

# Self-Calibrating 4D Novel View Synthesis from Monocular Videos Using Gaussian Splatting

Fang Li, Hao Zhang, Narendra Ahuja  
University of Illinois at Urbana-Champaign  
{fangli3, haoz19, n-ahuja}@illinois.edu

## Abstract

Gaussian Splatting (GS) has significantly elevated scene reconstruction efficiency and novel view synthesis (NVS) accuracy compared to Neural Radiance Fields (NeRF), particularly for dynamic scenes. However, current 4D NVS methods, whether based on GS or NeRF, primarily rely on camera parameters provided by COLMAP and even utilize sparse point clouds generated by COLMAP for initialization, which lack accuracy as well are time-consuming. This sometimes results in poor dynamic scene representation, especially in scenes with large object movements, or extreme camera conditions e.g. small translations combined with large rotations. Some studies simultaneously optimize the estimation of camera parameters and scenes, supervised by additional information like depth, optical flow, etc. obtained from off-the-shelf models. Using this unverified information as ground truth can reduce robustness and accuracy, which does frequently occur for long monocular videos (with e.g. > hundreds of frames). We propose a novel approach that learns a high-fidelity 4D GS scene representation with self-calibration of camera parameters. It includes the extraction of 2D point features that robustly represent 3D structure, and their use for subsequent joint optimization of camera parameters and 3D structure towards overall 4D scene optimization. We demonstrate the accuracy and time efficiency of our method through extensive quantitative and qualitative experimental results on several standard benchmarks. The results show significant improvements over state-of-the-art methods for 4D novel view synthesis. The source code will be released soon at <https://github.com/fangli333/SC-4DGS>.

## 1 Introduction

This paper is about joint optimization of camera parameters and high-fidelity dynamic scene representation for photorealistic 4D novel view synthesis (NVS). Neural Radiance Field (NeRF) [18] approaches have significantly advanced the performance of Novel View Synthesis (NVS). However, NeRF-based methods [18, 42, 12, 21, 22], with their reliance on ray-casting and point-sampling algorithms, still suffer from long preprocessing and training times, which highly restrict their applications in areas like AR/VR and 3D content generation. Recently, 3D-GS [11] has introduced explicit 3D representation and Differential-Gaussian-Rasterization in place of NeRF’s implicit neural representation and neural renderer. Such improvements [11, 39, 33] have significantly reduced the training time while maintaining high-granularity rendering of novel views. It is important to note that the effectiveness of current NVS methods highly relies on the accuracy of camera parameters, which are obtained either from COLMAP[26] or using self-calibration. The quality of these estimates affects the results obtained by various NVS methods.

Our goal is this paper is an NVS method that performs well with respect to eight properties. These are motivated by the limitations we have noted in the current methods. Specifically, we wish to overcome eight limitations (L1-L8) listed below. (L1): inaccurate feature extraction (L2): long

Table 1: **Some Common Limitations and Indications of Whether Different Approaches Overcome them (✓) or Not (X).** The total time shown is for camera & scene optimizations on the be11 data in NeRF-DS, using COLMAP docker [5] indicated by Co1.

Method	Dynamic Scene NVS	Optimize w/o G.T. Intrinsic	Extreme Geometry Scene	Long-video Camera estimation	Accurate Extracted Feature	$\leq 1$ more Models' Supervision	Optimize w/o 3D Priors	Total Optimization Time
NeRFmm [32]	X	✓	X	-	X	✓	✓	-
CF-3DGS [7]	X	X	-	-	-	✓	✓	-
Nope-NeRF [1]	X	X	-	-	-	✓	✓	-
LocalNeRF [17]	X	X	-	✓	-	✓	✓	-
InstantSplat [6]	X	✓	-	-	-	✓	X	-
FlowMap [27]	X	✓	-	✓	X	X	✓	-
Co1 + HyperNeRF [22]	✓	✓	X	✓	X	✓	✓	>64h
Co1 + Deform-3DGS [39]	✓	✓	X	✓	X	✓	✓	5.7h
Co1 + 4D-GS [33]	✓	✓	X	✓	X	✓	✓	-
RoDynRF [16]	✓	✓	✓	X	X	X	✓	30h
<b>Ours(SC-4DGS)</b>	✓	✓	✓	✓	✓	✓	✓	5h

training time; (L3): works only for static scenes, does not extend to dynamic scenes; (L4): does not work under extreme geometric conditions - having large object movements and camera rotations but small camera translations such as DAVIS[23] dataset; (L5): requires specification of camera intrinsics; (L6): requires specification of 3D prior model [31]; (L7): requires supervision from multiple models (to overcome their individual limitations); and (L8): does not work well on long videos.

Examples of past work whose limitations motivate specific properties are as follows. NeRFmm [32] requires the camera to forward-face the scenes and the rotation range of the camera be limited to  $\pm 20^\circ$  (L4). Nope-NeRF [1], CF-3DGS[7] and LocalNeRF[17] lack (L5) and also require monocular depth estimation from MiDaS [25]. Nope-NeRF also needs about 30 hours of optimization for each scene (L2). InstantSplat [6] cannot work without the 3D prior model DUST3R [31] in (L6). FlowMap [27] requires off-the-shelf models RAFT [28], MiDaS [25], and CoTracker [10] for computing optical flow, monocular depth estimation, and point tracking, and suffers from the inaccuracies arising from these models; Sec 4.3 shows failure cases. (L7) .

All the aforementioned methods fail on dynamic scenes (L3). 4D-GS [33], Deformable 3DGS [39] and HyperNeRF [22] methods address (L3) but depend on COLMAP [26] for camera parameter estimation in the absence of a good alternative for dynamic scenes. Even when they eliminate moving objects using motion masks, the performance of COLMAP and the 4D scene optimization models is still limited by (L1) and (L2), resulting in poor 4D NVS performance. Further, our experiments show that COLMAP [26] completely fails in the extreme geometry scenes with relatively large object movements and tiny camera translations but huge camera rotations such as seen in the DAVIS [23] dataset (L4), consistent with the findings in RodynRF [16] (Sec 4.1). As a result of incorporating supervision from MiDaS [25] and RAFT [28], RodynRF [16] is affected by (L2), (L7) and (L8), requiring over 28 hours' training for one monocular video with 50~80 frames and failing while optimizing long videos (e.g. > 800 frames). Tab 1 and Fig 5 lists these limitations and how different methods compare with respect to them, where ✓ indicates that the limitation is overcome (desirable).

In comparison to this SOTA, in this paper we propose a new method **SC-4DGS** which can robustly learn accurate camera parameters and reconstruct high-fidelity dynamic scene representations, free of the limitations above. SC-4DGS starts with our proposed Structural Points Extraction (SPE) algorithm (Sec 3.2), which can extract highly accurate 2D-3D mappings of structural points based only on CoTracker [10], for accurate camera estimation. We then jointly optimize camera parameters and 3D structural points supervised by the extracted 2D structural points and 2D-3D correspondence from SPE (Sec 3.3). Finally, with the optimized camera parameters and 3D structural points, optimal scene representations are optimized within the canonical field (to learn mean positions  $x$ , mean quaternions  $r$ , mean scaling  $s$ , and opacity  $\sigma$ ) and the deformation field defined (to learn  $\Delta x$ ,  $\Delta r$ ,  $\Delta s$ ) defined in Sec 3.4. We evaluate our approach on three standard public datasets including NeRF-DS [36], DAVIS [23], and Nvidia [41], and present quantitative and qualitative comparisons with existing methods.

Our **main contributions** are as follows:

- We introduce a new method **SC-4DGS** that possesses a number of desirable properties: it synthesizes high-fidelity novel views of dynamic scenes using Gaussian Splatting without

requiring camera priors and limitations on video length while taking less time than SOTA methods. Indeed, it overcomes all the limitations (L1-L8), as can be seen in Tab 1

- Towards the aforementioned performance, our SC-4DGS learns robust and accurate camera parameters, an ability whose lack has adversely affected the performances of many SOTA methods.
- We show that our method outperforms the current state-of-the-art methods in quantitative and qualitative terms on three standard benchmark datasets.

