., Xu, L., He, X., Yu, J.: Gnerf: Gan-based neural radiance field without posed camera. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 6351–6361 (2021)
34. Meuleman, A., Liu, Y.L., Gao, C., Huang, J.B., Kim, C., Kim, M.H., Kopf, J.: Progressively optimized local radiance fields for robust view synthesis. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 16539–16548 (2023)
35. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)* (2019)
36. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: *European Conference on Computer Vision*. pp. 405–421. Springer (2020)
37. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)* **41**(4), 1–15 (2022)
38. Niemeyer, M., Barron, J.T., Mildenhall, B., Sajjadi, M.S., Geiger, A., Radwan, N.: Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5480–5490 (2022)
39. Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5589–5599 (2021)
40. Oquab, M., Darcet, T., Moutakanni, T., Vo, H., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., et al.: Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193* (2023)
41. Park, K., Sinha, U., Barron, J.T., Bouaziz, S., Goldman, D.B., Seitz, S.M., Martin-Brualla, R.: Nerfies: Deformable neural radiance fields. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5865–5874 (2021)
42. Park, K., Sinha, U., Hedman, P., Barron, J.T., Bouaziz, S., Goldman, D.B., Martin-Brualla, R., Seitz, S.M.: Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228* (2021)

43. Peng, R., Gu, X., Tang, L., Shen, S., Yu, F., Wang, R.: Gens: Generalizable neural surface reconstruction from multi-view images. *Advances in Neural Information Processing Systems* **36**, 56932–56945 (2023)
44. Peng, R., Wang, R., Wang, Z., Lai, Y., Wang, R.: Rethinking depth estimation for multi-view stereo: A unified representation. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2022)
45. Pumarola, A., Corona, E., Pons-Moll, G., Moreno-Noguer, F.: D-nerf: Neural radiance fields for dynamic scenes. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 10318–10327 (2021)
46. Roessle, B., Barron, J.T., Mildenhall, B., Srinivasan, P.P., Nießner, M.: Dense depth priors for neural radiance fields from sparse input views. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 12892–12901 (2022)
47. Schonberger, J.L., Frahm, J.M.: Structure-from-motion revisited. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4104–4113 (2016)
48. Somraj, N., Karanayil, A., Soundararajan, R.: Simplenerf: Regularizing sparse input neural radiance fields with simpler solutions. In: *SIGGRAPH Asia 2023 Conference Papers*. pp. 1–11 (2023)
49. Song, J., Park, S., An, H., Cho, S., Kwak, M.S., Cho, S., Kim, S.: Därf: Boosting radiance fields from sparse input views with monocular depth adaptation. *Advances in Neural Information Processing Systems* **36** (2024)
50. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5459–5469 (2022)
51. Truong, P., Rakotosaona, M.J., Manhardt, F., Tombari, F.: Sparf: Neural radiance fields from sparse and noisy poses. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 4190–4200 (2023)
52. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689* (2021)
53. Wang, P., Tan, H., Bi, S., Xu, Y., Luan, F., Sunkavalli, K., Wang, W., Xu, Z., Zhang, K.: Pflrm: Pose-free large reconstruction model for joint pose and shape prediction. *arXiv preprint arXiv:2311.12024* (2023)
54. Wang, Z., Wu, S., Xie, W., Chen, M., Prisacariu, V.A.: NeRF—: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064* (2021)
55. Wei, Y., Liu, S., Rao, Y., Zhao, W., Lu, J., Zhou, J.: Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5610–5619 (2021)
56. Xia, Y., Tang, H., Timofte, R., Gool, L.V.: Sinerf: Sinusoidal neural radiance fields for joint pose estimation and scene reconstruction. In: *33rd British Machine Vision Conference* (2022)
57. Xu, D., Jiang, Y., Wang, P., Fan, Z., Shi, H., Wang, Z.: Sinnerf: Training neural radiance fields on complex scenes from a single image. In: *European Conference on Computer Vision*. pp. 736–753. Springer (2022)
58. Yen-Chen, L., Florence, P., Barron, J.T., Rodriguez, A., Isola, P., Lin, T.Y.: inerf: Inverting neural radiance fields for pose estimation. In: *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 1323–1330. IEEE (2021)
59. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 5752–5761 (2021)

- 60. Zhang, J., Zhan, F., Wu, R., Yu, Y., Zhang, W., Song, B., Zhang, X., Lu, S.: Vmrf: View matching neural radiance fields. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 6579–6587 (2022)
- 61. Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving neural radiance fields. arXiv preprint arXiv:2010.07492 (2020)
- 62. Zou, Z.X., Yu, Z., Guo, Y.C., Li, Y., Liang, D., Cao, Y.P., Zhang, S.H.: Triplane meets gaussian splatting: Fast and generalizable single-view 3d reconstruction with transformers. arXiv preprint arXiv:2312.09147 (2023)