

Modeling Ambient Scene Dynamics for Free-view Synthesis

Meng-Li Shih
University of Washington
Seattle, Washington, USA
mlshih@cs.washington.edu

Jia-Bin Huang
University of Maryland
College Park, Maryland, USA
jbhuang@umd.edu

Changil Kim
Meta
Seattle, Washington, USA
changil@meta.com

Rajvi Shah
Meta
Seattle, Washington, USA
rajvishah@meta.com

Johannes Kopf
Meta
Seattle, Washington, USA
jkopf@meta.com

Chen Gao
Meta
Seattle, Washington, USA
gaochen@meta.com

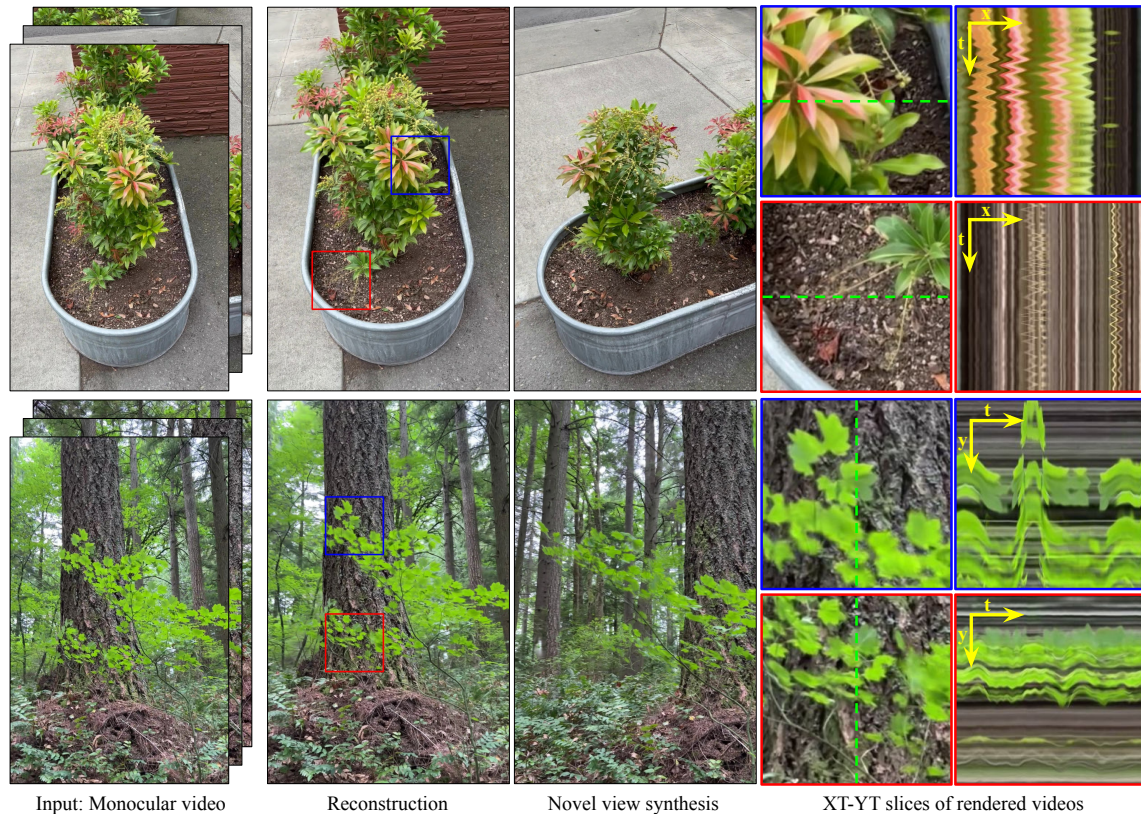


Figure 1: Dynamic free-view synthesis of scenes with ambient motion. Given a causally captured video (1st col) in a scene with ambient dynamics, e.g., swaying trees in a forest or oscillating leaves of a plant in an outdoor environment, we reconstruct highly detailed 3D scene and the scene dynamics. Our method produces sharp reconstruction (2nd col) and render high-fidelity novel views at a specific time step (3rd col). Our video rendering also capture complex ambient scene dynamics, providing immersive viewing experiences beyond navigating in purely static scenes (4th col).