## 2 Related Work

**Novel View Synthesis w/ COLMAP.** To reconstruct views of objects and scenes, existing methods employ different representations, including mesh representations [37, 43], planar representations [8, 9], point cloud representations [35, 44], neural field implicit representations [18, 17, 22], and the recently introduced explicit Gaussian representations [39, 33, 11]. Prior to 3D-GS [11], numerous NeRF-based enhancements [18] were made, including dynamic scene synthesis [24, 22, 21, 36], sparse-view scene reconstruction [13, 20, 34, 38], and high-fidelity mesh extraction [37, 29, 40]. However, NeRF-based methods share the limitation of long training time. To address this issue, the recently introduced 3D-GS [11] offers explicit 3D-GS representations and Differential-Gaussian-Rasterization rendering, implemented in CUDA [19]. This technique leverages learnable explicit 3D Gaussian ellipsoids, incorporating attributes like position, rotation, opacity, scale, and color for scene representations and reduce the time costs compared to NeRF-based methods. Several works [39, 33] have proposed its applications in dynamic scenes. Nonetheless, the efficiency and performance of both NeRF-based and 3D-GS-based methods are significantly hampered by the preprocessing time, and accuracy of camera parameters which are obtained from COLMAP [26].

**Novel View Synthesis w/o COLMAP.** Currently, COLMAP [26] is the most widely used method for camera parameter estimation, but its limitations have been a barrier. The inaccuracies of its SIFT [15] feature extraction, along with its time-consuming matching and reconstruction steps, have hindered its usage. Some methods [32, 14, 1, 7] attempt to jointly optimize camera poses and static scene representations, but they require camera intrinsics to be provided. InstantSplat [6] and FlowMap [27] have been introduced to address such limitations in static scenes. In dynamic environments, COLMAP [26] is widely used by state-of-the-art dynamic scene NVS methods [22, 39, 33] even though it is theoretically designed for static scenes - it optimizes camera parameters and 3D parameters in successive time steps and therefore cannot handle the changing 3D parameters in dynamic scenes. To overcome these challenges, RoDynRF [16] leverages supervision from monocular depth estimation [25] and optical flow estimation [28] in addition to RGB images. Unfortunately, RoDynRF struggles with long monocular videos and requires over 28 hours of training even for short videos. Compared with RoDynRF and other existing dynamic scene NVS methods utilizing COLMAP, our proposed method learns more accurate and robust camera parameters in less time without requiring any camera priors and produces comparable results.

## 3 Method

We present an overview of our method in Fig 1. Starting with a monocular video with  $N$  frames, we input the RGB frames  $\mathbf{F}_i^{rgb}$ , motion masks  $\mathbf{M}_i^{motion}$ , and frame times  $\mathbf{T}_i$  to our model,  $i \in [0, N - 1]$ . The steps in our method are presented in the subsections below. We first briefly review 3D Gaussian Splatting (3D-GS) in Sec 3.1. In Sec 3.2, we discuss our newly proposed Structural Point Extraction (SPE) algorithm, detailing how it can establish correspondences between 2D structural points in each frame and 3D structural points in world coordinates shared at successive frame times, and then extract them. Joint optimization of camera parameters and 3D structural points is presented in Sec 3.3. Finally, we discuss dynamic scene representation optimization in Sec 3.4.

### 3.1 Preliminaries

In significant contrast to the implicit scene representation in NeRF [18], 3D-GS [11] adopts a new way of representing the scene via explicit gaussian ellipsoids. Each 3D gaussian ellipsoid is parameterized by (a) center position  $\mu \in \mathbb{R}^3$  in world coordinates; (b) quaternion (rotation) matrix  $r \in \mathbb{R}^4$ ; (c) opacity scalar  $\alpha \in \mathbb{R}$ ; (d) scale factor  $\sigma \in \mathbb{R}^3$ , and (e) spherical harmonics (SH) coefficients  $c \in \mathbb{R}^k$

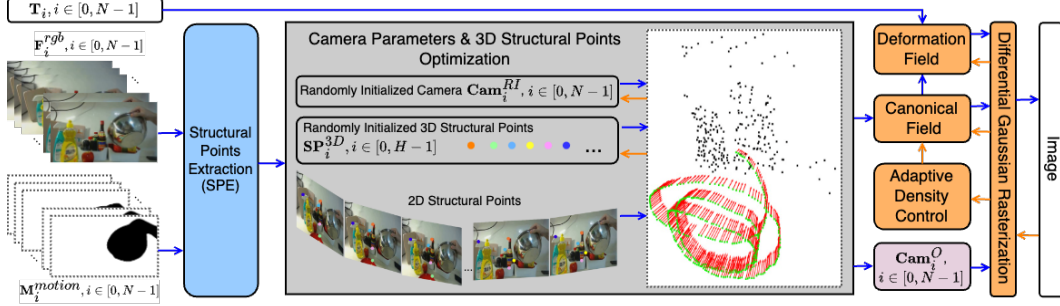


Figure 1: **Overview of SC-4DGS.** We take the basin data of NeRF-DS [36] dataset as the example. First, the SPE algorithm attempts to extract 2D structural points  $\mathbf{SP}^{2D}$  and 3D structural points  $\mathbf{SP}^{3D}$  through  $\mathbf{F}^{rgb}$  and  $\mathbf{M}^{motion}$ , and establish relationships among them. Then, the optimized cameras  $\mathbf{Cam}^O$  are learned in the second joint optimization module, starting with randomly initialized camera  $\mathbf{Cam}^{RI}$  and  $\mathbf{SP}^{3D}$ , supervised by the estimated  $\mathbf{SP}^{2D}$  from SPE. Finally, given  $\mathbf{T}$ , the optimized  $\mathbf{Cam}^O$  and  $\mathbf{SP}^{3D}$ , a Canonical Field and a Deformation Field (see text for details) are computed to optimize the mean representations and deformations of the scene, respectively, supervised by  $\mathbf{F}^{rgb}$ . In the middle of the figure, we show the learned camera positions  $\bullet$  and orientations  $\rightarrow$ , and the optimized  $\mathbf{SP}^{3D}$   $\bullet$ .  $\rightarrow$  and  $\leftarrow$  respectively represent operations flow and gradient flow.

( $k$  stands for the degree of freedom) representing the color by encoding the spatial distribution of light intensity across the surface of a sphere. The 3D gaussian ellipsoid  $G$  is computed from its covariance matrix  $\Sigma$  and its center  $\mu$ , and its 3D covariance matrix  $\Sigma$  is calculated from the scaling factor  $s$  and the quaternion matrix  $r$  as in the following Eq 1 and Eq 2:

$$G(x) = e^{-1/2(x-\mu)^T \Sigma^{-1}(x-\mu)} \quad (1)$$

$$\Sigma' = \mathbf{J}\mathbf{W}\Sigma\mathbf{W}^T\mathbf{J}^T, \quad \Sigma = \mathbf{R}\mathbf{S}\mathbf{S}^T\mathbf{R}^T \quad (2)$$

Given one 3D covariance matrix  $\Sigma$  in world coordinates and the world-to-camera transformation matrix  $\mathbf{W}$ , the corresponding 2D covariance matrix  $\Sigma'$  in camera coordinates is calculated as in Eq 2.  $\mathbf{J}$  is the Jacobian of the affine approximation of the projection transformation. The color of each pixel  $C(p)$  is rendered through the blending of the overlapped  $K$  2D gaussian ellipsoids on this pixel using volume rendering technique in Eq 3.

$$C(p) = \sum_{k \in K} c_k \alpha_k \Pi_k^{j-1} (1 - \alpha_k) \quad (3)$$

Here,  $c_k$  and  $\alpha_k$  represent the color and density of this point calculated from  $\Sigma$  multiplied by the opacity and SH color coefficients. Details can be seen in [11, 45].

