

TD-NeRF: Novel Truncated Depth Prior for Joint Camera Pose and Neural Radiance Field Optimization

Zhen Tan Zongtan Zhou Yangbing Ge Zi Wang Xieyuanli Chen* Dewen Hu*

Abstract—The reliance on accurate camera poses is a significant barrier to the widespread deployment of Neural Radiance Fields (NeRF) models for 3D reconstruction and SLAM tasks. The existing method introduces monocular depth priors to jointly optimize the camera poses and NeRF, which fails to fully exploit the depth priors and neglects the impact of their inherent noise. In this paper, we propose Truncated Depth NeRF (TD-NeRF), a novel approach that enables training NeRF from unknown camera poses - by jointly optimizing learnable parameters of the radiance field and camera poses. Our approach explicitly utilizes monocular depth priors through three key advancements: 1) we propose a novel depth-based ray sampling strategy based on the truncated normal distribution, which improves the convergence speed and accuracy of pose estimation; 2) to circumvent local minima and refine depth geometry, we introduce a coarse-to-fine training strategy that progressively improves the depth precision; 3) we propose a more robust inter-frame point constraint that enhances robustness against depth noise during training. The experimental results on three datasets demonstrate that TD-NeRF achieves superior performance in the joint optimization of camera pose and NeRF, surpassing prior works, and generates more accurate depth geometry. The implementation of our method has been released at <https://github.com/nubot-nudt/TD-NeRF>.

I. INTRODUCTION

Pose estimation and scene representation play crucial roles in 3D Reconstruction [1], [2] and Simultaneous Localization and Mapping (SLAM) [3]. In recent studies on scene representation methods [4]–[6], Neural Radiance Fields (NeRF) [7] have gained significant attention in the domains of robotics and autonomous driving, primarily because of their capacity to produce highly realistic images.

Given a set of posed images, NeRF [7] is capable of simulating radiance fields using neural networks. Most current NeRF methods [8]–[11] separate the pose estimation and reconstruction rendering processes. They use offline processing methods like Structure from Motion (SfM) [1], [12] to obtain camera poses from RGB images. Such images with poses are then fed into a radiance field network, known as posed NeRF. These loosely coupled approaches have multiple limitations. First, the accuracy of the rendering in NeRF relies on the accuracy of the camera poses. Second, the mutual correlation between the pose and the radiance field is ignored. In other words, camera poses can improve the accuracy of the

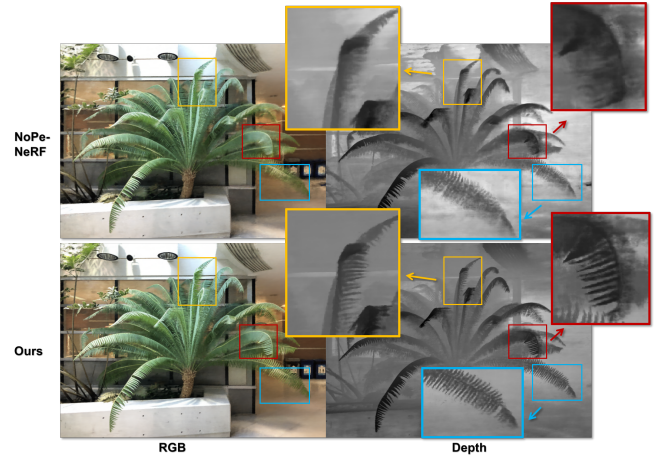


Fig. 1: Comparison with the state-of-the-art depth-based NeRF method NoPe-NeRF [16]. RGB images and depth images are rendered by NeRF with a coarse depth map.

radiance field, while radiance field information can assist in optimizing the camera pose. Third, failures occur when using methods such as COLMAP or SfM [1], [12] to generate the poses, which makes it impossible to carry out subsequent work such as radiance field reconstruction, and there is no remedy. Fourth, the loosely coupled approach is challenging to deploy in real-world scenarios for 3D reconstruction and SLAM tasks with large motion changes.

To reduce the dependence on pose in NeRF, some methods [11], [13]–[16] perform joint optimization of the poses and radiance fields. [13] adds a pose net for simultaneous optimization with the NeRF network, which is prone to fall into local optima or non-convergence; [14] and [15] use a priori information of the pose and refinement based on it, which is not effective when the pose is completely unknown; [11] is optimized for different camera models and does not discuss the case where the poses are unknown. Furthermore, [16] introduces depth priors for supervision, which fails to fully exploit the depth priors and neglects the impact of their inherent noise. [17] and [10] also utilize depth priors, but these depth priors are derived from sparse point clouds computed through SfM and are not used for joint optimization.

In this paper, to jointly optimize camera poses and NeRF for enhanced pose estimation accuracy, we leverage depth priors derived from monocular depth estimation networks. First, we introduce a sampling strategy named Truncated Depth-Based Sampling (TDBS), which utilizes depth priors informed by the truncated normal distribution to optimize

Z. Tan, Z. Zhou, Y. Ge, X. Chen, D. Hu are with the College of Intelligence Science and Technology, National University of Defense Technology, China. Z. Wang is with College of Aerospace Science and Engineering, National University of Defense Technology, China.

* indicates corresponding authors: X. Chen (xieyuanli.chen@nudt.edu.cn) and D. Hu (dwhu@nudt.edu.cn)

This work has partially been funded by the National Natural Science Foundation of China (12302252, U19A2083)

ray point sampling, as detailed in Sec. III-B. Specifically, we employ a coarse-to-fine training strategy in TDBS to obtain more accurate depth geometry (see Fig. 1). This strategy accelerates the convergence of pose optimization and enhances the precision of pose estimation. Second, we introduce a novel inter-frame point constraint (see Sec. III-C). In this constraint, we utilize the Gaussian kernel function to measure distances between inter-frame point clouds, thereby robustly handling depth noise and improving the accuracy of relative pose estimation and the quality of view synthesis.

In summary, the contributions of this paper include: (i) A robust depth-based NeRF method TD-NeRF is proposed, which jointly optimizes camera poses and radiance fields based on depth priors. Our approach can be applied to both indoor and outdoor scenes with large motion changes. (ii) We leverage depth priors for neural rendering and propose a truncated depth-based ray sampling strategy - TDBS. This strategy speeds up pose optimization convergence and improves the accuracy of pose estimation. (iii) We further propose the coarse-to-fine training strategy for the sampling strategy, effectively mitigating local minima of the model and enhancing depth geometry resolution. (iv) We propose a Gaussian point constraint that more robustly measures the distance between inter-frame point clouds, accounting for the depth noise.

