

CF-NeRF: Camera Parameter Free Neural Radiance Fields with Incremental Learning

Qingsong Yan¹, Qiang Wang^{2,*}, Kaiyong Zhao³, Jie Chen⁴,
Bo Li⁵, Xiaowen Chu^{6,5,*}, Fei Deng^{1,7}

¹ Wuhan University, Wuhan, China

² Harbin Institute of Technology (Shenzhen), Shenzhen, China

³ XGRIDS, Shenzhen, China

⁴ Hong Kong Baptist University, Hong Kong SAR, China

⁵ The Hong Kong University of Science and Technology, Hong Kong SAR, China

⁶ The Hong Kong University of Science and Technology (Guangzhou), Guangzhou, China

⁷ Hubei LuoJia Laboratory, Wuhan, China

yanqs_w hu@whu.edu.cn, qiang.wang@hit.edu.cn, kyzhao@xgrids.com, chenjie@comp.hkbu.edu.hk
bli@cse.ust.hk, xwchu@ust.hk, fdeng@sgg.whu.edu.cn

Abstract

Neural Radiance Fields (NeRF) have demonstrated impressive performance in novel view synthesis. However, NeRF and most of its variants still rely on traditional complex pipelines to provide extrinsic and intrinsic camera parameters, such as COLMAP. Recent works, like NeRFmm, BARF, and L2G-NeRF, directly treat camera parameters as learnable and estimate them through differential volume rendering. However, these methods work for forward-looking scenes with slight motions and fail to tackle the rotation scenario in practice. To overcome this limitation, we propose a novel camera parameter free neural radiance field (CF-NeRF), which incrementally reconstructs 3D representations and recovers the camera parameters inspired by incremental structure from motion (SfM). Given a sequence of images, CF-NeRF estimates the camera parameters of images one by one and reconstructs the scene through initialization, implicit localization, and implicit optimization. To evaluate our method, we use a challenging real-world dataset NeRF-Buster which provides 12 scenes under complex trajectories. Results demonstrate that CF-NeRF is robust to camera rotation and achieves state-of-the-art results without providing prior information and constraints.

Introduction

3D reconstruction is a hot topic in computer vision that aims to recover 3D geometry from RGB images. However, traditional methods contain lots of complex procedures, such as feature extraction and matching (Lowe 2004; Yi et al. 2016), sparse reconstruction (Agarwal et al. 2011; Wu 2013; Schonberger and Frahm 2016; Moulon et al. 2017), and dense reconstruction (Yao et al. 2018; Mi, Di, and Xu 2022; Yan et al. 2023). Consequently, traditional methods are not a differential end-to-end reconstruction pipeline and require high-quality results from each sub-module to achieve accurate results. When the quality of results is poor, it is challenging to identify which module is causing the problem.

*Corresponding author

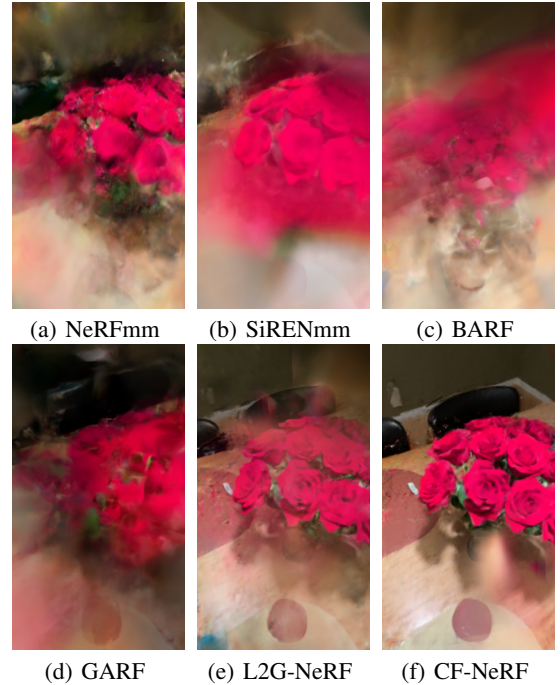


Figure 1: **Rendering results on the NeRF Buster** We select an image sequence from NeRF Buster (Warburg et al. 2023) and use novel views synthesis to compare the quality of camera parameters from NeRFmm (Wang et al. 2021b), SiRENmm (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022), L2G-NeRF (Chen et al. 2023) and our method CF-NeRF.

Recently, Neural Radiance Fields (NeRF) (Mildenhall et al. 2020; Müller et al. 2022; Yu et al. 2021a; Barron et al. 2021) has demonstrated a novel way to render highly realistic novel views with impressive quality. Without recovering 3D geometry, NeRF relies on multi-layer perception (MLP) to predict color and sigma for each point in the scene and samples several points along a ray to render a pixel through differential volume rendering. Unlike traditional 3D recon-

struction, NeRF simplifies the reconstruction into one step and implicitly represents the 3D scene. Benefiting from the excellent ability of NeRF, it has been further extended to dynamic scenes (Pumarola et al. 2021), large-scale (Turki, Ramanan, and Satyanarayanan 2022), and even surface (Wang et al. 2021a) and material reconstruction (Boss et al. 2021a).

Despite the remarkable performance of NeRF and its variants in novel view synthesis, they still require camera parameters before training. The most common processing pipeline is first recovering camera parameters using traditional complex methods like COLMAP (Schonberger and Frahm 2016) or OpenMVG (Moulon et al. 2017), and then training the NeRF through differential volume rendering. In other words, the differentiability of the whole reconstruction pipeline is destroyed and divided into two separate parts, resulting in the NeRF not being end-to-end and the reconstruction quality being unidirectionally dependent on traditional methods.

To unify camera parameter estimation and reconstruction, researchers have tried to recover or optimize camera parameters along with NeRF (Wang et al. 2021b; Lin et al. 2021; Chng et al. 2022; Kuang et al. 2022; Tancik et al. 2023). The straightforward idea is to treat camera parameters as learnable. NeRFmm (Wang et al. 2021b) proves that it is possible to estimate camera parameters by NeRF. BARF (Lin et al. 2021) recovers extrinsic camera parameters and the NeRF model by dynamically adjusting weights of different frequencies of positional encoding. GARF (Chng et al. 2022) replaces ReLU with Gaussian activations to obtain high-accuracy results. NeROIC (Kuang et al. 2022) and NeRF-Studio (Tancik et al. 2023) optimize camera parameters obtained from COLMAP and the NeRF simultaneously. However, these methods are only suitable for forward-looking scenes or scenes with initial camera parameters and cannot be directly used in the real world with complex movement.

This paper proposes a new end-to-end approach called camera parameter free NeRF (CF-NeRF) to address the limitations of existing NeRF-based methods in estimating camera parameters. Figure 7 compares rendered novel views by camera parameters estimated by NeRFmm (Wang et al. 2021b), SiRENmm (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022), L2G-NeRF (Chen et al. 2023), and our method CF-NeRF, where CF-NeRF is the only method that successfully reconstructs the 3D scene with rotation. Unlike other methods that simultaneously estimate all camera parameters, CF-NeRF inherits ideas from incremental structure from motion (SfM) and recovers camera parameters one by one. CF-NeRF contains three major components: initialization, implicit localization, and implicit optimization. CF-NeRF first uses initialization to recover the camera parameters and NeRF by a few images and then gradually estimates the camera parameters of other images through two steps: the implicit localization provides an initial camera parameter for the newly added image, and the implicit optimization optimizes the camera parameter of all images to reduce drift. Our contributions are as follows:

1. We propose a novel end-to-end method, CF-NeRF, that does not need prior information or constraints to recover the intrinsic and extrinsic camera parameters and the NeRF simultaneously.

2. We design an incremental training pipeline for the CF-NeRF, inspired by the incremental SfM, to avoid trapping to local minimal and is suitable for complex trajectories.
3. Experiments of our method achieve state-of-the-art results on the NeRFBuster dataset (Warburg et al. 2023) captured in the real world, proving that the CF-NeRF can estimate accurate camera parameters with the specifically designed training procedure.

Related Work

In this section, we first introduce the development of NeRF-related methods with known camera parameters. Then, we describe several camera parameter estimation methods using traditional SfM&SLAM (simultaneous localization and mapping) and the NeRF.

NeRF

NeRF (Mildenhall et al. 2020) uses the MLP to represent the 3D scene implicitly and can be trained through differential volume rendering from a set of images with known camera parameters. However, NeRF suffers from efficiency and needs around 1-2 days to train a scene and several minutes to render a novel view at the testing. Instant-NGP (Müller et al. 2022) builds a multi-resolution hash table to store space-aware feature vectors and reduces the complexity of the MLP network. Meanwhile, (Sun, Sun, and Chen 2022; Fridovich-Keil et al. 2022; Zhang et al. 2022a) try to use the coarse-to-fine strategy and (Yu et al. 2021a; Chen et al. 2022; Garbin et al. 2021; Wadhvani and Kojima 2022) update the network structure to speed up training or testing. Besides, NeRF faces another problem that it cannot work for large-scale, unbounded 3D scenes. NeRF++ (Zhang et al. 2020) and MipNeRF360 (Barron et al. 2021, 2022) utilize different sampling strategies for foreground and background to model unbounded 3D scenes by a finite volume. MegaNeRF (Turki, Ramanan, and Satyanarayanan 2022) and BlockNeRF (Tancik et al. 2022) split a large scene into multiple small regions and assign a network for each part. Moreover, (Martin-Brualla et al. 2021; Pumarola et al. 2021; Park et al. 2021; Attal et al. 2021) extend NeRF to dynamic scenes and (Jain, Tancik, and Abbeel 2021; Yu et al. 2021b; Niemeyer et al. 2022; Kim, Seo, and Han 2022) introduce context or geometry information into NeRF to suit scenes with sparse views. In addition to the advances in novel view synthesis, NeRF has made significant progress in geometric reconstruction (Yariv et al. 2021; Wang, Skorokhodov, and Wonka 2022; Darmon et al. 2022; Long et al. 2023, 2022; Fu et al. 2022). UniSURF (Oechsle, Peng, and Geiger 2021) and NeUS (Wang et al. 2021a) estimate the zero-level set of an implicit signed distance function instead of the space density. Furthermore, some work (Zhang et al. 2021; Verbin et al. 2022; Boss et al. 2021a; Kuang et al. 2022; Boss et al. 2021b, 2022) even combines BRDF and NeRF to decompose a scene into shape, reflectance, and illumination. However, all of these methods split the reconstruction into two steps and require traditional methods to provide camera parameters, which significantly limits the application of NeRF.