### 3.2 Structural Points Extraction (SPE)

**Terminology.** In this section, we define the variables used in the SPE algorithm to clarify the process.  $\mathbf{F}_i^{rgb}$  and  $\mathbf{M}_i^{motion}$  ( $i \in [0, N-1]$ ) represent the RGB frames and motion masks.  $N$  is the number of frames in one monocular RGB video.  $H$  is the total number of shared 3D structural points.  $\mathbf{SP}^{3D} \in \mathbb{R}^{H \times 3}$ ,  $\mathbf{SP}_i^{Pool} \in \mathbb{R}^{B \times 2}$ , and  $\mathbf{SP}_i^{2D} \in \mathbb{R}^{\tau \times 2}$ ,  $i \in [0, N-1]$  respectively represent the 3D structural points in the world coordinate, the potential 2D structural points, and the 2D structural points on each frame, with  $\tau$  and  $B$  are the number of required 2D structural points ( $\tau = 100$  by default) and the number of the potential 2D structural points on each frame.  $\mathbf{SP}_i^{Index}$ ,  $i \in [0, N-1]$  stores the mappings between 3D structural points and 2D structural points on each frame.  $\mathbf{Cam}_i^{RI}$  and  $\mathbf{Cam}_i^O$  are respectively the randomly initialized camera parameters and the optimized camera parameters.  $\mathbf{F}_i^{gray}$  is the grayscale mapping of  $\mathbf{F}_i^{rgb}$ .  $\mathbf{Grad}_i^{magn}$  is the gradient magnitude obtained by combining the gradient norms of all three color channels  $\mathbf{Grad}_i^r$ ,  $\mathbf{Grad}_i^g$ ,  $\mathbf{Grad}_i^b$  of  $\mathbf{F}_i^{rgb}$ .  $\mathcal{W}$  is the window size of the Maximum Filter [2].  $\mathbf{E}_i^{gray} \in \{0, 1\}$  represents the edge detection



information by Canny Edge Detector [3] on  $\mathbf{F}_i^{gray}$ .  $credit = 1^B$  is to mark the remaining  $\mathbf{SP}_i^{Pool}$  while processing frame by frame.  $\mathbf{P}^{pos}$  and  $\mathbf{P}^{index}$  are used to store the  $\mathbf{SP}_i^{2D}$  and  $\mathbf{SP}_i^{Index}$ , and they are all initialized with  $-1$  inside, representing 'TODO'.

**SPE algorithm 1** begins by initializing  $H = 0$  and progressively increases  $H$  by the number of newly introduced  $\mathbf{SP}^{3D}$ . While processing frames sequentially, given  $\tau$ , if  $\mathbf{P}_i^{index}$  is still somehow 'empty' ( $\mathbf{P}_i^{index}.any() == -1$ ), new  $\mathbf{SP}_i^{2D}$  is selected from  $\mathbf{SP}_i^{Pool}$  which stores the most 'representative' points on  $\mathbf{F}_i^{rgb}$  into the following process. The reasons why we do not randomly select points on  $\mathbf{F}_i^{rgb}$  and use CoTracker [10] to track it are discussed in Sec 4.3. For the most 'representative' points, we first calculate  $\mathbf{Grad}_i^{magni}$  of  $\mathbf{F}_i^{rgb}$ .  $\mathbf{Grad}_i^{magni}$  can represent the frequency of each pixel through the gradient magnitude of the sum of the gradient norm in each color channel. However, although the local maxima of  $\mathbf{Grad}_i^{magni}$  can tell us which pixels possess relatively high frequency, it still includes the pixels in the low-texture regions. Accurately tracking points on a low-texture surface is challenging for nearly all dense point-tracking models. To eliminate pixels in such low-texture regions, we use Canny Edge Detector to obtain  $\mathbf{E}_i^{gray}$  from the grayscale mapping  $\mathbf{F}_i^{gray}$ , according to the intensity across  $\mathbf{F}_i^{rgb}$ , and then do Maximum Filter( $\mathbf{E}_i^{gray} \cap \mathbf{Grad}_i^{magni}$ ,  $\mathcal{W}$ ) to obtain the most 'representative' points across  $\mathbf{F}_i^{rgb}$ . Then, we use  $\mathbf{M}_i^{motion}$  to mask out the potential 2D structural points on the deforming objects to acquire the final  $\mathbf{SP}_i^{Pool} \in \mathbb{R}^{B \times 2}$ , in that those are not appropriate for camera parameter estimation.  $B$  can be different according to different  $i$ .

Given  $\mathbf{SP}_i^{Pool}$ , we use CoTracker [10] to track the positions  $\mathbf{Pred}^{pos}$  and visibility  $\mathbf{Pred}^{vis}$  in the following frames. For  $\mathbf{SP}_i^{Pool}$  in  $\mathbf{F}_i^{rgb}$ , if either the corresponding  $\mathbf{Pred}_p^{vis} == 0$  or  $\mathbf{M}_p^{motion}[\mathbf{Pred}_p^{pos}] == 0$  ( $p > i$ ), we mark these points as missing. When the number of remaining points first gets to be below  $num$ , we randomly select  $num$  points from the remaining points of last frame using  $credit$ , and store the corresponding points before the current frame into  $\mathbf{P}^{pos}$ , then assign them index in  $\mathbf{P}_i^{index}$ . As the points continuously disappear in the following frames, we store only the existing points in each frame with their corresponding indexes. Our SPE algorithm performs such operations iteratively until all  $\mathbf{SP}_i^{2D}$  and  $\mathbf{SP}_i^{Index}$  are not 'empty'. The qualitative results of our SPE algorithm are shown in Fig 6 and Fig 7. More discussion on why we do not use the CoTracker output directly is in Sec 4.3.

### 3.3 Camera Parameters & 3D Structural Points Optimization

Given  $\mathbf{SP}^{2D}$ ,  $\mathbf{SP}^{Index}$  and  $\mathbf{SP}^{3D}$  from our SPE algorithm 1 in Sec 3.2, we now discuss how to conduct joint optimization of camera parameters  $\mathbf{Cam}$  and  $\mathbf{SP}^{3D}$ . We assume that the entire monocular video has a constant focal length  $f$ . We define the quaternion rotation matrix of  $\mathbf{F}_i^{rgb}$  as  $\mathbf{Quat}_i \in \mathbb{R}^4$ , and parameterize orientation  $\mathbf{R}_i \in \mathbb{R}^{3 \times 3}$ , translation  $\mathbf{T}_i \in \mathbb{R}^3$ , world-to-camera transformation matrix  $\mathbf{W2C}_i \in \mathbb{R}^{4 \times 4}$ , and perspective projection matrix  $\mathbf{PP} \in \mathbb{R}^{4 \times 4}$  following 3D-GS [11]. With the mappings from 3D world coordinates to the pixel locations on the images in Eq 4, we design the projection loss  $\mathcal{L}_{proj} = \sum_{i=0}^{N-1} \text{MSE}(\mathbf{SP}_i^{2Dproj} - \mathbf{SP}_i^{2D})$  between the 2D projected structural points  $\mathbf{SP}_i^{2Dproj}$  and  $\mathbf{SP}_i^{2D}$  from SPE 1, the distance error loss  $\mathcal{L}_{de} = \sum_{i=0}^{N-1} \text{MSE}(\text{Dist}(\mathbf{SP}_i^{2Dproj}) - \text{Dist}(\mathbf{SP}_i^{2D}))$  between the distances among  $\mathbf{SP}_i^{2Dproj}$  and the ones from  $\mathbf{SP}_i^{2D}$ , and the depth regularization loss  $\mathcal{L}_{dr} = \sum_{i=0}^{N-1} \text{ReLU}(-\mathbf{SP}_i^{Cam}[:, 3])$  for supervision in Eq 5.

$$\mathbf{SP}_i^{Cam} = \text{Homo}(\mathbf{SP}^{3D}[\mathbf{SP}_i^{Index}])\mathbf{W2C}^T\mathbf{PP}^T, \mathbf{SP}_i^{2Dproj} = \mathbf{SP}_i^{Cam}[:, :2]/\mathbf{SP}_i^{Cam}[:, 3] \quad (4)$$

$$\mathcal{L}_{Cali} = \mathcal{L}_{proj} + \mathcal{L}_{de} + \mathcal{L}_{dr} \quad (5)$$

Here,  $\mathbf{SP}_i^{Cam} \in \mathbb{R}^{\tau \times 4}$ ,  $\text{Dist}(\mathcal{X})$  represents the distance between each points  $\mathcal{X}$ . HOMO converts the 3D points in the world coordinates into homogeneous coordinates by concatenating 1. In total, we have  $7N + 1 + H$  parameters required to be optimized in this step. Such supervision is shown to be effective for camera parameter estimation in Fig 1 and the discussions in Sec 4.1.