II. RELATED WORK

Visual SfM and SLAM: Earlier work utilized SfM [1], [12], [18] and SLAM [3], [19], [20] systems to simultaneously reconstruct 3D structures and estimate sensor poses. These methods could be categorized into indirect methods and direct methods, where indirect methods [20], [21] employed keypoint detection and matching, while direct method [22] utilized photometric error. However, they suffered from lighting variations and textureless scenes. To address this issue, NeRF [7] represented 3D scenes using neural radiance fields and significantly improved the realism of scene representation. Given a set of posed images, NeRF [7] has demonstrated notable achievements in generating photo-realistic images and 3D reconstruction. Furthermore, its variants [23]–[25] successfully combined NeRF with SLAM methods in indoor environments.

NeRF and poses joint optimization: Recent studies have started to focus on the pose-unknown NeRF. NeRFmm [13] jointly optimized the camera poses and the neural radiance network for forward-facing scenes but was susceptible to local optima. SC-NeRF [11] proposed a joint optimization method that could be applied to different camera models. Inspired by SfM techniques, BARF [14] and GARF [26] introduced pose refinement methods. However, they relied on accurate initialization. To tackle this problem, GNeRF [15] used Generative Adversarial Networks [27], utilizing randomly initialized poses for complex outside-in scenarios. LU-NeRF [28] employed a connectivity graph-based approach for local-to-global pose estimation in 360° scenes. To achieve higher accuracy, [24], [25] utilized depth sensor information to assist NeRF and joint pose optimization. More

relevant to our work, NoPe-NeRF [16] used a monocular depth network DPT [29] to estimate coarse depth maps and leveraged depth generation to enforce pose relative constraints. However, it did not fully exploit the potential of monocular depth maps.

Different from existing works, we reassess the utilization of depth priors and propose a coarse-to-fine ray sampling strategy in Sec. III-B to efficiently optimize poses while enhancing pose estimation and novel view synthesis accuracy.

III. METHOD

The overview of our proposed method TD-NeRF is illustrated in Fig. 2. We tackle a crucial challenge in the current NeRF [7] research, which involves simultaneously optimizing the neural radiance fields and camera pose without given camera poses in both indoor and outdoor scenes (Sec. III-A). We first introduce a lightweight pre-trained depth estimation network, DPT [29] to obtain depth priors. To fully exploit depth priors and gain more accurate depth geometry, we propose a ray sampling strategy based on the truncated normal distribution [30] and depth priors called Truncated Depth-Based Sampling (TDBS). Further, to speed up the convergence during training, prevent falling into local minima, and make the rendered depth more accurate, we propose the coarse-to-fine training strategy for TDBS (Sec. III-B). Moreover, to enhance robustness against depth noise, we introduce a Gaussian Point Constraint (GPC) that measures the distance between inter-frame point clouds (Sec. III-C). Therefore, we effectively utilize the coarse depth map and seamlessly incorporate it into the synchronous training of NeRF and camera pose estimation.

A. Preliminary

a) *NeRF [7]:* It essentially models a mapping function that utilizes Multi-Layer Perceptron (MLP) layers, denoted as F_θ , which maps 3D points x and viewing directions \mathbf{d} to color c and volume density σ , i.e., $F_\theta(x, \mathbf{d}) \rightarrow (c, \sigma)$. The rendering process can be described as follows: 1) Given a sequence of images, the camera ray $r(t) = \mathbf{o} + t\mathbf{d}$ is cast into the scene, and a set of points \mathbf{x} is sampled along the ray; 2) The network F_θ is employed to estimate the density σ and color c of this set of points; 3) Volume rendering is used to integrate along the ray and obtain the color values $C(\mathbf{r})$ for synthesizing the final image. The entire integration process between $[t_n, t_f]$ is modeled as:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))c(\mathbf{r}(t), \mathbf{d})dt, \quad (1)$$

where $T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s))ds\right)$ is accumulated transmittance along a ray. More details could be referred to in [7].

b) *Joint Optimization of NeRF and Poses:* In the problem of pose-unknown NeRF, we simultaneously optimize the pose \mathbf{P} and parameters θ of NeRF. In existing methods, most approaches directly minimize the photometric error to obtain the final result. The mathematical expression is as follows:

$$\underset{\mathbf{P}, \theta}{\operatorname{argmin}} d(C(\mathbf{r}), \hat{C}(\mathbf{r})). \quad (2)$$

