# CT-NeRF: Incremental Optimizing Neural Radiance Field and Poses with Complex Trajectory

Yunlong Ran[*]
Zhejiang University
Hangzhou, China

Yanxu Li[*]
Zhejiang University
Hangzhou, China

Qi Ye[†]
Zhejiang University
Hangzhou, China

Yuchi Huo
Zhejiang University
Hangzhou, China

Zechun Bai
Shandong University
Jinan, China

Jiahao sun
Zhejiang University
Hangzhou, China

Jiming chen
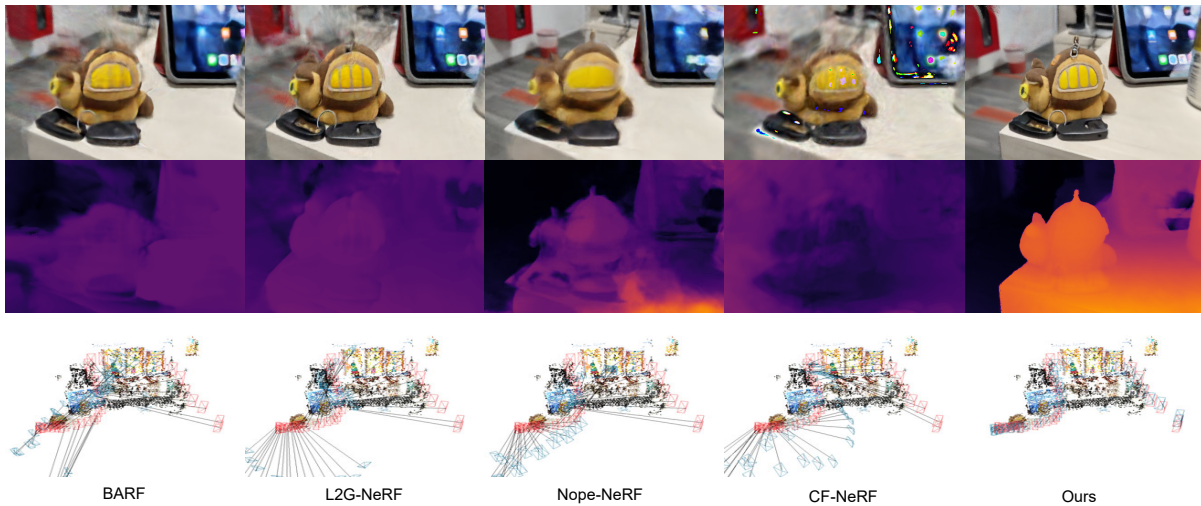Zhejiang University
Hangzhou, China

**Figure 1: Comparison on Free-dataset [36]. Top: novel view synthesis; Middle: depth maps; Bottom: the estimated trajectory errors (red rectangles for poses from COLMAP, blue rectangles for estimated, gray lines between them for errors). Our method enables more robust pose estimation, renders better novel views, and constructs better geometry than the state-of-the-arts.**

## ABSTRACT

Neural radiance field (NeRF) has achieved impressive results in high-quality 3D scene reconstruction. However, NeRF heavily relies on precise camera poses. While recent works like BARF have introduced camera pose optimization within NeRF, their applicability is limited to simple trajectory scenes. Existing methods struggle while tackling complex trajectories involving large rotations. To address this limitation, we propose CT-NeRF, an incremental reconstruction optimization pipeline using only RGB images without pose and depth input. In this pipeline, we first propose a local-global bundle adjustment under a pose graph connecting neighboring frames to enforce the consistency between poses to escape the local minima caused by only pose consistency with the scene structure. Further, we instantiate the consistency between poses as a reprojected geometric image distance constraint resulting from pixel-level correspondences between input image pairs. Through the incremental reconstruction, CT-NeRF enables the recovery of both camera poses and scene structure and is capable of handling scenes with complex trajectories. We evaluate the performance of CT-NeRF on two real-world datasets, NeRFBuster and Free-Dataset, which feature complex trajectories. Results show CT-NeRF outperforms existing methods in novel view synthesis and pose estimation accuracy.

## KEYWORDS

Pose estimation, Implicit representation, Structure from motion, SLAM

## 1 INTRODUCTION

Reconstructing high-fidelity, high-quality 3D scenes holds significant importance for the development of virtual reality / augmented reality, autonomous driving, and other domains. Recently, implicit

[*]Both authors contributed equally to the paper. Email: yunlong_ran@zju.edu.cn
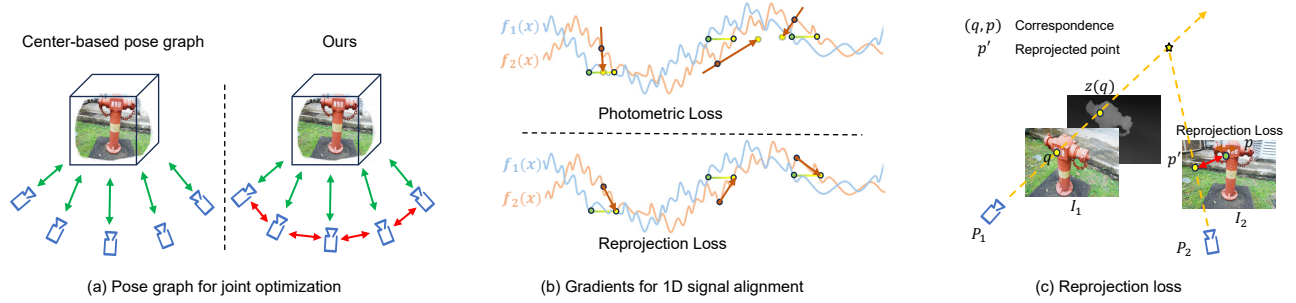[†]Corresponding author. Email: qi.ye@zju.edu.cn

**Figure 2: (a) Left: center-based pose graph to force pose consistent to the scene; Right: our pose graph to enable consistency between the camera poses in addition to the consistency to the scene. (b) The reprojection loss (bottom) provides an accurate gradient towards alignment, while the photometric loss (top) provides inconsistent gradients. (c) For a pair of correspondence $(q, p)$ between $I_1$ and $I_2$, $q$ is reprojected to the image plane of $I_2$ via the depth value of $q$. The reprojection loss is the geometric image distance between the reprojected point $p'$ and the ground truth corresponding point $p$.**

representations such as neural radiance fields (NeRF) [22] have achieved remarkable progress in reconstructing photo-realistic scenes given a sequence of RGB images and their corresponding camera poses. The camera poses for these high-fidelity reconstructions with implicit representations are primarily obtained through the structure from motion (SfM) methods, with the off-the-shelf tool COLMAP [26] being the most popular choice. These SfM methods estimate the camera poses through local registration between images and global bundle adjustment (BA) on both all camera poses and sparse 3D scene points. Therefore, accurate camera poses can only be acquired after all images are processed. Also, matching and registration of the methods are sensitive to image variations.

Recent works such as NeRFmm [37], SC-NeRF [17], BARF [20], GARF [8] and L2G-NeRF [6] tackle the dependency on the camera pose priors by treating the camera pose as learnable parameters and jointly optimizing poses and scenes offline. However, using images only to constrain the 3D space encounters many problems like wrong geometry, blurry textures, or floaters when very dense multiview images are not available; adding more freedoms for the camera poses to the optimization leads to worse results. Therefore, these methods often require initial camera parameters close to the ground truth poses in object-centered scenes with dense multiview observations, or small camera movements. Nope-NeRF [3] incorporates monocular depth to impose further constraints on adjacent images, enabling pose estimation for trajectories with relatively small camera motions and rotations while its initialization of all poses as identity matrices leads to local optima when facing complex trajectories.

On the other hand, following the classic SfM pipelines, CF-NeRF [41] adds images incrementally, initializes the pose for a newly added image with the pose for the previous one, and optimizes the poses (and the scene). With the incremental strategy, the method is capable of reconstructing the real-world scene under complex camera trajectories. However, it still suffers from large pose errors and inferior reconstruction quality as shown in Fig. 1, Fig. 4 and Fig. 5. The reason can be attributed to two aspects. First, the bundle adjustment constructs a center-based graph as shown in Fig. 2 (a) left, optimizing only the consistency between the camera poses

and the center implicit global scene while neglecting the consistency between the pose and the multiview images. When the global structure falls into local minima, the camera poses cannot be recovered; in turn, the structure cannot find a way to escape the local minima as the poses and structure are optimized jointly. Secondly, the method only uses the visual difference between the rendering images and raw images whereas BARF [20] observes that as natural images are typically complex signals, gradient-based registration with pixel value differences is susceptible to suboptimal solutions if poorly initialized, as shown in the top figure of Fig. 2 (b). The coarse-to-fine pose estimation proposed by BARF [20] mitigates this issue but requires good initialization or dense forward-facing images. In complex trajectories, the camera typically exhibits large motions, and views covering a region are much sparser than the forward-facing scenario.

To tackle the issues and enable accurate pose estimation and reconstruction for complex trajectories, we propose a novel incremental joint optimization method for implicit radiance fields and camera poses named CT-NeRF. For the first issue, as shown in Fig. 2 (a) right, we propose a joint incremental reconstruction and pose estimation pipeline with pose graphs connecting edges between camera poses upon the center-based pose graphs for a local-global bundle adjustment (BA). The graph forms many subgraphs and forces consistency between the camera poses, which helps to recover the poses when the scene and poses are consistent but the scene is actual in a local minimum during BA. For the second issue and also instantiating the pose consistency between the pose edges, we introduce a geometric image instance, *i.e.* the reprojected Euclidean distance between the correspondences of two input images. In addition to providing consistency constraints for pose edges, the reprojected distance benefits the pose and scene optimization in three aspects: 1) it provides direct direction to align the poses, whereas the pixel value differences do not necessarily correlate to the pose error: as shown in Fig. 2 (b) gradients based on the pixel value difference are not consistent while the gradients from the reprojected geometric image distance are; 2) the reprojected distance requires the depth of the scene to warp a pixel in an image to the other one as shown in Fig. 2 (c) and therefore, the gradient can help the convergence of the geometry of the scene directly; 3)

correspondence learning networks typically leverage large scale pair-wise image datasets and the reprojected loss based on the correspondences is robust to occlusion, lighting variation, textureless and large motions compared to raw image losses.

In summary, our main contributions are threefold:

- We design an incremental reconstruction pipeline for neural radiance fields using only RGB images under complex camera trajectories, without pose and depth input.
- We propose to construct pose graphs with in-between pose consistency edges for BA and instantiate the consistency as a reprojected geometric image distance constraint from the learned correspondences between input images for robust pose and scene optimization.
- We achieve significant improvements in pose estimation accuracy and reconstruction quality compared to state-of-the-art methods in complex trajectories.

## 2 RELATED WORK

**SFM and SLAM** In the field of computer vision, given a set of input images, SfM and SLAM aim to concurrently estimate camera poses and reconstruct the scene. The distinction lies in the fact that SLAM operates online, emphasizing runtime performance, while SfM does not require online operation but demands higher accuracy. SfM methods can be categorized into incremental [26, 27, 40], global [9, 18], and hierarchical [14] approaches: The incremental approach initializes with two images and progressively registers and reconstructs additional images one by one. The global approach registers and reconstructs all images simultaneously. The hierarchical approach first groups images, performs registration and reconstruction for each group, and then conducts a global optimization. SLAM methods are primarily divided into filter-based [1, 2, 5] and graph optimization-based [4, 12, 23] approaches. Filter-based methods mainly utilize state estimation strategies such as Kalman filtering and particle filtering to incrementally estimate the posterior distributions of camera poses and key point locations. Graph optimization-based methods abstract camera poses at different times as nodes and the observation constraints at different robot locations as edges connecting the nodes, then employ bundle adjustment (BA) algorithms for global optimization. Our proposed incremental pipeline is inspired by the incremental SfM and SLAM approaches.