### 3.4 Dynamic Scene Representations Optimization

Unlike the existing 3D-GS-based methods [11, 33, 39] which set up the initial point clouds either as the sparse point clouds from COLMAP [26] or as a cube full of random dense points, our SC-4DGS

Table 2: **Camera Parameter Prediction Errors for NeRF-DS Dataset.** We use the COLMAP camera parameters as ground truth and show the error measures ATE↓/RPR trans↓/EPR rot↓ for each method.

Method	bell	as	basin	plate	press	cup	sieve
RoDynRF	.14/.18/11.5	.13/.17/10.8	.13/.15/11.8	.13/.18/10.3	.13/.18/10.5	.12/.16/11.8	.14/.18/15.7
Ours	.02/.03/2.05	.03/.04/3.68	.02/.03/1.94	.07/.09/6.86	.03/.03/3.68	.01/.01/1.42	.02/.03/2.57

takes the optimized 3D structural points  $\mathbf{SP}_i^{3DO}$  from Sec 3.3 into the following dynamic scene representation optimization with Adaptive Density Control [39]. We implement a Canonical Field  $\mathcal{G}_C$  [39] and a Deformation Field  $\mathcal{G}_D$  [33, 39] to learn the canonical scene representations  $x, r, s, \sigma$  and deformation scene representations  $\Delta x, \Delta r, \Delta s$  respectively as Eq 6.

$$\begin{aligned} \Delta x_i, \Delta r_i, \Delta s_i &= \mathcal{G}_D(\mathcal{P}(\mathcal{X}_{\mathcal{G}_C}), \mathcal{P}(\mathbf{T}_i)) \\ \mathbf{I}_i^{\text{Render}} &= \mathcal{R}(x + \Delta x_i, r + \Delta r_i, s + \Delta s_i, \sigma) \end{aligned} \quad (6)$$

In Eq 6,  $\mathcal{X}_{\mathcal{G}_C}$  represents the learnt gaussian centers from  $\mathcal{G}_C$ ,  $\mathcal{X}_{\mathcal{G}_C} \leftarrow \mathcal{G}_C(\mathbf{SP}^{3D})$ ;  $\mathcal{P}$  stands for the positional encoding on  $\mathcal{X}_{\mathcal{G}_C}$  and time  $\mathbf{T}_i$ , following [18, 39]. Besides, for simplicity, we use  $\mathcal{R}$  standing for the Differential Gaussian Rasterization [11] to render the image with main parameters, and using RGB loss [11]  $\mathcal{L}_{RGB} = (1 - \lambda)\mathcal{L}_1 + \lambda\mathcal{L}_{D-SSIM}$ ,  $\lambda = 0.2$  for supervision.

## 4 Experiments

COLMAP [26] and the recently released FlowMap [27] are the dominant approaches to static scene camera parameter estimation, and RoDynRF [16] is the state-of-the-art method for estimating camera parameters for dynamic scenes. Since COLMAP [26] and FlowMap [27] are both fundamentally developed for static scenes, and as discussed in the limitation section of FlowMap [27] paper, the estimated camera parameters of FlowMap are less accurate than the ones from COLMAP [26], in the following subsections we use COLMAP [26] and RoDynRF [16] as two baselines for camera parameters estimation; the evaluation setups are discussed in Appendix Sec A.2.

### 4.1 Evaluation of Estimated Camera Parameters

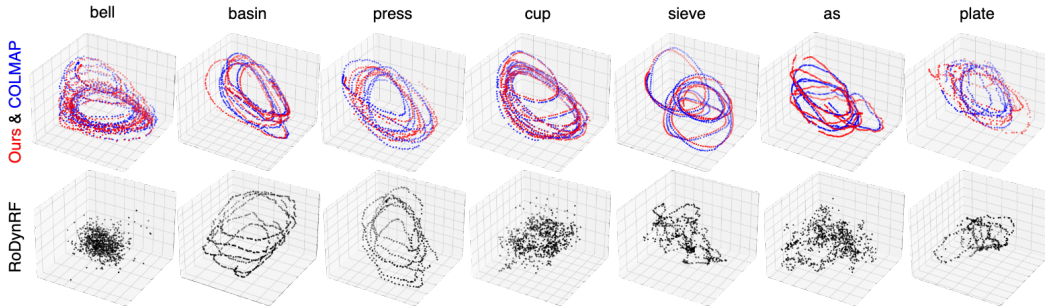


Figure 2: **Visual Camera Comparisons on NeRF-DS.** The red •, blue •, and black • bullets respectively represent the estimated camera poses by our approach, COLMAP [26], and RodynRF [16].

We first present quantitative and qualitative evaluations of camera parameter estimates. For quantitative evaluation, we use Absolute Trajectory Error (ATE), Relative Pose Error for Translation (RPE Trans), and Relative Pose Error for Rotation (RPE rot), which represent the global discrepancy between two trajectories, errors in translation between consecutive imaging instants (frames), and errors in orientation between consecutive imaging instants (frames). (See Appendix Sec. A.3 for details.) For the NeRF-DS [36] dataset, containing **long monocular videos** and **small object movements** with relatively **large camera movements**, we show quantitative and visual comparisons in Table 2 and Fig 2. Our method obtains estimates comparable to COLMAP’s [26], but RoDynRF [16] fails here, showing our method is more accurate than RoDynRF [16]. Furthermore, we show the optimized point cloud comparisons between Deformable-3DGS [39] with COLMAP cameras and our SC-4DGS with our estimated cameras in Fig 3. The optimized point clouds of Deformable-3DGS [39] contain

points that appear to be floating and not a part of scene geometry plate, e.g., the ones in the colored boxes. In contrast, the optimized point clouds from our method are spatially more concentrated and geometrically contiguous. We can use such comparisons to see the misalignments between camera pose estimates from COLMAP and our method and their relative accuracies (Fig 2). For example, we can select the correspondences of a point across three frames and do ray casting from each camera center going through the corresponding frame point. If the camera parameters are accurate, these three rays should intersect at the same 3D point. When the camera parameter estimates have errors, the resulting lack of triangulation relation among these three rays, the RGB loss minimization during rendering will encourage the Adaptive Density Control to add more floating points in each frame. This phenomenon becomes more pronounced as the distance to the camera increases. We present more such comparisons in Appendix Sec A.6 and Fig 9. Further, we note that in Tab. 2, the EPR rot errors seem relatively bigger than ATE errors and RPE trans errors. The reason is that while calculating EPR rot, the differences between the scales of our learned camera coordinates and COLMAP camera coordinates amplify the EPR rot errors, without affecting the accuracy of our optimized camera parameters. The ATE and RPE trans metrics in Tab. 2, visual camera results in Fig 2, and the rendering results in Fig 5 help bring out the effectiveness of our method.



Figure 3: **Optimized Point Cloud Comparisons.** We take the plate scene in the NeRF-DS [36] dataset as the example here and show more in Appendix Sec A.6. The boxes and the corresponding viewpoints are color-coded. The dense points due to the back wall plane formed using our estimated camera parameters, shown in the **green boxes** and **blue boxes**, can be seen to be more reasonable, in comparison with the scattered points from the same back wall formed using COLMAP camera parameters. Similar comments apply to the **red boxes** corresponding to the window points.