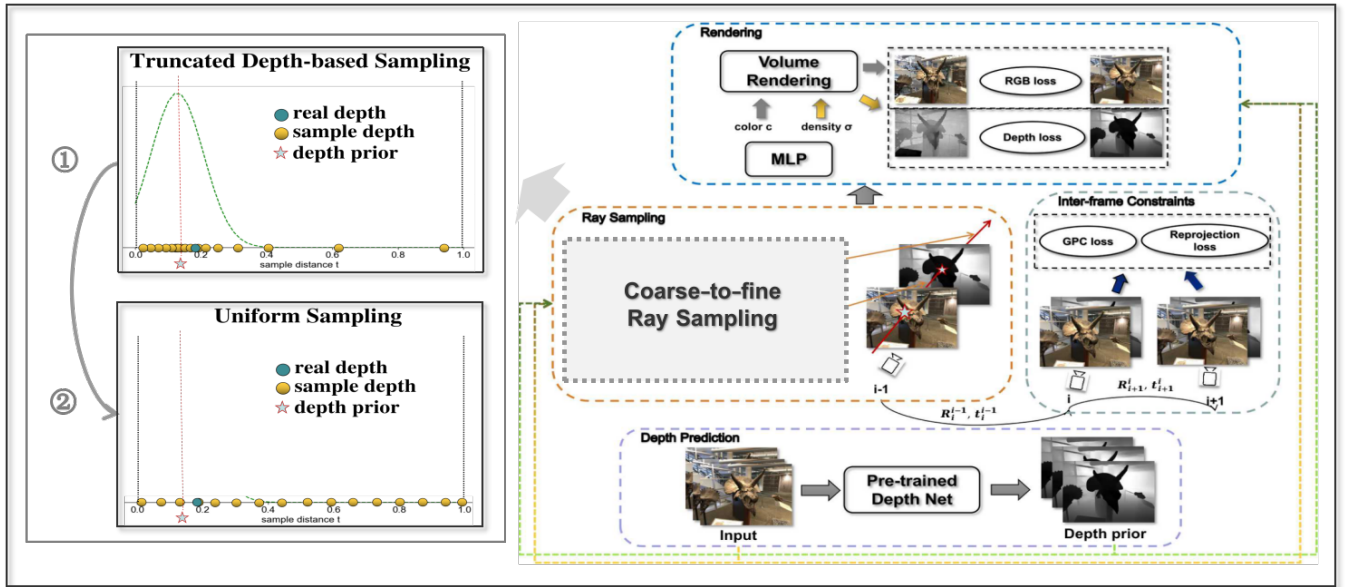


Fig. 2: Overview of our method. The inputs are RGB images without poses, and RGB images are first processed by a pre-trained depth network to obtain depth priors. Then, we employ a truncated normal distribution to optimize the ray sampling of each pixel based on the depth priors with a coarse-to-fine training strategy (①: coarse step, ②: fine step). Subsequently, the sampled points are fed into an MLP to estimate the color c and the density σ . Next, RGB and depth images are integrated by color c and σ by utilizing volume rendering. Finally, the radiance field is optimized by supervising depth and RGB. Additionally, we incorporate depth information to calculate GPC and reprojection loss between point clouds, providing constraints for inter-frame pose optimization and refinement.

In NoPe-NeRF [16], $\mathcal{L}_{self-depth}$, \mathcal{L}_{pc} , and \mathcal{L}_{rgb-s} are defined by incorporating the depth map to self-supervise depth and constrain inter-frame variations. The integration of the depth map $D_i^{nerf}(\mathbf{r})$ can be achieved through point sampling, similar to the volume integration process for color. Mathematically, this can be expressed as follows:

$$D_i^{nerf}(\mathbf{r}) = \int_{t_n}^{t_f} T(t)\sigma(\mathbf{r}(t))dt, \quad (3)$$

where t_n and t_f denote the farthest and nearest distances of the sampling points, respectively. Further, we use a self-supervised method to compute the loss between the rendered depth map $D_i^{nerf}(\mathbf{r})$ and the priori depth map \hat{D}_i^{dpt} . Since the a priori depth map obtained by the pre-trained monocular depth network does not have multi-view consistency, to recover a sequence of multi-view consistent depth maps, we use a linear transformation $s_i\hat{D}_i^{dpt} + k_i$ to undistort depth maps:

$$\mathcal{L}_{self-depth} = \sum_i^N \left\| (s_i\hat{D}_i^{dpt} + k_i) - D_i^{nerf} \right\|. \quad (4)$$

where $\{(s_i, k_i) | i = 0, \dots, N\}$ are learnable parameters that will be optimized along with the NeRF network. From Eq. (3), the sampling points are sampled uniformly in that interval $[t_n, t_f]$, however, this ignores the role of the depth priors. In the following Sec. III-B, we leverage the depth priors and improve the integration of $D_i^{nerf}(\mathbf{r})$ by proposing TDBS to make the depth self-supervising more robust.

B. Truncated Depth-Based Sampling (TDBS)

a) *full-stage*: In our approach, we utilize a monocular depth network, DPT [29], to estimate the depth of each

image. However, the depth map estimated by the pre-trained network is a coarse depth map, where the depth is generally accurate in most regions but fails to capture the fine details of certain objects. The previous state-of-the-art method, NoPe-NeRF [16], considers overall depth distortion, which does not fully exploit the depth information. Therefore, we rethink the role of the monocular depth priors in rendering and propose that a coarse depth map can provide a prior for ray sampling, assisting in sampling by assuming that the real surface is near the estimated depth. To achieve this goal, we design the TDBS ray sampling strategy based on the truncated normal distribution [30]. This strategy enables sampling from a truncated normal distribution $\psi(\bar{\mu}, \bar{\sigma}, a, b; x)$ with mean $\bar{\mu}$ (depth value) and variance $\bar{\sigma}$, truncated within the interval $[a, b]$. The mathematical expression of the probability density function $\psi(\cdot)$ is as follows:

$$\psi(\bar{\mu}, \bar{\sigma}, a, b; x) = \begin{cases} 0 & x \leq a; \\ \frac{\phi(\bar{\mu}, \bar{\sigma}^2; x)}{\Phi(\bar{\mu}, \bar{\sigma}^2; b) - \Phi(\bar{\mu}, \bar{\sigma}^2; a)} & a < x < b; \\ 0 & b \leq x, \end{cases} \quad (5)$$

where $\phi(\cdot)$ and $\Phi(\cdot)$ denote the probability density and distribution function of the standard normal distribution, respectively. We use inverse distribution sampling [31] to obtain sampling points along the ray. The number of sampling points is set to 128. In practice, the values of a , b , and σ are set to 0.001, 1.0, and 0.1, respectively.

b) *Coarse-to-fine Training Strategy*: Relying on the aforementioned strategy, we can predict the depth of the majority of pixels with relative accuracy. However, the depth estimated by DPT is inaccurate or even incorrect in many details. If we use full-stage TDBS, these erroneous depth

values tend to degrade the rendering quality. Therefore, to avoid the model from falling into a local minima during the training, we propose a coarse-to-fine training strategy for TDBS. In the coarse stage, we utilize TDBS to converge the density values of the majority of points to the optimal values. In the fine stage, we adopt Uniform Sampling. The purpose of this stage is to recompute the erroneous depth values based on the already optimized majority of points, aiming to optimize the noisy and erroneous points. In summary, the coarse-to-fine TDBS strategy can be defined as follows:

$$S(\mathbf{r}) = \begin{cases} S_{\text{TDBS}}(\mathbf{r}) & \text{epoch} < T_s \\ S_{\text{uni}}(\mathbf{r}) & \text{epoch} \geq T_s \end{cases}, \quad (6)$$

where $S_{\text{TDBS}}(\mathbf{r})$ and $S_{\text{uni}}(\mathbf{r})$ denote TDBS and Uniform Sampling, respectively. T_s is a threshold controlling the sampling method. In our approach, T_s corresponds to the epoch at which the learning rate begins to schedule. In this way, according to Eq. (3), a more robust and accurate initial value of $D_i^{\text{nerf}}(\mathbf{r})$ can be obtained, which is self-supervised with $D_i^{\text{dpt}}(\mathbf{r})$, resulting in a fast and robust convergence of the model. We empirically set T_s to 1000.

