# Correspondence-Guided SfM-Free 3D Gaussian Splatting for NVS

**Wei Sun[1], Xiaosong Zhang[2], Fang Wan[1], Yanzhao Zhou[1], Yuan Li[1*], Qixiang Ye[1], Jianbin Jiao[1]**

[1]School of Electronic, Electrical and Communication Engineering, University of Chinese Academy of Sciences, China
[2]Beijing Academy of Artificial Intelligence, China
sunwei162@mails.ucas.ac.cn, xszhang@baai.ac.cn, {wanfang, zhouyanzhao, liyuan23, qxye, jiaojb}@ucas.ac.cn

## Abstract

Novel View Synthesis (NVS) without Structure-from-Motion (SfM) pre-processed camera poses—referred to as SfM-free methods—is crucial for promoting rapid response capabilities and enhancing robustness against variable operating conditions.

Recent SfM-free methods have integrated pose optimization, designing end-to-end frameworks for joint camera pose estimation and NVS. However, most existing works rely on per-pixel image loss functions, such as L2 loss. In SfM-free methods, inaccurate initial poses lead to misalignment issue, which, under the constraints of per-pixel image loss functions, results in excessive gradients, causing unstable optimization and poor convergence for NVS.

In this study, we propose a correspondence-guided SfM-free 3D Gaussian splatting for NVS. We use correspondences between the target and the rendered result to achieve better pixel alignment, facilitating the optimization of relative poses between frames. We then apply the learned poses to optimize the entire scene. Each 2D screen-space pixel is associated with its corresponding 3D Gaussians through approximated surface rendering to facilitate gradient back-propagation. Experimental results underline the superior performance and time efficiency of the proposed approach compared to the state-of-the-art baselines.

## Introduction

Novel-view synthesis serves as a fundamental objective within the realm of computer vision. The recent surge in NVS popularity is largely attributable to the success of Neural Radiance Fields (NeRFs) (Mildenhall et al. 2021) and 3D Gaussian Splatting (3DGS) (Kerbl et al. 2023). However, these methods require densely captured views with accurately labeled camera poses, which is often not feasible in practical scenarios. Often, camera poses are obtained from SfM methods like COLMAP (Schonberger and Frahm 2016) as a pre-processing step to NeRF and 3DGS, which is not only time-consuming but also prone to fail due to its sensitivity to feature extraction errors and difficulties in handling textureless or repetitive regions.

Recent studies (Bian et al. 2023; Lin et al. 2021; Wang et al. 2021; Fu et al. 2024; Jiang et al. 2024) have focused on reducing the reliance on SfM by integrating pose estimation directly within the NVS framework. However, we would like to note that existing approaches typically rely on per-pixel image loss functions (such as L2 loss) from a pair of RGB images and compute per-pixel color derivatives with respect to desired scene parameters. The rendered result and the target do not sufficiently overlap or align because the camera pose is inaccurate at the initial stage of optimization. This problem is exacerbated when there is significant camera movement between consecutive views, at which point achieving perfect per-pixel alignment between the rendered result and the target becomes even more challenging. This misalignment issue, under the constraints of per-pixel image loss, often results in excessive gradients, leading to instability in the optimization process and difficulty in convergence.

To address this problem, we introduce a **C**orrespondence-**G**uided SfM-free **3D G**aussian **S**platting for NVS (CG-3DGS), a novel approach that integrates 2D correspondence detection (Sun et al. 2021; Tang et al. 2022), and computes derivatives on associated points instead of on a fixed grid of pixels. Specifically, we detect the 2d correspondence to find a pixel matching between rendered and target images and design a novel loss function based on the pixel matching. We then develop an approximated surface rendering pipeline for 3D Gaussians, which propagates disturbances from the 2D screen space to the parameters of the 3D Gaussians for differentiable scene optimization. Our derivatives are dense and could account for long-range object motions through the correspondence-based loss function, naturally leading to better robustness in optimization.

Inspired by but fundamentally distinct from CF-3DGS (Fu et al. 2024), we construct a two-step optimization pipeline: (i) We initialize an auxiliary 3D Gaussian set given frame t with depth back-projection, and we sample the next nearby frame t+1. Our goal is to learn an affine transformation that can transform the 3D Gaussians in frame t to render the pixels in frame t+1. Correspondence-based loss function provides the gradients for optimizing the affine transformation, which is essentially optimizing the relative camera pose between frames t − 1 and t. This process continues iteratively until we obtain all the relative poses between frames 0, 1, ..., t. (ii). We initialize another 3D Gaussians set, where we perform scene optimization with all the frames and their corresponding learned camera poses.

---
*Corresponding author.

This paper primarily contributes the following:

- We introduce the correspondence-guided SfM-free 3D Gaussian Splatting for NVS, minimizing the impact of pixel misalignment.
- We integrate the 3DGS framework with effective correspondence supervision without time-consuming surface rendering, offering a unified differentiable pipeline for NVS without SfM pre-processing.
- Our method boosts time efficiency, and delivers superior results compared to the state-of-the-art methods.

## Related Work

### Novel View Synthesis

Various 3D scene representations are utilized to produce realistic images from new viewpoints, including planes (Horry, Anjyo, and Arai 1997; Hoiem, Efros, and Hebert 2005), meshes (Hu et al. 2020; Riegler and Koltun 2020, 2021), point clouds (Xu et al. 2022; Zhang et al. 2022), and multi-plane images (Tucker and Snavely 2020; Zhou et al. 2018; Li et al. 2021). NeRFs (Mildenhall et al. 2021) have recently become prominent for their superior photorealistic rendering capabilities, with numerous enhancements such as better anti-aliasing (Barron et al. 2021, 2022, 2023; Zhang et al. 2020) and improved reflectance (Verbin et al. 2022; Attal et al. 2023).

More recently, the use of point-cloud-based representations has surged due to its rendering efficiency (Xu et al. 2022; Zhang et al. 2022; Kerbl et al. 2023; Luiten et al. 2023; Kopanas et al. 2022; Yifan et al. 2019). For example, Zhang (Zhang et al. 2022) introduce a method to learn the per-point position and view-dependent appearance through a differentiable splat-based renderer initialized from object masks. Furthermore, 3DGS (Kerbl et al. 2023) facilitates real-time rendering of novel views using its explicit representation combined with a novel differential point-based splatting technique. Nevertheless, these methods typically depend on pre-computed camera parameters derived from SfM techniques (Hartley and Zisserman 2003; Schonberger and Frahm 2016; Mur-Artal, Montiel, and Tardos 2015; Taketomi, Uchiyama, and Ikeda 2017).