**NeRF-based SFM and SLAM** Implicit neural representations have gained prominence since 2019 [25]. Compared to traditional explicit representations that store geometric information in a relatively fixed and simple manner, implicit neural representations can better handle complex topological structures and geometric details. The classic algorithm for implicit neural representations, Vanilla NeRF [22] (Neural Radiance Fields), is based on the theory of volume rendering and utilizes a Multi-Layer Perceptron (MLP) to learn the implicit neural representation of a static scene, achieving high-quality novel view synthesis. Consequently, some researchers have considered combining SfM and SLAM with NeRF, not only to reduce NeRF's dependence on the accuracy of input image poses but also to enhance the scene representation capability of SfM and SLAM. BARF [20] was the first to integrate the core bundle adjustment (BA) algorithm from SfM with NeRF and adopted a coarse-to-fine reconstruction strategy, progressively aligning camera poses during the reconstruction process. To address the issue of camera poses being prone to local optima in BARF, L2G-NeRF [6] proposed a Local-to-Global alignment strategy, allowing camera poses to converge more easily to the global optimum. NoPe-NeRF [3] introduced monocular depth information and key point matching information, respectively, to constrain the relative camera pose relationships, ensuring global consistency of camera poses. LocalRF [21] proposed a progressive strategy based on video sequences to gradually optimize local regions. CF-NeRF [41] employed an incremental learning approach to enable reconstruction under complex trajectories. LU-NeRF [7] introduced a Local-to-Global pose estimation strategy, enabling pose estimation and scene reconstruction from datasets with completely unknown camera poses. As for NeRF-based SLAM methods, iMAP [29] adopted two threads: Tracking and Mapping. The Tracking thread optimizes the camera pose using the current model and performs key frame selection, while the Mapping thread jointly optimizes the poses of the keyframes and the model. iMAP uses a single MLP to represent the entire scene, limiting its scalability. Nice-SLAM [42] improved upon iMAP by combining Hierarchical Feature Grids and MLP as the scene representation, enabling the application to large-scale scenes. Co-SLAM [35] and e-SLAM [19] introduced multi-resolution hash encoding and tri-plane representations, respectively, to improve system frame rate and scene representation capability, building upon Nice-SLAM. Regrettably, current NeRF-based SLAM methods necessitate dense image sequences, while NeRF-based SfM techniques face difficulties in accommodating complex camera trajectories. To tackle these limitations, we introduce CT-NeRF, an incremental optimization framework that leverages additional correspondence constraints. CT-NeRF employs an incremental optimization process, iteratively refining the reconstruction as new images are integrated.

**Correspondence in pose-estimate** Local feature matching plays a crucial role in SfM and SLAM. The traditional feature matching pipeline consists of three main steps: feature detection, feature descriptor computation, and feature matching. Through feature detectors, the search space for matching can be effectively reduced, and the generated sparse matches are often sufficient to handle most tasks. However, in low-texture regions or repetitive patterns, these methods often fail due to the inability to detect sufficient feature points. With the flourishing development of deep learning, some researchers have started to leverage data-driven dense feature matching to enhance the accuracy and robustness of SfM and SLAM. For instance, Droid-SLAM [32] utilizes the dense optical flow learned by RAFT [31] for feature matching, achieving higher accuracy and robustness compared to traditional SLAM, and rarely failing in experimental scenarios. Detector-free SfM [16] leverages the matching strategy of Loftr [30] without feature detectors, exhibiting significant advantages in low-texture regions and winning multiple competitions. Droid-SLAM and Detector-free SfM focus on pose estimation and reconstructing the scene with sparse points. On the other hand, SPARF [33] utilizes the correspondences from DKM matching for implicit neural reconstruction under noisy poses with several views. Different from these works, our method focuses on an incremental pose and implicit scene joint optimization pipeline with complex trajectories.

# 3 METHOD

We first present the formulation of incremental scene reconstruction and pose estimation: given a set of sequential images $\mathcal{I} = \{I_1, I_2, ..., I_n\}$, where $n$ represents the $n^{th}$ frame captured in a camera trajectory, we aim to jointly optimize the poses $\mathcal{P} = \{P_1, P_2, ..., P_n\}$ for the images and a neural radiance field model $\Theta$ representing the 3D scene captured by the images by adding one image at a time sequentially. To achieve the goal, we design an incremental reconstruction pipeline for the neural radiance field without pose priors. Our pipeline consists of five parts as shown in Fig. 3. The scene is initialized using a small set of images. Subsequently, for each input image, tracking is applied to estimate the rough camera pose for a new image. A window optimization is followed to refine the poses of images within the window and also reconstruct the local structural components. To further incorporate the consistency of all visited camera poses, global optimization (bundle adjustment) is performed on all images to optimize the global camera poses and the overall scene structure. Tracking, window, and global optimization repeat until all images are added. After all images are added, post optimization iteratively refines the entire scene and all camera poses until convergence.

## 3.1 Preliminary: NeRF with Pose Optimization

We define the camera projection function $\pi$ that projects a space point $\mathbf{x} \in \mathcal{R}^3$ to a pixel $p \in \mathcal{R}^2$ as

$$p = \pi(\mathbf{x}, P, K), \tag{1}$$

where camera pose $P$ is the camera-to-world transformation of image $I$ and $K$ is the intrinsic (we assume all images share the same intrinsic in a trajectory). $P = [R, t]$, where $R \in SO(3)$ represents rotation and $t \in \mathcal{R}^3$ translation. The homogenization operations are omitted for clarity. The backprojection function $\pi^{-1}$ projects the pixel coordinate location $p$ into a space point with depth $z$

$$x = \pi^{-1}(p, P, K, z). \tag{2}$$

NeRF maps a 3D location $\mathbf{x} \in \mathbb{R}^3$ and a view direction $\mathbf{d} \in \mathbb{R}^3$ to a radiance color $\mathbf{c} \in \mathbb{R}^3$ and volume density $\sigma \in \mathbb{R}$ with an MLP parameterized by $\Theta$. It optimizes the model $\Theta$ and camera poses $\mathcal{P}$ by minimizing photometric loss between rendered images $\hat{\mathcal{I}}$ and input images $\mathcal{I}$

$$\Theta^*, P^* = \arg\min_{\Theta, P} \mathcal{L}_{photo}(\hat{\mathcal{I}}, \mathcal{P} \mid \mathcal{I}), \tag{3}$$

where $\mathcal{L}_{photo} = \frac{1}{n}\sum_1^n ||\hat{I}_i(P_i, \Theta) - I_i||_2^2$ and $\hat{I}_i$ can be obtained through volume rendering. For each pixel $q$ in image $\hat{I}_i$, its color is rendered by aggregating predicted colors $\mathbf{c}$ and densities $\sigma$ alone the ray $r = \mathbf{o} + \mathbf{d}s$ where $\mathbf{o}$ and $\mathbf{d}$ can be obtained by $\mathbf{o}, \mathbf{d} = \pi^{-1}(p, P, K)$ and $s \in (s_{near}, s_{far}]$ represents sample distance.

$$\hat{I}_i(q) = \int_{s_{near}}^{s_{far}} T(s)\sigma(r(s))\mathbf{c}(r(s), \mathbf{d})ds, \tag{4}$$

where $T(s) = \exp(-\int_{s_{near}}^{s} \sigma(r(h))dh)$ indicates how much light is transmitted on ray up to $s$. In the same way, depth map $\hat{D}_i$ can be rendered.

$$\hat{D}_i(q) = \int_{s_{near}}^{s_{far}} T(s)\sigma(r(s))sds. \tag{5}$$

## 3.2 Reprojected Geometric Image Distance

As aforementioned, we aim to incorporate the geometric distance constraint from correspondences between images to jointly optimize camera poses and 3D scene model parameters under complex trajectories. For correspondence generation, many existing correspondence learning networks can be exploited. In this work we choose DKM [11] as it is the current state-of-the-art work in dense correspondence matching, achieving high matching accuracy and strong robustness. Given a pair of adjacent images $I_1, I_2$ in the input trajectory $\mathcal{I}$, a pre-trained DKM model can generate full-resolution pixel-level correspondences between them while predicting the confidence of each pixel match $\alpha$. We use a set $M$ to represent the output correspondences for two input images

$$M = \{m = (q, p, \alpha) \mid q \in I_1, p \in I_2, \alpha \in (0, 1]\}. \tag{6}$$

According to the multiview geometry theory [15], four pairs of correspondences can solve the relative pose between the images and using the triangulation technique, the 3D scene points for the pairs can be acquired. Though the correspondences only produce sparse 3D points and our representation is implicit, the theory provides important information for our problem: 1) correspondences provide a way to solve the pose estimation problem without knowing a 3D scene; 2) the pose estimation requires only sparse correspondences for the pose estimation; 3) the correspondences embed the 3D information.

We define our geometric image distance for the pose estimation between two images under implicit scene representation as follows. As correspondence in practice often contains noise, we randomly sample $N_m$ correspondences from the correspondence set $M$ with confidence above a threshold $t_m$, to serve as a correspondences set $M_{sparse}$ for pose estimation,

$$M_{sparse} = \{m_1, m_2, ..., m_{N_m}\} \sim \{m \mid m \in M, \alpha > t_m\}. \tag{7}$$

In classic SfM pipelines, algorithms like RANSAC [13] are exploited to remove outliers and get robust poses and 3D points. Though differentiable RANSAC [39] may be deployed in our implicit joint optimization, we choose a simpler strategy: 1) weighting the reprojection error according to the confidence of correspondence estimation; 2) randomly sampling a small set of samples in each iteration. In every training iteration for the scene model $\Theta$, $N_s$ pairs of correspondences are randomly fetched from $M_{sparse}$. Sampled correspondence sets with the confidence from multiple different iterations to optimize the pose and/or scene parameters serve as the similar purpose of RANSAC.

For each correspondence $(q, p, \alpha)$, its depth value $z(q; \Theta)$ can be rendered by Eq. (5), and $q$ can be backprojected into a 3D point and reprojected to $I_2$. Then we have the geometric image distance

$$\mathcal{L}_{rp}(I_1, I_2) = \frac{1}{N_s}\sum_1^{N_s} \alpha|\pi((\pi^{-1}(q, P_1, K, z(q; \Theta))), P_2, K) - p|, \tag{8}$$

which is reprojection error. $P_1, P_2$ are the poses to be estimated for images $I_1, I_2$ respectively. The gradient with respect the camera poses $P_1, P_2$, rendered depths and further the model $\Theta$ can be obtained from Eq. (8), indicating the gradients can help the depth and scene reconstruction in addition to pose estimation.

The reprojection error is referred as *Gold Standard* [15]. The error is a quadratic convex function of pose parameters, pointing the right direction for pose optimization without local minima. However, the pixel value difference does not necessarily correlate to the
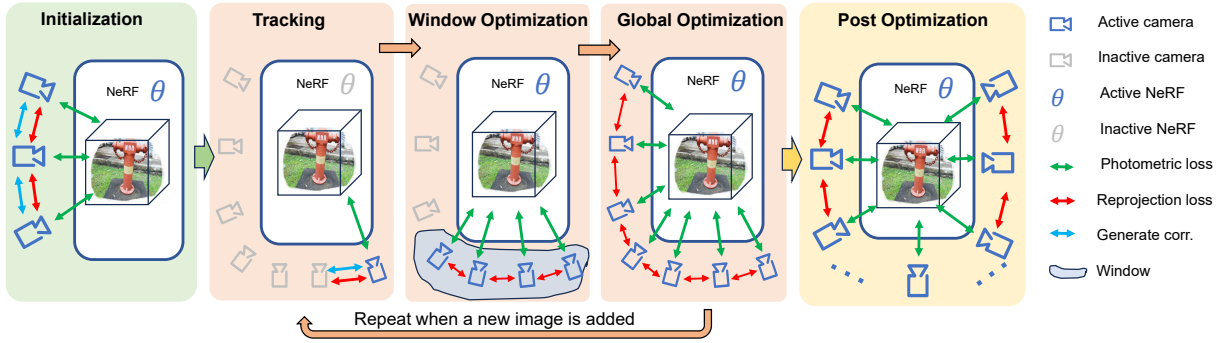
**Figure 3: Our incremental optimization pipeline for neural radiance fields and pose estimation.**

correct direction for the pose optimization as shown in BARF [20]. Compared with the pixel value difference, the reprojection error also provides more robust gradients for the scene estimation as the pixel value difference is susceptible to the lighting, occlusion, sparse views, etc, while the correspondence learning network for the reprojection error learns these factors during training.

### 3.3 Tracking

When a new frame $I_i$ is added to the training process, tracking provides a rough estimate of the camera pose, which becomes particularly crucial when there exist violent changes in camera motion. The pose of the new frame is initialized based on the previous frame $P_i = P_{i-1}$. Tracking performs pose estimation through the reprojection error with adjacent frames and photometric loss with the scene. The loss function can be formulated as:

$$\mathcal{L}(\mathcal{E}) = \mathcal{L}_{nrp}(\mathcal{E}) + \mathcal{L}_{photo}(\mathcal{E}), \tag{9}$$

$$\mathcal{L}_{nrp}(\mathcal{E}) = \frac{1}{N_e * 2 - 2} \sum_1^{N_e} (\mathcal{L}_{rp}(I_i, I_{i-1}) + \mathcal{L}_{rp}(I_i, I_{i+1})), \tag{10}$$

where $\mathcal{L}_{nrp}$ is reprojection loss for paris of neighboring images and $\mathcal{E} = \{I_1, I_2, ..., I_{N_e}\}$ is an optimization frame set with $N_e$ frames.

For the tracking, $\mathcal{E}_{tracking} = \{I_{i-1}, I_i\}$ consists of a new frame and the preceding frame ($N_e = 2$) and we optimize only the pose of the newly added frame $P_i$, keeping other optimizable pose parameters and the network parameters $\Theta$ fixed.