C. Gaussian Point Constraint (GPC)

Given the depth point clouds of a scene, a method [16] utilizes the Chamfer distance [32] to measure the distance between the inter-frame point clouds as a constraint for pose estimation. However, the Chamfer distance can be sensitive to noise in the estimated depth map, since it imposes a strong constraint on point clouds and can penalize even small deviations between the estimated and ground truth point clouds.

Therefore, we propose a point cloud constraint based on the Gaussian kernel function [33], which is expressed as:

$$w_{ij} = \exp\left(-\frac{\|\mathbf{p}_i^m - \mathbf{p}_j^{m+1}\|^2}{2\sigma_{pc}^2}\right), \quad (7)$$

where w_{ij} denotes the weight between points \mathbf{p}_i^m and \mathbf{p}_j^{m+1} , \mathbf{p}_i^m represents the i_{th} point in the m_{th} frame of the point cloud, σ_{pc} denotes the standard deviation of the Gaussian kernel function, and $\|\mathbf{p}_i - \mathbf{p}_j\|^2$ represents the Euclidean distance between points \mathbf{p}_i and \mathbf{p}_j . Here, σ_{pc} is empirically set to 1.

By applying the Gaussian kernel function to the distance matrix, the corresponding weight values can be determined based on the proximity of points. Points that are closer together have higher weights, while points that are farther apart have lower weights. This allows for greater emphasis on the distances between adjacent points, thus more accurately reflecting the relationships between point clouds. We can incorporate this into the regularization term \mathcal{L}_{GPC} :

$$D_{ij} = w_{ij} \cdot \|\mathbf{p}_i^m - \mathbf{p}_j^{m+1}\|, \quad (8)$$

$$\mathcal{L}_{GPC} = \sum D_{ij}, \quad (9)$$

where D_{ij} is the distance between points \mathbf{p}_i^m and \mathbf{p}_j^{m+1} .

D. Reprojection Loss

The GPC mentioned above establishes associations between 3D points across frames. Additionally, the photometric error between pixel correspondences across frames is also minimized. Hence, the reprojected photometric loss is defined as:

$$\mathcal{L}_{reproj} = \sum_{(m,n)} \|I_m(K_m P_m), I_n(K_n(R_m^n P_m + t_m^n))\|, \quad (10)$$

where (R_m^n, t_m^n) denote the rotation matrix and translation from the m_{th} camera coordinate to the n_{th} camera coordinate, P_m represents the point cloud derived from the m_{th} depth map and RGB image, K_m denotes a projection matrix for the m_{th} camera, and $I(\cdot)$ represents the pixel value of the point in the image.

E. Overall Training Loss

Assembling all loss terms, we get the overall loss function:

$$\mathcal{L} = \mathcal{L}_{rgb} + \lambda_1 \mathcal{L}_{self-depth} + \lambda_2 \mathcal{L}_{GPC} + \lambda_3 \mathcal{L}_{reproj}, \quad (11)$$

where $\lambda_1, \lambda_2, \lambda_3$ represent weights assigned to different loss terms. In practice, we set $\lambda_1 = 0.04, \lambda_2 = 1.0, \lambda_3 = 1.0$.

IV. EXPERIMENTAL EVALUATION

The main focus of this work is how the depth prior can be leveraged to achieve simultaneous optimization of NeRF and camera pose. We compare our method with pose-unknown methods: Nope-NeRF [16] and NeRFmm [13].

We present our experiments to show the capabilities of our method. The results of our experiments also support our key claims, which are: (i) Our proposed method TD-NeRF can optimize the camera pose and radiance field simultaneously based on the depth prior and can be applied to both indoor and outdoor scenes. A specific large motion changes dataset validates that our method can be applied in the presence of large motion changes as well. (ii) The proposed TDBS sampling strategy based on the depth prior not only speeds up the convergence of training but also yields better pose estimation results. (iii) Further, the experimental results verify that our proposed coarse-to-fine training strategy approach is effective in avoiding falling into local optima and obtaining more accurate depth geometry. (iv) Finally, the experiments verify that the Gaussian point constraint is more robust compared to other inter-frame point constraints.

A. Experimental Setup

a) Datasets: We conduct experiments on three datasets to evaluate the performance of our method, including LLFF [34], Tanks and Temples [35], and BLEFF [13].

LLFF [34]: We first conduct experiments on the forward-facing dataset as that in NeRFmm [13] and NeRF [7], which has 8 scenes containing 20-62 images. The resolution of the training images is 756×1008 , and every 8th image is used for novel view synthesis.

Tanks and Temples [35]: 8 scenes are used to evaluate the quality of novel view synthesis and camera pose estimation. To emulate scenarios of bandwidth limitations in