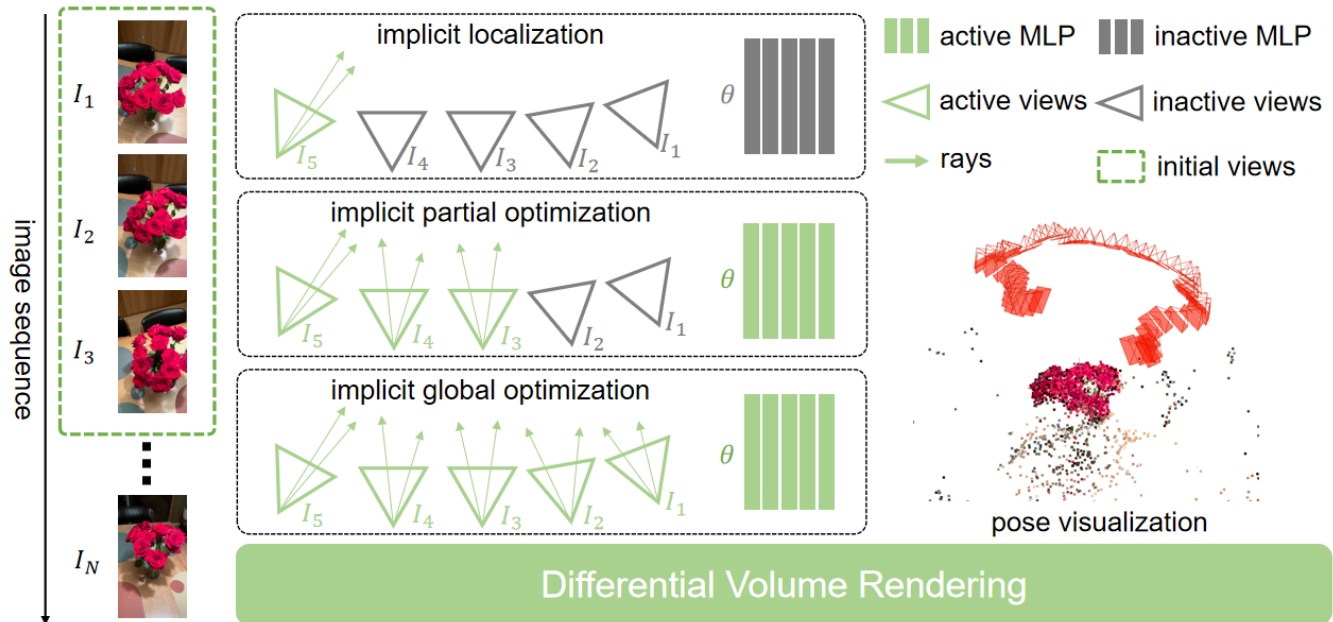


Figure 2: **The pipeline of CF-NeRF** CF-NeRF can estimate the weight θ of NeRF \mathcal{F} and the camera parameter δ . After initializing through a few selected images, CF-NeRF recovers δ of the image one by one through implicit localization that only optimizes the newly added image and implicit optimization that refines θ and δ . Implicit optimization can be divided into partial and global optimization depending on the number of images used. We visualize δ reconstructed by CF-NeRF and sparse points from COLMAP (Schonberger and Frahm 2016) to show that CF-NeRF can reconstruct rotation in image sequences.

Camera Parameter Estimation

Traditional SfM (Agarwal et al. 2011; Wu 2013; Chatterjee and Govindu 2013; Moulon, Monasse, and Marlet 2013; Schonberger and Frahm 2016; Moulon et al. 2017) and SLAM (Mur-Artal, Montiel, and Tardos 2015; Engel, Schöps, and Cremers 2014; Sumikura, Shibuya, and Sakurada 2019; Engel, Koltun, and Cremers 2017) can estimate camera parameters for given images. However, these methods divide the reconstruction pipeline into several different modules that need hand-crafted features (Lowe 2004; Bay, Tuytelaars, and Van Gool 2006) or learning-based methods (Yi et al. 2016; DeTone, Malisiewicz, and Rabinovich 2018) or the optical flow (Lucas and Kanade 1981; Teed and Deng 2020) to establish image correspondences, and then reconstruct a sparse scene and camera parameters through several separate non-differentiable steps.

In light of these limitations, it is worth exploring to estimate camera parameters during the training process of NeRF. The most direct attempt to utilize NeRF is the visual localization, where iNeRF (Yen-Chen et al. 2021), NeDDF (Ueda et al. 2022), and PNeRF (Lin et al. 2023) try to estimate the extrinsic camera parameter of a new image by a pre-trained NeRF model. Then, NeRFmm (Wang et al. 2021b) and SiRENmm (Guo and Sherwood 2021) take the NeRF and camera parameters as learnable and prove that it is possible to train the NeRF model from scratch without camera parameters, but they only work for forward-looking scenes. To further enhance accuracy in forward-looking or rotation scenes with initial camera parameters, BARF (Lin et al. 2021) dynamically adjusts the weight of the positional

encoding, GARF (Chng et al. 2022) replaces the ReLU activate function with the Gaussian activation function, and L2G-NeRF (Chen et al. 2023) introduces a local-to-global registration. Interestingly, GNeRF (Meng et al. 2021) and VMRF (Zhang et al. 2022b) assume there is a prior known distribution of camera parameters to decrease the freedom of camera parameters during training the NeRF model. Meanwhile, other researchers try to add different external restrictions to guide the camera parameter estimation. SCNeRF (Jeong et al. 2021) and Level- S^2 fM (Xiao et al. 2023) rely on feature matches to guide camera parameters estimation. NoPe-NeRF (Bian et al. 2023), iMap (Sucar et al. 2021), NeRF-SLAM (Rosinol, Leonard, and Carlone 2022), Nice-SLAM (Zhu et al. 2022), and Nicer-SLAM (Zhu et al. 2023) integrate depth maps from active sensors or CNN networks to tune the NeRF. Additionally, LocalLR (Meuleman et al. 2023) combines depth maps and optical flow to train NeRF.

Regrettably, images acquired from real-world scenarios often exhibit a multitude of challenges. These challenges include rotations and the absence of prior information of camera parameters. Furthermore, the introduction of external constraints can augment the intricacy and unpredictability of the reconstruction process. To solve these problems, we propose CF-NeRF inspired by the traditional incremental SfM, which does not require any prior information or external constraints while reconstructing the 3D scene and camera parameters end-to-end from image sequences, demonstrating the powerful reconstruction capability of the NeRF after using a specific training strategy.

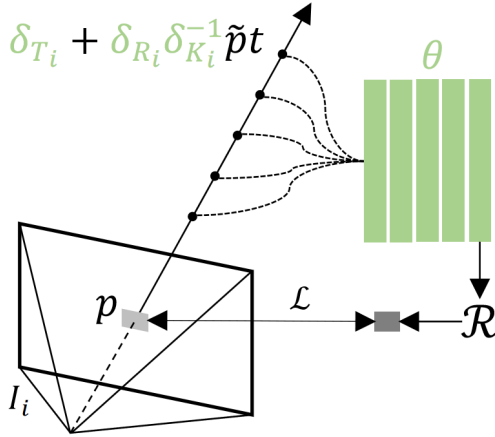


Figure 3: **Estimated Parameters** CF-NeRF estimates the weight θ of NeRF model and the camera parameter δ , which include the camera rotation δ_R , the camera translation δ_T , and camera intrinsic parameter δ_K .

Method

In this section, we provide an overview of the proposed method. Firstly, we introduce the preliminary background of the NeRF and the traditional incremental SfM. Then, we explain the details of CF-NeRF that can recover camera parameters from image sequences.