ABSTRACT

We introduce a novel method for dynamic free-view synthesis of an ambient scenes from a monocular capture bringing a immersive quality to the viewing experience. Our method builds upon the recent advancements in 3D Gaussian Splatting (3DGS) that can faithfully reconstruct complex static scenes. Previous attempts to extend 3DGS to represent dynamics have been confined to bounded scenes or require multi-camera captures, and often fail to generalize to unseen motions, limiting their practical application. Our

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
SIGGRAPH Conference Papers '24, July 27-August 1, 2024, Denver, CO, USA
© 2024 Copyright held by the owner/author(s).
ACM ISBN 979-8-4007-0525-0/24/07
<https://doi.org/10.1145/3641519.3657488>

approach overcomes these constraints by leveraging the periodicity of ambient motions to learn the motion trajectory model, coupled with careful regularization. We also propose important practical strategies to improve the visual quality of the baseline 3DGS static reconstructions and to improve memory efficiency critical for GPU-memory intensive learning. We demonstrate high-quality photorealistic novel view synthesis of several ambient natural scenes with intricate textures and fine structural elements. We show that our method significantly outperforms prior methods both qualitatively and quantitatively. Project page: <https://ambientGaussian.github.io/>

CCS CONCEPTS

• **Computing methodologies** → **Rendering**; Point-based models.

KEYWORDS

novel view synthesis, 3D gaussians, ambient motion

ACM Reference Format:

Meng-Li Shih, Jia-Bin Huang, Changil Kim, Rajvi Shah, Johannes Kopf, and Chen Gao. 2024. Modeling Ambient Scene Dynamics for Free-view Synthesis. In *Special Interest Group on Computer Graphics and Interactive Techniques Conference Conference Papers '24 (SIGGRAPH Conference Papers '24)*, July 27-August 1, 2024, Denver, CO, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3641519.3657488>

1 INTRODUCTION

We live in a dynamic world with rich and detailed textures and motions, where elements like leaves swaying rhythmically with the subtle influences of the wind and blades of grass are in constant, gentle motion. The *ambiance* makes the viewing experience immersive. Our work aims to achieve high-quality dynamic free-view synthesis from a *monocular* capture of an ambient (plant) scene, providing a life-like viewing experience. This has significant implications in areas such as virtual 3D teleportation [Orts-Escolano et al. 2016] and virtual touring [Broxton et al. 2020].

Dynamic view synthesis has seen considerable exploration through NeRF-based methods. These vary from learning static templates in canonical spaces [Park et al. 2021a,b] to modeling motion fields across different timestamps [Gao et al. 2021; Li et al. 2021; Liu et al. 2023]. However, these methods are limited by their training speed and overall reconstruction quality. Recently, 3D Gaussian Splatting (3D-GS) [Kerbl et al. 2023a] has emerged as a powerful tool for reconstructing complex static scenes with high fidelity and fast training speeds. It represents the scene with 3D Gaussians. To render an image, we project the 3D Gaussians onto a 2D plane (splatting). The positions, rotations, sizes, colors, and opacities of these 3D Gaussians can be optimized by minimizing the photometric loss between the rendered image and the corresponding input image. The extension of Gaussian Splatting to 4D has opened new possibilities for dynamic reconstruction, either through modeling the scene with 4D Gaussian primitives [Yang et al. 2023b] or a set of 3D Gaussians [Luiten et al. 2024]. However, most of these methods depend on multi-camera captures [Luiten et al. 2024], a requirement impractical for most users. A few methods can produce dynamic view synthesis from a monocular capture but are limited to bounded scenes [Kratimenos et al. 2023; Wu et al. 2023]. Additionally, these

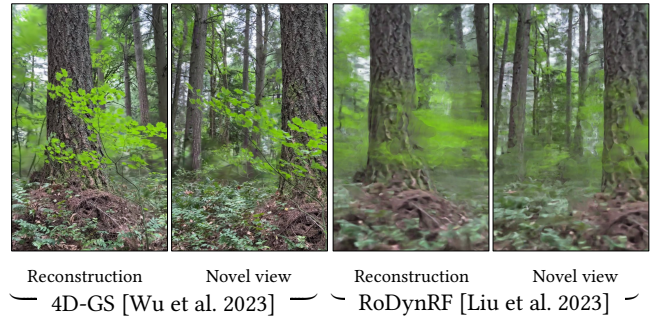


Figure 2: Limitations of existing methods. Here we highlight the limitations of state-of-the-art dynamic radiance fields in addressing the ambient dynamics in an unbounded scene. RoDynRF suffers from severe blurriness due to the use of voxel-grid representation. 4D-GS can recover some spatial details for contents close to the cameras, but struggles with handling ambient motion, resulting in unstable foreground motion and inaccuracies in background motion. Please refer to the supplementary video for comparison.

methods assume that all motion is observed; thus, they are unable to generalize to unseen motions.