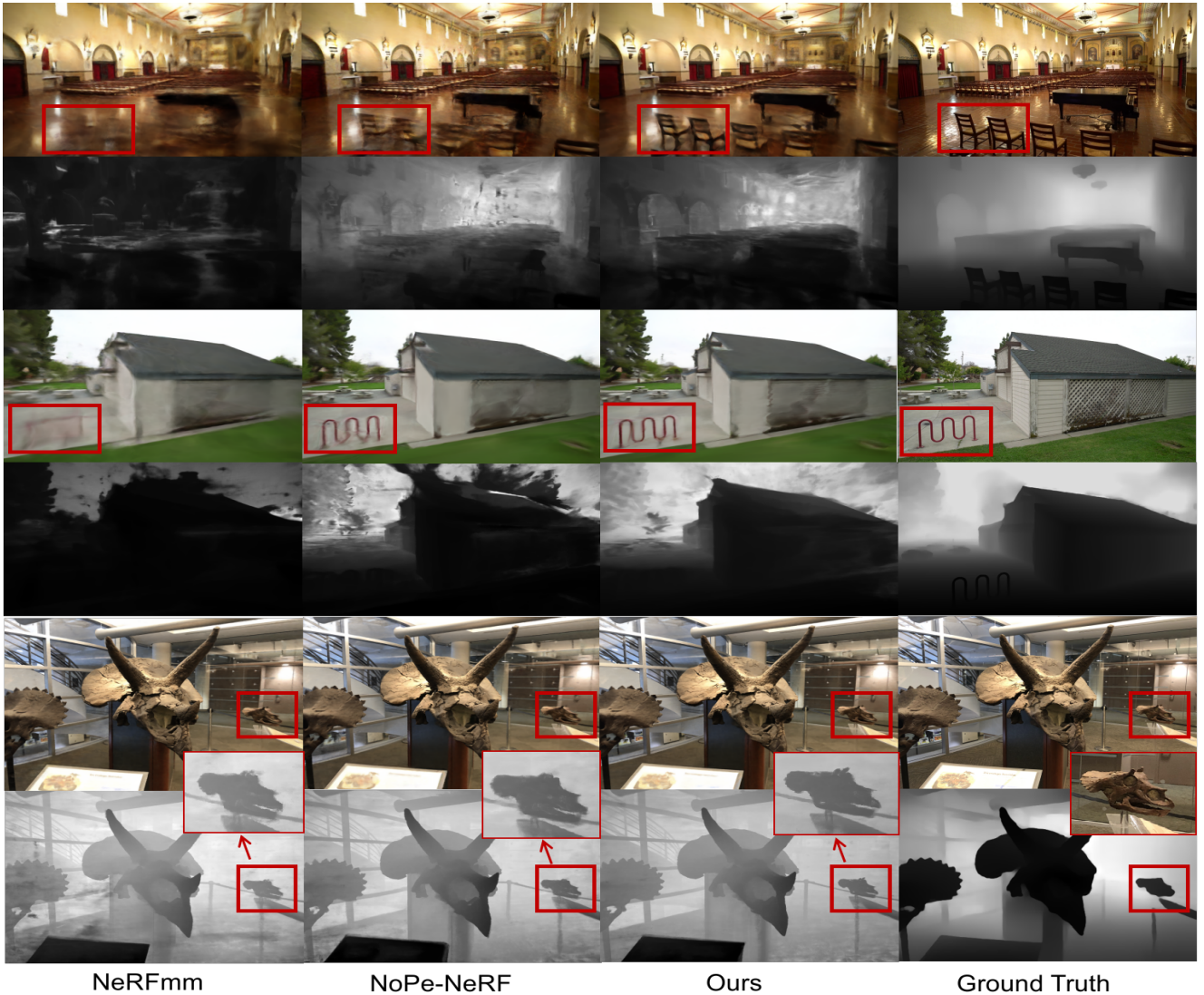


Fig. 3: Qualitative Comparison of Novel View Synthesis on Tanks & Temples (top: 4 rows) and LLFF (bottom: last 2 rows) dataset. The rendered RGB and depth images are visualized above. TD-NeRF is able to recover better details for both RGB and depth geometry, as shown in the red box. (The ground truth of depth is generated by DPT [29].)

data transmission and large motion changes, we employed a sampling approach on the dataset provided by Nope-NeRF [16]. Specifically, we selected every 5th image, resulting in a reduced frame rate of 6 fps (origin: 30 fps). During training, we utilized images with a resolution of 960×540 . Similarly, for the Family scene, we divided the frames equally, allocating half for training and the remaining half for testing.

BLEFF [13]: This dataset proposed by NeRFmm has the ground truth of the pose and can be used specifically to check the accuracy of the pose estimation. To simulate the case of large motion changes, we select 5 scenes of subset $t_{010}r_{010}$ with large rotation and translation to evaluate pose accuracy and novel view synthesis. Each scene consists of 31 images, all of which are sampled to a resolution of 780×520 .

b) Metrics: To evaluate our proposed method, we consider two aspects of test views, including the quality of novel view synthesis and camera pose estimation.

For novel view synthesis: following previous methods [7]–[9], we use (1) Peak signal-to-noise ratio (PSNR); (2) Structural Similarity Index (SSIM); and (3) Learned Perceptual Image Patch Similarity (LPIPS).

For camera pose evaluation: we use visual odometry metrics [36]; (4) Relative Pose Error (RPE); and (5) Absolute Trajectory Error (ATE).

c) Implementation Details: We implement our approach in PyTorch, Ubuntu20.04, RTX4090. We adopt the multi-network structure based on Nope-NeRF [16] without any modification. The activation function is Softplus. We sample 128 points along each ray using a truncated normal distribution combined with depth priors (Sec. III-B) with noise within a predefined range. For LLFF, the predefined range is (0.0, 1.0), and for Tanks and Temples, the range is (0.1, 10). The initial learning rate for NeRF is 0.001 and for pose and distortion is 0.0005. The training process is divided into 3 phases: the first phase is the simultaneous optimization

TABLE I: Camera pose estimation results on LLFF.

Scenes	Ours			NoPe-NeRF			NeRFmm		
	RPE _t	RPE _r	ATE	RPE _t	RPE _r	ATE	RPE _t	RPE _r	ATE
fern	0.088	0.1413	0.0008	0.326	1.2690	0.0040	\	1.780	0.0290
flower	0.036	0.0482	0.0005	0.055	0.2361	0.0009	\	4.840	0.0160
fortress	0.066	0.2310	0.0010	0.104	0.4550	0.0020	\	1.360	0.0250
horns	0.177	0.3320	0.0030	0.264	0.5440	0.0050	\	5.550	0.0440
leaves	0.140	0.0270	0.0013	0.147	0.0474	0.0014	\	3.900	0.0160
orchids	0.135	0.1121	0.0016	0.313	0.9708	0.0051	\	4.960	0.0510
room	0.082	0.2846	0.0014	0.329	1.0339	0.0054	\	2.770	0.0300
trex	0.432	0.6034	0.0068	0.557	0.7339	0.0089	\	4.670	0.0360
average	0.145	0.2225	0.0021	0.262	0.6613	0.0041	\	3.729	0.0309

 TABLE II: Quantitative evaluation of novel view synthesis and camera pose estimation on Tanks and Temples. We adopt the COLMAP pose as the ground truth. Unlike the dataset provided by NoPe-NeRF [16], to simulate the **bandwidth-constrained scenario**, the original data is sampled at 1 frame every 5 frames, and the frequency is reduced from 30 fps to 6 fps. (†: indicates our re-implementation in 6 fps.)