Preliminary Background

NeRF NeRF can generate realistic images from a set of images $I = (I_1, I_2, \dots, I_N)$ from N different places without explicitly reconstructing. However, NeRF needs associated camera parameters δ , including camera rotation $\delta_R = (\delta_{R_1}, \delta_{R_2}, \dots, \delta_{R_N})$, camera translation $\delta_T = (\delta_{T_1}, \delta_{T_2}, \dots, \delta_{T_N})$, and intrinsic camera parameter δ_K . Given a NeRF model \mathcal{F} and corresponding weight θ , it can estimate color c and density σ through an implicit function $c(x, \vec{d}), \sigma(x) = \mathcal{F}_\theta(x, \vec{d})$ with a point x and a view direction \vec{d} . To render a pixel p , NeRF needs to sample several points $x_p(t) = o + \vec{d}t$ along a ray shooting from the view position o and generate the color c_p by the volume rendering function \mathcal{R} as Eq. 1 shows, where $\mathcal{T}(t) = \exp(-\int_{t_n}^t \sigma(x_p(s)) ds)$ indicates the accumulated transmittance along the ray. t_n and t_f are the near and far bounds of the ray.

$$c_p = \mathcal{R}(p|\theta) = \int_{t_n}^{t_f} \mathcal{T}(t) \sigma(x_p(t)) c(x_p(t), \vec{d}) dt \quad (1)$$

Benefiting from the differential property of the volume rendering, NeRF can be trained end to end by minimizing the difference between c_p and observed color $I(p)$ as Eq. 2 shows, where \mathcal{L} is the loss function. To be noted, NeRF only estimates θ and borrows δ from traditional SfM methods. However, NeRFmm (Wang et al. 2021b) prove that it is possible to estimate θ and δ simultaneously under the forward-looking situation.

$$\arg \min_{\theta} \left\{ \sum_{I_i \in I} \sum_{p \in I_i} \mathcal{L}(\mathcal{R}(p|\theta), I_i(p)) \right\} \quad (2)$$

Incremental SfM Given a set of images, the incremental SfM can recover δ one by one in a linear time (Wu 2013) and contains four steps (Schonberger and Frahm 2016):

Initialization The selection of an initial two-view is essential because a suitable initial two-view improves the robustness and quality of the reconstruction. With a given two-view and its matched features, incremental SfM computes the relative pose by multi-view geometry (MVG) and triangulates 3D points to initial the scene.

Image Registration After initialization, incremental SfM adds images to the scene in order. Given a new image, incremental SfM builds the 2D-3D relationship by matching its features with images in the scene and recovers the camera parameter by Perspective-n-Point (PnP).

Triangulation As a newly added image observes additional information that can extend the scale of the scene, incremental SfM triangulates more 3D points based on the new image and matched features.

Bundle Adjustment Adding new images and 3D points without refinement leads to drift. Therefore, it is essential to apply bundle adjustment (BA) by minimizing the re-projection error. In terms of efficiency, incremental SfM proposes partial BA that refines only a subset of images, and global BA that optimizes all images.

CF-NeRF

Fusing the differentiability of NeRF and the reconstruction strategy of SfM, we propose CF-NeRF, which is capable of estimating the camera parameter under complex movement from sequential images. CF-NeRF consists of three modules: initialization, implicit localization, and implicit optimization, as Figure 2 shows. To convenient later introduction, we define the set of images we have completed estimating the camera parameter as E , which starts from \emptyset .

Parameter CF-NeRF estimates the camera parameter δ and the weight θ of NeRF, as Figure 3 shows. The δ includes δ_R , δ_T , and δ_K . During the differential volume rendering, we calculate the ray $\vec{r}_p(t) = \delta_{T_i} + \delta_{R_i} \delta_{K_i}^{-1} \vec{p} t$ of pixel p in image $I_i \in I$, where \vec{p} is the homogeneous expression of p . Following NeRFmm (Wang et al. 2021b), we use the axis-angle to represent δ_R and assume all images have the same camera intrinsic parameter without distortion so that δ_K only contains the focal length. We initialize δ_R and δ_T to zero, and set δ_K to 53° by a common field of view. The θ is the weight of \mathcal{F} that implicitly represents the scene. The activation function determines how to initialize θ . NeRF using ReLU are initialized according to NeRF (Mildenhall et al. 2020), while NeRF using sine are initialized according to SIREN (Sitzmann et al. 2020).

Initialization Similar to incremental SfM, CF-NeRF requires initialize θ , δ_{R_1} , δ_{T_1} , and δ_K before adding images to E . We select the first N_{init} images I_{init} from I to optimise these parameters by Eq. 3 with ξ_{init} iterations. Since the rotation between adjacent images is not large and NeRF is hard

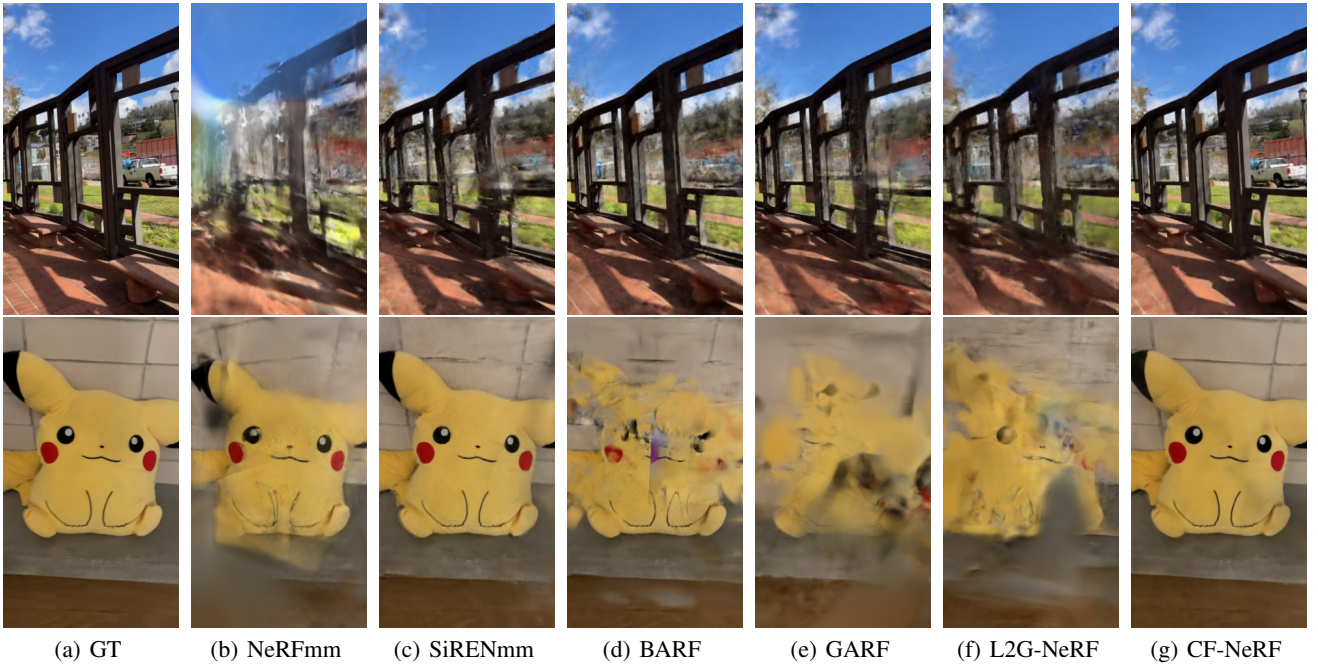


Figure 4: **Qualitative results on NeRFBuster.** We select two sequences from NeRFBuster (Warburg et al. 2023) and render novel views to compare the quality of camera parameters. Our method CF-NeRF generate high-quality images, while results of NeRFmm (Wang et al. 2021b), SiRENmm (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022) and L2G-NeRF (Chen et al. 2023) contain lots of noise.

to estimate rotation (Lin et al. 2021; Chng et al. 2022; Chen et al. 2023), we do not estimate the rotation in the initialization to reduce the freedom. After initialization, we add I_1 to E and keep θ , δ_{R_1} , δ_{T_1} , and δ_K but discard other camera parameters. Note that, unlike the initialization in the previous section, the initialization here is data-specific, similar to the warm-up procedure.

$$\arg \min_{\theta, \delta_T, \delta_K} \left\{ \sum_{I_i \in I_{init}} \sum_{p \in I_i} \mathcal{L}(\mathcal{R}(p|\theta, \delta_{T_i}, \delta_K), I_i(p)) \right\} \quad (3)$$

Implicit Localization After initialization, CF-NeRF estimates the camera parameter of the remaining images one by one and determines δ_{R_n} and δ_{T_n} for each new image I_n by localization. Specifically, we first initialize δ_{R_n} and δ_{T_n} by $\delta_{R_{n-1}}$ and $\delta_{T_{n-1}}$, and then optimize them by minimizing Eq. 4 with fixed θ through ξ_{loc} iterations. The localization is similar to iNeRF (Yen-Chen et al. 2021), but CF-NeRF does not have a pre-trained \mathcal{F} .

$$\arg \min_{\delta_{R_n}, \delta_{T_n}} \left\{ \sum_{p \in I_n} \mathcal{L}(\mathcal{R}(p|\delta_{R_n}, \delta_{T_n}), I_n(p)) \right\} \quad (4)$$

Implicit Optimization Although implicit localization can roughly determine δ_{R_n} and δ_{T_n} , it faces two problems: the observation from I_n is not added to NeRF, and the localization does not take the multi-view consistency into account to reduce drift. Incremental SfM solves these problems using two separate steps: triangulation and BA, while CF-NeRF benefits from the volume rendering and deals with these problems together. However, it is time-consuming to

optimize all images in E every time a new image is added. Therefore, CF-NeRF splits optimization into implicit partial optimization and implicit global optimization.