### SfM-Free Modeling for Novel View Synthesis

Initial efforts in SfM-free novel view synthesis include iNeRF (Yen-Chen et al. 2021), which employs key-point matching to estimate camera poses. NeRFmm (Wang et al. 2021) introduces a method for joint optimization of camera pose and NeRF itself. Techniques such as those proposed in BARF (Lin et al. 2021) focus on learning neural 3D representations and aligning camera frames using hierarchical positional encodings. The approach in Nope-NeRF (Bian et al. 2023) incorporates monocular depth priors to simultaneously capture relative poses and synthesize new views. (Meuleman et al. 2023) uses a combination of pre-trained depth and optical-flow priors to refine blockwise NeRFs, which helps in the sequential recovery of camera poses.

In more generalizable settings, methods like SRT (Sajjadi et al. 2022), VideoAE (Lai et al. 2021), RUST (Sajjadi et al. 2023), MonoNeRF (Tian, Du, and Duan 2023), DBARF (Chen and Lee 2023), and FlowCam (Smith et al. 2023) aim to learn a scene representation from unposed videos using the implicit framework of NeRF. Despite these efforts, they often fail to achieve satisfactory view synthesis without specific scene optimization and share NeRF's original limitations, such as the inability to render explicit primitives in real time.

The inherent complexity of NeRF's implicit modeling often complicates the simultaneous optimization of scene and camera poses. Recent advancements like 3DGS, with its explicit point-based scene representation, facilitate real-time rendering and efficient optimization. New developments, such as those in CF-3DGS (Fu et al. 2024), continue to push the limits of simultaneous scene and pose optimization, CF-3DGS (Fu et al. 2024) employs a progressive training strategy to reduce the cumulative noise associated with the pose optimization process. However, these methods consistently rely on per-pixel image loss, which always results in excessive gradients and unstable optimization when the optimization target and rendering result deviate significantly from perfect per-pixel alignment. This is a common issue in SfM-free scenarios due to the inaccurate initial pose estimation, which leads us to explore the integration of correspondence in the simultaneous scene and pose optimization.

## Method

In this paper, we leverage 3D Gaussians to reconstruct photo-realistic scenes from sequential frames of a video stream. Given a sequence of unposed images $\{I_1, \ldots, I_K\}$ with camera intrinsics, our goal is to better reconstruct the complete scene via a joint optimization of the camera poses and the 3D representation (i.e. 3D Gaussians). We detail our method in the following sections, starting from a brief review of the representation and rendering process of 3D Gaussians (Sec 3.1). Then, we propose a correspondence-guided pose optimization, a simple yet effective method to estimate the relative camera pose from each pair of nearby frames (Sec 3.2). Finally, we briefly introduce how to reconstruct scenes using the estimated poses (Sec 3.3).

### Revisiting 3D Gaussian Splatting

3D Gaussian Splatting (3DGS) is a point-based novel view synthesis technique that uses 3D Gaussians to model the scene. The Gaussian attributes are optimized based on a set of input training views denoted by ground truth images $\mathcal{I}_{gt} = \{I_i \in \mathbb{R}^{H \times W}\}_{i=1}^{K}$ and associated camera poses $\mathcal{P}_{gt} = \{W_i \in \mathbb{R}^{3 \times 4}\}_{i=1}^{K}$. The Gaussian initialization is derived from a sparse point cloud created via SfM across the training views. To increase the number of Gaussians in areas where small-scale geometry is insufficiently reconstructed, a Gaussian densification process is periodically applied during training.

Each Gaussian $G_i$ in the scene is described by several parameters: its position $x_i \in \mathbb{R}^3$, scale $s_i \in \mathbb{R}^3$, rotation $r_i \in \mathbb{R}^4$, base color $c_i \in \mathbb{R}^3$, view-dependent spherical har-
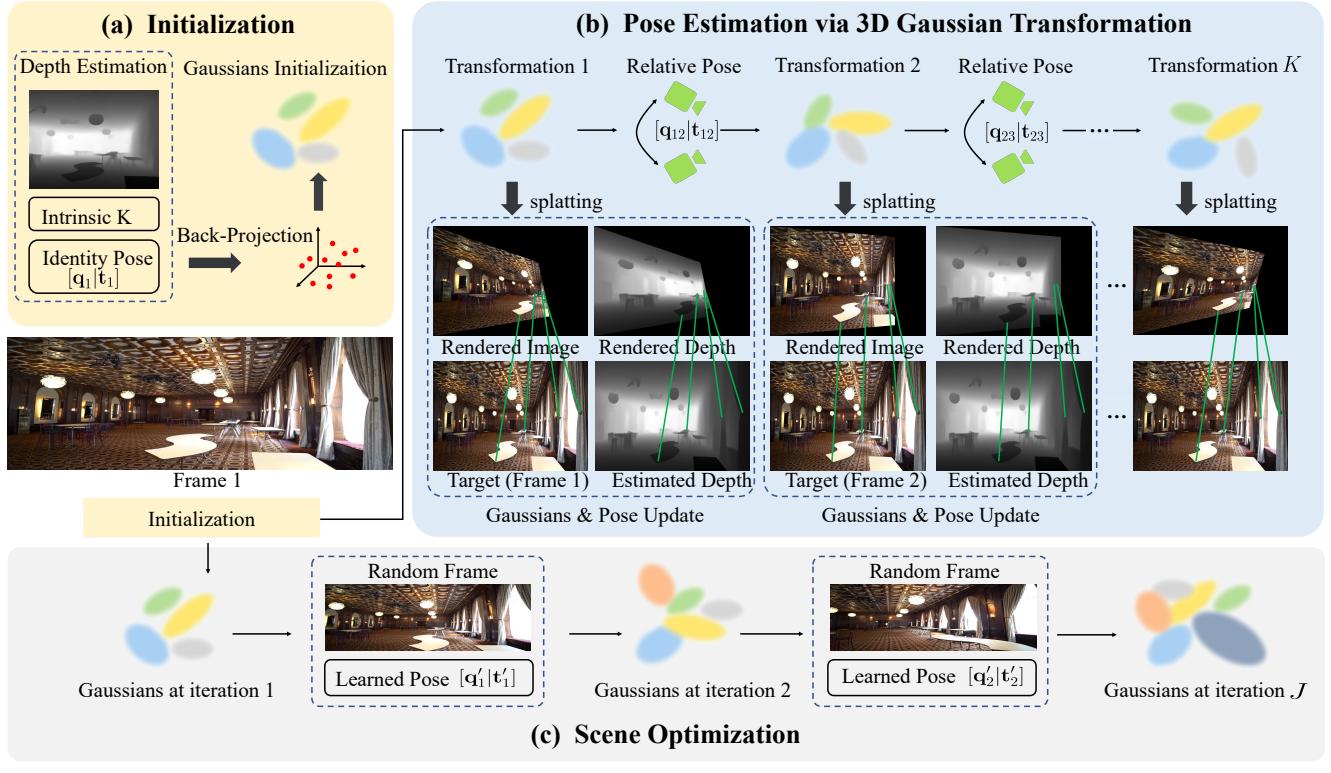
Figure 1: Overview of our CG-3DGS. (a) We utilize camera intrinsics and the identity pose to back-project depth estimate into a point cloud, initializing a set of 3D Gaussians. (b) These 3D Gaussians are used to simulate camera pose changes between adjacent frames through SE-3 transformations. First, we update the parameters of $G_t$ based on the rendering results on frame $t$, and use the SE-3 transformed $G_t$ as $G_{t+1}$ to render frame $t+1$. At this point, we freeze $G_t$ and only update the parameters of the SE-3 transformation. This iterative process continues until the relative poses between all adjacent frames in a video sequence are estimated. The optimization is based on the correspondence between the rendered result and the ground truth. (c) After pose estimation, the same point cloud is also used to initialize a set of 3D Gaussians used for rendering the scene, and frames with estimated poses are randomly sampled for conventional training of the 3D Gaussians.