We next demonstrate the better performance of our method over COLMAP and RoDynRF [16] for the case of dynamic scenes with relatively **large object movements** and **tiny camera translations** but **huge camera rotations**. In our experiments on the DAVIS [23] dataset, we find that COLMAP [26] fails on more than 80% of the scenes in DAVIS [23] even when we provide the ground truth motion mask. This finding is also consistent with the observations made in RoDynRF [16]. In investigating why COLMAP fails, we find that it successfully extracts features in the first step but the failures start from the exhaustive feature matching step.

As shown in Fig. 4, COLMAP [26] fails on such scenes, so the methods relying on COLMAP like Deformable-3DGS [39] also fail. By contrast, RodynRF [16] and our method can provide reasonable camera parameter results and good rendering performance. In our experiments, we found that in most cases, the camera poses estimated by our method and RoDynRF are similar and lead to high-quality renderings, demonstrating that both learn good scene representations. However, in some cases, there are obvious differences between the camera poses estimated by our SC-4DGS and RoDynRF. In these cases, we observe that our renderings exhibit higher fidelity and detail (shown in the middle of Fig 4). This suggests that our method can estimate camera poses more reliably.

## 4.2 Rendering Evaluation

We evaluate the performance of novel view synthesis through quantitative results and qualitative comparisons in Tab 3 and Fig 5. Due to the inability to estimate accurate camera parameters from long videos, RoDynRF [16] cannot learn good dynamic scene representations. Therefore, it cannot render high-quality images from novel views. Its PSNR and SSIM are significantly lower than the others' by approximately 10.00 and 0.20, representing low image-reconstruct quality. Its LPIPS is higher

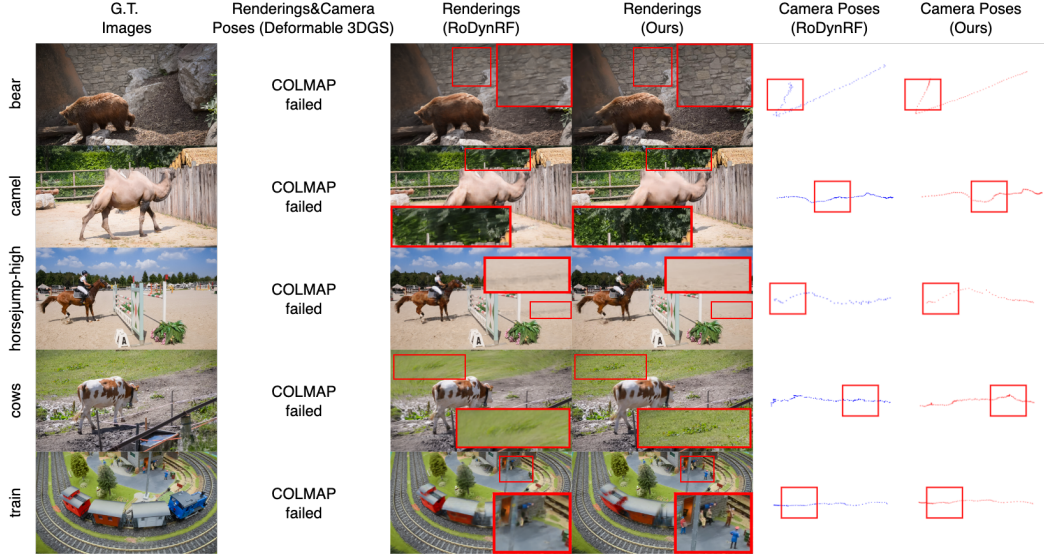


Figure 4: **Rendering & Camera Pose Comparisons on DAVIS.** For each scene, we show the camera pose comparisons and rendering comparisons among Deformable-3DGS [39], RoDynRF [16] and ours, marking the relatively large pose or rendering differences with **red boxes**.



Figure 5: **Visual Novel View Synthesis Results on NeRF-DS.**

than the others' by around 0.35, meaning more noise. The rendered novel view frames in Fig 5 often have blur, floating points, or other noise. Regarding the rendering comparisons between ours and Deformable-3DGS [39], we use the output from COLMAP [26] in all experiments of Deformable-3DGS, and our calibration results in all experiments of our method. Under our self-calibration, our SC-4DGS can achieve comparable PSNR, SSIM, and LPIPS with the Deformable-3DGS [39], and even have high-quality rendering details in some rendering frames. For example, although the PSNR of Deformable-3DGS on the plate scene is 0.005 better than ours, as shown in the fifth column of Fig 5, our method can render more detailed lighting and shadow effects on the plate. The reason behind this is that due to the less inaccurate camera parameters from COLMAP [26], the Adaptive Density Control adds more points to the scene to adapt to the RGB loss, as shown in Fig 3. Such floating points might help improve NVS performance at the viewpoints near the training views; however, because of the wrong geometry, the NVS performance at the viewpoints far from the training views is extremely poor.

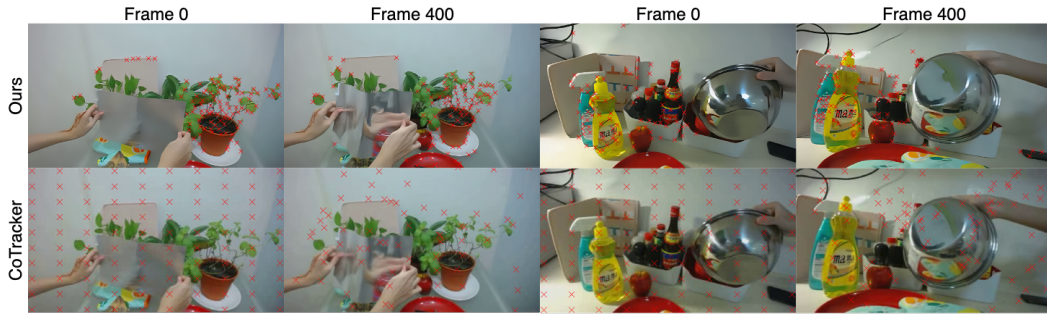
#### 4.3 Why we do not use the CoTracker output directly - Ablation Study

Our SPE algorithm is based on the point tracking (PT) method used in CoTracker [10]. This tracking helps obtain correspondence between  $\mathbf{SP}_i^{3D}$  and  $\mathbf{SP}_i^{2D}$ . Since it is difficult for almost all state-of-the-art dense PT models, e.g., [10, 30], to accurately track every point in every frame in each scene, we do not fully trust the results from any off-the-shelf models unlike CF-3DGS [7], RoDynRF [16]



Table 3: **Rendering Results: PSNR  $\uparrow$ /SSIM $\uparrow$ /LPIPS $\downarrow$  for NeRF-DS.**

Method	Metric	bell	as	basin	plate	press	cup	sieve
Deformable 3DGS [39]	PSNR	31.9745	36.7205	33.1977	30.0201	37.1836	36.4271	37.2177
	SSIM	0.9297	0.9629	0.9450	0.9141	0.9642	0.9621	0.9692
	LPIPS	0.1174	0.0927	0.0974	0.1304	0.0980	0.0854	0.0801
RoDynRF [16]	PSNR	22.7290	20.9097	20.3676	25.4060	22.5965	20.9977	28.2272
	SSIM	0.7018	0.7409	0.6845	0.8230	0.7612	0.6791	0.8537
	LPIPS	0.3959	0.4411	0.4201	0.2819	0.3971	0.4891	0.2695
Ours	PSNR	31.2146	36.4721	32.9480	30.0175	35.8331	35.9480	36.9229
	SSIM	0.9240	0.9611	0.9416	0.9140	0.9535	0.9590	0.9668
	LPIPS	0.1231	0.0879	0.1001	0.1318	0.1007	0.0841	0.0755

Figure 6: **Quality of point tracking by SPE (Ours) vs CoTracker for NeRF-DS.**

and FlowMap [27]. In Fig 6, we show that the direct implementation of CoTracker [10] results in major point-tracking errors. The red  $\times$  denotes the location of the selected structural point in frame 0, and of the same structural point after tracking in frame 400. Even the state-of-the-art dense point-tracking model CoTracker cannot correctly track points in low-texture background regions like walls, due to their highly similar features. Also in high-texture foreground regions such as leaves, numerous points exhibit confusingly similar features. Alternatively, our results in Fig 6 show that the SPE algorithm can filter out the most reliable point-tracking results as 2D structural points. If some points become invisible during tracking, the automatic structural point adaptation mechanism introduces new structural points for tracking to continue estimation of frame-to-frame relationships. For space reasons, here we show only qualitative results for two dynamic scenes in the NeRF-DS dataset; Appendix Sec A.1 contains more comparisons.