In this paper, we introduce a new algorithm for dynamic free-view synthesis from a monocular capture of a plant scene that addresses these limitations. More specifically, we learn a set of template 3D Gaussians in canonical space and then learn how to deform these 3D Gaussians to model the time-varying volume density and appearance of the plant scene. We fix the sizes, colors and opacities of the 3D Gaussians while learning time-dependent positions and rotations. To model the motion, we propose to fit the observed motion of each 3D Gaussian to a *trajectory*. We find that various forms of periodic motion, especially those observed in nature, such as the motion of plants, can be effectively modeled using combinations of periodic functions and their associated coefficients. Inspired by DCT-NeRF [Wang et al. 2021], we utilize basis functions derived from the discrete cosine transform (DCT) to represent the time-variant deformation. For each 3D Gaussian, we use a small MLP to predict the coefficients for each basis, with the final trajectory being the linear combination of all these bases. This method enables us to *extrapolate* motion trajectories beyond the observed data, thereby broadening the scope of our synthesis capabilities.

However, one of the primary challenges in this process is the network’s ability to learn plausible deformations from just a single 2D image observation at each time step. While infinitely many solutions can correctly render the input video, not all of them will result in plausible, photorealistic dynamic novel view synthesis. We introduce rigidity regularization [Luiten et al. 2024] to assist in learning a more consistent deformation in both spatial and temporal dimensions.

We validate our method’s dynamic free-view synthesis performance on our proposed real-world plant scene dataset, demonstrating its effectiveness in synthesizing unseen motion. Additionally, we showcase the capability of our method to edit motion, highlighting its versatility and practical applications.

Our *contributions* are summarized as follows:

- We introduce a high-quality dynamic free-view synthesis method from a monocular capture of unbounded ambient scenes.
- We leverage motion trajectory to generalize to unseen motions.
- We provide a real-world plant scene dataset featuring ambient motion. We demonstrate our method’s capabilities in both motion editing and motion synthesis.

2 RELATED WORK

Novel view synthesis. Novel view synthesis has been studied since decades ago. Earlier research relies [Buehler et al. 2001; Cayon et al. 2015; Debevec et al. 1996; Gortler et al. 1996; Hanrahan 1996; Heigl et al. 1999; Kopf et al. 2014] on proxy geometry such as depth map to blend multiple source view and project to the target view. Beyond that, [Hedman et al. 2018] encodes deep features to the proxy geometry and decodes it with a learning-based encoder to provide better rendering quality. Besides, [Flynn et al. 2019; Srinivasan et al. 2019; Zhou et al. 2018] exploit multi-plane images with different learning strategies, resulting in efficient and high-quality view synthesis in forward-facing scenarios. A few years ago, NeRF [Mildenhall et al. 2020] proposed a way based on implicit field representation and used differentiable volumetric rendering to render images and update the implicit field. Since then, numerous works have been proposed to improve the view-synthesis technique in every aspect. For example, [Barron et al. 2021, 2022, 2023; Xu et al. 2022] are proposed to improve the quality in bounded and unbounded scenes and [Fridovich-Keil et al. 2022; Müller et al. 2022; Sun et al. 2022; Yu et al. 2021] are proposed to improve training and inference speed. Recently, 3D Gaussian Splatting [Kerbl et al. 2023b] introduces a novel approach to scene representation. This method leverages anisotropic 3D Gaussians as a representation for scenes and introduces an efficient differentiable rasterizer that exploit the splatting of these Gaussians onto the image plane. The key advantage of this technique is its ability to achieve fast, high-resolution rendering. However, these approach are limited to reconstruct and rendering static scene.