	View Synthesis Quality			Camera Pose Estimation		
	PSNR↑	SSIM↑	LPIPS↓	RPE _t	RPE _r	ATE
NeRFmm†	21.26	0.57	0.56	6.6800	1.6345	0.1663
NoPe-NeRF†	23.63	0.67	0.46	3.6065	0.4991	0.0909
ours	23.76	0.67	0.45	3.1759	0.4728	0.0762

of all losses; the second phase gradually reduces all weights to 0 except RGB loss; and in the third phase, the network is finely optimized for a single loss. All models are trained for 12000 epochs unless otherwise specified. The initial weights of the different losses are $\lambda_1 = 0.04$, $\lambda_2 = 1.0$, $\lambda_3 = 1.0$.

B. Performance

1) *On Camera Pose Estimation*: Our proposed method significantly outperforms other baselines in all metrics for camera pose estimation. NoPe-NeRF [16] is the state-of-the-art NeRF designed for pose-unknown novel view synthesis. NeRFmm [13] is one of the key benchmarks in the pose-unknown NeRF approach. We evaluated its performance on three datasets. For LLFF [34] and Tanks and Temples [35], where the ground truth pose is not available, we utilized the pose provided by COLMAP [18] as the reference. As shown in Tabs. I and II. Our method demonstrates a significant reduction in errors on the LLFF, achieving improvements of 44.8%, 66.4%, and 49.8% compared to the state-of-the-art methods. Similarly, on the Tanks and Temples, our method achieves reductions in errors of 8.88%, 5.27%, and 10.02%, respectively. To further validate the robustness of our approach, we conducted experiments on the BLEFF [13], which contains the ground truth of the camera pose. Our method outperforms the current state-of-the-art method by a considerable margin in terms of error reduction. We provide quantitative results and a qualitative visualization in Tab. III and Fig. 4, respectively. Our method achieves better quantitative and qualitative results than previous state-of-the-art methods. This experiment demonstrates that our method can optimize the camera pose simultaneously in both indoor and outdoor scenes with large motion changes.

 TABLE III: Quantitative evaluation of camera pose estimation on BLEFF (subset: $t_{010}r_{010}$). The ground truth of the pose is provided by BLEFF dataset.

Scenes	Ours			NoPe-NeRF			NeRFmm		
	RPE _t	RPE _r	ATE	RPE _t	RPE _r	ATE	RPE _t	RPE _r	ATE
airplane	0.2786	0.0904	0.0028	23.0097	2.5415	0.3661	0.3206	0.1004	0.0029
bed	0.2650	0.02391	0.0025	0.9045	0.1023	0.0091	0.2709	0.0246	0.0025
classroom	0.4783	0.0479	0.0053	1.9869	0.2140	0.0177	1.3417	0.1324	0.0115
halloween	0.2691	0.1165	0.0028	5.6785	1.2731	0.1140	10.7105	1.0577	0.2024
castle	0.4380	0.0943	0.0038	1.4011	0.3094	0.0123	0.5669	0.1220	0.0050
average	0.3458	0.0746	0.0034	6.5961	0.8881	0.1038	2.6421	0.2874	0.0449



Fig. 4: Pose Estimation Comparison. We visualize the camera poses on LLFF (scene: fern). red: ground truth; blue: predicted pose

2) *On Novel View Synthesis*: The second experiment evaluates the quality of novel view synthesis and illustrates that our approach is capable of further improving the quality of rendering on different datasets. We visualize the results on Tanks & Temples and LLFF datasets, as shown in Fig. 3. Our method can better reproduce the rendering details in both indoor and outdoor scenes. Especially in outdoor scenes, our method renders depth geometry well. Whereas previous methods tend to perform poorly in such open scenes. The quantitative results are summarised in Tab. II and Tab. IV.

3) *Effectiveness of coarse-to-fine TDBS*: The third experiment is to evaluate the effectiveness of our sampling strategy and coarse-to-fine training strategy: the accuracy of depth estimation and camera pose estimation and faster training convergence. As shown in Fig. 6, our TDBS strategy generates more accurate depth geometry than the state-of-the-art method. Specifically, the full-stage TDBS performs better at capturing details in nearby objects, while the innovative coarse-to-fine TDBS reliably reconstructs depth maps with realistic details at any distance. Furthermore, coarse-to-fine TDBS can also achieve better results in terms of novel view synthesis. Moreover, the effectiveness of our coarse-to-fine training strategy is substantiated by quantitative experiments, as indicated in the results shown in Tab. V. Compared with the uniform and full-stage training strategies, the average errors of the coarse-to-fine training strategy in camera pose

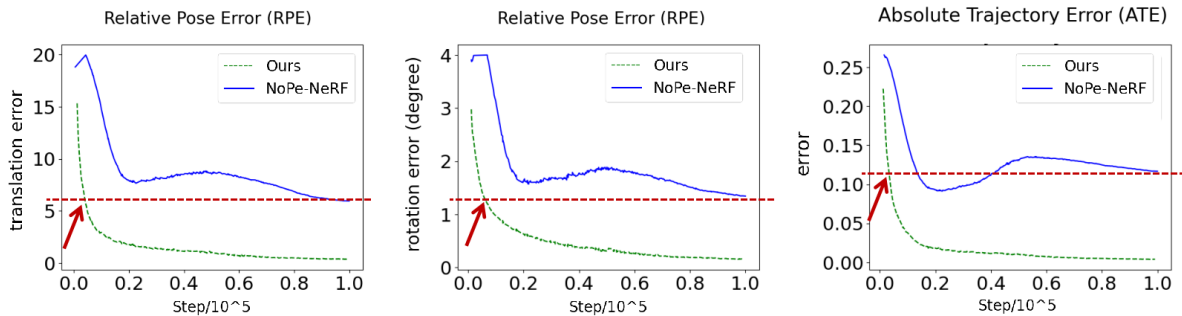


Fig. 5: Visualization of convergence. The experiment is conducted on the dataset BLEFF (scene: bed1). The blue and green colors denote NoPe-NeRF and ours, respectively. At the red arrow, the error of our method already reaches the final convergence result of NoPe-NeRF.

TABLE IV: Novel View Synthesis results on LLFF dataset.