monics $h_i \in \mathbb{R}^{15 \times 3}$, and opacity $\alpha_i \in \mathbb{R}$. Collectively, these parameters are grouped as:

$$\mathcal{G} = \{G_i = \{x_i, s_i, r_i, c_i, h_i, \alpha_i\}\}_{i=1}^{N}, \qquad (1)$$

where $N$ denotes the total count of Gaussians.

During synthesis, scaling and rotation parameters are translated into matrices $S_i$ and $R_i$. The Gaussian $G_i$ is spatially characterized in the 3D scene by its center point, or mean position, $x_i$ and a decomposable covariance matrix $\Sigma_i$:

$$G_i(x_i) = e^{-\frac{1}{2} x_i^T \Sigma_i^{-1} x_i}, \Sigma_i = R_i S_i S_i^T R_i^T. \qquad (2)$$

To facilitate the differentiation of 3D Gaussian rendering, the Gaussian projection process is applied from a specific camera pose $W$, approximating the splatting of a 3D Gaussian onto the 2D image plane:

$$\Sigma^{2D} = J W \Sigma W^\top J^\top \qquad (3)$$

where $J$ represents the Jacobian of the affine approximation of the projective transformation.

For each pixel, the final rendered color and depth can be formulated as the alpha-blending of $N$ ordered Gaussians that overlap the pixel:

$$\hat{C} = \sum_i^N \mathbf{c}_i \alpha_i \prod_j^{i-1}(1 - \alpha_j),$$
$$\hat{D} = \sum_i^N \mathbf{d}_i \alpha_i \prod_j^{i-1}(1 - \alpha_j), \qquad (4)$$

with $\mathbf{c}_i$, $\alpha_i$, and $\mathbf{d}_i$ representing the color, opacity, and depth of the Gaussians, respectively.

The optimization of the 3DGS model relies on minimizing a composite loss function using stochastic gradient descent:

$$L(G|W, I) = ||\hat{I} - I||_1 + L_{SSIM}(\hat{I}, I), \qquad (5)$$

where $\hat{I}$ is the rendering result and $I$ is the ground truth image. The overall loss combines L1 loss for residual minimization and SSIM loss for structural similarity.

### Correspondence-guided Pose Optimization

**Initialization from Monocular Depth.** As shown in Fig. 1 (a), for the initial frame $I_1$, which is at timestep 1,

we apply a standard monocular depth network to produce a depth map, represented as $D_1$. We then construct the point cloud $P$ by back-projecting the depth map $D_1$ using the default identity camera pose (orthogonal projection) and camera intrinsics, and use $P$ to initialize 3D Gaussians instead of relying on SfM-derived points. Following this initialization, we optimize a set of 3D Gaussians $G_1$, adjusting all attributes to reduce the correspondence-based loss between the rendered image and the ground truth $I_1$,

$$G_1{}^\star = \arg \min_{c_1, r_1, s_1, \alpha_1} \mathcal{L}_{cor}(\mathcal{R}(G_1), I_1), \qquad (6)$$

where $\mathcal{R}$ denotes the rendering operation of 3DGS. The correspondence-based loss $\mathcal{L}_{cor}$ is detailed in Sec. 3.2.3.

**Pose Estimation via 3D Gaussians Transformation.** The problem of camera pose estimation is addressed by predicting the transformation of 3D Gaussians as discussed in CF-3DGS (Fu et al. 2024). Starting with the Gaussian center's position $\mu$, we project it into the 2D camera plane with camera pose $W$ as $\mu_{2D} = K \frac{W\mu}{(W\mu)_z}$. Hence, estimating the camera pose effectively involves determining the transformations of these 3D Gaussians.

For the relative camera pose estimation, we apply a learnable SE-3 affine transformation $T_t$ to the pretrained 3D Gaussians $G_t{}^*$, transforming it into frame $t + 1$, represented as $G_{t+1} = T_t \odot G_t$. This transformation $T_t$ is refined by minimizing the photometric loss between the rendered images and the subsequent frame $I_{t+1}$:

$$T_t{}^* = \arg \min_{T_t} \mathcal{L}_{cor}(\mathcal{R}(T_t \odot G_t), I_{t+1}), \qquad (7)$$

During this optimization phase, we preserve the attributes of the pretrained 3D Gaussians $G_t{}^*$ unchanged to distinctly separate the effects of camera motion from changes such as deformation, densification, pruning, or self-rotation of the Gaussians. The transformation matrix $T$, comprising quaternion rotations $\mathbf{q} \in \mathfrak{so}(3)$ and translation vectors $\mathbf{t} \in \mathbb{R}^3$, facilitates the estimation of relative camera poses between consecutive frames. After this, we have estimated the relative camera pose between frames $I_t$ and $I_{t+1}$. As the next frame $I_{t+2}$ becomes available, this process is repeated: we optimize the 3D Gaussians to obtain $G_{t+1}^*$, similar to what is described at the end of Sec. 3.2.1; we then optimize the relative pose between $I_{t+1}$ and $I_{t+2}$, and could subsequently infer the relative pose between $I_t$ and $I_{t+2}$.

**Correspondence-based Loss** We utilize off-the-shelf detectors (Tang et al. 2022; Sun et al. 2021) to establish the 2D correspondences between ground truth image $I$ and rendered result $\hat{I}(W)$ for the pose optimization. The 2D screen-space coordinates in $I$ are represented as $\mathcal{K} = \{\mathbf{k}^{(1)}, \mathbf{k}^{(2)}, ..., \mathbf{k}^{(M)}\}$, where $M$ represents the total number of points. Correspondingly, the 2D screen-space coordinates in $\hat{I}(W)$ are $\mathcal{K}' = \{\mathbf{k}'^{(1)}, \mathbf{k}'^{(2)}, ..., \mathbf{k}'^{(M)}\}$. The optimization objective is to align $\mathbf{k}$ with $\mathbf{k}'$, visualized in Fig. 1 (b).