View synthesis of dynamic scenes. View synthesis of dynamic scenes is a challenging problem as it needs to deal with the reconstruction of geometry, appearance, and scene representation all at once. Earlier works [Bansal et al. 2020; Collet et al. 2015; Jain and Wakimoto 1995; Kanade et al. 1997; Li et al. 2017, 2018; Ma et al. 2023; Yang et al. 2002] tackle this problem by employing multiple-view camera setting to reconstruct the proxy geometry better and bake appearance on it. Beyond this, [Thonat et al. 2021] uses a multiscale representation of motion that allows for looping and blending. Since the emergence of NeRF [Mildenhall et al. 2020], numerous related papers [Du et al. 2021; Fridovich-Keil et al. 2023; Işık et al. 2023; Li et al. 2022a,b; Lin et al. 2023b; Shao et al. 2023b; Song et al. 2023; Wang et al. 2023, 2022] make huge breakthrough on this domain. In addition, the requirement of dynamic view synthesis using only monocular video as input started to be fulfilled. For example, [Gao et al. 2021; Li et al. 2021; Liu et al. 2023] rely on frame-to-frame correspondence to supervise the motion, [Fang et al. 2022; Guo et al. 2022; Park et al. 2021a,b; Petitjean et al. 2023] relies on the mapping from 3D space to a canonical space at each time-step to better fuse

the information from multiple frames, and [Li et al. 2023c; Wang et al. 2021] use trajectories to leverage the advantages of both sides. In the past few months, there are a bunch of works [Das et al. 2023; Huang et al. 2023; Katsumata et al. 2023; Kratimenos et al. 2023; Li et al. 2023a; Liang et al. 2023; Lin et al. 2023a; Luiten et al. 2024; Shao et al. 2023a; Wu et al. 2023; Yang et al. 2023a,c; Yu et al. 2023] start to exploit 3D Gaussian Splatting and applied it to the dynamic scene scenario. These works reach high-quality results in rendered images and fast inference time. However, challenges still exist, particularly when the input is monocular video instead of multiple view.

Modeling ambient motion. People have used various techniques in the frequency domain to model natural, oscillatory 3D motion, such as the swaying of trees in the wind. For example, [Davis et al. 2015] shows how to use modal analysis in the frequency domain to describe the scene from input video and even simulate its motion. [Petitjean et al. 2023] shows that such technique can be adapted to 3D scene. [Li et al. 2023b] shows that one can generate various plausible periodic motions by leveraging diffusion models.

3 METHOD

We first introduce the background of Gaussian Splatting in Section 3.1. Following that, we provide an outline of our method in Section 3.3 with details. Subsequently, we delve into the specifics of modeling periodic motion in Section 3.4. Finally, in Section 3.5, we discuss strategies for mitigating memory consumption during training with many Gaussians along with other pertinent implementation details.

3.1 3D Gaussian Splatting

3D Gaussian Splatting (3D-GS)[Kerbl et al. 2023a] is a recently proposed explicit representation for scene reconstruction. A set of 3D Gaussians represents the scene. Each Gaussian G is defined by a center point $\mu \in \mathbb{R}^3$ and a covariance matrix $\Sigma \in \mathbb{R}^{3 \times 3}$, composed of a rotation matrix $R \in \mathbb{R}^{3 \times 3}$ and a 3-D scaling diagonal matrix $S \in \mathbb{R}^{3 \times 3}$:

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

$$\Sigma = RSS^T R^T \quad (2)$$

Usually, the rotation matrix R and scaling matrix S are represented by the rotation unit quaternion q and scaling factors $s \in \mathbb{R}^3$, respectively, for easier optimization. In addition to Eq. (1) and (2) which define the basic properties of a 3D Gaussian, additional attributes are incorporated to enable photo-realistic rendering. These include color $c \in \mathbb{R}^3$ defined by coefficients of spherical harmonics (SH), and opacity $\alpha \in \mathbb{R}$. To render an image, 3D-GS [Kerbl et al. 2023a] first arrange these Gaussians based on their proximity to a designated viewpoint, then exploit the over-composite blending function to aggregate them for rendering:

$$C = \sum_{i \in N} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (3)$$

The positions, rotations, sizes, colors, and opacities of these 3D Gaussians can be optimized by minimizing the photometric loss(i.e. L1 loss, SSIM loss) between the rendered image and the corresponding input image.