## 5 Conclusions & Limitations

**Conclusions.** In this paper, we propose **SC-4DGS** for 4D novel view synthesis w/o camera priors. Our experiments demonstrate that our approach yields more reliable and robust estimates of camera parameters than the state-of-the-art for videos of varying lengths and scenarios, particularly for extreme geometry scenes; obtains optimal dynamic scene representations, and synthesizes high-fidelity RGB images from novel views. We believe that our camera parameter estimation algorithm may also benefit other tasks requiring camera self-calibration.

**Limitations.** A limitation that needs to be overcome is our underlying assumption that the focal length remains constant across frames; if overcome, it would allow self-calibrating 4D NVS with variable zoom effects. Another limitation is that our method requires ground truth motion masks as input which becomes difficult to specify for scenes containing areas of high-speed fluid motion.

## References

- [1] Wenjing Bian, Zirui Wang, Kejie Li, Jia-Wang Bian, and Victor Adrian Prisacariu. Nope-nerf: Optimising neural radiance field with no pose prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4160–4169, 2023.
- [2] Gary Bradski. The opencv library. *Dr. Dobb’s Journal: Software Tools for the Professional Programmer*, 25(11):120–123, 2000.
- [3] John Canny. A computational approach to edge detection. *IEEE Transactions on pattern analysis and machine intelligence*, (6):679–698, 1986.
- [4] Kelvin CK Chan, Shangchen Zhou, Xiangyu Xu, and Chen Change Loy. Investigating tradeoffs in real-world video super-resolution. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5962–5971, 2022.
- [5] dgrnwd. colmap-docker. <https://hub.docker.com/r/dgrnwd/colmap-docker>, 2019. Docker image.
- [6] Zhiwen Fan, Wenyan Cong, Kairun Wen, Kevin Wang, Jian Zhang, Xinghao Ding, Danfei Xu, Boris Ivanovic, Marco Pavone, Georgios Pavlakos, et al. Instantsplat: Unbounded sparse-view pose-free gaussian splatting in 40 seconds. *arXiv preprint arXiv:2403.20309*, 2024.
- [7] Yang Fu, Sifei Liu, Amey Kulkarni, Jan Kautz, Alexei A Efros, and Xiaolong Wang. Colmap-free 3d gaussian splatting. *arXiv preprint arXiv:2312.07504*, 2023.
- [8] Derek Hoiem, Alexei A Efros, and Martial Hebert. Automatic photo pop-up. In *ACM SIGGRAPH 2005 Papers*, pages 577–584. 2005.
- [9] Youichi Horry, Ken-Ichi Anjyo, and Kiyoshi Arai. 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*, pages 225–232, 1997.
- [10] Nikita Karaev, Ignacio Rocco, Benjamin Graham, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Cotracker: It is better to track together. *arXiv preprint arXiv:2307.07635*, 2023.
- [11] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4):1–14, 2023.
- [12] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lerf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023.
- [13] Mijeong Kim, Seonguk Seo, and Bohyung Han. Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022.
- [14] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021.
- [15] Tony Lindeberg. Scale invariant feature transform. 2012.
- [16] Yu-Lun Liu, Chen Gao, Andreas Meuleman, Hung-Yu Tseng, Ayush Saraf, Changil Kim, Yung-Yu Chuang, Johannes Kopf, and Jia-Bin Huang. Robust dynamic radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13–23, 2023.
- [17] Andreas Meuleman, Yu-Lun Liu, Chen Gao, Jia-Bin Huang, Changil Kim, Min H Kim, and Johannes Kopf. Progressively optimized local radiance fields for robust view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16539–16548, 2023.

- [18] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021.
- [19] John Nickolls, Ian Buck, Michael Garland, and Kevin Skadron. Scalable parallel programming with cuda: Is cuda the parallel programming model that application developers have been waiting for? *Queue*, 6(2):40–53, 2008.
- [20] Michael Niemeyer, Jonathan T Barron, Ben Mildenhall, Mehdi SM Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5480–5490, 2022.
- [21] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021.
- [22] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021.
- [23] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 724–732, 2016.
- [24] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021.
- [25] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020.
- [26] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4104–4113, 2016.
- [27] Cameron Smith, David Charatan, Ayush Tewari, and Vincent Sitzmann. Flowmap: High-quality camera poses, intrinsics, and depth via gradient descent. *arXiv preprint arXiv:2404.15259*, 2024.
- [28] Zachary Teed and Jia Deng. 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*, pages 402–419. Springer, 2020.
- [29] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021.
- [30] Qianqian Wang, Yen-Yu Chang, Ruojin Cai, Zhengqi Li, Bharath Hariharan, Aleksander Holynski, and Noah Snavely. Tracking everything everywhere all at once. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19795–19806, 2023.
- [31] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. *arXiv preprint arXiv:2312.14132*, 2023.
- [32] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf-: Neural radiance fields without known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021.
- [33] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint arXiv:2310.08528*, 2023.

- [34] Dejia Xu, Yifan Jiang, Peihao Wang, Zhiwen Fan, Humphrey Shi, and Zhangyang Wang. Sinnerf: Training neural radiance fields on complex scenes from a single image. In *European Conference on Computer Vision*, pages 736–753. Springer, 2022.
- [35] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5438–5448, 2022.
- [36] Zhiwen Yan, Chen Li, and Gim Hee Lee. Nerf-ds: Neural radiance fields for dynamic specular objects. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8285–8295, 2023.
- [37] Gengshan Yang, Minh Vo, Natalia Neverova, Deva Ramanan, Andrea Vedaldi, and Hanbyul Joo. Banmo: Building animatable 3d neural models from many casual videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2863–2873, 2022.
- [38] Jiawei Yang, Marco Pavone, and Yue Wang. Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8254–8263, 2023.
- [39] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101*, 2023.
- [40] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. *Advances in Neural Information Processing Systems*, 34:4805–4815, 2021.
- [41] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5336–5345, 2020.
- [42] Hao Zhang, Fang Li, and Narendra Ahuja. Open-nerf: Towards open vocabulary nerf decomposition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3456–3465, 2024.
- [43] Hao Zhang, Fang Li, Samyak Rawlekar, and Narendra Ahuja. Learning implicit representation for reconstructing articulated objects. *arXiv preprint arXiv:2401.08809*, 2024.
- [44] Qiang Zhang, Seung-Hwan Baek, Szymon Rusinkiewicz, and Felix Heide. Differentiable point-based radiance fields for efficient view synthesis. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–12, 2022.
- [45] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization, 2001. VIS’01.*, pages 29–538. IEEE, 2001.



## A Appendix

### A.1 Datasets

As the novel view synthesis tasks highly rely on the estimated camera parameters from COLMAP [26], the NVS datasets will be released only if the contributors can make sure that COLMAP [26] can estimate the relatively accurate camera parameters in each scene. Due to such reasons, besides NeRF-DS [36] and Nvidia [41], following RoDynRF [16], we also evaluate our results on the DAVIS [23] dataset which includes more general wild scenarios with large movements and deformation. In summary, we test our methods on challenging datasets with large object movements and large camera movements. The scenarios include indoor scenes, outdoor urban scenes, and outdoor wild scenes.