Each time localizing a new image I_n , CF-NeRF performs implicit partial optimization. We select I_n and previous $N_{part} - 1$ images to construct the partial image set I_{part} , then optimizes them with ξ_{part} iterations, as Eq.5 shows.

$$\arg \min_{\theta, \delta_R, \delta_T} \left\{ \sum_{I_i \in I_{part}} \sum_{p \in I_i} \mathcal{L}(\mathcal{R}(p|\theta, \delta_{R_i}, \delta_{T_i}), I_i(p)) \right\} \quad (5)$$

When the number of images in E can be evenly divided by N_{glob} , CF-NeRF employs implicit global optimization for θ and all images in E to enhance the overall accuracy and reduce drifts with ξ_{glob} iterations, as Eq. 6 shows.

$$\arg \min_{\theta, \delta_R, \delta_T, \delta_K} \left\{ \sum_{I_i \in I_E} \sum_{p \in I_i} \mathcal{L}(\mathcal{R}(p|\theta, \delta_{R_i}, \delta_{T_i}, \delta_K), I_i(p)) \right\} \quad (6)$$

Coarse-to-Fine CF-NeRF uses a coarse-to-fine strategy to improve robustness. CF-NeRF first constructs a Gaussian pyramid with depth d_G , then recovers all parameters at a low-resolution image through the incremental pipeline. Finally, CF-NeRF directly performs implicit global optimization with a higher resolution in each scale of the Gaussian pyramid with ξ_G iterations.

Loss Function To improve robustness, we employ the Smooth-L1 loss function, as Eq. 7 shows, where gt represents the ground truth, pr is the estimated value, and β is the set to 1.0 by default.

Table 1: We quantitatively evaluated on 12 scenes from the NeRFBuster dataset (Warburg et al. 2023), which is captured in the real world with complex trajectories, including rotations and forward-looking. CF-NeRF achieves state-of-the-art in ΔR , ΔT , PSNR, and LPIPS compared to NeRFmm (Wang et al. 2021b), SiRENmm (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022), L2G-NeRF (Chen et al. 2023).

		NeRFBuster											
		aloe	art	car	century	flowers	garbage	picnic	pikachu	pipe	plant	roses	table
$\Delta R \downarrow$	NeRFmm	159.9731	177.5907	129.5799	119.6261	106.9195	150.8233	154.7781	113.7001	164.8207	165.0298	102.2748	115.2991
	SiRENmm	155.1513	177.3642	127.2673	89.0172	103.8741	82.9375	44.3671	25.3603	159.7565	114.0761	132.5378	93.2612
	BARF	158.6689	59.1868	133.4527	101.6013	88.6842	88.7832	69.7201	41.0302	64.5250	143.1975	133.7586	111.2884
	GARF	125.9797	171.9169	153.5585	105.1868	106.0595	84.4992	49.8503	32.7285	126.9602	156.6058	118.9748	164.6562
	L2G-NeRF	124.2368	24.3753	55.7291	131.9683	96.3498	110.4782	146.3119	116.8911	70.9131	70.7562	95.9982	116.4972
	CF-NeRF	12.1226	19.2496	17.5570	9.6811	8.2556	9.7658	12.6501	11.3067	19.9926	4.8968	5.1229	4.5837
	NeRFmm	11.5935	15.0762	23.9514	24.5934	12.8753	16.3842	12.9675	25.6841	19.3563	23.6613	8.0367	13.3849
SiRENmm	11.4912	14.8720	27.7235	28.8582	15.4841	13.3099	8.3607	15.8052	20.4572	31.2943	9.0498	13.8350	
BARF	9.6196	17.3299	36.0351	26.0549	15.3166	17.1936	9.6679	18.7184	18.8936	31.8629	8.5266	17.5453	
GARF	13.9100	17.1527	26.5184	26.0096	14.1459	16.1054	11.7021	17.4975	18.8366	32.3170	11.2654	15.0367	
L2G-NeRF	14.4012	13.4240	20.5634	23.2650	7.4559	17.7167	12.9408	32.4048	11.4012	18.5012	10.7110	12.8061	
CF-NeRF	3.3788	2.2821	6.5452	2.7383	2.7026	4.0535	1.2833	4.0586	9.4491	3.4346	1.1945	1.2127	
PSNR \uparrow	NeRFmm	20.6912	16.8220	17.7504	16.0326	18.4377	17.7229	18.4300	25.3819	20.0978	21.1558	17.7735	14.1651
	SiRENmm	22.7462	20.3890	22.0268	18.0252	19.4640	17.2283	21.5628	27.8706	20.9538	23.4980	16.7480	18.8135
	BARF	22.4366	21.1947	16.7665	15.3436	17.8350	15.9065	19.1846	23.0386	19.9728	25.5135	13.6741	13.8227
	GARF	19.0241	19.3556	15.4460	14.4117	16.2955	15.3383	15.4035	20.9663	18.5371	20.5600	13.1274	12.6677
	L2G-NeRF	21.3398	19.8099	17.3255	16.6476	18.0016	13.6077	18.5268	22.4939	18.1787	19.0160	17.2614	15.5658
	CF-NeRF	26.9367	26.5293	22.4654	21.7072	21.6950	22.4736	22.5475	32.3661	22.2719	25.7312	24.3918	26.8491
	NeRFmm	0.5560	0.4954	0.5991	0.5793	0.5778	0.5661	0.6113	0.3683	0.5614	0.4927	0.5371	0.6073
SiRENmm	0.4508	0.4034	0.4450	0.4785	0.5048	0.5193	0.5227	0.2883	0.5170	0.3256	0.5333	0.4659	
BARF	0.3328	0.3511	0.5361	0.5394	0.5552	0.5480	0.5358	0.3440	0.5198	0.3217	0.6138	0.5913	
GARF	0.5257	0.4055	0.5984	0.5845	0.6158	0.5931	0.6086	0.3987	0.5688	0.4345	0.6189	0.6356	
L2G-NeRF	0.4620	0.4186	0.5409	0.5116	0.5466	0.6016	0.5530	0.4051	0.4741	0.3840	0.4788	0.5309	
CF-NeRF	0.1939	0.2316	0.3983	0.3627	0.3983	0.3859	0.4686	0.1679	0.4453	0.2594	0.2831	0.3011	

$$\mathcal{L}(pr, gt) = \begin{cases} 0.5 * (gt - pr)^2 / \beta & \text{if } |gt - pr| < \beta \\ |gt - pr| - 0.5 * \beta & \text{otherwise} \end{cases} \quad (7)$$

Experiments

Dataset

Instead of using the forward-looking dataset LLFF (Mildenhall et al. 2020) as usually conducted in the previous studies, we use a real-world dataset from NeRFBuster (Warburg et al. 2023), mainly rotating around an object, to evaluate our method. We sample around 50 frames for each scene and re-size all images to 480×270 . We use COLMAP (Schonberger and Frahm 2016) to estimate ground truth (GT) camera parameters for all scenes. More details of the dataset can be found in the supplementary material.

Implementation

CF-NeRF is implemented using PyTorch. Similar to NeRFmm (Wang et al. 2021b), CF-NeRF does not have hierarchical sampling and uses the coarse network, which has eight layers and the dimension of the hidden layers is set to 128. And we use the sine activation function instead of the ReLU, as SiRENmm (Guo and Sherwood 2021) is more robust than NeRFmm. We utilize the Adam optimizer (Kingma and Ba 2014) to optimize all learnable parameters. Specifically, we set the learning rate of θ to 0.001, which undergoes a decay of 0.9954 every 200 epochs. Similarly, the learning rate of δ is set to 0.001 and undergoes a decay of 0.9000 every 2000 epochs. Here we describe how to set the hyper-parameters in CF-NeRF. We set N_{init} and N_{part}

to 3 to meet the minimum requirements that can filter outliers based on MVG. To balance drift and efficiency, we set N_{glob} to 5. Considering the input image resolution, we set d_G to 3 to reconstruct all parameters by coarse-to-fine strategy. The most important parameter in CF-NeRF is iteration, which is the epoch number for each image. During initialization, we set ξ_{init} to 3000 to guarantee that θ and δ can be correctly initialized with fewer images. Subsequently, during the incremental training, we maintain a consistent value of ξ , setting $\xi = \xi_{loc} = \xi_{part} = \xi_{glob} = \xi_G$ to 900, thus reconstructing the scene from images one by one. Throughout all our experiments, we use the NVIDIA RTX3090.

Evaluation

To demonstrate the performance of the proposed method, we conduct a comprehensive comparison between CF-NeRF and several state-of-the-art models, including NeRFmm (Wang et al. 2021b) SiRENmm (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022), and L2G-NeRF (Chen et al. 2023). We use all images for camera parameter estimation without employing a train/test split. To evaluate the quality of the camera parameters, we calculate the average translation error ΔT and the average rotation error ΔR by aligning the estimated camera parameters δ_R and δ_T with COLMAP using a similarity transformation Sim(3) (Lin et al. 2021). It is worth noting that δ_T represents a relative translation error rather than an absolute measurement, as COLMAP can not reconstruct an absolute scale of the scene. We further evaluate the estimated camera parameters through a novel view synthesis by PSNR and LPIPS. To ensure a fair comparison and avoid the influence of varying

Table 2: **Ablation experiments.** We compare the accuracy of camera parameter estimation of CF-NeRF under different hyper-parameter settings, including the iteration ξ , the global optimization frequency N_{glob} and the coarse-to-fine strategy, where C means the coarse stage and F means the fine stage.