**Initialization** is crucial as it affects the subsequent tracking performance and the final pose estimation and reconstruction quality. In our approach, we select the first $N_{init}$ images in the sequence as the initialization images set $\mathcal{E}_{init} = \{I_1, I_2, ..., I_{N_{init}}\}$. The initialization is achieved by minimizing Eq. (9). Due to the difficulty in optimizing the rotation $R$ parameters in the initial stages, we fix the rotation parameters and do not optimize them during the initialization.

### 3.4 Joint Optimization

Using the center-based graph and photo loss in equation Eq. (3) for the joint optimization or bundle adjustment is susceptible to local minima when only RGB inputs are available. It only constrains all the poses in consistency with the scene while reconstructing the scene with implicit neural radiance fields from only

RGB images (especially sparse images) are prone to converge to wrong geometry, which is demonstrated in many previous works [10, 24, 28, 34]. If the scene is stuck in a local minimum, the bundle adjustment can not get the pose right as long as the poses conform with the twisted scene. As shown in Table 1, though the poses from BARF[20], L2G-NeRF[6], Nope-NeRF[3] are fairly deviated, visual metrics like PSNR for the scene maintains high, causing the pose not able to escape the local minima. In contrast, we construct pose graphs with constraints between pose edges for the joint optimization. The subgraphs formed between the poses and between the pose and the scene (shown in Fig. 2 (a) right) enable consistency of all the in-between poses in the image set. The reprojection error forces the consistency of the adjacent poses with the correspondences of input images. Further, we design a combination of local window and global BA strategy to balance the integration of new information and consistency with existing estimation.

Our joint optimization uses the same equation Eq. (10) in tracking. However, the scene model $\Theta$ is learnable, and different frame sets $\mathcal{E}$ for the optimization are maintained.

**Window optimization** To joint optimize the scene and the camera pose for a newly added image, window optimization selects the most recent $N_{window}$ frames as the optimization set $\mathcal{E}_{window} = \{I_{i-N_{window}+1}, ..., I_{i-1}, I_i\}$. The window optimization process fixes the camera poses outside the window and optimizes the camera poses within the window and network $\Theta$ using Eq. (9) by performing local bundle adjustment. Different from the existing implicit SfM and SLAM methods which consist of tracking for the newly added image and global bundle adjustment for all the input images, the extra window optimization improves the pose estimation accuracy by enhancing consistency with the near previous frames and makes the network learn faster from the information of the new frame by leaving older frames out.

**Global optimization** Relying solely on tracking and window optimization can lead to cumulative errors and even failure in pose estimation. To address this issue, global optimization incorporates all frames currently added to the training process into the optimization set $\mathcal{E}_{global} = \{I_1, I_2, ..., I_i\}$. By applying Eq. (9) to optimize all frames simultaneously, global optimization significantly enhances the robustness and accuracy of pose estimation.

**Post optimization** Before all frames are added, the learning rate of the network $lr_\Theta$ and poses $lr_{\hat{\mathcal{P}}}$ are fixed. The positional encoding

**Figure 4: Qualitative Comparison on Free-Dataset [36]. Rendered views and depths (top left corner of each image)**

control parameter $\alpha_{pe}$ is also fixed. After the whole frames are added, the post optimization process gradually reduces the learning rate and increases the frequency of positional encoding to iteratively refine the entire scene and camera poses through Eq. (9), ultimately obtaining the final results.

### 3.5 Training procedure

**Training Pipeline** The pipeline is initialized with $N_{init}$ frames. The stage is optimized for $\beta_{init}$ iterations. Afterwards, for each subsequent frame added, tracking is performed for $\beta_{tracking}$ iterations, window optimization for $\beta_{window}$ iterations, and global optimization for $\beta_{global}$ iterations. These three stages repeat with a new frame added and continue until all frames have been added. Finally, the post optimization stage consisting of $\beta_{post}$ iterations is conducted to further refine the reconstruction.

**Positional Encoding** The coarse-to-fine positional encoding plays an important role in accurate pose estimation [20], as excessively high frequencies can hinder this process. To address this, we employ the BARF [20] positional encoding frequency control method. Specifically, before the post-optimization stage, we ensure a low-frequency setting for the positional encoding control parameter $\alpha_{pe}$. During post optimization, we keep the same coarse-to-fine strategy as the BARF.

## 4 EXPERIMENTS

### 4.1 Experiment Settings

**Dataset** We evaluated our method on the challenging datasets with complex trajectories, **NeRFBuster** [38] and **Free-Dataset** [36]. **NeRFBuster** consists of a total of 12 scenes, with most trajectories revolving around a central object. We employ sequences selected by CF-NERF [41], with approximately 50 images per scene. We chose every 8th image from each sequence for novel view synthesis as the test set. All images are downsampled to a resolution of $480 \times 270$. Ground truth poses are estimated using COLMAP, as provided by

CF-NeRF. **Free-Dataset** comprises 7 scenes with arbitrary trajectories, predominantly in outdoor environments characterized by highly dynamic camera motions. We select 50 images per scene in sequential order, and every 8th image is designated as the test set. The images are downsampled to a resolution of $312 \times 487$, and the ground truth poses are obtained through COLMAP [26].

**Implementation Details** Our approach is implemented based on the BARF [20] framework. The majority of the hyperparameters in our network model align with the BARF Real-World Scenes setting, including the network learning rate $lr_\Theta$ decay from $1 \times 10^{-3}$ to $1 \times 10^{-4}$, pose learning rate $lr_\mathcal{P}$ decay from $3 \times 10^{-3}$ to $1 \times 10^{-5}$, inverse sampling of 128 points along each ray with an inverse range of $[1, 0)$, a batch size of 1024, and linearly adjust $\alpha$ for post optimization phase from iteration 20K to 100K. We randomly select $N_m = 10000$ correspondences with confidence scores higher than $t_m = 0.2$ from dense correspondences as sparse correspondences. During each iteration, $N_s = \frac{256}{len(\mathcal{E})}$ correspondences are randomly chosen from this set for reprojection loss. For NeRFBuster scenes, we set $N_{init} = 3$, $N_{window} = 4$, $\beta_{init} = 2000$, $\beta_{tracking} = 100$, $\beta_{window} = 200$, $\beta_{global} = 500$, and $\beta_{post} = 200K$. As the frames in in Free-Dataset exhibit larger camera motions and smaller overlap, the network requires more iterations to estimate poses accurately and achieve convergence. For Free-Dataset scenes, we set $N_{init} = 3$, $N_{window} = 4$, $\beta_{init} = 4000$, $\beta_{tracking} = 200$, $\beta_{window} = 400$, $\beta_{global} = 900$, and $\beta_{post} = 200K$.

**Metrics** We primarily evaluate our method by assessing the quality of novel view synthesis and the accuracy of pose estimation. As the variables for the scene and the cameras are up to a 3D similarity transformation, existing work [3, 6, 20] aligns the optimized poses to the ground truth using Sim(3) with Procrusters analysis on the camera locations for pose error computation and test pose initialization (termed as Sim(3) below) and then runs an additional test-time optimization on the trained model to reduce the pose error that may influence the view synthesis quality. Since all existing

**Table 1: Evaluations of the pose accuracy (top 2 rows) and the novel view quality (bottom 3 rows) on NeRFBuster [38]. $\Delta T$ is the transition error in ground truth scale and $\Delta R$ is rotation error in degree.**

| Metrics | Method | aloe | art | car | century | garbage | flowers | picnic | pikachu | pipe | plant | roses | table | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\Delta R\downarrow$ | BARF [20] | 128.599 | 45.739 | 97.148 | 114.261 | 103.259 | 79.702 | 81.418 | 112.996 | 166.701 | 140.270 | 125.974 | 139.675 | 111.312 |
| | L2G-NeRF [6] | 117.475 | 28.247 | 161.862 | 59.730 | 90.889 | 88.750 | 99.076 | 123.543 | 106.838 | 71.551 | 139.057 | 144.848 | 102.656 |
| | Nope-NeRF [3] | 101.589 | 32.345 | 113.063 | 150.253 | 149.459 | 161.859 | 148.710 | 158.059 | 99.836 | 138.816 | 150.050 | 114.783 | 126.569 |
| | CF-NeRF [41] | 6.703 | 76.306 | 29.079 | 11.013 | 74.163 | 10.672 | 109.868 | 13.243 | 122.345 | 18.664 | 3.903 | 3.835 | 39.983 |
| | Ours | **3.163** | **3.151** | **0.701** | **2.343** | **0.902** | **0.481** | **1.938** | **7.708** | **2.302** | **6.302** | **0.570** | **1.154** | **2.560** |
| $\Delta T\downarrow$ | BARF | 6.039 | 4.040 | 5.043 | 5.434 | 4.663 | 4.693 | 3.007 | 3.772 | 3.763 | 5.865 | 4.952 | 4.076 | 4.612 |
| | L2G-NeRF | 4.986 | 4.402 | 4.764 | 5.895 | 4.926 | 4.214 | 6.451 | 5.592 | 2.764 | 5.055 | 4.199 | | 4.795 |
| | Nope-NeRF | 5.151 | 5.302 | 5.401 | 3.202 | 5.571 | 4.742 | 4.819 | 3.757 | 4.983 | 5.896 | 5.399 | 5.817 | 5.004 |
| | CF-NeRF | 0.637 | 1.549 | 1.621 | 0.497 | 0.548 | 0.745 | 1.285 | 0.879 | 5.757 | 0.685 | 0.182 | 0.274 | 1.222 |
| | Ours | **0.168** | **0.030** | **0.035** | **0.134** | **0.039** | **0.039** | **0.106** | **0.548** | **0.164** | **0.225** | **0.038** | **0.045** | **0.131** |
| $PSNR\uparrow$ | BARF | 23.56 | 20.55 | 23.69 | 18.73 | 19.92 | 23.14 | 22.91 | 31.58 | **23.43** | 29.38 | 21.87 | 25.88 | 23.72 |
| | L2G-NeRF | 23.47 | 22.58 | 23.98 | 19.32 | 20.36 | 24.52 | 22.18 | **33.66** | 21.99 | **29.63** | 21.74 | 25.60 | 24.09 |
| | Nope-NeRF | 22.42 | 21.53 | 22.62 | 19.55 | 20.59 | 21.91 | 22.97 | 27.39 | 21.33 | 25.69 | 19.91 | 26.83 | 22.73 |
| | CF-NeRF | 23.32 | 23.50 | 22.04 | 21.35 | 21.30 | 23.91 | **23.31** | 31.51 | 22.24 | 25.89 | 23.42 | 26.71 | 24.04 |
| | Ours neighbor | 23.17 | 25.76 | 24.90 | 21.64 | 21.62 | 26.14 | 23.04 | 30.50 | 23.02 | 27.19 | 22.14 | 30.85 | 25.00 |
| | Ours sim(3) | **24.36** | **26.73** | **27.41** | **22.56** | **22.69** | **27.37** | 23.04 | 22.91 | 23.13 | 22.64 | **29.63** | **32.73** | **25.43** |
| $SSIM\uparrow$ | BARF | 0.59 | 0.68 | 0.72 | 0.52 | 0.54 | 0.72 | **0.54** | 0.92 | **0.62** | 0.85 | 0.69 | 0.84 | 0.69 |
| | L2G-NeRF | 0.58 | 0.74 | 0.74 | 0.56 | 0.53 | 0.76 | 0.50 | **0.94** | 0.55 | **0.87** | 0.69 | 0.83 | 0.69 |
| | Nope-NeRF | 0.52 | 0.73 | 0.69 | 0.57 | 0.56 | 0.67 | 0.53 | 0.88 | 0.53 | 0.80 | 0.64 | 0.86 | 0.67 |
| | CF-NeRF | 0.56 | 0.74 | 0.66 | 0.63 | 0.57 | 0.72 | 0.53 | 0.93 | 0.54 | 0.80 | 0.72 | 0.85 | 0.69 |
| | Ours neighbor | 0.56 | 0.82 | 0.74 | 0.63 | 0.58 | 0.78 | 0.52 | 0.91 | 0.58 | 0.83 | 0.71 | 0.89 | 0.71 |
| | Ours sim(3) | **0.61** | **0.83** | **0.79** | **0.65** | **0.63** | **0.81** | 0.52 | 0.76 | 0.59 | 0.71 | **0.88** | **0.91** | **0.72** |
| $LPIPS\downarrow$ | BARF | 0.36 | 0.20 | 0.29 | 0.40 | 0.38 | 0.27 | **0.44** | **0.09** | 0.33 | 0.14 | 0.24 | 0.23 | **0.28** |
| | L2G-NeRF | 0.37 | 0.18 | 0.28 | **0.35** | 0.42 | **0.14** | 0.49 | 0.10 | 0.39 | **0.14** | 0.24 | 0.24 | 0.29 |
| | Nope-NeRF | 0.51 | 0.33 | 0.43 | 0.48 | 0.50 | 0.47 | 0.54 | 0.30 | 0.52 | 0.34 | 0.48 | 0.35 | 0.44 |
| | CF-NeRF | 0.43 | 0.24 | 0.41 | 0.37 | 0.46 | 0.33 | 0.55 | 0.11 | 0.50 | 0.20 | 0.21 | 0.23 | 0.34 |
| | Ours neighbor | 0.36 | **0.14** | 0.26 | 0.50 | 0.44 | 0.22 | 0.49 | 0.11 | 0.41 | 0.22 | 0.14 | **0.17** | 0.29 |
| | Ours sim(3) | **0.35** | **0.14** | **0.25** | 0.50 | 0.43 | 0.22 | 0.49 | 0.29 | 0.40 | 0.30 | **0.09** | **0.17** | 0.30 |

methods compared below except ours struggle to obtain reasonable initial test poses through Sim(3) and fail to perform the test-time optimization under the complex trajectories, we adopt the approach of Nope-NeRF [3] for these methods, *i.e.* initializing a test image pose with the estimated pose of the training frame that is closest to it (termed as neighbor below). For our method, we provide results using both initialization methods. We report PSNR, SSIM, and LPIPS for the view synthesis, and rotation and translation errors for the pose estimation.