Scenes	Ours			Nope-NeRF			NeRFmm			COLMAP		
	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS
fern	20.08	0.61	0.40	19.79	0.59	0.44	21.67	0.61	0.50	22.22	0.64	0.47
flower	28.88	0.85	0.19	29.04	0.85	0.19	25.34	0.71	0.37	25.25	0.71	0.36
fortress	28.67	0.79	0.25	27.46	0.74	0.26	26.20	0.63	0.49	27.60	0.73	0.38
horns	25.80	0.76	0.34	25.00	0.73	0.37	22.53	0.61	0.50	24.25	0.68	0.44
leaves	20.23	0.64	0.37	19.97	0.63	0.38	18.88	0.53	0.47	18.81	0.52	0.47
orchids	17.90	0.51	0.41	18.03	0.49	0.43	16.73	0.55	0.39	19.09	0.51	0.46
room	28.02	0.90	0.26	27.90	0.89	0.29	25.84	0.84	0.44	27.77	0.87	0.40
trex	25.40	0.83	0.30	24.89	0.81	0.32	22.67	0.72	0.44	23.19	0.74	0.41
average	24.37	0.74	0.32	24.01	0.72	0.34	22.48	0.65	0.45	23.52	0.68	0.42

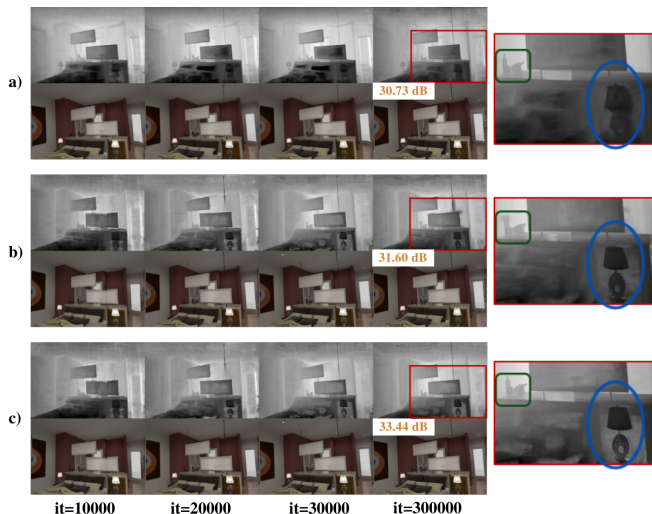


Fig. 6: Visual comparison of TDBS strategy over training time. We compare a baseline method with two different training strategies on the BLEFF dataset (scene: bed1). a: NoPe-NeRF [16], b: the full-step TDBS, and c: the coarse-to-fine TDBS.

estimation are reduced by 73.3% and 46%, respectively. Additionally, convergence comparisons between TDBS and the original method, illustrated in Fig. 5, reveal that TDBS not only accelerates the convergence speed during training but also improves the accuracy of pose estimation. Most importantly, TDBS achieves the previously established minimal error within an estimated 1000 epochs, a fraction (one-tenth) of the epochs previously needed, while successfully avoiding local optima, thereby allowing for continued optimization. This advancement is a testament to the robustness of our

TABLE V: Ablation study results of different training strategies on BLEFF (scene: bed1)

	View Synthesis Quality			Camera Pose Estimation		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	RPE $_t$	RPE $_r$	ATE
uniform	30.73	0.95	0.12	0.9045	0.1023	0.0091
full-step TDBS	31.60	0.94	0.14	0.4900	0.0487	0.0045
coarse-to-fine TDBS	33.44	0.95	0.11	0.2650	0.0239	0.0025

TABLE VI: Ablation study results of different strategies on LLFF (scene: fern).

	View Synthesis Quality			Camera Pose Estimation		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	RPE $_t$	RPE $_r$	ATE
ours w/o GPC	20.16	0.60	0.42	0.092	0.1780	0.0009
ours w/o TDBS	20.08	0.61	0.40	0.125	0.3650	0.0010
ours w/o TDBS + GPC	19.99	0.60	0.42	0.207	0.7643	0.0024
baseline (NoPe-NeRF)	19.79	0.59	0.44	0.326	1.2690	0.0040
ours	20.20	0.61	0.39	0.092	0.1450	0.0009

TABLE VII: Ablation study results of different point constraints on LLFF (scene: fern).

	View Synthesis Quality			Camera Pose Estimation		
	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	RPE $_t$	RPE $_r$	ATE
distribution	19.75	0.58	0.45	0.398	1.3081	0.0040
KL divergence	19.95	0.59	0.43	0.251	0.9587	0.0029
ICP	19.79	0.59	0.44	0.326	1.2690	0.0040
GPC (ours)	20.08	0.61	0.40	0.125	0.3650	0.0010

strategy in elevating model training efficiency and precision.

4) *Ablation Study*: In this section, we study the impact of different components of our algorithm and the importance of Gaussian Point Constraint (GPC).

Different Strategies. We consider two variants of our algorithm: TDBS and GPC. Tab. VI illustrates the results. When the GPC component is removed, the results of novel view synthesis and pose estimation do not decrease much, but if the TDBS component is removed, the results of pose estimation have a relatively large effect. This suggests that the TDBS strategy is a determining factor in the improvement of the camera pose estimation. It is because that TDBS obtain better depth geometry, pose estimation based on better depth geometry will naturally have a smaller error. Besides, TDBS and GPC have a superimposed effect on pose estimation.

Inter-frame Point Constraint. In addition, we compare different inter-frame point constraint methods, including ICP,

point cloud distribution, and KL divergence methods. From the quantitative results in Tab. VII, along with the improved quality of the view synthesis, our proposed constraint is significantly reduced by 59% compared to both second-best constraints in terms of the camera pose estimation error.

V. CONCLUSION

In this paper, we present a novel approach, TD-NeRF, for joint optimization of camera poses and radiance fields. Our approach operates the coarse depth map obtained from the pre-trained depth network for depth self-supervision and proposes a coarse-to-fine sampling strategy TDBS based on the truncated normal distribution, which improves the quality of the rendered depth and speeds up the optimization. Furthermore, our method proposes an inter-frame point cloud constraint. This allows us to successfully improve the accuracy of the pose estimation significantly and the quality of the novel view synthesis both indoors and outdoors with large motion changes. We implemented and evaluated our approach on different datasets, provided comparisons to other existing methods, and supported all claims made in this paper. Our experiments illustrate the generalizability of our proposed method on different datasets.