To enable gradient back-propagation from the matching of $\mathbf{k}$ and $\mathbf{k}'$ to the 3D Gaussians shaping the surface, we employ a differentiable approximate surface renderer, described in Sec. 3.2.4, to render the screen-space coordinates
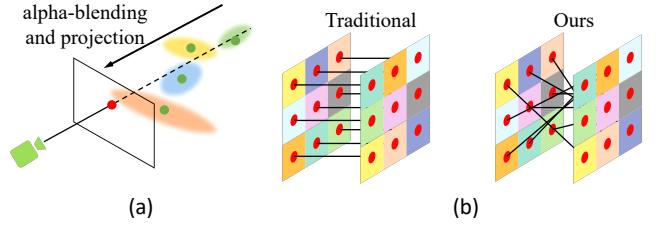


Figure 2: (a) By performing alpha-blending on the center coordinates of the 3D Gaussians, an approximate 3D surface point is generated and projected onto the 2D screen. (b) The comparison between traditional methods and our method. The fundamental difference is the technique used for aligning pixels.

at $\mathbf{k}'^{(i)}, i = 1, 2, ..., M$ as $q(\mathbf{k}'^{(i)})$. The resulting loss function is expressed as:

$$\mathcal{L}_{\text{cor-rgb}} = \sum_{i=1}^{M} ||q(\mathbf{k}'^{(i)}) - \mathbf{k}^{(i)}||_1. \qquad (8)$$

Notably, $q(\mathbf{k}'^{(i)})$ numerically matches $\mathbf{k}'^{(i)}$, yet it creates a pathway for gradients to flow back to the underlying 3D representation without altering the original 3DGS.

Moreover, incorporating short-range relations through pixel-wise supervision can assist in stabilizing the optimization process. This loss is formulated as:

$$\mathcal{L}_{\text{pix-rgb}} = ||I - \hat{I}(W)||_1. \qquad (9)$$

Furthermore, the depth matching process involves equating the monocular depth at $\mathbf{k}$, denoted as $d(\mathbf{k})$, with the rendered depth at $\mathbf{k}'$, denoted as $\hat{d}(\mathbf{k}')$. The corresponding loss term is defined as:

$$\mathcal{L}_{\text{cor-depth}} = \sum_{i=1}^{M} ||\hat{d}(\mathbf{k}'^{(i)}) - d(\mathbf{k}^{(i)})||_1. \qquad (10)$$

The correspondence-based loss consolidates these components:

$$\mathcal{L}_{\text{cor}} = \lambda_1 \mathcal{L}_{\text{cor-rgb}} + \lambda_2 \mathcal{L}_{\text{pix-rgb}} + \lambda_3 \mathcal{L}_{\text{cor-depth}}, \qquad (11)$$

where $\lambda_1 = 10, \lambda_2 = 1, \lambda_3 = 1$.

**Approximated Surface Rendering** Our aim in correspondence-based optimization is to transmit gradient information from a 2D screen-space location to its associated 3D surface location. Essentially, we seek to link disturbances at a 2D screen-space location with those at its 3D surface counterpart.

Given the volumetric nature of the 3D Gaussian representation, explicit surfaces are not present. However, reconstructing an explicit surface is extremely time-consuming (Park et al. 2019), and modifying the rendering logic of the 3D Gaussian to obtain surfaces would also significantly increase the training duration (Jiang et al. 2024). Instead, according to previous studies (Keetha et al. 2024; Chung, Oh, and Lee 2024; Fu et al. 2024; Yan et al. 2024),

| scenes | Ours | | | CF-3DGS | | | Nope-NeRF | | | BARF | | | NeRFmm | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS |
| Church | **32.14** | **0.96** | **0.08** | 30.23 | 0.93 | 0.11 | 25.17 | 0.73 | 0.39 | 23.17 | 0.62 | 0.52 | 21.64 | 0.58 | 0.54 |
| Barn | **33.19** | **0.94** | **0.07** | 31.23 | 0.90 | 0.10 | 26.35 | 0.69 | 0.44 | 25.28 | 0.64 | 0.48 | 23.21 | 0.61 | 0.53 |
| Museum | **31.62** | **0.94** | **0.08** | 29.91 | 0.91 | 0.11 | 26.77 | 0.76 | 0.35 | 23.58 | 0.61 | 0.55 | 22.37 | 0.61 | 0.53 |
| Family | **34.80** | **0.97** | **0.04** | 31.27 | 0.94 | 0.07 | 26.01 | 0.74 | 0.41 | 23.04 | 0.61 | 0.56 | 23.04 | 0.58 | 0.56 |
| Horse | **35.45** | **0.97** | **0.04** | 33.94 | 0.96 | 0.05 | 27.64 | 0.84 | 0.26 | 24.09 | 0.72 | 0.41 | 23.12 | 0.70 | 0.43 |
| Ballroom | **33.91** | **0.97** | **0.04** | 32.47 | 0.96 | 0.07 | 25.33 | 0.72 | 0.38 | 20.66 | 0.50 | 0.60 | 20.03 | 0.48 | 0.57 |
| Francis | **33.80** | **0.92** | **0.13** | 32.72 | 0.91 | 0.14 | 29.48 | 0.80 | 0.38 | 25.85 | 0.69 | 0.57 | 25.40 | 00.69 | 0.52 |
| Ignatius | **31.14** | **0.94** | **0.06** | 28.43 | 0.90 | 0.09 | 23.96 | 0.61 | 0.47 | 21.78 | 0.47 | 0.60 | 21.16 | 0.45 | 0.60 |
| mean | **33.26** | **0.95** | **0.07** | 31.28 | 0.93 | 0.09 | 26.34 | 0.74 | 0.39 | 23.42 | 0.61 | 0.54 | 22.50 | 0.59 | 0.54 |

Table 1: **Novel view synthesis results on Tanks and Temples**. Each baseline method is trained with its public code under the original settings and evaluated with the same evaluation protocol. The best results are highlighted in bold.

the depth of an expected 3D surface point $D(\mathbf{k})$ relative to a 2D screen-space point $\mathbf{k}$ is computed as follows:

$$D(\mathbf{k}) = \sum_i d_i \alpha_i(\mathbf{k}) \prod_{j=1}^{i-1} (1 - \alpha_j(\mathbf{k})), \qquad (12)$$

where $d_i$ indicates the $z$-axis position of the Gaussian centers within the camera coordinate system, and $\alpha_i$ and $\alpha_j$ represent the alpha-blending coefficients for the $i^{\text{th}}$ and $j^{\text{th}}$ Gaussian, respectively.