## 4.2 Comparison with Pose-Unknown Methods

We compare with the state-of-the-art methods for joint optimization of scenes and poses from RGB images, *i.e.* BARF [20], L2G-NeRF [6], Nope-NeRF [3], and CF-NeRF [41].

**Results on the Object Centered Dataset** Table 1 presents the pose evaluation results on the **NeRFBuster** dataset. BARF, L2G-NeRF, and Nope-NeRF initialize all poses as identity matrices and then perform bundle adjustment to jointly optimize poses and the scene. With pose initialization far from the actual poses and sparse views not able to effectively constrain the scenes, these methods frequently fail to recover the camera poses and geometry. CF-NeRF manages to estimate poses but still suffers from larger errors as CF-NeRF only constrains poses through photo loss.

Our method achieves significantly smaller errors compared to these methods.

Despite the significant pose errors and poor structural quality of methods such as BARF, they still achieve surprisingly "good" results on PSNR, SSIM, and LPIPS in Table 1, almost on-par with our results on the view synthesis. We attribute this to overfitting both in the training stage and test-time optimization. In Fig. 5, we visualize the trajectory of the NeRFbuster **garbage** scene and select two

frames with an abrupt change in the estimated camera trajectory. We then render novel views by interpolating between these poses. As shown in Fig. 5 (b), the rendering results are unreasonable, while CF-NeRF and our method can render smooth view transitions (In the supplementary video, we further show BARF, L2G-NeRF, Nope-NeRF, and CF-NeRF renders view inconsistent effects). The other evidence is that these methods struggle to reconstruct the geometry as shown Fig. 6 and Fig. 4. During test-time optimization, as the estimated trajectories of these methods diverge far from the ground truth and they fail to perform test-time optimization using Sim(3), the poses of the neighboring frames are used to initialize the test frames. This initialization causes the pose of a test frame to converge to a "pseudo " pose close to the estimated pose of the closest neighboring train frame. Then the network for the scene further overfits to the "pseudo" ground truth test pose and image ($P, I$ for example) pairs and renders an image $I'$ using $P$ to calculate visual metrics with $I$.

**Results on the Free Trajectory Dataset** We also conduct our experiments on the **Free-Dataset**, which consists of more challenging scenarios with arbitrary trajectory variations and reduced frame overlap (please refer to the supplementary material for visualization of the sequences). In Table 2, we report the results on the Free-Dataset, which demonstrates that our method exhibits more significant advantages under more challenging scenes. Fig. 4 shows that in addition to the superior quality of the novel view synthesis, our method can produce depths of good quality while most existing methods fail to.

## 4.3 Ablation Study

In this subsection, we conduct ablation studies to investigate the impact of various components in our method. We ablate **projection**
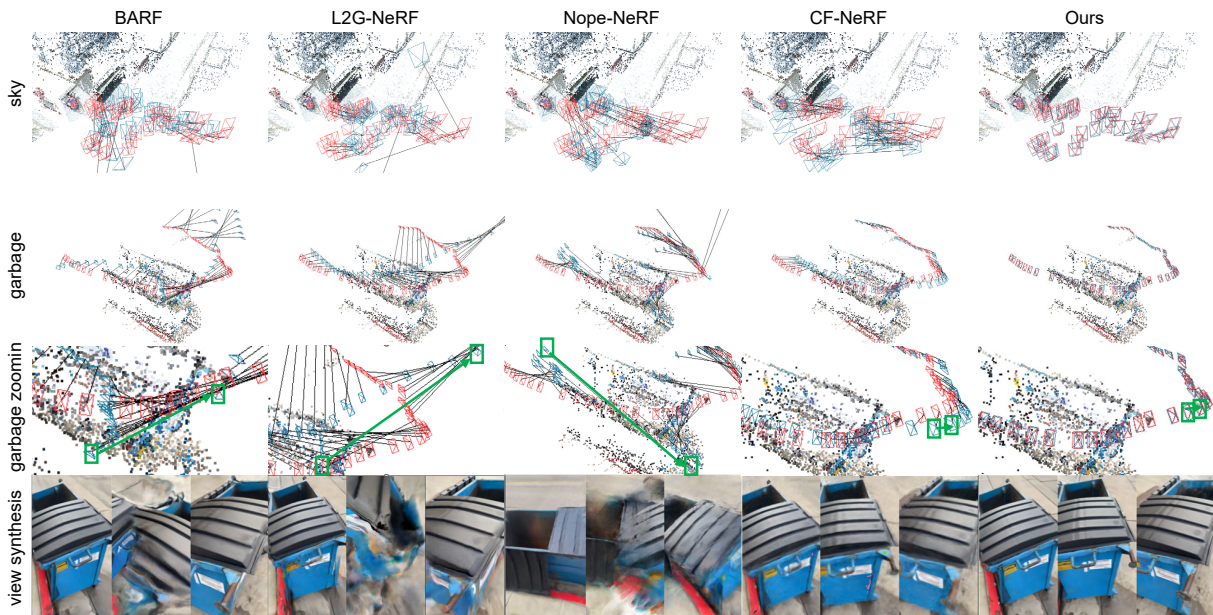
**Figure 5: Trajectory comparison.** We visualize camera poses of both estimated (blue) and COLMAP (red). Sparse 3D points for the scenes are from COLMAP. While there are abrupt changes in the trajectories of BARF, L2G-NeRF, and Nope-NeRF, the changes are steady along the trajectories of CF-NeRF and ours. The bottom row shows rendered interframes between two frames of abrupt changes denoted by green rectangles.
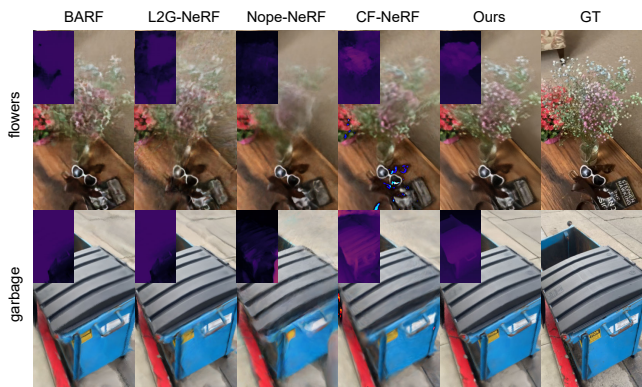


**Figure 6: Qualitative Comparison on NeRFbuster [38].** Rendered views and depths (top left corner of each image).

**Table 2: Evaluations of the pose accuracy and the novel view quality on Free-Dataset [36].** $\Delta T$ is the transition error in ground truth scale and $\Delta R$ is rotation error in degree.

| Method | $\Delta R\downarrow$ | $\Delta T\downarrow$ | $PSNR\uparrow$ | $SSIM\uparrow$ | $LPIPS\downarrow$ |
|---|---|---|---|---|---|
| BARF | 61.098 | 3.498 | 19.56 | 0.52 | 0.45 |
| L2G-NeRF | 110.303 | 6.587 | 19.95 | 0.54 | 0.45 |
| Nope-NeRF | 144.202 | 4.693 | 18.67 | 0.51 | 0.66 |
| CF-NeRF | 55.329 | 2.385 | 18.30 | 0.42 | 0.72 |
| Ours neighbor | **2.805** | **0.161** | 18.69 | 0.49 | 0.49 |
| Ours sim(3) | **2.805** | **0.161** | 22.46 | 0.59 | **0.43** |

**loss**, **tracking**, **window optimization**, and **global optimization** components individually. Table 3 shows that removing tracking, window, or global optimization leads to performance degradation, but the method remains functional. However, removing reprojection loss leads to dramatic pose errors. We refer readers to supplementary material for more results for the ablation.

**Table 3: Ablation study on reprojection loss, tracking, window optimization, and global optimization.**

| Method | $\Delta R\downarrow$ | $\Delta T\downarrow$ | $PSNR\uparrow$ | $SSIM\uparrow$ | $LPIPS\downarrow$ |
|---|---|---|---|---|---|
| Ours w/o reproj. loss | 56.040 | 1.904 | 16.18 | 0.44 | 0.63 |
| Ours w/o tracking | 5.302 | 0.280 | 23.66 | 0.67 | 0.35 |
| Ours w/o window opt. | 3.562 | 0.182 | 24.57 | 0.70 | 0.33 |
| Ours w/o global opt. | 6.189 | 0.234 | 22.66 | 0.64 | 0.40 |
| Ours | **2.560** | **0.131** | **25.43** | **0.72** | **0.30** |

## 5 CONCLUSION

We present CT-NeRF, a method capable of recovering poses and reconstructing scenes from image sequences captured along complex trajectories. We first introduce correspondence and reprojected geometric image distance to impose extra constraints on the optimization graph, enabling robust and accurate pose estimation and scene structure reconstruction. Subsequently, we detail our incremental learning process for pose recovery, including initialization, tracking, window optimization, and global optimization. Through comparative and ablation experiments, we demonstrate the superiority of our method and the necessity of its individual components. Although our method enables joint pose estimation

and reconstruction under complex camera trajectories, we only explore simple pose graphs. More sophisticated graph optimization is required for very long trajectories. Also, evaluation datasets, protocols, and metrics are required for complex camera trajectories as discussed in the paper, the current visual metrics can not fully reflect the reconstruction quality.

# REFERENCES

[1] Yassin Abdelrasoul, Abu Bakar Sayuti HM Saman, and Patrick Sebastian. 2016. A quantitative study of tuning ROS gmapping parameters and their effect on performing indoor 2D SLAM. In *2016 2nd IEEE international symposium on robotics and manufacturing automation (ROMA)*. IEEE, 1–6.

[2] Tim Bailey, Juan Nieto, Jose Guivant, Michael Stevens, and Eduardo Nebot. 2006. Consistency of the EKF-SLAM algorithm. In *2006 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 3562–3568.

[3] Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. 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.

[4] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. 2021. Orb-slam3: An accurate open-source library for visual, visual–inertial, and multimap slam. *IEEE Transactions on Robotics* 37, 6 (2021), 1874–1890.

[5] José A Castellanos, Ruben Martinez-Cantin, Juan D Tardós, and José Neira. 2007. Robocentric map joining: Improving the consistency of EKF-SLAM. *Robotics and autonomous systems* 55, 1 (2007), 21–29.

[6] Yue Chen, Xingyu Chen, Xuan Wang, Qi Zhang, Yu Guo, Ying Shan, and Fei Wang. 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.

[7] Zezhou Cheng, Carlos Esteves, Varun Jampani, Abhishek Kar, Subhransu Maji, and Ameesh Makadia. 2023. LU-NeRF: Scene and pose estimation by synchronizing local unposed nerfs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 18312–18321.

[8] Shin-Fang Chng, Sameera Ramasinghe, Jamie Sherrah, and Simon Lucey. 2022. Gaussian activated neural radiance fields for high fidelity reconstruction and pose estimation. In *European Conference on Computer Vision*. Springer, 264–280.

[9] Zhaopeng Cui and Ping Tan. 2015. Global structure-from-motion by similarity averaging. In *Proceedings of the IEEE International Conference on Computer Vision*. 864–872.

[10] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. 2022. Depth-supervised nerf: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12882–12891.

[11] Johan Edstedt, Ioannis Athanasiadis, Mårten Wadenbäck, and Michael Felsberg. 2023. DKM: Dense kernelized feature matching for geometry estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17765–17775.

[12] Jakob Engel, Thomas Schöps, and Daniel Cremers. 2014. LSD-SLAM: Large-scale direct monocular SLAM. In *European conference on computer vision*. Springer, 834–849.