**NeRF-DS.** The NeRF-DS [36] dataset comprises seven monocular long videos from seven distinct dynamic scenes, each containing between 400 and 800 frames. Every scene includes at least one specular object and features a mix of low-texture and high-texture backgrounds. Furthermore, the dataset exhibits significant scene and camera movements. Due to some blurriness in the provided frames, we have applied the RealBasicVSR [4] model for deblurring. To ensure fairness, deblurred frames are used consistently across all experiments and ablation studies involving the NeRF-DS [36] dataset. The ground truth RGB images and motion masks are provided. Camera parameters are determined using COLMAP [26]. In line with other studies [39, 36], we utilize the highest resolution images available ( $480 \times 270$ ) for scene optimization.

**DAVIS.** DAVIS [23] dataset has 40 short monocular video sequences of 40 different dynamic large scenes including animals, humans, vehicles, etc. Each video sequence contains 50 - 80 frames with at least one moving object with the ground truth RGB images and motion masks. Following RoDynRF [16], we choose the challenging sequences with relatively large camera and objects' movements for comparisons and use the largest resolution ( $1920 \times 1080$ ) for optimization.

**Nvidia.** Nvidia [41] dataset contains 9 dynamic scenes. Each scene is recorded by 12 cameras at 12 timestamps and at each timestamp, each camera takes one image. For each scene, the input monocular video is made up of selecting one image taken by one camera at one timestamp without duplications. More detail can be referred to [41]. In the Nvidia dataset, the camera parameters are preprocessed by COLMAP [26], and the foreground masks are given. We use the default resolution ( $960 \times 540$ ) for optimization.

### A.2 Implementation Details

**Device.** All experiments in this paper are conducted on one NVIDIA A100 40GB GPU.

**Model Setup.** During the camera parameters & 3D structural points optimization period, we parameterize quaternion, translation, and focal length for each camera, and  $H$  learnable 3D structural points. We also adopt the constant learning rates equaling 0.01, 0.01, 1.0, and 0.01 for quaternion, translation, focal length, and the positions of 3D structural points. In the overall detailed dynamic scene optimization step, we follow the same model settings with Deformable-3DGS [39].

**Evaluation Setup.** For camera pose evaluations, following RoDynRF [16], we take COLMAP [26] as the baseline, although sometimes COLMAP [26] cannot work well in the scenes with large movements. In our evaluations of novel view synthesis, we adapt our assessment approaches to accommodate the diverse configurations of various datasets. For example, the NeRF-DS [36] dataset provides training videos and testing videos, however, the estimated camera parameters from different works on different videos can be in different coordinates, and the direct alignment will cause large errors of aligned rotations. Under such circumstances, we split the frames of the training videos in each scene into training sets and testing sets. For every two frames, we incorporate the first frame into the training set and the second frame into the testing set. Besides, since COLMAP [26] always fails on DAVIS [23], following the setup in RoDynRF [16], we only compare the scene representations and camera poses between ours and RoDynRF since no testing sets are available. Besides, the tremendous movements of objects across frames make Although our approach is developed for monocular video input, we also show the camera pose comparisons between ours and COLMAP [26] on the non-video Nvidia [41] dataset.

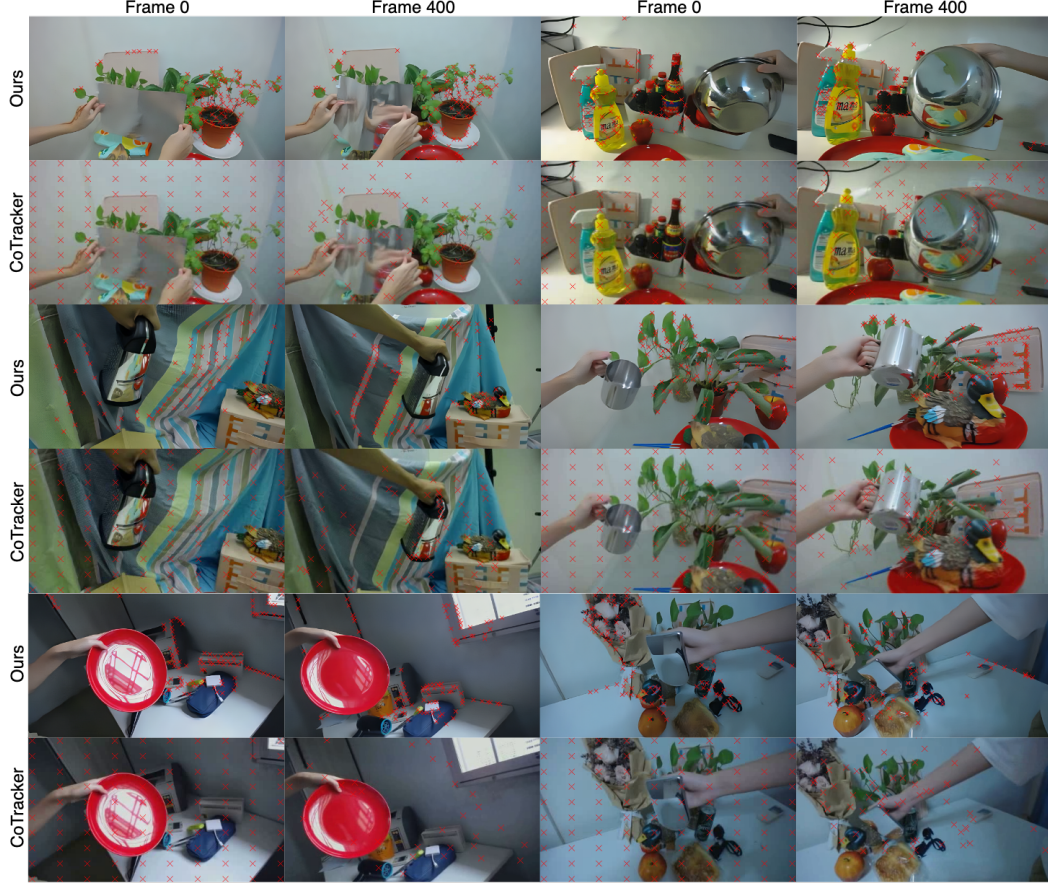


Figure 7: More Ablation Results on SPE v.s. CoTracker.

### A.3 Evaluation Metrics

We follow the same camera and rendering evaluation metrics as the existing works [16].

#### A.3.1 Rendering Evaluation Metrics

**PSNR.** PSNR is a widely used metric for measuring the quality of reconstructed images compared to ground truth images. It is expressed in decibels (dB). The higher the PSNR, the better the quality of the reconstructed image. PSNR is calculated using the mean squared error (MSE) between the original and the reconstructed image:

$$PSNR = 10 \cdot \log_{10} \left( \frac{MAX^2}{MSE(Image_{Recon}, Image_{G.T.})} \right), \quad (7)$$

where  $MAX$  is the maximum possible pixel value of the image.

**SSIM.** SSIM is designed for assessing the perceived quality of digital images and videos. SSIM focuses on the changes in structural information, luminance, and contrast. The values of SSIM range from -1 to 1, where 1 represents perfect. The formula is defined as:

$$SSIM = \frac{(2\mu_x\mu_y + c_1)(2\sigma_{xy} + c_2)}{(\mu_x^2 + \mu_y^2 + c_1)(\sigma_x^2 + \sigma_y^2 + c_2)}, \quad (8)$$

where  $x$  and  $y$  represent two images,  $\mu_x$  and  $\mu_y$  are the average of  $x$  and  $y$ .  $\sigma_x^2$  and  $\sigma_y^2$  are variance of  $x$  and  $y$ .  $\sigma_{xy}$  stands for the covariance of  $x$  and  $y$ , and  $c_1, c_2$  are the regularization factors.

**LPIPS.** LPIPS is a recent metric measuring the perceptual difference between two images as perceived by a trained neural network. In this paper, we use the same trained network with existing works [11, 16, 39]. LPIPS quantifies the difference in the feature representations of images within the neural network, suggesting that a higher LPIPS score indicates a greater perceptual difference.

### A.3.2 Camera Evaluation Metrics

**Absolute Trajectory Error (ATE).** ATE measures the discrepancy between the true trajectory and the estimated trajectory that a robot or a camera follows over a period of time. It provides a global error measurement over the entire trajectory. ATE is computed by aligning the estimated trajectory with the ground truth trajectory and then computing the Euclidean distances between corresponding points on the aligned trajectories.