As illustrated in Fig. 2, the corresponding expected 3D surface point $\Psi(\mathbf{k})$ at $\mathbf{k}$ could then be defined by:

$$\Psi(\mathbf{k}) = \sum_i \boldsymbol{\mu}_i \alpha_i(\mathbf{k}) \prod_{j=1}^{i-1} (1 - \alpha_j(\mathbf{k})), \qquad (13)$$

where $\boldsymbol{\mu}_i \in \mathbb{R}^3$ represents the center position of the $i^{\text{th}}$ Gaussian. This formula provides an approximation of the 3D surface point without relying on time-consuming surface reconstruction method like signal distance function (SDF) (Park et al. 2019).

**Correspondence Cache Mechanism.** In our pose optimization process, correspondences between the rendered views and the reference views are stored in a cache. By reusing these correspondences in subsequent iterations, we achieve a significant reduction in computation time without significantly degrading performance, as adjacent frames in continuous video tend to exhibit similar features and poses. Concretely, rather than recalculating correspondence points for each image pair in every iteration, we strategically update these correspondences every $H$ iterations—where $H$ is empirically set to 50.

### Scene Optimization

Following the camera pose optimization, we proceed to optimize a new set of 3D Gaussians that ultimately represent the scene. Similarly to the pose optimization phase, we start by generating a set of initialized 3D Gaussians using the depth estimation from the first frame $I_1$. Here, we keep the camera pose fixed and focus solely on minimizing the photometric loss as in the vanilla 3DGS (Kerbl et al. 2023). During the optimization, we randomly sample frames from the training set of the scene and utilize the associated optimized poses for training, as shown in Fig. 1 (c).

## Experiments

### Experimental Setup

**Datasets.** We conduct extensive experiments on several real-world datasets, including Tanks and Temples (Knapitsch et al. 2017) and CO3D-V2 (Reizenstein et al. 2021). **Tanks and Temples:** Following the methodology in Nope-NeRF (Bian et al. 2023), we assess the quality of novel view synthesis and the accuracy of pose estimation across eight diverse scenes that encompass both indoor and outdoor environments. We select seven images from each 8-frame sequence for training and evaluate the novel view synthesis on the remaining frame. Camera poses are estimated and assessed on all training images following alignment according to Umeyama's method (Umeyama 1991). **CO3D-V2:** This dataset comprises thousands of object-centric videos, maintaining view of the full object while the camera moves in a complete circle around it. Deriving camera poses from CO3D videos is more challenging compared to Tanks and Temples due to the large and complex camera movements. We randomly select four scenes from different object categories and follow the same protocol as CF-3DGS (Fu et al. 2024) to divide the training and testing sets.

**Metrics.** We assess the performance of novel view synthesis and camera pose estimation tasks. For the former, we evaluate using standard metrics such as Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index (SSIM) (Wang et al. 2004), and Learned Perceptual Image Patch Similarity (LPIPS) (Zhang et al. 2018). For the latter, we utilize established visual odometry metrics such as Absolute Trajectory Error (ATE) and Relative Pose Error (RPE).

**Implementation Details.** Our implementation leverages the PyTorch framework (Paszke et al. 2017) and adheres to the optimization parameters specified in 3DGS (Kerbl et al. 2023), unless noted otherwise. Importantly, we continuously adjust the opacity throughout the training process to effectively limit the unchecked growth of Gaussian components caused by inaccuracies in pose estimation. For the Tanks and Temples and CO3D V2 datasets, the off-the-shelf monocular depth networks used are DPT (Ranftl, Bochkovskiy, and Koltun 2021) and Zoe (Bhat et al. 2023), respectively. The initial learning rate is set at $10^{-5}$ and is progressively reduced to $10^{-6}$ until the model converges. All experiments are performed on a single RTX 3090 GPU.

| scenes | Ours | | | CF-3DGS | | | Nope-NeRF | | | BARF | | | NeRFmm | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $RPE_t \downarrow$ | $RPE_r \downarrow$ | $ATE \downarrow$ | $RPE_t$ | $RPE_r$ | $ATE$ | $RPE_t$ | $RPE_r$ | $ATE$ | $RPE_t$ | $RPE_r$ | $ATE$ | $RPE_t$ | $RPE_r$ | $ATE$ |
| Church | 0.006 | 0.016 | 0.001 | 0.008 | 0.018 | 0.002 | 0.034 | 0.008 | 0.008 | 0.114 | 0.038 | 0.052 | 0.626 | 0.127 | 0.065 |
| Barn | 0.029 | 0.030 | 0.002 | 0.034 | 0.034 | 0.003 | 0.046 | 0.032 | 0.004 | 0.314 | 0.265 | 0.050 | 1.629 | 0.494 | 0.159 |
| Museum | 0.047 | 0.203 | 0.004 | 0.052 | 0.215 | 0.005 | 0.207 | 0.202 | 0.020 | 3.442 | 1.128 | 0.263 | 4.134 | 1.051 | 0.346 |
| Family | 0.024 | 0.020 | 0.001 | 0.022 | 0.024 | 0.002 | 0.047 | 0.015 | 0.001 | 1.371 | 0.591 | 0.115 | 2.743 | 0.537 | 0.120 |
| Horse | 0.109 | 0.053 | 0.003 | 0.112 | 0.057 | 0.003 | 0.179 | 0.017 | 0.003 | 1.333 | 0.394 | 0.014 | 1.349 | 0.434 | 0.018 |
| Ballroom | 0.033 | 0.020 | 0.003 | 0.037 | 0.024 | 0.003 | 0.041 | 0.018 | 0.002 | 0.531 | 0.228 | 0.018 | 0.449 | 0.177 | 0.031 |
| Francis | 0.026 | 0.147 | 0.005 | 0.029 | 0.154 | 0.006 | 0.057 | 0.009 | 0.005 | 1.321 | 0.558 | 0.082 | 1.647 | 0.618 | 0.207 |
| Ignatius | 0.027 | 0.012 | 0.003 | 0.033 | 0.032 | 0.005 | 0.026 | 0.005 | 0.002 | 0.736 | 0.324 | 0.029 | 1.302 | 0.379 | 0.041 |
| mean | **0.037** | 0.063 | **0.003** | 0.041 | 0.069 | 0.004 | 0.080 | **0.038** | 0.006 | 1.046 | 0.441 | 0.078 | 1.735 | 0.477 | 0.123 |