[13] Martin A Fischler and Robert C Bolles. 1981. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Commun. ACM* 24, 6 (1981), 381–395.

[14] Riccardo Gherardi, Michela Farenzena, and Andrea Fusiello. 2010. Improving the efficiency of hierarchical structure-and-motion. In *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 1594–1600.

[15] Richard Hartley and Andrew Zisserman. 2003. *Multiple view geometry in computer vision*. Cambridge university press.

[16] Xingyi He, Jiaming Sun, Yifan Wang, Sida Peng, Qixing Huang, Hujun Bao, and Xiaowei Zhou. 2023. Detector-Free Structure from Motion. In *arxiv*.

[17] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. 2021. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5846–5854.

[18] Nianjuan Jiang, Zhaopeng Cui, and Ping Tan. 2013. A global linear method for camera pose registration. In *Proceedings of the IEEE international conference on computer vision*. 481–488.

[19] Mohammad Mahdi Johari, Camilla Carta, and François Fleuret. 2023. Eslam: Efficient dense slam system based on hybrid representation of signed distance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 17408–17419.

[20] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. 2021. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5741–5751.

[21] Andreas Meuleman, Yu-Lun Liu, Chen Gao, Jia-Bin Huang, Changil Kim, Min H Kim, and Johannes Kopf. 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.

[22] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.

[23] Raul Mur-Artal, Jose Maria Martinez Montiel, and Juan D Tardos. 2015. ORB-SLAM: a versatile and accurate monocular SLAM system. *IEEE transactions on robotics* 31, 5 (2015), 1147–1163.

[24] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. 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.

[25] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. 2019. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 165–174.

[26] Johannes L Schonberger and Jan-Michael Frahm. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4104–4113.

[27] Noah Snavely, Steven M Seitz, and Richard Szeliski. 2008. Modeling the world from internet photo collections. *International journal of computer vision* 80 (2008), 189–210.

[28] Jiuhn Song, Seonghoon Park, Honggyu An, Seokju Cho, Min-Seop Kwak, Sungjin Cho, and Seungryong Kim. 2024. DäRF: Boosting Radiance Fields from Sparse Input Views with Monocular Depth Adaptation. *Advances in Neural Information Processing Systems* 36 (2024).

[29] Edgar Sucar, Shikun Liu, Joseph Ortiz, and Andrew J Davison. 2021. imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 6229–6238.

[30] Jiaming Sun, Zehong Shen, Yuang Wang, Hujun Bao, and Xiaowei Zhou. 2021. LoFTR: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8922–8931.

[31] Zachary Teed and Jia Deng. 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*. Springer, 402–419.

[32] Zachary Teed and Jia Deng. 2021. Droid-slam: Deep visual slam for monocular, stereo, and rgb-d cameras. *Advances in neural information processing systems* 34 (2021), 16558–16569.

[33] Prune Truong, Marie-Julie Rakotosaona, Fabian Manhardt, and Federico Tombari. 2023. Sparf: Neural radiance fields from sparse and noisy poses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4190–4200.

[34] Chen Wang, Jiadai Sun, Lina Liu, Chenming Wu, Zhelun Shen, Dayan Wu, Yuchao Dai, and Liangjun Zhang. 2023. Digging into depth priors for outdoor neural radiance fields. In *Proceedings of the 31st ACM International Conference on Multimedia*. 1221–1230.

[35] Hengyi Wang, Jingwen Wang, and Lourdes Agapito. 2023. Co-slam: Joint coordinate and sparse parametric encodings for neural real-time slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13293–13302.

[36] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and Wenping Wang. 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.

[37] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. 2021. NeRF−−: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064* (2021).

[38] Frederik Warburg, Ethan Weber, Matthew Tancik, Aleksander Holynski, and Angjoo Kanazawa. 2023. Nerfbusters: Removing Ghostly Artifacts from Casually Captured NeRFs. arXiv:2304.10532 [cs.CV]

[39] Tong Wei, Yash Patel, Alexander Shekhovtsov, Jiří Matas, and Daniel Barath. 2023. Generalized Differentiable RANSAC. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. 17603–17614. https://doi.org/10.1109/ICCV51070.2023.01618

[40] Changchang Wu. 2013. Towards linear-time incremental structure from motion. In *2013 International Conference on 3D Vision-3DV 2013*. IEEE, 127–134.

[41] Qingsong Yan, Qiang Wang, Kaiyong Zhao, Jie Chen, Bo Li, Xiaowen Chu, and Fei Deng. 2023. CF-NeRF: Camera Parameter Free Neural Radiance Fields with Incremental Learning. arXiv:2312.08760 [cs.CV]

[42] Zihan Zhu, Songyou Peng, Viktor Larsson, Weiwei Xu, Hujun Bao, Zhaopeng Cui, Martin R Oswald, and Marc Pollefeys. 2022. Nice-slam: Neural implicit scalable encoding for slam. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12786–12796.

## A  MORE IMPLEMENTATION DETAILS

Our network architecture follows the BARF [20] approach, utilizing a single 8-layer MLP network with a width of 128. All SOTA methods employ their official open-source implementations. For test-optimization, NoPe-NeRF adopts its official implementation, while all other methods undergo 100 iterations of test-optimization after per-image neighbor initialization before evaluation. The Sim(3) alignment approach is also derived from the official open-source version of BARF.

**Dataset** We choose two datasets NeRFBuster [38] which used in CF-NeRF [41] and Free-Dataset [36] which consists of more challenging scenarios with arbitrary trajectory variations and reduced frame overlap as shown in Fig. 7. We utilize the NeRFBuster sequences processed by CF-NeRF. For each scene, CF-NeRF selects approximately 50 images based on their overlap, ordered sequentially. Regarding the Free-Dataset, the **sky** scene comprises images with indexes from 50 to 100, while all other scenes consist of images with indexes from 0 to 50. All selected sequences present considerable challenges.

## B  TESTING METHODS

As mentioned in the main paper, to calculate the metrics for test images, two sequential steps during testing are required: alignment of trajectories for pose quality assessment and test-time optimization for view synthesis quality assessment.

**Alignment** A 3D similarity transformation Sim(3) for the scene and the cameras can be obtained through different methods.

- **Sim(3)** BARF [20] and L2G-NeRF [6] align estimated poses to the ground truth through Sim(3) obtained by Procrustes analysis on the camera pose locations.
- **Sim(3) with rotation** CF-NeRF [41] finds that the Procrustes analysis used for Sim(3) is unreliable when all cameras lie in a line or the camera translation contains noise. To overcome the problem, CF-NeRF adds a virtual point $(0, 0, 1)$ in the camera coordinate of each image and uses the camera parameter to transform it to the world coordinate, then uses the camera rotation during the alignment process (termed as rotation). However, we find the approach of CF-NeRF will cause more transition errors.

We list the results of pose error on **buster** both aligned by the approach of **Sim(3)** and **Sim(3) with rotation** in Table 4. The results show that the approach of **Sim(3) with rotation** can reduce rotation errors while causing more transition errors. When the accuracy of poses is high, **Sim(3) with rotation** takes rare benefits on $\Delta R$ but harms to $\Delta T$. As a result, we employ the **Sim(3)** approach in the main paper for all methods to align two trajectories and then calculate pose errors.

**Test-time optimization** Here we outline previous testing methods with different combinations of initialization and test-time optimization.

- **Sim(3) + opt.** In BARF [20], the poses are first initialized using Sim(3) alignment with Procrustes analysis on the camera pose locations. Then, an additional test-time optimization is used to further adjust the test poses. This initialization works well when the estimated poses can be aligned precisely to COLMAP poses. However, incorrect pose estimations can affect the Sim(3) alignment.
- **Estimated + no opt.** CF-NeRF [41] recovers all poses without employing a test/train split and then tests every 8th image. However, such an approach leads to results indistinguishable whether the rendered results are due to overfitting or successful reconstruction.
- **Neighbor + opt.** Nope-NeRF [3] initializes the test image pose with the estimated pose of the training frame that is closest to it. Neighbor initialization works well when the framerate is high and the test pose is near the neighbor pose. Facing complex trajectories and reduced overlap it struggles to supply a good initialization as shown in Table 1 in the main paper and Table 6.

Due to the substantial alignment errors, all methods except ours struggle to obtain reasonable initial test poses through Sim(3). In our main paper for testing results, we adopt **Neighbor + opt.** for all the methods  and also provide results using  **Sim(3) + opt.** .

**Overfitting** As described in Table 1 and Section 4.2 of the main text and Table 6, although methods like BARF converge to significant pose errors and poor structural quality, they still achieve comparable novel view synthesis metrics to our method. We attribute this to the network converging to local optima. The left part of Fig. 8 illustrates the poses estimated by BARF, where the three green boxes indicate three pose segments after fitting. To render a video, we fit B-spline functions to the estimated poses to get a smooth camera trajectory. The right part visualizes novel views synthesized for poses within the segments (a,c,e) and between segments (b, d) on the B-spline trajectory. The novel views synthesized for poses in the three segments (a, c, e) appear normal. However, the visualization results for the interpolated poses (b,d) between segments are unreasonable. It seems that each segment fits a sub-scene and the images for the interpolated poses between segments stitch different scenes together. Fig. 9 shows more results of this issue. These methods do not recover correct poses and scene geometry. During test-time optimization, poses for testing images are initialized with the estimated poses of neighboring training images. With the twisted scene and fragmented pose trajectories after training, the test-time optimization results in the pose of a test frame converging to a "pseudo " pose close to the estimated pose of the closest neighboring training frame. Then the network for the scene further overfits to the "pseudo" ground truth test pose and image ($P, I$ for example) pairs and renders an image $I'$ using $P$ to calculate visual metrics with $I$, leading to high view synthesis metrics. Notice that during the process, $P$ diverges far from its true pose to the direction minimizing $I'$ and $I$ when the scene geometry and camera poses exhibit large errors.

## C  COMPARISON TO NERF WITH COLMAP POSE

We additionally compare the novel view synthesis quality of the NeRF model trained with our estimated pose and COLMAP pose (we use it as GT) to demonstrate the pose accuracy estimated. On average, **NeRF + Our pose** achieves novel view quality close to that of **NeRF + COLMAP pose**. In some scenes, our poses have large estimation error, like **Pikachu** and **plant**. Both pose error

**Figure 7: Consecutive frames in the NeRFBuster dataset (top) and the Free Dataset (bottom). The Free Dataset exhibits more pronounced camera motion, posing greater challenges.**

**Table 4: Pose accuracy aligned by approaches of BARF and CF-NeR. $\Delta T$ is the transition error in ground truth scale and $\Delta R$ is rotation error in degree. Sorted in descending order by $\Delta R$.**

| Metrics | Method | aloe | art | car | century | garbage | flowers | picnic | pikachu | pipe | plant | roses | table | mean |
|---------|--------|------|-----|-----|---------|---------|---------|--------|---------|------|-------|-------|-------|------|
| $\Delta R\downarrow$ | CF-NeRF [41] Sim(3) with rotation | 21.918 | 25.702 | 22.653 | 11.245 | 9.061 | 9.915 | 13.489 | 12.046 | 173.343 | 11.056 | 7.002 | 3.837 | 26.772 |
| | CF-NeRF Sim(3) | 6.703 | 76.306 | 29.079 | 11.013 | 74.163 | 10.672 | 109.868 | 13.243 | 122.345 | 18.664 | 3.903 | 3.835 | 39.983 |
| | Ours Sim(3) with rotation | 3.618 | **0.469** | **0.545** | **2.237** | 0.921 | 0.596 | 2.118 | 7.698 | 2.320 | **5.212** | 1.919 | 1.223 | **2.406** |
| | Ours Sim(3) | **3.163** | 3.151 | 0.701 | 2.343 | **0.902** | **0.481** | **1.938** | **7.708** | **2.302** | 6.008 | **0.570** | **1.154** | 2.560 |
| $\Delta T\downarrow$ | CF-NeRF Sim(3) with rotation | 3.858 | 5.064 | 8.423 | 3.655 | 4.018 | 4.305 | 3.372 | 4.949 | 36.930 | 5.154 | 1.130 | 2.200 | 6.921 |
| | CF-NeRF Sim(3) | 0.637 | 1.549 | 1.621 | 0.497 | 0.548 | 0.745 | 1.285 | 0.879 | 5.757 | 0.685 | 0.182 | 0.274 | 1.222 |
| | Ours Sim(3) with rotation | 0.701 | 0.237 | 0.086 | 0.517 | 0.256 | 0.215 | 0.347 | 0.983 | 0.515 | 0.519 | 0.371 | 0.247 | 0.416 |
| | Ours Sim(3) | **0.168** | **0.030** | **0.035** | **0.134** | **0.039** | **0.039** | **0.106** | **0.548** | **0.164** | **0.225** | **0.038** | **0.045** | **0.131** |