		NeRFBuster											
		aloe	art	car	century	flowers	garbage	picnic	pikachu	pipe	plant	roses	table
$\Delta R \downarrow$	$F, \xi = 600, N_{glob} = 10$	17.8029	24.3389	17.1692	11.5924	11.6163	9.1240	14.6452	13.0037	19.0749	5.3354	5.9091	6.4731
	$F, \xi = 900, N_{glob} = 10$	14.8730	22.8142	17.8879	11.1201	10.4707	8.6973	11.4209	12.0625	18.4305	4.7303	5.8538	6.8481
	$C, \xi = 900, N_{glob} = 5$	12.4862	19.1647	17.4755	9.7177	8.4555	9.6460	12.3162	10.9802	19.9855	5.1579	5.5133	5.2821
	$F, \xi = 900, N_{glob} = 5$	12.1226	19.2496	17.5570	9.6811	8.2556	9.7658	12.6501	11.3067	19.9926	4.8968	5.1229	4.5837
$\Delta T \downarrow$	$F, \xi = 600, N_{glob} = 10$	4.5457	5.9307	7.5697	2.9652	3.6234	4.5340	2.6677	5.3105	9.3544	4.4384	1.3324	1.9535
	$F, \xi = 900, N_{glob} = 10$	3.9111	6.1190	7.3752	3.6834	3.3956	4.3080	3.4918	3.2682	8.4666	3.6109	1.3013	2.3973
	$C, \xi = 900, N_{glob} = 5$	3.4681	2.2770	6.6250	2.8224	2.7405	4.1085	1.2886	4.2462	9.6998	3.5309	1.2182	1.2232
	$F, \xi = 900, N_{glob} = 5$	3.3788	2.2821	6.5452	2.7383	2.7026	4.0535	1.2833	4.0586	9.4491	3.4346	1.1945	1.2127

network backbones across different methods, we uniformly use the NerfAcc (Li, Tancik, and Kanazawa 2022), where we select one image for testing in every eight images and the remaining is for training.

Results

We performed qualitative and quantitative evaluations of these methods on 12 scenes of the NeRFBuster (Warburg et al. 2023) dataset. Notably, BARF (Lin et al. 2021), GARF (Chng et al. 2022), and L2G-NeRF (Chen et al. 2023) require manual setting the focal length. In contrast, NeRFmm (Wang et al. 2021b), SiRENmm (Guo and Sherwood 2021), and CF-NeRF have the ability to estimate the focal length.

Table 3 shows the results of qualitative experiments. Our method obtains the highest accuracy camera parameters, while all other methods fail outright. It is important to understand that ΔR and ΔT are calculated by aligning the camera positions with Sim(3) and that a slight difference in camera position can lead to huge errors. The rotation error ΔR of our method CF-NeRF is roughly around 10° , while the other methods are around 100° . Moreover, the translation error δ_T of CF-NeRF is approximately about 4, while all other methods are around 15. Although NeRFmm, SiRENmm, BARF, GARF, and L2G-NeRF claim high accuracy on forward-looking scenes from scratch, they are unsuitable for scenes with rotation and are prone to be trapped in a local minimum. In contrast, CF-NeRF recovers the camera parameters sequentially and can effectively handle image sequences with complex trajectories. Furthermore, SiRENmm outperforms NeRFmm in camera parameter estimation, which is why CF-NeRF uses the sine activate function.

Table 3 also shows the quality of the novel view synthesis, which serves as an additional evaluation criterion for the quality of camera parameters. CF-NeRF achieves state-of-the-art results on PSNR and LPIPS, whereas other methods struggle to generate high-quality results across multiple novel views. Interestingly, the reconstruction results of other methods appear reasonable compared to their poor camera parameters, mainly due to the high over-fitting ability of NeRF and the fact that partial camera parameters are correctly reconstructed. We further visualize the rendering results of three scenes from different methods in Figure 7 and Figure 4. CF-NeRF can generate high-quality results, while other methods have lots of noise in their results due to their inability to provide accurate camera parameters for all im-

ages. More visualization results are in supplementary.

Ablation Experiments

To understand the influence of hyper-parameters in CF-NeRF, we conduct several ablation experiments on the iteration ξ , the global optimization frequency N_{glob} , and the coarse-to-fine strategy. These experiments are performed on the NeRFBuster dataset (Warburg et al. 2023), and results are presented in Table 2.

The iteration ξ The iteration ξ is the most important hyper-parameter in our method, determining how many times to optimize the camera parameter for each image. We compare two configurations: $F, \xi = 600, N_{glob} = 10$ and $F, \xi = 900, N_{glob} = 10$. Table 2 reveals that increasing ξ improves the final results for almost all scenes. This observation aligns with NeRF (Mildenhall et al. 2020) and iNeRF (Yen-Chen et al. 2021), where NeRF requires a large number of iterations to converge, and iNeRF enhances the quality of camera parameters through more iterations.

The global optimization frequency N_{glob} To mitigate drift while maintaining efficiency, CF-NeRF employs the implicit global optimization when every N_{glob} image is added E . We conduct two experiments $F, \xi = 900, N_{glob} = 10$ and $F, \xi = 900, N_{glob} = 5$ to find out the influence of N_{glob} . As highlighted in Table 2, reducing N_{glob} yields improved final results, which can be attributed to the fact that global optimization ensures global consistency to avoid NeRF trap into a local minimum.

The coarse-to-fine strategy CF-NeRF adopts a coarse-to-fine strategy to avoid directly estimating camera parameters on high-resolution images, where the fine stage refines initial results from the coarse stage. We conduct two experiments $C, \xi = 900, N_{glob} = 5$ and $F, \xi = 900, N_{glob} = 5$, where C means the coarse stage and F means the fine stage. Results in Table 2 demonstrate that the fine stage outperforms the coarse stage across almost all scenes. Adopting the coarse-to-fine strategy facilitates the training process of CF-NeRF, as the pixel gradient is smoother at the coarse stage and has less RGB information to learn.

Limitation

Although CF-NeRF achieves state-of-the-art results in camera parameter estimation, surpassing other NeRF-based methods, there are still some gaps between CF-NeRF and

COLMAP (Schonberger and Frahm 2016), and the accuracy can be further improved through the adjustment of the sample space (Wang et al. 2023) or the utilization of a more robust function (Sabour et al. 2023).

Conclusion

This paper presents CF-NeRF, a novel end-to-end method that does not require prior camera parameters to deal with image sequences with complex trajectories. Following the pipeline of incremental SfM, CF-NeRF contains three major sub-modules: initialization, implicit localization, and implicit optimization. Experiments on the NeRFBuster dataset demonstrate that CF-NeRF achieves state-of-the-art results, while NeRFmm, SiRENmm, BARF, GARF, and L2G-NeRF only work for forward-looking scenes and get trapped in the local minimum on the NeRFBuster dataset. More importantly, CF-NeRF highlights the unlimited potential of NeRF and differential volume rendering, showing that NeRF has impressive reconstruction capabilities and can also be used to estimate camera parameters in complex trajectories.

Acknowledgments The research was supported in part by a RGC RIF grant under the contract R6021-20, RGC CRF grants under the contracts C7004-22G and C1029-22G, and RGC GRF grants under the contracts 16209120, 16200221, and 16207922. This research was also supported by the National Natural Science Foundation of China (No. 62302126), and the Shenzhen Science and Technology Program (No. RCBS20221008093125065, No. JCYJ20220818102414030).

Supplementary material

How to use COLMAP

COLMAP (Schonberger and Frahm 2016) is a widely used 3D reconstruction software that can produce stable reconstruction results, providing the intrinsic and extrinsic camera parameters and the sparse 3D point clouds. In this paper, we use COLMAP to process all datasets and use these results as the ground truth (GT) to evaluate camera parameters estimated by different methods. However, we make some modifications to the default settings of COLMAP to align with the unique characteristics of the datasets.

1. We assume no distortions in images and use the pinhole camera model that only considers the focal length.
2. We upsample the image twice during feature extraction to ensure that COLMAP can extract enough feature points with traditional feature extraction algorithms.
3. We use the sequential matching method to match images, which is much faster than the exhaustive matching method and more suitable for the incremental structure from motion (SfM) for an image sequence.

Details of the NeRFBuster dataset

In all experiments, we use the challenging NeRFBuster dataset (Warburg et al. 2023), which is captured in the real world. This dataset contains 12 scenes, including aloe, art, car, century, flowers, garbage, picnic, pikachu, pipe, plant, flowers, and table, as Figure 5 shows. Furthermore, we present camera parameters and sparse point clouds results obtained from COLMAP in Figure 6. Unlike the LLFF dataset (Mildenhall et al. 2020), which is limited to forward-looking motion, the NeRFBuster dataset is more challenging for camera parameters estimation, because it is captured under complex movement and contains a variety of objects.

How to use NeRFacc

Instead of only evaluating the camera parameter by comparing the camera translation and the camera rotation, we also indirectly evaluate the camera parameter by novel view synthesis. Because we use all images during camera parameter estimation and different methods use different backbones, we uniformly use NeRFacc (Li, Tancik, and Kanazawa 2022) to measure Peak Signal-to-Noise Ratio (PSNR) and Learned Perceptual Image Patch Similarity (LPIPS). For each scene in NeRFBuster, we follow these steps:

1. We align the estimated camera parameters to the GT camera parameters provided by COLMAP through a similarity transformation $\text{Sim}(3)$ (Lin et al. 2021).
2. We split the images in each scene into a training/split set, with one-eighth of the images as test images and the rest as training images.
3. We use the bounding box of sparse points provided by the COLMAP and utilize NeRFacc to reconstruct the scene with the aligned camera parameters implicitly.
4. We train the model for 20000 epochs on the train split and evaluate PSNR and LPIPS on the test split.