Table 2: **Pose accuracy on Tanks and Temples**. Note that we use COLMAP poses in Tanks and Temples as the "ground truth". The unit of $RPE_r$ is in degrees, ATE is in the ground truth scale and $RPE_t$ is scaled by 100. The best results of means are highlighted in bold.
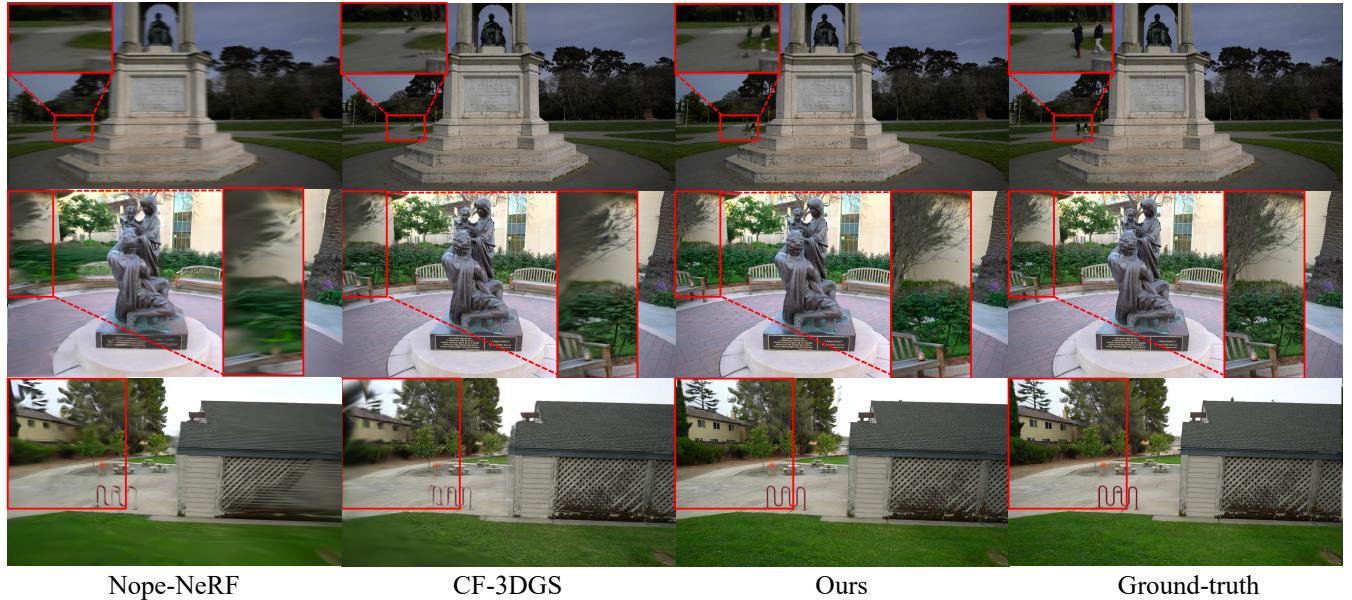


Figure 3: **Qualitative comparison for novel view synthesis on Tanks and Temples.** Our approach produces more realistic rendering results than other baselines. **Better viewed when zoomed in.**

## Comparing with SfM-Free Methods

In this subsection, we compare our method with several baselines including CF-3DGS (Fu et al. 2024), Nope-NeRF (Bian et al. 2023), BARF (Lin et al. 2021) and NeRFmm (Wang et al. 2021) on both novel view synthesis and camera pose estimation.

**Novel View Synthesis.** In contrast to conventional approaches where camera poses for testing views are provided, we need to additionally ascertain the camera poses of test views. We utilize the same protocol as outlined in CF-3DGS to optimize the camera poses for these testing views. This identical procedure is applied across all baseline models to maintain a consistent basis for comparison.

We present the comparative analysis on the Tanks and Temples dataset in Table 1. Our approach consistently surpasses competing methods across all evaluated metrics. Remarkably, our direct training strategy achieves superior PSNR values even compared to CF-3DGS, which leverages carefully designed progressive 3D Gaussians training strategy, with a notable increase of 3.5 points in the Family scene.

Qualitative results are shown in Fig. 3. The images generated using our method are distinctly sharper and could retain small objects within the scene, such as the walking person in the first scene shown in Fig. 3, which correlates with the significantly improved scores for SSIM and LPIPS, as detailed in Table 1.

**Camera Pose Estimation.** The learned camera poses are post-processed by Procrustes analysis as in (Lin et al. 2021; Bian et al. 2023) and compared with the ground-truth poses of training views. The quantitative results of camera pose estimation are summarized in Table 2. Our approach achieves comparable performance with the current state-of-the-art results. We hypothesize that the relatively poorer performance in terms of $RPE_r$ may be attributed to relying solely on photometric loss for relative pose estimation. In contrast, Nope-NeRF incorporates additional constraints on relative poses beyond photometric loss, including the chamfer distance be-

| Method | Times ↓ | 46_2587_7531 | | | 407_54965_106262 | | | 429_60388_117059 | | | 437_62482_122880 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR ↑ | SSIM ↑ | LPIPS ↓ | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS |
| Nope-NeRF | ∼30 h | 25.3 | 0.73 | 0.46 | 25.53 | 0.83 | 0.58 | 22.19 | 0.62 | 0.56 | 20.81 | 0.59 | 0.51 |
| CF-3DGS | ∼2 h | 25.44 | 0.80 | 0.21 | 27.80 | 0.84 | 0.35 | 24.44 | 0.68 | 0.36 | 22.95 | 0.66 | 0.41 |
| Ours | **∼1.5 h** | **26.43** | **0.85** | **0.15** | **28.46** | **0.88** | **0.27** | **25.72** | **0.74** | **0.29** | **24.32** | **0.69** | **0.32** |

Table 3: **Novel view synthesis results on CO3D V2**. Each baseline method is trained with its public code under the original settings and evaluated with the same evaluation protocol. The best results are highlighted in bold.

| Method | Times ↓ | 46_2587_7531 | | | 407_54965_106262 | | | 429_60388_117059 | | | 437_62482_122880 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | RPE$_t$ ↓ | RPE$_r$ ↓ | ATE ↓ | RPE$_t$ | RPE$_r$ | ATE | RPE$_t$ | RPE$_r$ | ATE | RPE$_t$ | RPE$_r$ | ATE |
| Nope-NeRF | ∼30 h | 0.426 | 4.226 | 0.023 | 0.553 | 4.685 | 0.057 | 0.398 | 2.914 | 0.055 | 0.591 | 2.014 | 0.041 |
| CF-3DGS | ∼2 h | 0.095 | 0.447 | 0.009 | 0.31 | 0.243 | 0.008 | 0.134 | 0.542 | 0.018 | 0.252 | 0.493 | 0.018 |
| Ours | **∼1.5 h** | **0.041** | **0.274** | **0.005** | **0.14** | **0.182** | **0.003** | **0.092** | **0.239** | **0.008** | **0.116** | **0.284** | **0.009** |

Table 4: **Pose accuracy on CO3D V2**. Note that the camera poses provided by CO3D as the "ground truth". The unit of RPE$_r$ is in degrees, ATE is in the ground truth scale and RPE$_t$ is scaled by 100. The best results are highlighted in bold.