**Relative Pose Error for Translation (RPE Trans).** RPE Trans measures the error in the translation part of the pose between consecutive poses or over a fixed time/distance interval. Unlike ATE, RPE focuses on the local accuracy of the motion estimation, examining how well the system preserves the relative motion between two points in time or space.

**Relative Pose Error for Rotation (RPE Rot).** RPE Rot measures the error in the orientation between estimated poses relative to the true orientation. This metric is computed by determining the difference in orientation between the estimated and ground truth poses over short sequences and is typically expressed in angular units (like degrees or radians).

### A.4 More Ablation Study Results

**SPE v.s. CoTracker.** Here in Fig 7 we show more comparisons between the results from the direct implementation of CoTracker [10] and the results from our introduced SPE algorithm. The comparisons are conducted on the NeRF-DS [36] dataset. As there are two scenes with completely the same background, we only show one of them here.

### A.5 Camera Comparisons on Nvidia dataset

Our SC-4DGS model requires the monocular videos as input, and one set of individual images as input might somehow degrade the performance of our method. Despite this, we still test our method and COLMAP [26] on the non-video Nvidia [41] dataset, as shown in Fig 8. In most scenes, our method can still produce comparable results with COLMAP [26].

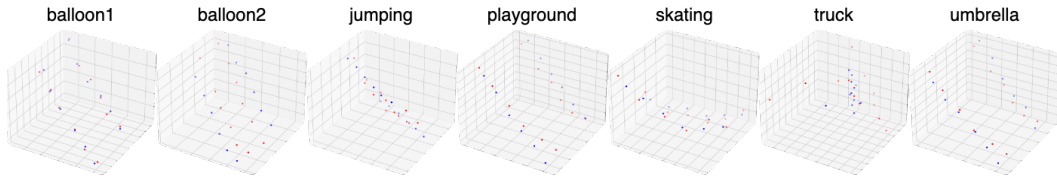


Figure 8: **Visual Camera Comparisons on Nvidia.** The red • and blue • respectively represent the estimated camera poses from our SC-4DGS and COLMAP [26]

### A.6 More Optimized Scene Point Clouds Comparisons

In Fig 9, we show more optimized point cloud results of the scenes in the NeRF-DS [36] dataset. Among these comparisons, Deformable-3DGS [39] prefers to add more floating points to adapt to the RGB loss. However, these floating points cannot represent the real geometry of the scenes. This can promote the accuracy of the novel view synthesis from easy test viewpoints like the ones around the training viewpoints. However, it will destroy the real geometry of scenes, leading to poor NVS performance at the views that are not around the training viewpoints.

### A.7 Structural Point Extraction

We show the detailed SPE algorithm in Alg 1.

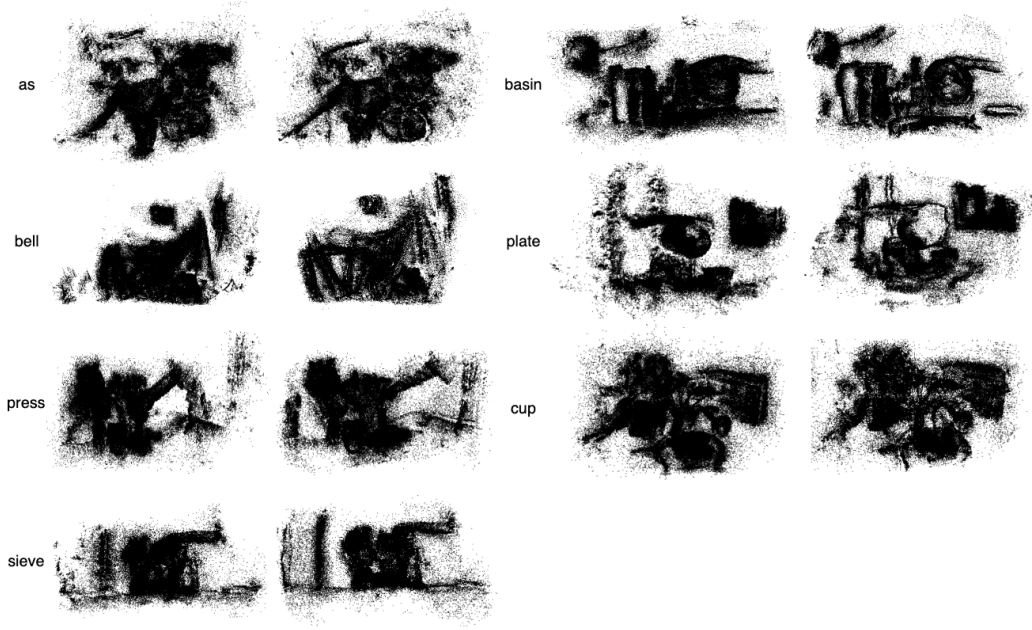


Figure 9: **Optimized Scene Point Clouds on NeRF-DS.** For each scene, the left figure is from Deformable 3DGS [39], the right figure is from Our method.

---

**Algorithm 1** Structural Points Extraction (SPE)

---

```

1: Initialize  $H = 0$ ,  $\mathbf{P}^{pos} \in -\mathbb{1}^{N \times \tau \times 2}$ ,  $\mathbf{P}^{index} \in -\mathbb{1}^{N \times \tau}$ .
2: for  $i = 0, 1, 2, \dots, N - 1$  do
3:   if  $-1$  in  $\mathbf{P}_i^{index}$  then:
4:      $\mathbf{E}_i^{gray} \leftarrow \text{Canny Edge Detector}(\mathbf{F}_i^{gray})$ 
5:      $\text{Grad}_i^{magni} \leftarrow \text{Sqrt}(\text{Sum}(\text{Grad}_i^r, \text{Grad}_i^g, \text{Grad}_i^b))$ 
6:      $\text{SP}_i^{Pool} \leftarrow \text{Maximum Filter}(\mathbf{E}_i^{gray} \cap \text{Grad}_i^{magni}, \mathcal{W}) \cap \mathbf{M}_i^{motion}$ 
7:      $credit = \mathbb{1}^B, p = i + 1$ 
8:      $\text{Pred}^{pos}, \text{Pred}^{vis} = \text{CoTracker}(\mathbf{F}_{j, (j>i)}^{rgb}, \text{SP}_i^{Pool})$ 
9:     while  $p \leq N - 1$  do
10:       $credit[\text{Pred}_p^{vis} == 0 \cup \mathbf{M}_p^{motion}[\text{Pred}_p^{pos} == 0]] = 0$ 
11:      if  $\sum(credit_{p-1}) > num \ \& \ \sum(credit_p) < num$  then
12:         $\mathbf{P}_m^{index}, m \in [i, p - 1] \leftarrow [H, H + num]$   $\triangleright$  Assign index
13:         $\mathbf{P}_m^{pos}, m \in [i, p - 1] \leftarrow \text{Pred}_m^{pos}[\text{Random}(\text{where}(credit_{p-1} == 1), num)]$ 
14:      end if
15:      if  $\sum(credit_{p-1}) < num \ \& \ \sum(credit_p) < num$  then
16:         $\mathbf{P}_p^{index} \leftarrow [H, H + num][credit_p[\text{Random}(\text{where}(credit_{p-1} == 1), num)] == 1]$ 
17:         $\mathbf{P}_p^{pos} \leftarrow \text{Pred}_p^{pos}[\text{where}(credit_p == 1)]$ 
18:      end if
19:       $p += 1, H += num$ 
20:    end while
21:  else PASS
22: end if
23: end for
24:  $\text{SP}_i^{2D} \leftarrow \mathbf{P}_i^{pos}, \text{SP}_i^{Index} \leftarrow \mathbf{P}_i^{index} (i \in [0, N - 1]), \text{SP}^{3D} \leftarrow \frac{1}{2} \mathbb{1}^{H \times 3}$ 

```

---