During the training process, we employ the default settings from NeRFacc and utilize the MipNeRF (Barron et al. 2021) as the backbone, which is motivated by the abundance of open-space scenes present in the NeRFBuster dataset.

How to align camera parameters

We follow BARF (Lin et al. 2021) to align two camera parameters through $\text{Sim}(3)$, where only the camera translation is used. However, we find this method is unreliable when all cameras lie in a line or the camera translation contains noise. To overcome this problem and improve robustness, we add a virtual point $(0, 0, 1)$ in the camera coordinate of each image and use the camera parameter to transform it to the world coordinate, thus using the camera rotation during the alignment process.

Visualization Results

We qualitatively compare the camera parameters of our method CF-NeRF with state-of-the-art methods NeRFmm (Wang et al. 2021b), SiRENmm (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022), and L2G-NeRF (Chen et al. 2023) on the NeRFBuster dataset

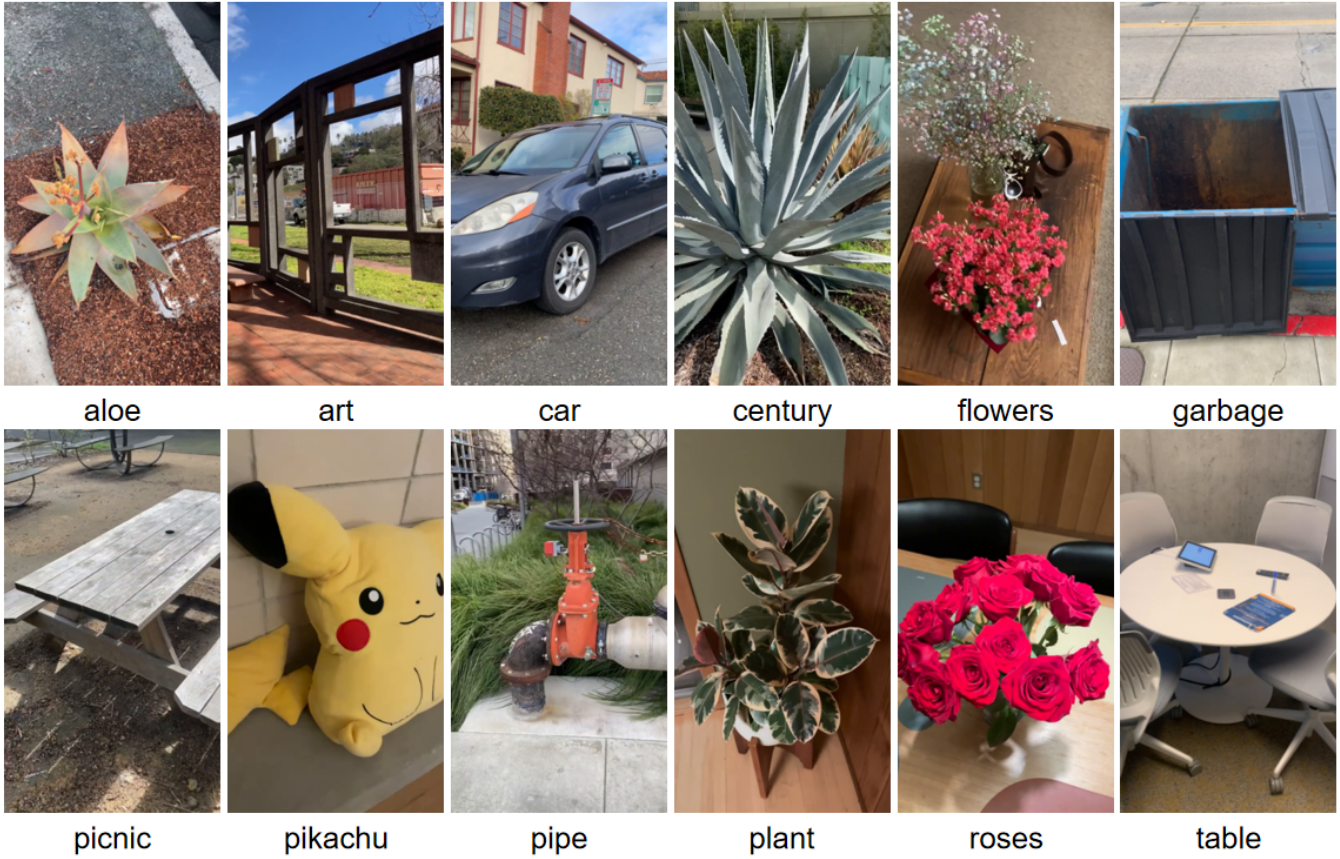


Figure 5: **12 scenes from NeRF Buster** We show one image from each scene in the NeRF Buster dataset (Warburg et al. 2023), which contains various objects from indoor or outdoor environments.

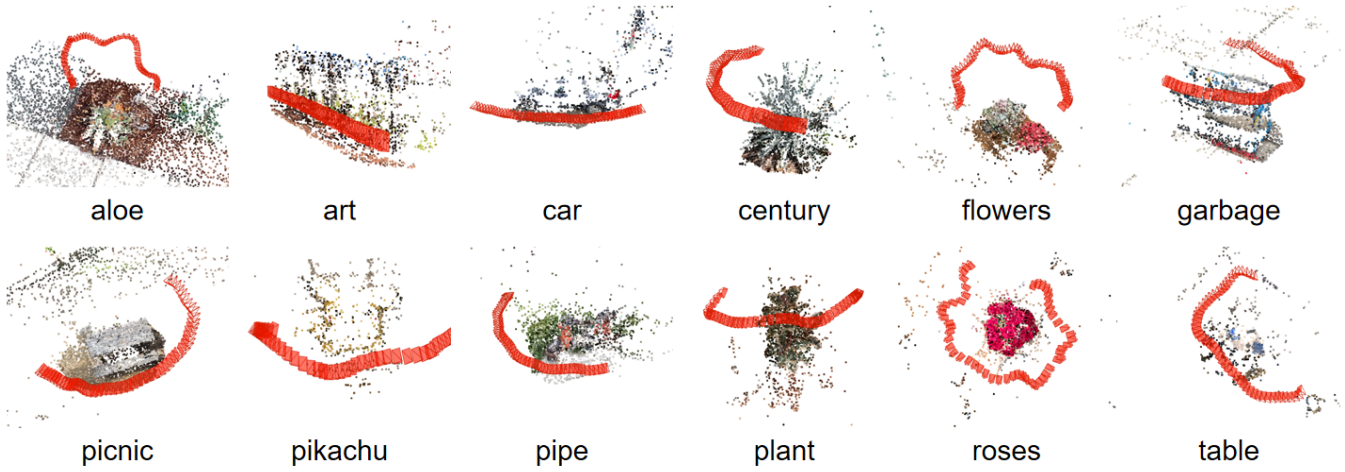


Figure 6: **Sparse reconstruction results of the NeRF Buster** The NeRF Buster dataset (Warburg et al. 2023) is captured under complex trajectories and mainly rotates around objects.

(Warburg et al. 2023). Figure 7 shows the camera parameter estimation results of the different methods as well as the sparse points of COLMAP, and it is not difficult to notice that only our CF-NeRF generates correct results, while the other methods generate entirely wrong results. For a full comparison, we also provide a video in the supplementary, using NeRFacc to render the video according to the camera parameters of the different methods.

Table 3: **Quantitative evaluation results.** We evaluate our method on LLFF and Real360 from NeRF. BARF (Lin et al. 2021) needs a pre-calculated focal length from COLMAP. NeRFmm (Wang et al. 2021b) and CF-NeRF (ours) do not need prior camera parameters.

scene	rotation↓			translation↓			PSNR↑		
	NF-	ours	BARF	NF-	ours	BARF	NF-	ours	BARF
fern	1.2724	0.8669	0.2244	0.2469	0.2218	0.0811	23.71	26.38	27.19
flower	4.5220	10.7385	1.1812	0.6570	0.5587	0.0482	23.04	26.80	29.06
fortress	174.296	3.0570	0.2691	6.0263	1.3652	0.0972	9.62	27.19	26.92
horns	4.3958	3.3078	0.2477	0.7677	1.1096	0.0644	18.74	26.71	29.29
leaves	0.4504	11.2940	1.7657	0.2216	0.8991	0.1823	18.94	20.88	20.12
orchids	4.7278	0.1951	0.2309	0.2728	0.0703	0.0641	14.79	22.26	22.21
room	3.1178	3.7610	0.2552	0.3987	2.2731	0.1712	24.20	28.76	32.63
trex	5.0095	0.6298	1.0239	0.3608	0.4072	0.3525	23.61	28.37	28.86
vasedec	148.9928	6.0560	165.0392	18.6165	0.3143	12.2520	16.64	23.63	14.74
pinecone	122.9956	0.4887	42.6464	29.5298	0.1387	8.2233	10.66	23.80	10.88

Experiments on LLFF

Throughout our paper, we use the NeRFBuster, since it is a real-world dataset captured by mobile phones and contains quite complex trajectories, whereas the commonly used LLFF contains only forward views with simple trajectories. However, to fully compare with other methods, we conduct experiments on LLFF as well. As shown in Table 3, our method achieves competitive results to NeRFmm and BARF (with pre-calculated focal length) in the eight forward-looking scenes of LLFF. However, for the two scenes with rotation in NeRF360, our method achieves much better results than NeRFmm and BARF.