| scenes | w.o. correspondence | | | | Ours | | | |
|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | RPE$_t$ | RPE$_r$ | PSNR | SSIM | RPE$_t$ | RPE$_r$ |
| Church | 27.95 | 0.88 | 0.031 | 0.089 | 32.14 | 0.96 | 0.006 | 0.016 |
| Barn | 28.20 | 0.89 | 0.127 | 0.194 | 33.19 | 0.94 | 0.029 | 0.030 |
| Museum | 27.95 | 0.83 | 0.074 | 0.212 | 31.62 | 0.94 | 0.047 | 0.203 |
| Family | 29.12 | 0.83 | 0.051 | 0.028 | 34.80 | 0.97 | 0.024 | 0.020 |
| Horse | 29.43 | 0.87 | 0.135 | 0.061 | 35.45 | 0.97 | 0.109 | 0.053 |
| Ballroom | 28.19 | 0.84 | 0.056 | 0.064 | 33.91 | 0.97 | 0.033 | 0.020 |
| Francis | 28.57 | 0.79 | 0.103 | 0.159 | 33.80 | 0.92 | 0.026 | 0.147 |
| Ignatius | 26.66 | 0.76 | 0.150 | 0.044 | 31.14 | 0.94 | 0.027 | 0.012 |
| mean | 28.26 | 0.84 | 0.091 | 0.106 | **33.26** | **0.95** | **0.037** | **0.063** |

Table 5: **Ablation for Correspondence on Tanks and Temples.** Performance on both novel view synthesis and camera pose estimation. The best results of means are highlighted in bold.

| scenes | Ours | | | COLMAP + 3DGS | | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS |
| Church | 32.14 | 0.96 | 0.08 | 29.93 | 0.93 | 0.09 |
| Barn | 33.19 | 0.94 | 0.07 | 31.08 | 0.95 | 0.07 |
| Museum | 31.62 | 0.94 | 0.08 | 34.47 | 0.96 | 0.05 |
| Family | 34.80 | 0.97 | 0.04 | 27.93 | 0.92 | 0.11 |
| Horse | 35.45 | 0.97 | 0.04 | 20.91 | 0.77 | 0.23 |
| Ballroom | 33.91 | 0.97 | 0.04 | 34.48 | 0.96 | 0.04 |
| Francis | 33.80 | 0.92 | 0.13 | 32.64 | 0.92 | 0.15 |
| Ignatius | 31.14 | 0.94 | 0.06 | 30.20 | 0.93 | 0.08 |
| mean | **33.26** | **0.95** | **0.07** | 30.20 | 0.92 | 0.10 |

Table 6: **Comparison to 3DGS trained with SfM poses**. We report the performance of novel view synthesis using ours and vanilla 3DGS. The best results of means are highlighted in bold.

tween two point clouds. As indicated in (Bian et al. 2023), omitting the point cloud loss leads to a significant decrease in pose accuracy.

### Performance in Complex Camera Motions

While the camera motions involved in the Tanks and Temples dataset are relatively minor, we extend the validation of our method's robustness to the CO3D videos, which feature more intricate and demanding camera movements.

As shown in Table 3 and Table 4, our approach not only excels in novel view synthesis but also clearly surpasses CF-3DGS in pose estimation, reinforcing the findings from the Tanks and Temples experiments and underscoring the precision and robustness of our proposed method in scenarios characterized by complex camera motions.

### Ablation Study

**Effectiveness of Correspondence.** We assess the impact of correspondence-guided optimization by substituting it for traditional pixel-wise supervision. Performance metrics for novel view synthesis and camera pose estimation with and without correspondence-guided optimization are detailed in Table 5. Our observations confirm that correspondence plays a crucial role in enhancing both novel view synthesis and pose estimation accuracy. In the absence of correspondence,

inaccurate initial poses lead to significant deviations between the screen space coordinates of objects in the rendered images and those in the GT images, resulting in poor gradients quality and unstable optimization of the 3D Gaussians model.

**Comparison with 3DGS with SfM Poses.** Our analysis extends to comparing the quality of novel view synthesis of our method with that of the conventional 3DGS model (Kerbl et al. 2023), which utilizes poses derived from SfM technique on the Tanks and Temples dataset. As shown in Table 6, our integrated optimization framework delivers performance on par with the 3DGS model that incorporates SfM-derived poses. In scenes where SfM pose estimation is challenging, there is a significant improvement in performance, as observed in the Horse scene.

## Conclusion

We introduce a novel correspondence-guided SfM-free 3D Gaussian splatting for NVS method that enhances novel-view synthesis by avoiding SfM pre-processing. Our approach effectively optimizes relative poses between frames through correspondence estimation and achieves a differentiable pipeline using our proposed approximated surface rendering technique. Experimental results confirm the superiority of our method in terms of quality and efficiency.

# References

Attal, B.; Huang, J.-B.; Richardt, C.; Zollhöfer, M.; Kopf, J.; O'Toole, M.; and Kim, C. 2023. HyperReel: High-Fidelity 6-DoF Video With Ray-Conditioned Sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 16610–16620.

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

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

Barron, J. T.; Mildenhall, B.; Verbin, D.; Srinivasan, P. P.; and Hedman, P. 2023. Zip-NeRF: Anti-aliased grid-based neural radiance fields. *arXiv preprint arXiv:2304.06706*.

Bhat, S. F.; Birkl, R.; Wofk, D.; Wonka, P.; and Müller, M. 2023. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*.

Bian, W.; Wang, Z.; Li, K.; Bian, J.-W.; and Prisacariu, V. A. 2023. Nope-nerf: Optimising neural radiance field with no pose prior. 4160–4169.

Chen, Y.; and Lee, G. H. 2023. DBARF: Deep Bundle-Adjusting Generalizable Neural Radiance Fields. 24–34.

Chung, J.; Oh, J.; and Lee, K. M. 2024. Depth-regularized optimization for 3d gaussian splatting in few-shot images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 811–820.

Fu, Y.; Liu, S.; Kulkarni, A.; Kautz, J.; Efros, A. A.; and Wang, X. 2024. Colmap-free 3d gaussian splatting.

Hartley, R.; and Zisserman, A. 2003. *Multiple view geometry in computer vision*.

Hoiem, D.; Efros, A. A.; and Hebert, M. 2005. Automatic photo pop-up. In *ACM SIGGRAPH 2005 Papers*, 577–584.

Horry, Y.; Anjyo, K.-I.; and Arai, K. 1997. Tour into the picture: using a spidery mesh interface to make animation from a single image. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, 225–232.

Hu, R.; Ravi, N.; Berg, A.; and Pathak, D. 2020. Worldsheet: Wrapping the World in a 3D Sheet for View Synthesis from a Single Image. In *ICCV*.

Jiang, K.; Fu, Y.; Varma T, M.; Belhe, Y.; Wang, X.; Su, H.; and Ramamoorthi, R. 2024. A Construct-Optimize Approach to Sparse View Synthesis without Camera Pose. In *ACM SIGGRAPH 2024 Conference Papers*, 1–11.