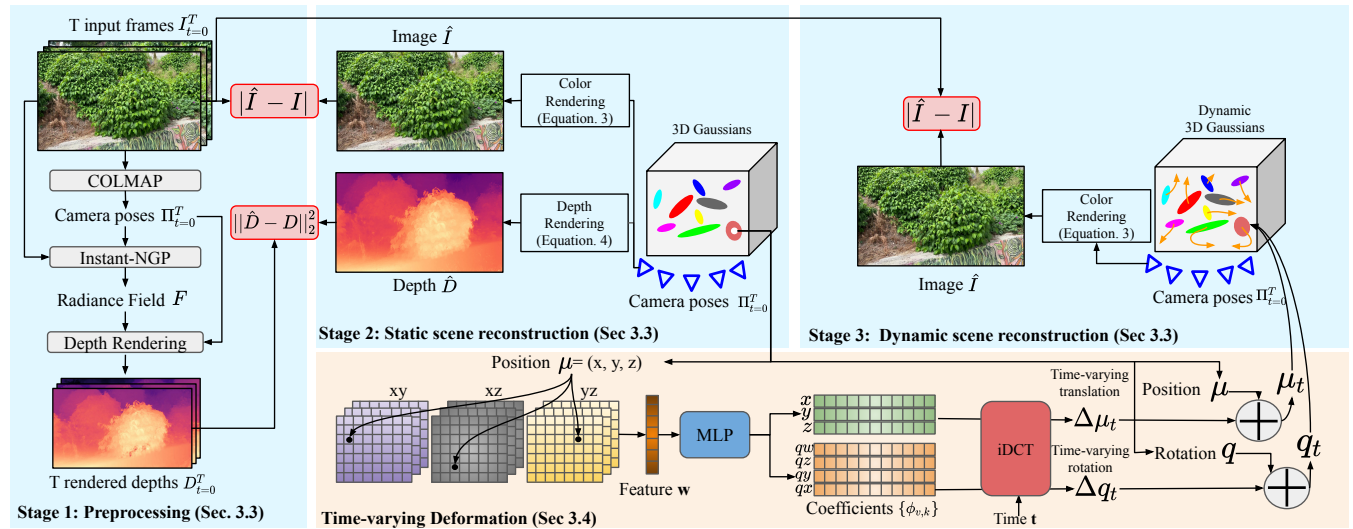


Figure 3: Method Overview. Our method comprises three stages: 1) pre-processing, 2) static scene reconstruction, and 3) dynamic scene reconstruction. In pre-processing stage, we extract the rendered depth map using Instant-NGP [Müller et al. 2022]. The rendered depth from radiance field provide essential depth regularization for unbounded scenes due to poor 3D point cloud recovery in distant regions. In the static scene reconstruction phase, we leverage both photometric and depth loss information obtained from the captured photos and reconstructed radiance field (from Instant-NGP), respectively. This stage allows us to produce high-quality static 3D Gaussian representations of the scene. However, the resulting representations do not model time-varying components like the ambient scene motion. The dynamic regions inevitably are blurry (due to inconsistent photometric losses across different frames). In the dynamic scene reconstruction stage, we introduce temporal parameters to explicitly model the dynamics of each individual 3D Gaussians. We do so by using a triplane-based representation to predict the DCT coefficients for each point (the center position of each 3D Gaussian). Using the predicted coefficients, we can recover the time-varying *translation* and *rotation* given any time step t . We supervise the motion representation using photometric loss and rigidity regularization. The resulting representation allow us to model the time-varying 3D Gaussians and thus render high-quality frames from novel view and time.

3.2 Problem Statement

We study the problem of dynamic view synthesis from monocular captures of ambient scenes. Specifically, we focus on plant scenes, where the magnitude of motion is typically confined to a specific range, such as the swaying motion of trees. While the design of 3D Gaussian splatting handles static scenes effectively, further adaptations are necessary to accommodate the temporal dynamics of scenes. Specifically, properties such as the center position μ and rotation q must be made time-dependent and are denoted as μ_t, q_t , respectively (see 3.3 for detail). Here, the subscript t is the time index.

3.3 Method Overview

We show our proposed framework in Fig. 3. Given an input video sequence of T frames and the corresponding camera parameters, we designed a three-stage approach to reconstruct the dynamic ambient scene using a collection of temporal 3D Gaussians. These Gaussians serve as comprehensive representations, capturing the scene’s geometry, appearance, and dynamic motion. The three stages of our approach are **pre-processing**, **static scene reconstruction**, and **dynamic scene reconstruction**.

Pre-processing. We first run COLMAP on the video frames $I_{t=0}^T$ to obtain the corresponding camera poses $\Pi_{t=0}^T$ and a sparse point cloud. These points help initialize a set of 3D Gaussians. We have observed a strong correlation between the quality of 3D Gaussian splatting and the adequacy of the initial point coverage. This connection stems from the nature of 3D Gaussian splatting, which relies on an explicit representation. If COLMAP fails to extract feature points in certain regions, it results in a lack of Gaussians to begin with. To mitigate this issue, 3D-GS [Kerbl et al. 2023a] introduced split and clone strategies, gradually enabling Gaussians to expand into poorly covered areas. While this approach shows promise in reconstructing content close to the initial points, it faces challenges in reconstructing content farther from them. The absence of Gaussians in areas with insufficient initial point coverage leads to noticeable artifacts and incorrect geometry (Fig. 4), especially in the large unbounded scenes that are the focus of this paper.