REFERENCES

- [1] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 4104–4113.
- [2] O. Özyeşil, V. Voroninski, R. Basri, and A. Singer, "A survey of structure from motion*," *Acta Numerica*, vol. 26, pp. 305–364, 2017.
- [3] T. Bailey and H. Durrant-Whyte, "Simultaneous localization and mapping (slam): Part ii," *IEEE Robotics and Automation Magazine*, vol. 13, no. 3, pp. 108–117, 2006.
- [4] G. Riegler and V. Koltun, "Free view synthesis," in *Proc. of the Europ. Conf. on Computer Vision (ECCV)*, 2020, pp. 623–640.
- [5] R. Tucker and N. Snavely, "Single-view view synthesis with multiple images," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 551–560.
- [6] R. B. Rusu and S. Cousins, "3d is here: Point cloud library (pcl)," in *Proc. of the IEEE Intl. Conf. on Robotics & Automation (ICRA)*, 2011, pp. 1–4.
- [7] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.
- [8] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-nerf 360: Unbounded anti-aliased neural radiance fields," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 5470–5479.
- [9] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Trans. on Graphics (TOG)*, vol. 41, no. 4, pp. 1–15, 2022.
- [10] B. Roessle, J. T. Barron, B. Mildenhall, P. P. Srinivasan, and M. Nießner, "Dense depth priors for neural radiance fields from sparse input views," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 12 892–12 901.
- [11] Y. Jeong, S. Ahn, C. Choy, A. Anandkumar, M. Cho, and J. Park, "Self-calibrating neural radiance fields," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021, pp. 5846–5854.
- [12] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm, "Pixel-wise view selection for unstructured multi-view stereo," in *Proc. of the Europ. Conf. on Computer Vision (ECCV)*, 2016.
- [13] Z. Wang, S. Wu, W. Xie, M. Chen, and V. A. Prisacariu, "Nerf: Neural radiance fields without known camera parameters," *arXiv preprint arXiv:2102.07064*, 2021.
- [14] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey, "Barf: Bundle-adjusting neural radiance fields," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021, pp. 5741–5751.
- [15] Q. Meng, A. Chen, H. Luo, M. Wu, H. Su, L. Xu, X. He, and J. Yu, "Gnerf: Gan-based neural radiance field without posed camera," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021, pp. 6351–6361.
- [16] W. Bian, Z. Wang, K. Li, J.-W. Bian, and V. A. Prisacariu, "Nope-nerf: Optimising neural radiance field with no pose prior," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 4160–4169.
- [17] K. Deng, A. Liu, J.-Y. Zhu, and D. Ramanan, "Depth-supervised nerf: Fewer views and faster training for free," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 12 882–12 891.
- [18] M. Pollefeys, R. Koch, and L. V. Gool, "Self-calibration and metric reconstruction inspite of varying and unknown intrinsic camera parameters," *Intl. Journal of Computer Vision (IJCV)*, vol. 32, no. 1, pp. 7–25, 1999.
- [19] C. Stachniss, J. Leonard, and S. Thrun, *Springer Handbook of Robotics, 2nd edition*. Springer Verlag, 2016, ch. Chapt. 46: Simultaneous Localization and Mapping.
- [20] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Trans. on Robotics (TRO)*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [21] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, "Monoslam: Real-time single camera slam," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 29, no. 6, pp. 1052–1067, 2007.
- [22] H. Alismail, B. Browning, and S. Lucey, "Photometric bundle adjustment for vision-based slam," in *Proc. of the Asian Conf. on Computer Vision (ACCV)*, 2017, pp. 324–341.
- [23] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-slam: Neural implicit scalable encoding for slam," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 12 786–12 796.
- [24] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, "imap: Implicit mapping and positioning in real-time," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021, pp. 6229–6238.
- [25] L. Yen-Chen, P. Florence, J. T. Barron, A. Rodriguez, P. Isola, and T.-Y. Lin, "Inerf: Inverting neural radiance fields for pose estimation," in *Proc. of the IEEE/RISJ Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2021, pp. 1323–1330.
- [26] S.-F. Chng, S. Ramasinghe, J. Sherrah, and S. Lucey, "Garf: gaussian activated radiance fields for high fidelity reconstruction and pose estimation," *arXiv preprint arXiv:2204.05735*, 2022.
- [27] A. Creswell, T. White, V. Dumoulin, K. Arulkumaran, B. Sengupta, and A. A. Bharath, "Generative adversarial networks: An overview," *IEEE Signal Processing Magazine*, vol. 35, no. 1, pp. 53–65, 2018.
- [28] Z. Cheng, C. Esteves, V. Jampani, A. Kar, S. Maji, and A. Makadia, "Lu-nerf: Scene and pose estimation by synchronizing local unposed nerfs," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2023.
- [29] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," in *Proc. of the IEEE/CVF Intl. Conf. on Computer Vision (ICCV)*, 2021, pp. 12 179–12 188.
- [30] J. Burkardt, "The truncated normal distribution," *Department of Scientific Computing Website, Florida State University*, vol. 1, p. 35, 2014.
- [31] K. Mosegaard and A. Tarantola, "Monte carlo sampling of solutions to inverse problems," *Journal of Geophysical Research: Solid Earth*, vol. 100, no. B7, pp. 12 431–12 447, 1995.
- [32] M. A. Butt and P. Maragos, "Optimum design of chamfer distance transforms," *IEEE Trans. on Image Processing*, vol. 7, no. 10, pp. 1477–1484, 1998.
- [33] S. S. Keerthi and C.-J. Lin, "Asymptotic behaviors of support vector machines with gaussian kernel," *Neural computation*, vol. 15, no. 7, pp. 1667–1689, 2003.
- [34] B. Mildenhall, P. P. Srinivasan, R. Ortiz-Cayon, N. K. Kalantari, R. Ramamoorthi, R. Ng, and A. Kar, "Local light field fusion: Practical view synthesis with prescriptive sampling guidelines," *ACM Trans. on Graphics (TOG)*, vol. 38, no. 4, pp. 1–14, 2019.
- [35] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and temples: Benchmarking large-scale scene reconstruction," *ACM Trans. on Graphics (TOG)*, vol. 36, no. 4, pp. 1–13, 2017.
- [36] D. Nistér, O. Naroditsky, and J. Bergen, "Visual odometry," in *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2004, pp. 1–1.