**Table 5: Comparison to NeRF with COLMAP(GT) pose.**

| scenes | Ours | | | | | NeRF + Our pose | | | NeRF + COLMAP pose | | |
|--------|------|------|------|------|------|------|------|------|------|------|------|
| | $\Delta R\downarrow$ | $\Delta T\downarrow$ | $PSNR\uparrow$ | $SSIM\uparrow$ | $LPIPS\downarrow$ | $PSNR\uparrow$ | $SSIM\uparrow$ | $LPIPS\downarrow$ | $PSNR\uparrow$ | $SSIM\uparrow$ | $LPIPS\downarrow$ |
| pikachu | 7.708 | 0.548 | 22.91 | 0.76 | 0.29 | 23.62 | 0.79 | 0.28 | **37.06** | **0.97** | **0.05** |
| plant | 6.302 | 0.225 | 22.64 | 0.71 | 0.30 | 20.51 | 0.63 | 0.39 | **28.27** | **0.85** | **0.24** |
| aloe | 3.163 | 0.168 | 24.36 | 0.61 | 0.35 | **25.51** | **0.68** | **0.26** | 24.04 | 0.58 | 0.40 |
| art | 3.151 | 0.030 | 26.73 | 0.83 | 0.14 | **13.77** | **0.34** | **0.56** | 12.90 | 0.30 | 0.60 |
| century | 2.343 | 0.134 | 22.56 | 0.65 | 0.50 | 14.08 | 0.25 | 0.69 | **14.28** | **0.28** | **0.68** |
| pipe | 2.302 | 0.164 | 23.13 | 0.59 | 0.40 | 21.90 | 0.55 | 0.43 | **23.05** | **0.63** | **0.37** |
| picnic | 1.938 | 0.106 | 23.04 | 0.52 | 0.49 | 22.52 | 0.51 | 0.45 | **25.25** | **0.67** | **0.31** |
| table | 1.154 | 0.045 | 32.73 | 0.91 | 0.17 | **25.64** | **0.84** | **0.28** | 23.82 | 0.82 | 0.27 |
| flowers | 0.902 | 0.039 | 22.69 | 0.63 | 0.43 | **16.29** | **0.30** | **0.66** | 15.38 | 0.27 | 0.67 |
| car | 0.701 | 0.035 | 27.41 | 0.79 | 0.25 | **21.52** | **0.67** | **0.33** | 18.93 | 0.60 | 0.40 |
| roses | 0.570 | 0.038 | 29.63 | 0.88 | 0.09 | **30.09** | **0.90** | **0.09** | 27.77 | 0.84 | 0.16 |
| garbage | 0.481 | 0.039 | 27.37 | 0.81 | 0.22 | **18.34** | **0.59** | **0.41** | 13.53 | 0.36 | 0.67 |
| mean | 2.560 | 0.131 | 25.43 | 0.72 | 0.30 | 21.14 | 0.59 | **0.40** | 22.02 | **0.60** | **0.40** |

during training and test pose misalignment lead to worse view quality of **NeRF + Our pose** in these scenes. In scenes with small pose error, **NeRF + Our pose** gains similar view quality with **NeRF + COLMAP pose**, even better in many scenes. In many scenes, our methods achieve better view quality than **NeRF + Our pose** and **NeRF + COLMAP pose**. We attribute it to coarse to fine positional encoding, reprojection loss, and joint optimization of poses and scenes.

## D MORE RESULTS

Detailed results of Free-Dataset are shown in Table 6. We provide more qualitative comparisons with state-of-the-art works. Fig. 13, Fig. 14, Fig. 15 and Fig. 10 illustrates a comparison of novel view synthesis quality. Fig. 16, Fig. 17, Fig. 18 and Fig. 11 demonstrates a comparison of depth map rendering quality. Fig. 19, Fig. 20 and Fig. 12 present a comparison of the reconstructed trajectories obtained by various methods, where the red boxes represent the COLMAP
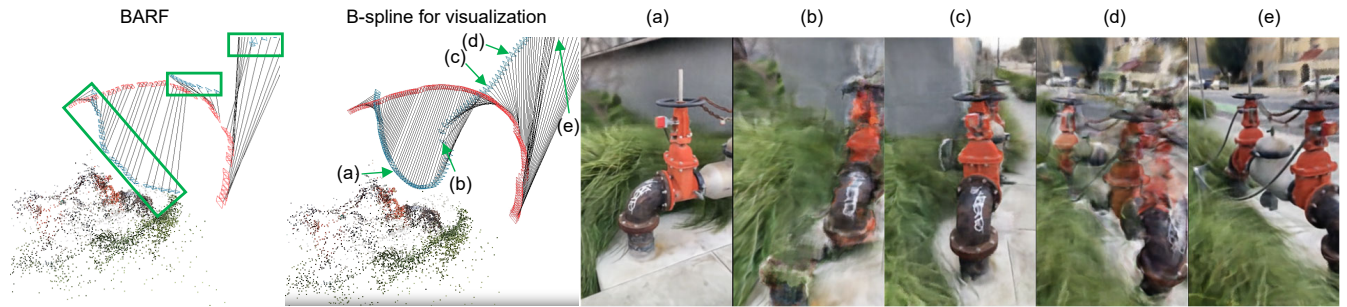
**Figure 8: Visualization of segments and interpolated views. (a), (c), and (e) are novel views in segments. (b) and (d) are interpolated novel views**

**Table 6: Evaluations of the pose accuracy (top 2 rows) and the novel view quality (bottom 3 rows) on Free-Dataset [36]. $\Delta T$ is the transition error in ground truth scale and $\Delta R$ is rotation error in degree.**

| Metrics | Method | grass | hydrant | lab | pillar | road | sky | stair | mean |
|---|---|---|---|---|---|---|---|---|---|
| | BARF [20] | 124.875 | 74.091 | 124.754 | 16.908 | 64.433 | 22.197 | 0.425 | 61.098 |
| | L2G-NeRF [6] | 114.356 | 170.250 | 56.227 | 131.588 | 109.558 | 27.245 | 162.898 | 110.303 |
| $\Delta R\downarrow$ | Nope-NeRF [3] | 158.408 | 140.245 | 165.086 | 153.613 | 144.478 | 67.749 | 179.836 | 144.202 |
| | CF-NeRF [41] | 36.875 | 36.129 | 150.882 | 18.282 | 49.790 | 94.082 | 1.260 | 55.329 |
| | Ours | **7.785** | **0.454** | **6.126** | **0.124** | **3.001** | **2.054** | **0.089** | **2.805** |
| | BARF | 7.273 | 3.890 | 6.675 | 1.475 | 4.587 | 0.583 | **0.004** | 3.498 |
| | L2G-NeRF | 7.962 | 7.203 | 4.786 | 7.707 | 7.107 | 0.849 | 10.498 | 6.587 |
| $\Delta T\downarrow$ | Nope-NeRF | 2.920 | 4.904 | 2.062 | 4.044 | 7.201 | 1.156 | 10.564 | 4.693 |
| | CF-NeRF | 2.850 | 2.018 | 5.998 | 1.382 | 1.808 | 2.518 | 0.121 | 2.385 |
| | Ours | **0.304** | **0.032** | **0.526** | **0.008** | **0.201** | **0.046** | 0.007 | **0.161** |
| | BARF | 18.00 | 17.33 | 18.73 | 20.17 | 19.01 | 15.66 | 28.00 | 19.56 |
| | L2G-NeRF | **18.29** | 17.46 | **21.18** | 19.93 | 20.49 | 17.90 | 24.41 | 19.95 |
| $PSNR\uparrow$ | Nope-NeRF | 17.02 | 18.33 | 17.55 | 18.99 | 19.08 | 15.39 | 24.35 | 18.67 |
| | CF-NeRF | 18.15 | 17.85 | 16.25 | 20.25 | 18.85 | 15.23 | 21.51 | 18.30 |
| | Ours neighbor | 17.57 | 17.95 | 17.94 | 21.91 | 19.30 | 15.12 | 21.07 | 18.69 |
| | Ours Sim(3) | 16.96 | **22.54** | 14.88 | **26.23** | **24.06** | **24.37** | **28.16** | **22.46** |
| | BARF | 0.40 | 0.33 | 0.58 | 0.52 | 0.47 | 0.49 | 0.83 | 0.52 |
| | L2G-NeRF | **0.42** | 0.32 | **0.67** | 0.53 | 0.52 | 0.55 | 0.74 | 0.54 |
| $SSIM\uparrow$ | Nope-NeRF | 0.40 | 0.37 | 0.63 | 0.47 | 0.44 | 0.60 | 0.66 | 0.51 |
| | CF-NeRF | 0.34 | 0.30 | 0.44 | 0.44 | 0.40 | 0.43 | 0.56 | 0.42 |
| | Ours neighbor | 0.40 | 0.35 | 0.53 | 0.56 | 0.49 | 0.45 | 0.64 | 0.49 |
| | Ours Sim(3) | 0.36 | **0.50** | 0.41 | **0.67** | **0.61** | **0.77** | **0.83** | **0.59** |
| | BARF | **0.51** | 0.60 | 0.38 | 0.50 | 0.51 | 0.48 | **0.18** | 0.45 |
| | L2G-NeRF | **0.51** | 0.61 | **0.26** | 0.51 | 0.57 | 0.41 | 0.25 | 0.45 |
| $LPIPS\downarrow$ | Nope-NeRF | 0.75 | 0.69 | 0.56 | 0.68 | 0.79 | 0.63 | 0.52 | 0.66 |
| | CF-NeRF | 0.77 | 0.82 | 0.57 | 0.73 | 0.83 | 0.70 | 0.59 | 0.72 |
| | Ours neighbor | 0.59 | 0.57 | 0.39 | 0.42 | 0.65 | 0.50 | 0.30 | 0.49 |
| | Ours Sim(3) | 0.61 | **0.50** | 0.55 | **0.37** | **0.48** | **0.31** | 0.21 | **0.43** |