Keetha, N.; Karhade, J.; Jatavallabhula, K. M.; Yang, G.; Scherer, S.; Ramanan, D.; and Luiten, J. 2024. SplaTAM: Splat Track & Map 3D Gaussians for Dense RGB-D SLAM. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 21357–21366.

Kerbl, B.; Kopanas, G.; Leimkühler, T.; and Drettakis, G. 2023. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4): 1–14.

Knapitsch, A.; Park, J.; Zhou, Q.-Y.; and Koltun, V. 2017. Tanks and Temples: Benchmarking Large-Scale Scene Reconstruction. *ACM Transactions on Graphics*.

Kopanas, G.; Leimkühler, T.; Rainer, G.; Jambon, C.; and Drettakis, G. 2022. Neural point catacaustics for novel-view synthesis of reflections. *ACM Transactions on Graphics (TOG)*, 41(6): 1–15.

Lai, Z.; Liu, S.; Efros, A. A.; and Wang, X. 2021. Video autoencoder: self-supervised disentanglement of static 3d structure and motion. 9730–9740.

Li, J.; Feng, Z.; She, Q.; Ding, H.; Wang, C.; and Lee, G. H. 2021. Mine: Towards continuous depth mpi with nerf for novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 12578–12588.

Lin, C.-H.; Ma, W.-C.; Torralba, A.; and Lucey, S. 2021. Barf: Bundle-adjusting neural radiance fields. 5741–5751.

Luiten, J.; Kopanas, G.; Leibe, B.; and Ramanan, D. 2023. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*.

Meuleman, A.; Liu, Y.-L.; Gao, C.; Huang, J.-B.; Kim, C.; Kim, M. H.; and Kopf, J. 2023. Progressively optimized local radiance fields for robust view synthesis. 16539–16548.

Mildenhall, B.; Srinivasan, P. P.; Tancik, M.; Barron, J. T.; Ramamoorthi, R.; and Ng, R. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*.

Mur-Artal, R.; Montiel, J. M. M.; and Tardos, J. D. 2015. ORB-SLAM: a versatile and accurate monocular SLAM system. *IEEE transactions on robotics*.

Park, J. J.; Florence, P.; Straub, J.; Newcombe, R.; and Lovegrove, S. 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.

Paszke, A.; Gross, S.; Chintala, S.; Chanan, G.; Yang, E.; DeVito, Z.; Lin, Z.; Desmaison, A.; Antiga, L.; and Lerer, A. 2017. Automatic differentiation in pytorch.

Ranftl, R.; Bochkovskiy, A.; and Koltun, V. 2021. Vision transformers for dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, 12179–12188.

Reizenstein, J.; Shapovalov, R.; Henzler, P.; Sbordone, L.; Labatut, P.; and Novotny, D. 2021. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 10901–10911.

Riegler, G.; and Koltun, V. 2020. Free View Synthesis. In *ECCV*.

Riegler, G.; and Koltun, V. 2021. Stable view synthesis. In *CVPR*.

Sajjadi, M. S.; Mahendran, A.; Kipf, T.; Pot, E.; Duckworth, D.; Lučić, M.; and Greff, K. 2023. RUST: Latent Neural Scene Representations from Unposed Imagery. 17297–17306.

Sajjadi, M. S. M.; Meyer, H.; Pot, E.; Bergmann, U.; Greff, K.; Radwan, N.; Vora, S.; Lucic, M.; Duckworth, D.; Dosovitskiy, A.; Uszkoreit, J.; Funkhouser, T.; and Tagliasacchi, A. 2022. Scene Representation Transformer: Geometry-Free Novel View Synthesis Through Set-Latent Scene Representations.

Schonberger, J. L.; and Frahm, J.-M. 2016. Structure-from-motion revisited. In *CVPR*.

Smith, C.; Du, Y.; Tewari, A.; and Sitzmann, V. 2023. Flow-Cam: Training Generalizable 3D Radiance Fields without Camera Poses via Pixel-Aligned Scene Flow.

Sun, J.; Shen, Z.; Wang, Y.; Bao, H.; and Zhou, X. 2021. LoFTR: Detector-free local feature matching with transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 8922–8931.

Taketomi, T.; Uchiyama, H.; and Ikeda, S. 2017. Visual SLAM algorithms: A survey from 2010 to 2016. *IPSJ Transactions on Computer Vision and Applications*, 9(1): 1–11.

Tang, S.; Zhang, J.; Zhu, S.; and Tan, P. 2022. Quadtree attention for vision transformers. *arXiv preprint arXiv:2201.02767*.

Tian, F.; Du, S.; and Duan, Y. 2023. Mononerf: Learning a generalizable dynamic radiance field from monocular videos. 17903–17913.

Tucker, R.; and Snavely, N. 2020. Single-view view synthesis with multiplane images. In *CVPR*.

Umeyama, S. 1991. Least-squares estimation of transformation parameters between two point patterns. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 13(04): 376–380.

Verbin, D.; Hedman, P.; Mildenhall, B.; Zickler, T.; Barron, J. T.; and Srinivasan, P. P. 2022. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 5481–5490. IEEE.

Wang, Z.; Bovik, A. C.; Sheikh, H. R.; and Simoncelli, E. P. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*.

Wang, Z.; Wu, S.; Xie, W.; Chen, M.; and Prisacariu, V. A. 2021. NeRF–: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*.

Xu, Q.; Xu, Z.; Philip, J.; Bi, S.; Shu, Z.; Sunkavalli, K.; and Neumann, U. 2022. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 5438–5448.

Yan, C.; Qu, D.; Xu, D.; Zhao, B.; Wang, Z.; Wang, D.; and Li, X. 2024. Gs-slam: Dense visual slam with 3d gaussian splatting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 19595–19604.

Yen-Chen, L.; Florence, P.; Barron, J. T.; Rodriguez, A.; Isola, P.; and Lin, T.-Y. 2021. iNeRF: Inverting neural radiance fields for pose estimation. In *IROS*, 1323–1330. IEEE.

Yifan, W.; Serena, F.; Wu, S.; Öztireli, C.; and Sorkine-Hornung, O. 2019. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6): 1–14.

Zhang, K.; Riegler, G.; Snavely, N.; and Koltun, V. 2020. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*.

Zhang, Q.; Baek, S.-H.; Rusinkiewicz, S.; and Heide, F. 2022. Differentiable point-based radiance fields for efficient view synthesis. In *SIGGRAPH Asia 2022 Conference Papers*, 1–12.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*.

Zhou, T.; Tucker, R.; Flynn, J.; Fyffe, G.; and Snavely, N. 2018. Stereo magnification: Learning view synthesis using multiplane images.