References

Agarwal, S.; Furukawa, Y.; Snavely, N.; Simon, I.; Curless, B.; Seitz, S. M.; and Szeliski, R. 2011. Building rome in a day. *Communications of the ACM*, 54(10): 105–112.

Attal, B.; Laidlaw, E.; Gokaslan, A.; Kim, C.; Richardt, C.; Tompkin, J.; and O’Toole, M. 2021. Törf: Time-of-flight radiance fields for dynamic scene view synthesis. *Advances in neural information processing systems*, 34: 26289–26301.

Barron, J. T.; Mildenhall, B.; Tancik, M.; Hedman, P.; Martin-Brualla, R.; and Srinivasan, P. P. 2021. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5855–5864.

Barron, J. T.; Mildenhall, B.; Verbin, D.; Srinivasan, P. P.; and Hedman, P. 2022. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5470–5479.

Bay, H.; Tuytelaars, T.; and Van Gool, L. 2006. Surf: Speeded up robust features. *Lecture notes in computer science*, 3951: 404–417.

Bian, W.; Wang, Z.; Li, K.; Bian, J.-W.; and Prisacariu, V. A. 2023. Nope-nerf: Optimising neural radiance field with no pose prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4160–4169.

Boss, M.; Braun, R.; Jampani, V.; Barron, J. T.; Liu, C.; and Lensch, H. 2021a. Nerf: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 12684–12694.

Boss, M.; Engelhardt, A.; Kar, A.; Li, Y.; Sun, D.; Barron, J. T.; Lensch, H.; and Jampani, V. 2022. SAMURAI: Shape and material from unconstrained real-world arbitrary image collections. *arXiv preprint arXiv:2205.15768*.

Boss, M.; Jampani, V.; Braun, R.; Liu, C.; Barron, J.; and Lensch, H. 2021b. Neural-pil: Neural pre-integrated lighting for reflectance decomposition. *Advances in Neural Information Processing Systems*, 34: 10691–10704.

Chatterjee, A.; and Govindu, V. M. 2013. Efficient and robust large-scale rotation averaging. In *Proceedings of the IEEE International Conference on Computer Vision*, 521–528.

Chen, A.; Xu, Z.; Geiger, A.; Yu, J.; and Su, H. 2022. Tensorf: Tensorial radiance fields. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, 333–350. Springer.

Chen, Y.; Chen, X.; Wang, X.; Zhang, Q.; Guo, Y.; Shan, Y.; and Wang, F. 2023. Local-to-global registration for bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8264–8273.

Chng, S.-F.; Ramasinghe, S.; Sherrah, J.; and Lucey, S. 2022. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *European Conference on Computer Vision*, 264–280. Springer.

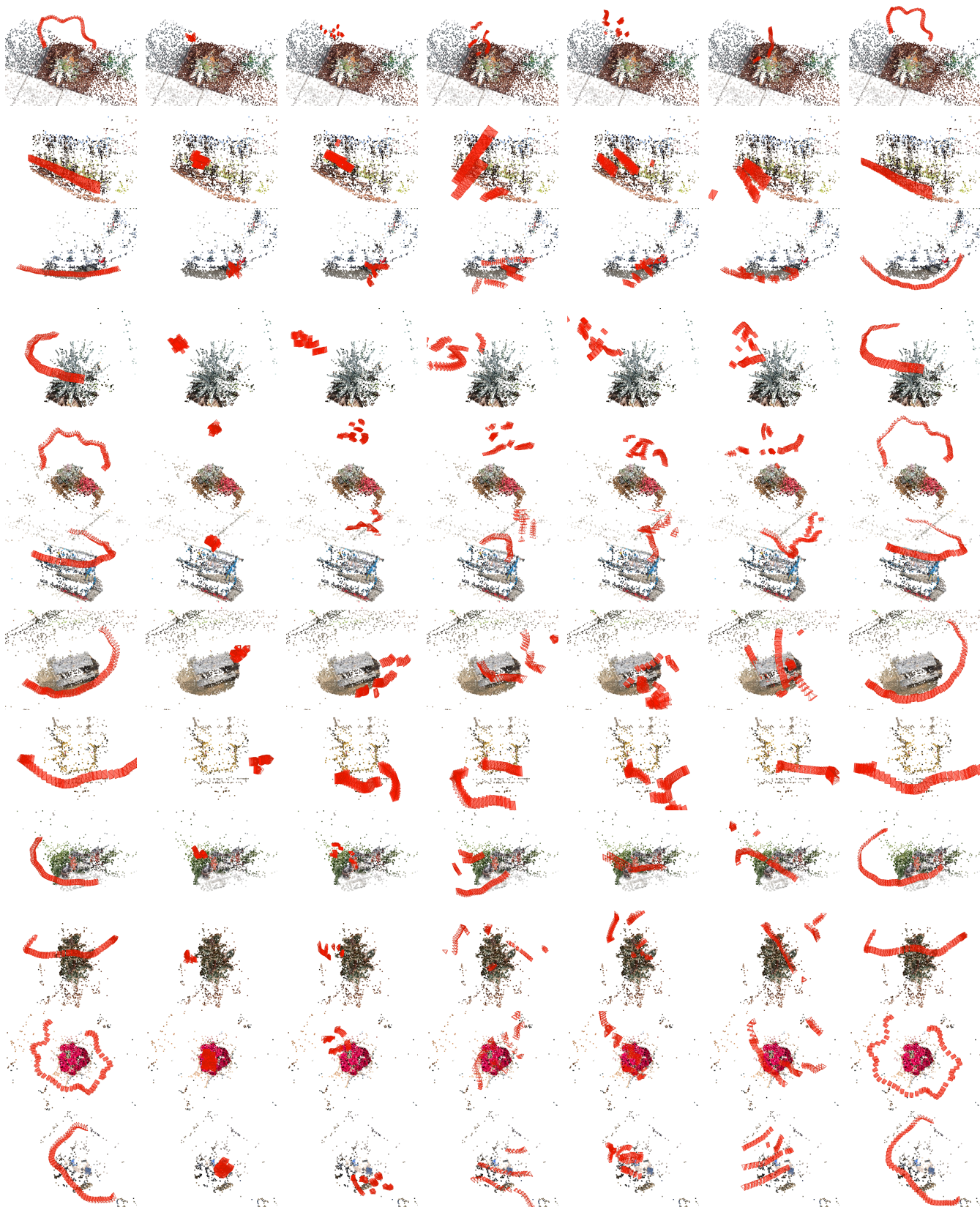
Darmon, F.; Basclé, B.; Devaux, J.-C.; Monasse, P.; and Aubry, M. 2022. Improving neural implicit surfaces geometry with patch warping. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 6260–6269.

DeTone, D.; Malisiewicz, T.; and Rabinovich, A. 2018. Superpoint: Self-supervised interest point detection and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 224–236.

Engel, J.; Koltun, V.; and Cremers, D. 2017. Direct sparse odometry. *IEEE transactions on pattern analysis and machine intelligence*, 40(3): 611–625.

Engel, J.; Schöps, T.; and Cremers, D. 2014. LSD-SLAM: Large-scale direct monocular SLAM. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part II 13*, 834–849. Springer.

Fridovich-Keil, S.; Yu, A.; Tancik, M.; Chen, Q.; Recht, B.; and Kanazawa, A. 2022. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5501–5510.



(a) GT (b) NeRFmm (c) SiRENmm (d) BARF (e) GARF (f) L2G-NeRF (g) CF-NeRF

Figure 7: **Comparison of camera parameters.** We compare camera parameters of our method CF-NeRF with NeRFmm (Wang et al. 2021b), SiNeRF (Guo and Sherwood 2021), BARF (Lin et al. 2021), GARF (Chng et al. 2022), and L2G-NeRF (Chen et al. 2023) along with sparse points from COLMAP.