poses, and the blue boxes depict the estimated poses. The point cloud, processed from the COLMAP output, serves as a reference for relative positioning.
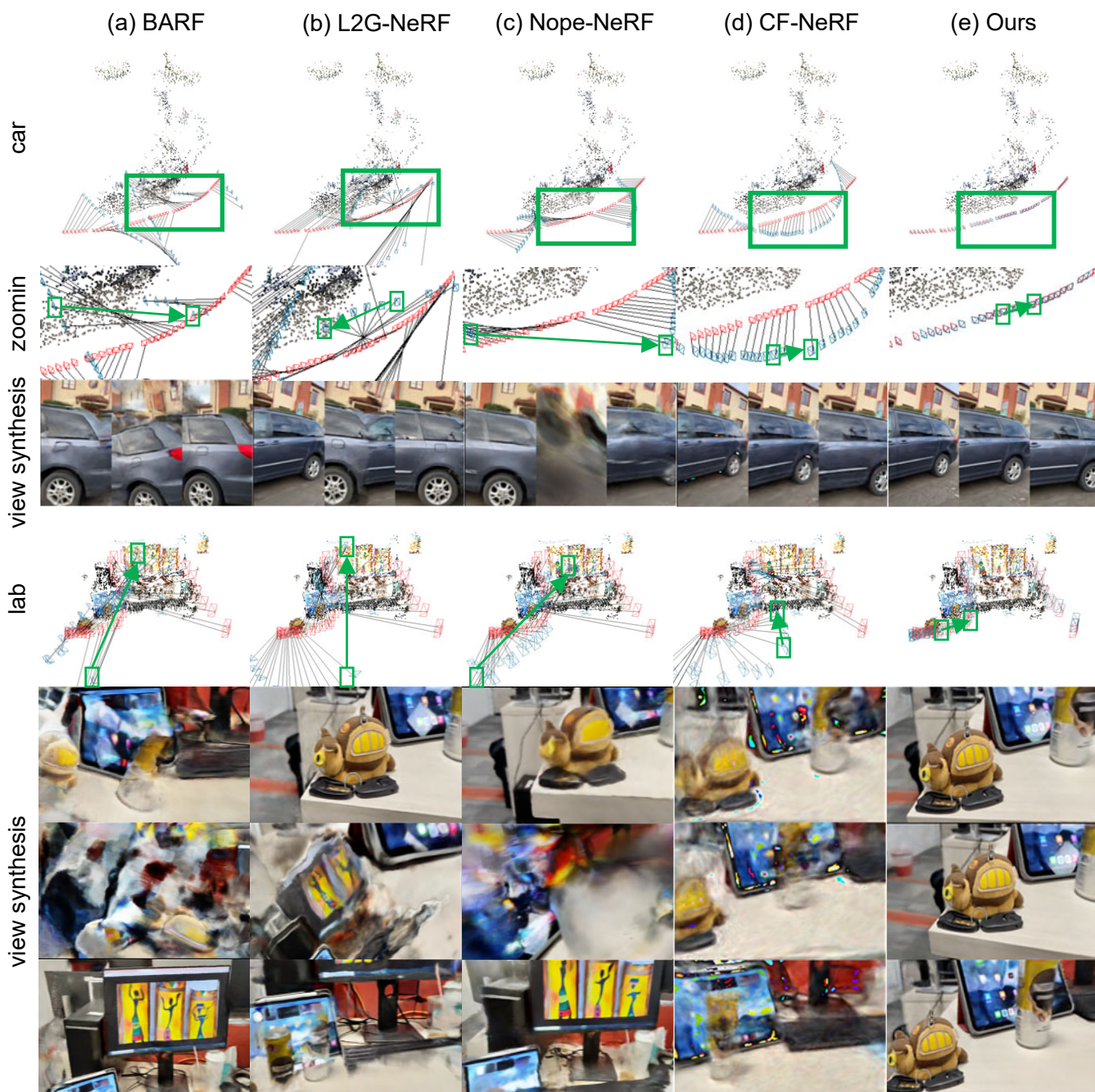
**Figure 9: Trajectory comparison.** We visualize camera poses of both estimated (blue) and COLMAP (red). Sparse 3D points for the scenes are from COLMAP. While there are abrupt changes in the trajectories of BARF, L2G-NeRF, and Nope-NeRF, the changes are steady along the trajectories of CF-NeRF and ours. The bottom row shows rendered interframes between two frames of abrupt changes denoted by green rectangles.
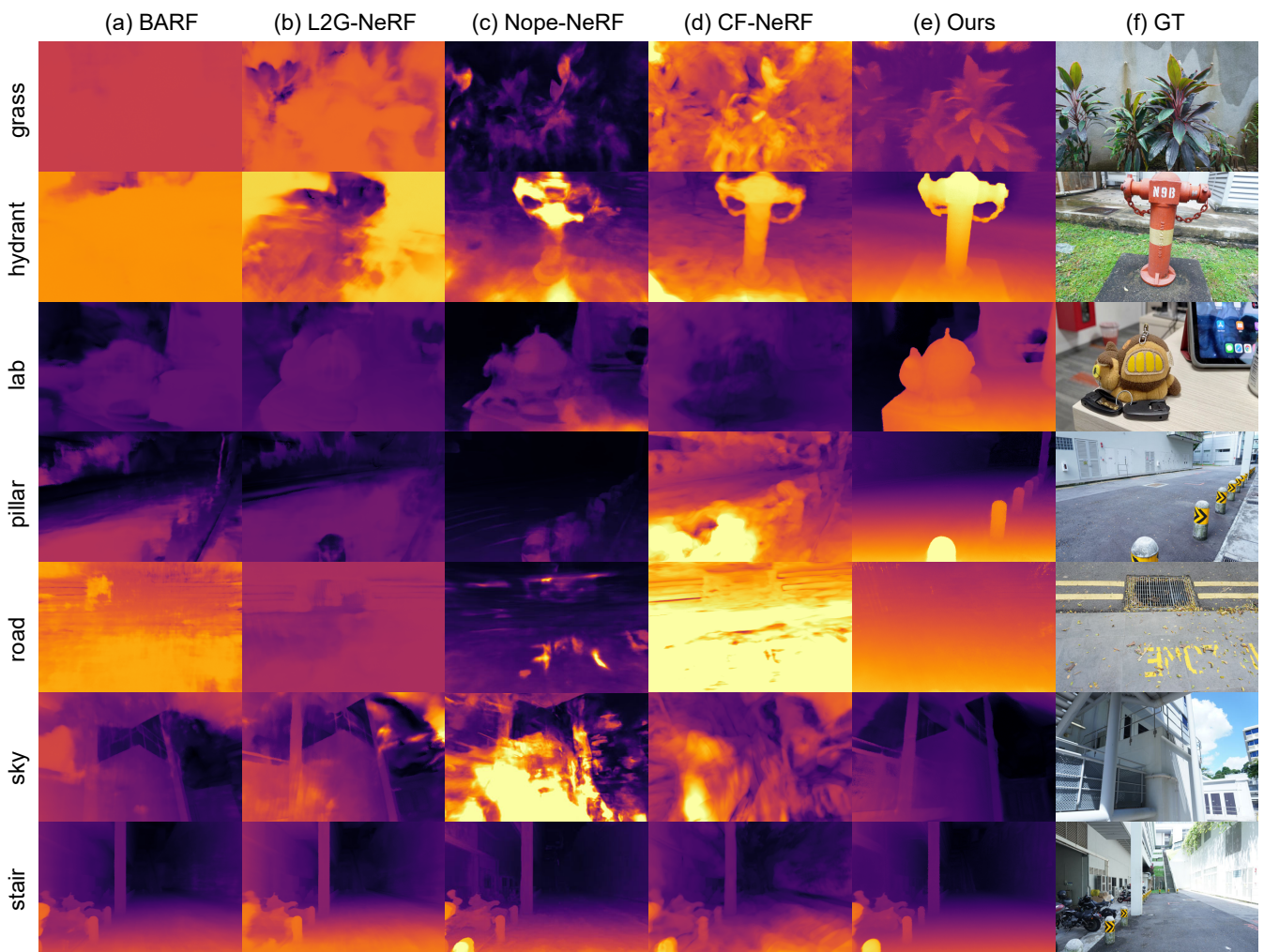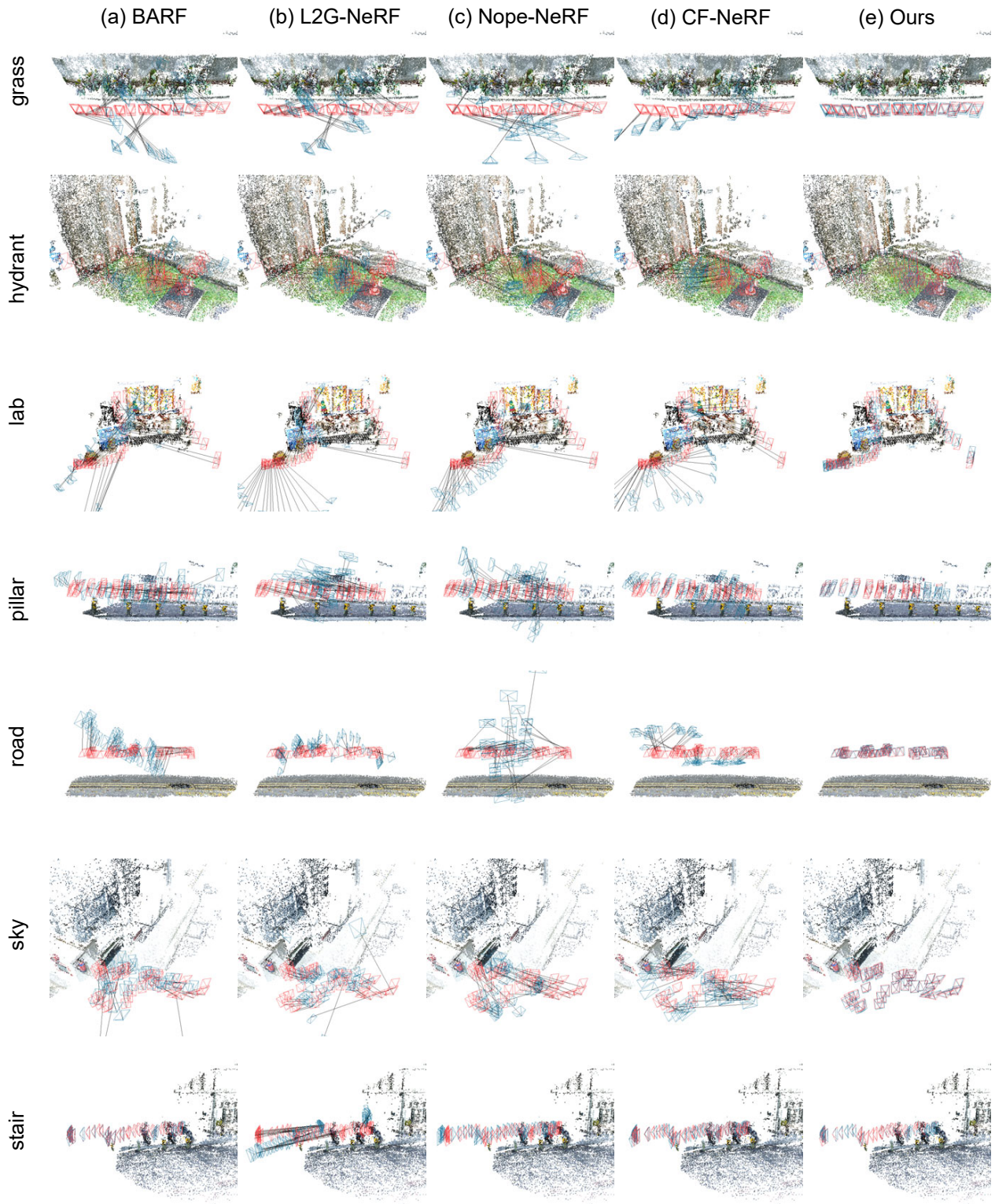
**Figure 10: Comparison of novel views on Free-Dataset [36].**

**Figure 11: Comparison of rendered depths on Free-Dataset [36].**

**Figure 12: Trajectory comparison on Free-Dataset [36]. We visualize camera poses of both estimated (blue) and COLMAP (red). Sparse 3D points for the scenes are from COLMAP.**
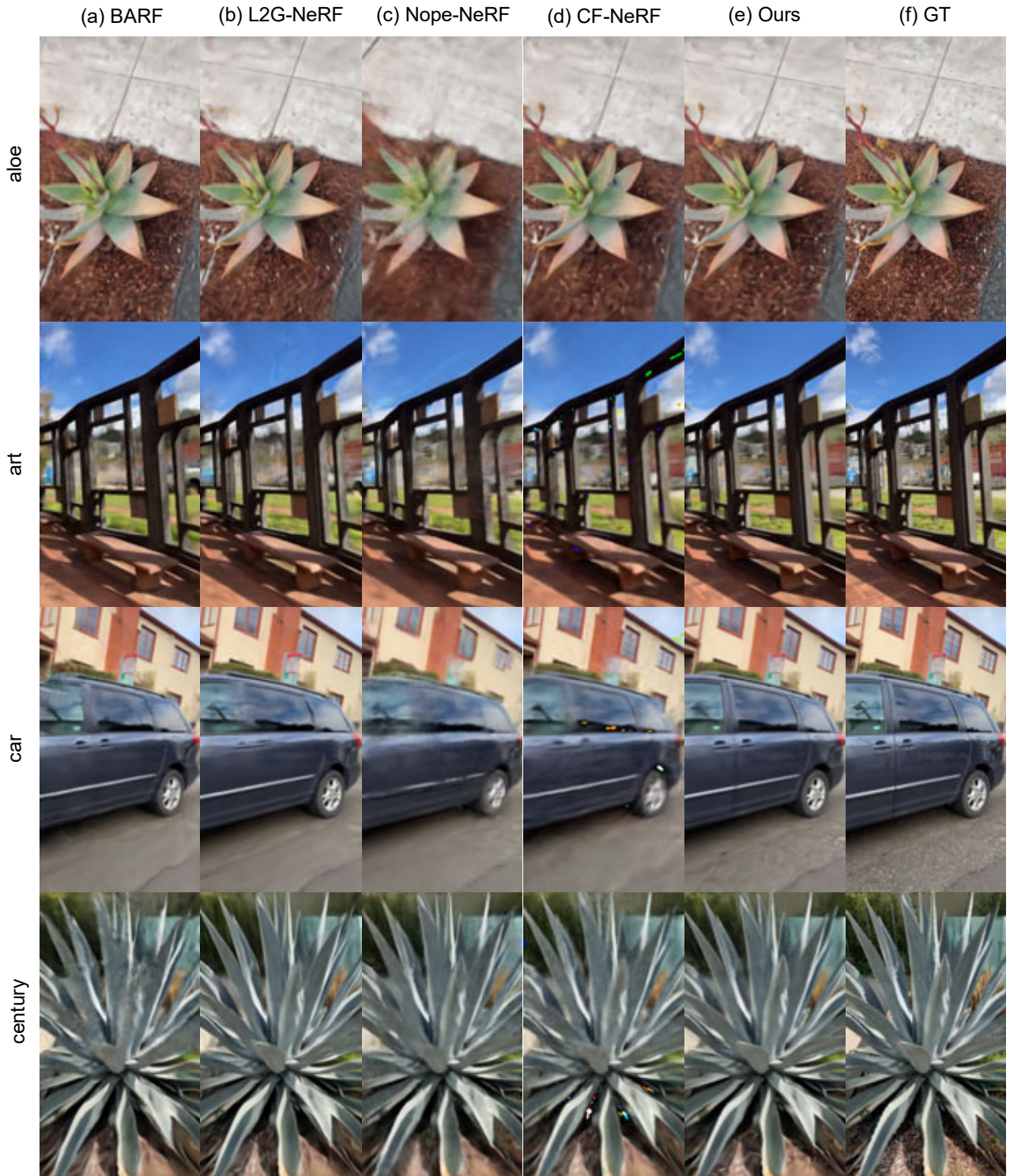
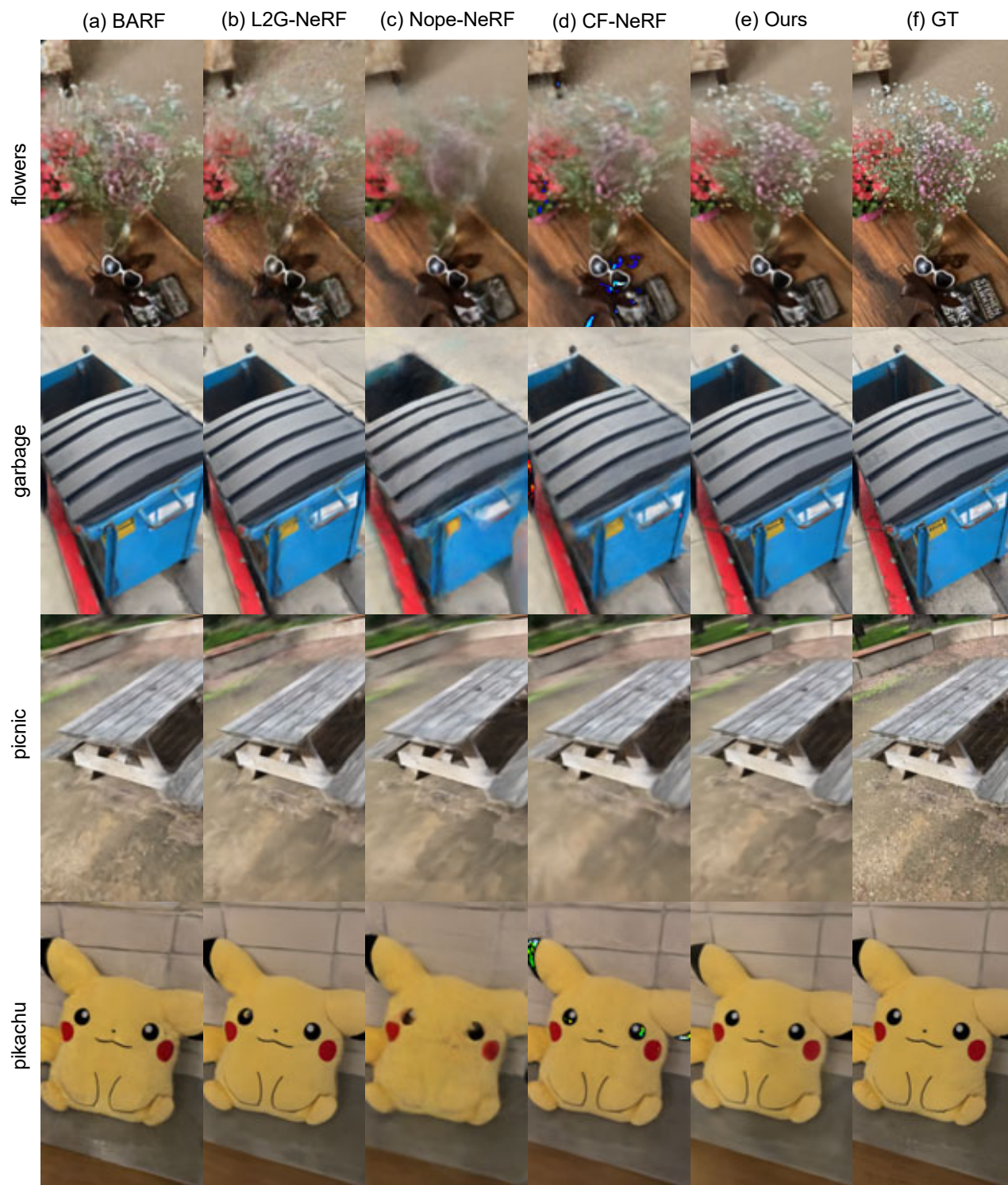**Figure 13: Comparison of novel views on NeRFBuster [38]. Part one.**

**Figure 14: Comparison of novel views on NeRFBuster [38]. Part two.**

|  | (a) BARF | (b) L2G-NeRF | (c) Nope-NeRF | (d) CF-NeRF | (e) Ours | (f) GT |
|---|---|---|---|---|---|---|



**Figure 15: Comparison of novel views on NeRFBuster [38]. Part three.**

Figure 16: Comparison of rendered depths on NeRFBuster [38]. Part one.

**Figure 17: Comparison of rendered depths on NeRFBuster [38]. Rendered depths. Part two.**

**Figure 18: Comparison of rendered depths on NeRFBuster [38]. Rendered depths. Part three.**
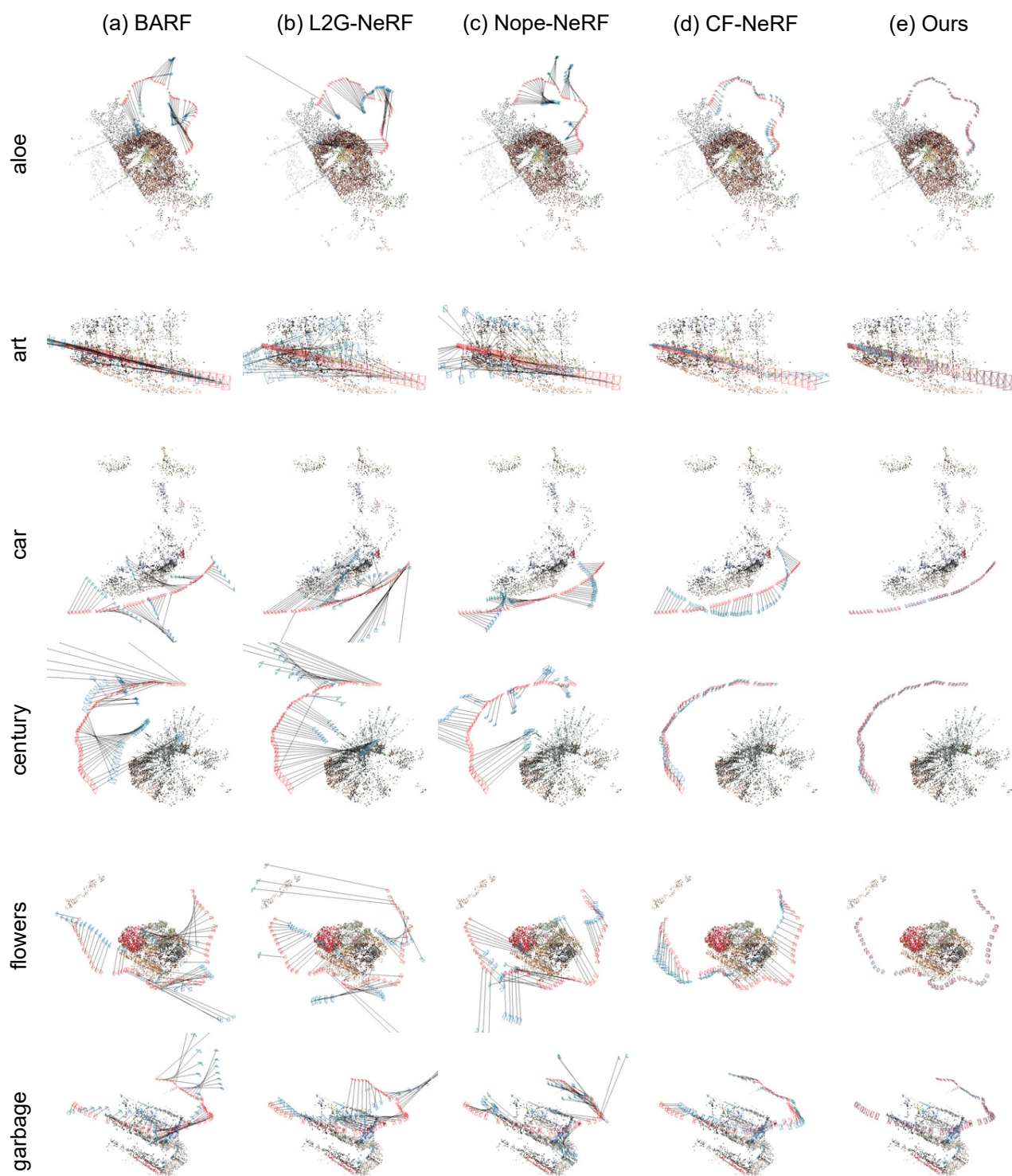
**Figure 19: Trajectory comparison on NeRFBuster [38]. We visualize camera poses of both estimated (blue) and COLMAP (red). Sparse 3D points for the scenes are from COLMAP. Part one.**
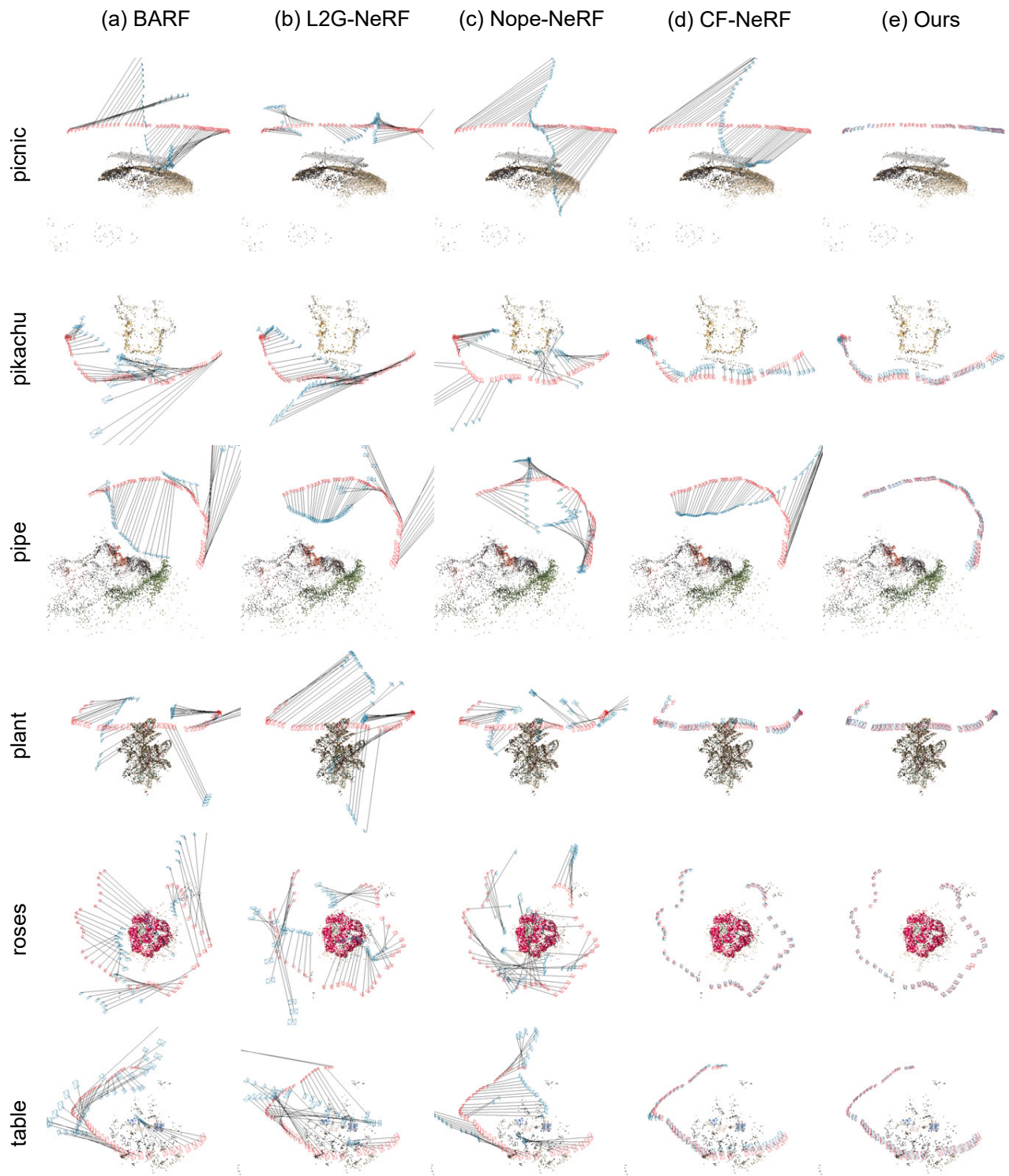
**Figure 20: Trajectory comparison on NeRFBuster [38]. We visualize camera poses of both estimated (blue) and COLMAP (red). Sparse 3D points for the scenes are from COLMAP. Part two.**