- Fu, Q.; Xu, Q.; Ong, Y. S.; and Tao, W. 2022. Geo-neus: Geometry-consistent neural implicit surfaces learning for multi-view reconstruction. *Advances in Neural Information Processing Systems*, 35: 3403–3416.
- Garbin, S. J.; Kowalski, M.; Johnson, M.; Shotton, J.; and Valentin, J. 2021. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 14346–14355.
- Guo, J.; and Sherwood, A. 2021. improved-nerfmm. <https://github.com/ventusff/improved-nerfmm>.
- Jain, A.; Tancik, M.; and Abbeel, P. 2021. Putting nerf on a diet: Semantically consistent few-shot view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5885–5894.
- Jeong, Y.; Ahn, S.; Choy, C.; Anandkumar, A.; Cho, M.; and Park, J. 2021. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5846–5854.
- Kim, M.; Seo, S.; and Han, B. 2022. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12912–12921.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kuang, Z.; Olszewski, K.; Chai, M.; Huang, Z.; Achlioptas, P.; and Tulyakov, S. 2022. NeROIC: neural rendering of objects from online image collections. *ACM Transactions on Graphics (TOG)*, 41(4): 1–12.
- Li, R.; Tancik, M.; and Kanazawa, A. 2022. NerfAcc: A General NeRF Acceleration Toolbox. *arXiv preprint arXiv:2210.04847*.
- Lin, C.-H.; Ma, W.-C.; Torralba, A.; and Lucey, S. 2021. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5741–5751.
- Lin, Y.; Müller, T.; Tremblay, J.; Wen, B.; Tyree, S.; Evans, A.; Vela, P. A.; and Birchfield, S. 2023. Parallel inversion of neural radiance fields for robust pose estimation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 9377–9384. IEEE.
- Long, X.; Lin, C.; Liu, L.; Liu, Y.; Wang, P.; Theobalt, C.; Komura, T.; and Wang, W. 2023. Neuraludf: Learning unsigned distance fields for multi-view reconstruction of surfaces with arbitrary topologies. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20834–20843.
- Long, X.; Lin, C.; Wang, P.; Komura, T.; and Wang, W. 2022. Sparseneus: Fast generalizable neural surface reconstruction from sparse views. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, 210–227. Springer.
- Lowe, D. G. 2004. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60: 91–110.
- Lucas, B. D.; and Kanade, T. 1981. An iterative image registration technique with an application to stereo vision. In *IJCAI’81: 7th international joint conference on Artificial intelligence*, volume 2, 674–679.
- Martin-Brualla, R.; Radwan, N.; Sajjadi, M. S.; Barron, J. T.; Dosovitskiy, A.; and Duckworth, D. 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.
- Meng, Q.; Chen, A.; Luo, H.; Wu, M.; Su, H.; Xu, L.; He, X.; and Yu, J. 2021. Gnerf: Gan-based neural radiance field without posed camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6351–6361.
- Meuleman, A.; Liu, Y.-L.; Gao, C.; Huang, J.-B.; Kim, C.; Kim, M. H.; and Kopf, J. 2023. Progressively optimized local radiance fields for robust view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16539–16548.
- Mi, Z.; Di, C.; and Xu, D. 2022. Generalized binary search network for highly-efficient multi-view stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12991–13000.
- Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I*, 405–421.
- Moulon, P.; Monasse, P.; and Marlet, R. 2013. Global fusion of relative motions for robust, accurate and scalable structure from motion. In *Proceedings of the IEEE international conference on computer vision*, 3248–3255.
- Moulon, P.; Monasse, P.; Perrot, R.; and Marlet, R. 2017. Openmvg: Open multiple view geometry. In *Reproducible Research in Pattern Recognition: First International Workshop, RRRP 2016, Cancún, Mexico, December 4, 2016, Revised Selected Papers I*, 60–74. Springer.
- Müller, T.; Evans, A.; Schied, C.; and Keller, A. 2022. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4): 1–15.
- Mur-Artal, R.; Montiel, J. M. M.; and Tardos, J. D. 2015. ORB-SLAM: a versatile and accurate monocular SLAM system. *IEEE transactions on robotics*, 31(5): 1147–1163.
- Niemeyer, M.; Barron, J. T.; Mildenhall, B.; Sajjadi, M. S.; Geiger, A.; and Radwan, N. 2022. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5480–5490.
- Oechsle, M.; Peng, S.; and Geiger, A. 2021. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5589–5599.
- Park, K.; Sinha, U.; Hedman, P.; Barron, J. T.; Bouaziz, S.; Goldman, D. B.; Martin-Brualla, R.; and Seitz, S. M. 2021. HyperNeRF: a higher-dimensional representation for topologically varying neural radiance fields. *ACM Transactions on Graphics (TOG)*, 40(6): 1–12.

- Pumarola, A.; Corona, E.; Pons-Moll, G.; and Moreno-Noguer, F. 2021. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10318–10327.
- Rosinol, A.; Leonard, J. J.; and Carlone, L. 2022. NeRF-SLAM: Real-Time Dense Monocular SLAM with Neural Radiance Fields. *arXiv preprint arXiv:2210.13641*.
- Sabour, S.; Vora, S.; Duckworth, D.; Krasin, I.; Fleet, D. J.; and Tagliasacchi, A. 2023. RobustNeRF: Ignoring Distractors with Robust Losses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 20626–20636.
- Schonberger, J. L.; and Frahm, J.-M. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4104–4113.
- Sitzmann, V.; Martel, J.; Bergman, A.; Lindell, D.; and Wetzstein, G. 2020. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33: 7462–7473.
- Sucar, E.; Liu, S.; Ortiz, J.; and Davison, A. J. 2021. iMAP: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 6229–6238.
- Sumikura, S.; Shibuya, M.; and Sakurada, K. 2019. Open-VSLAM: A versatile visual SLAM framework. In *Proceedings of the 27th ACM International Conference on Multimedia*, 2292–2295.
- Sun, C.; Sun, M.; and Chen, H.-T. 2022. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5459–5469.
- Tancik, M.; Casser, V.; Yan, X.; Pradhan, S.; Mildenhall, B.; Srinivasan, P. P.; Barron, J. T.; and Kretzschmar, H. 2022. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8248–8258.
- Tancik, M.; Weber, E.; Ng, E.; Li, R.; Yi, B.; Wang, T.; Kristoffersen, A.; Austin, J.; Salahi, K.; Ahuja, A.; et al. 2023. Nerfstudio: A modular framework for neural radiance field development. In *ACM SIGGRAPH 2023 Conference Proceedings*, 1–12.
- Teed, Z.; and Deng, J. 2020. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, 402–419. Springer.
- Turki, H.; Ramanan, D.; and Satyanarayanan, M. 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.
- Ueda, I.; Fukuhara, Y.; Kataoka, H.; Aizawa, H.; Shishido, H.; and Kitahara, I. 2022. Neural Density-Distance Fields. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*, 53–68. Springer.
- Verbin, D.; Hedman, P.; Mildenhall, B.; Zickler, T.; Barron, J. T.; and Srinivasan, P. P. 2022. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 5481–5490. IEEE.
- Wadhvani, K.; and Kojima, T. 2022. SqueezeNeRF: Further factorized FastNeRF for memory-efficient inference. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2717–2725.
- Wang, P.; Liu, L.; Liu, Y.; Theobalt, C.; Komura, T.; and Wang, W. 2021a. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*.
- Wang, P.; Liu, Y.; Chen, Z.; Liu, L.; Liu, Z.; Komura, T.; Theobalt, C.; and Wang, W. 2023. F2-NeRF: Fast Neural Radiance Field Training with Free Camera Trajectories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4150–4159.
- Wang, Y.; Skorokhodov, I.; and Wonka, P. 2022. HF-NeuS: Improved Surface Reconstruction Using High-Frequency Details. In *Advances in Neural Information Processing Systems*.
- Wang, Z.; Wu, S.; Xie, W.; Chen, M.; and Prisacariu, V. A. 2021b. NeRF-: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*.
- Warburg, F.; Weber, E.; Tancik, M.; Holynski, A.; and Kanazawa, A. 2023. Nerfbusters: Removing Ghostly Artifacts from Casually Captured NeRFs. *arXiv preprint arXiv:2304.10532*.
- Wu, C. 2013. Towards linear-time incremental structure from motion. In *2013 International Conference on 3D Vision-3DV 2013*, 127–134. IEEE.
- Xiao, Y.; Xue, N.; Wu, T.; and Xia, G.-S. 2023. Level- S^2 fm: Structure From Motion on Neural Level Set of Implicit Surfaces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 17205–17214.
- Yan, Q.; Wang, Q.; Zhao, K.; Li, B.; Chu, X.; and Deng, F. 2023. Rethinking Disparity: A Depth Range Free Multi-View Stereo Based on Disparity. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 3091–3099.
- Yao, Y.; Luo, Z.; Li, S.; Fang, T.; and Quan, L. 2018. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, 767–783.
- Yariv, L.; Gu, J.; Kasten, Y.; and Lipman, Y. 2021. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34: 4805–4815.
- Yen-Chen, L.; Florence, P.; Barron, J. T.; Rodriguez, A.; Isola, P.; and Lin, T.-Y. 2021. inerf: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 1323–1330. IEEE.
- Yi, K. M.; Trulls, E.; Lepetit, V.; and Fua, P. 2016. Lift: Learned invariant feature transform. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VI 14*, 467–483. Springer.

- Yu, A.; Li, R.; Tancik, M.; Li, H.; Ng, R.; and Kanazawa, A. 2021a. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 5752–5761.
- Yu, A.; Ye, V.; Tancik, M.; and Kanazawa, A. 2021b. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 4578–4587.
- Zhang, J.; Huang, J.; Cai, B.; Fu, H.; Gong, M.; Wang, C.; Wang, J.; Luo, H.; Jia, R.; Zhao, B.; et al. 2022a. Digging into Radiance Grid for Real-Time View Synthesis with Detail Preservation. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, 724–740. Springer.
- Zhang, J.; Zhan, F.; Wu, R.; Yu, Y.; Zhang, W.; Song, B.; Zhang, X.; and Lu, S. 2022b. Vmrf: View matching neural radiance fields. In *Proceedings of the 30th ACM International Conference on Multimedia*, 6579–6587.
- Zhang, K.; Riegler, G.; Snavely, N.; and Koltun, V. 2020. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*.
- Zhang, X.; Srinivasan, P. P.; Deng, B.; Debevec, P.; Freeman, W. T.; and Barron, J. T. 2021. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)*, 40(6): 1–18.
- Zhu, Z.; Peng, S.; Larsson, V.; Cui, Z.; Oswald, M. R.; Geiger, A.; and Pollefeys, M. 2023. Nicer-slam: Neural implicit scene encoding for rgb slam. *arXiv preprint arXiv:2302.03594*.
- Zhu, Z.; Peng, S.; Larsson, V.; Xu, W.; Bao, H.; Cui, Z.; Oswald, M. R.; and Pollefeys, M. 2022. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 12786–12796.