Thus, we propose to leverage an implicit representation, specifically NeRF [Mildenhall et al. 2020], to constrain the geometry further. Unlike 3D Gaussian splatting, NeRF and its subsequent developments [Barron et al. 2021, 2022, 2023; Müller et al. 2022] utilize an implicit representation, allowing for arbitrary sampling of points across the 3D space. As a result, it achieves plausible depth

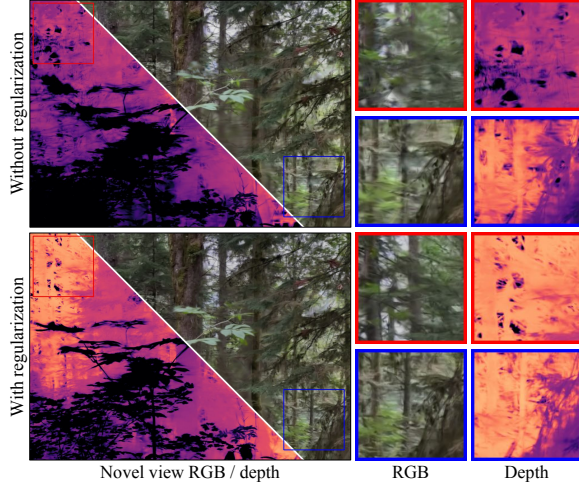


Figure 4: The importance of depth regularization. The quality of 3D-GS depends heavily on the accurate 3D point cloud initialization. In unbounded scenes, however, the geometry of scene elements far away from the camera cannot be reliably reconstructed with structure from motion algorithms (due to small motion parallax). Consequently, 3D-GS tends to predict incorrect geometry in the background and render blurry image due to the lack of initial Gaussians (top). We address this challenge by applying the depth regularization. With the regularization, we observe more accurate and detailed appearance in our rendering (bottom).

estimation in the under-covered region in the unbounded scene. Leveraging the capabilities of NeRF, we fit an NGP [Müller et al. 2022] to the unbounded scene in just 10 minutes and subsequently render depth maps D for each input frame.

One way to use the depth D is to sample additional 3D points through depth unprojection and use them to initialize more 3D Gaussians. However, this can lead to excessive points for rendering, causing computational strain and out-of-memory issues. Therefore, we propose a more effective depth regularization. In the following paragraph, we will detail how these depth maps guide the Gaussians, thus extending their influence into regions considerably distant from the initial points.

Static scene reconstruction. We focus on plant scenes, where the magnitude of motion is typically confined to a specific range, such as the swaying motion of trees. Moreover, a significant portion of our scene remains static, e.g., rocks, land, and tree trunks. Therefore, we begin by assuming the absence of motion within the scene, intending to introduce motion later in the process. This approach offers the advantage of breaking down the complex task of simultaneously reconstructing geometry, appearance, and motion into two more manageable sub-problems: first reconstructing geometry and appearance in the canonical space, and subsequently addressing the reconstruction of motion. As mentioned previously, 3D-GS heavily depends on COLMAP 3D point cloud initialization. The absence of 3D points in under-covered regions will lead to incorrect geometries. To solve this issue, we introduce a depth regularization. By replacing color c_i with depth value of Gaussian, z_i in Eq. 3, we can

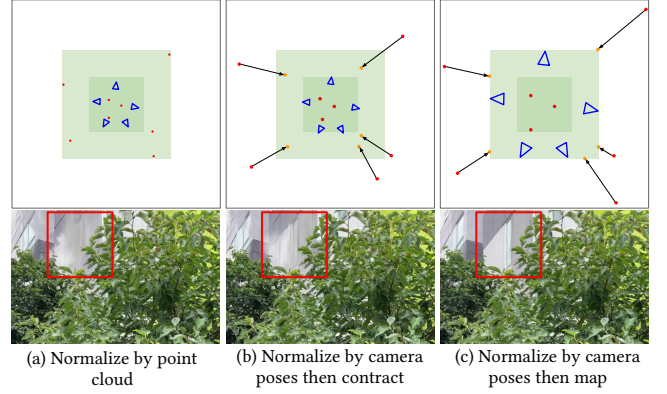


Figure 5: Effect of different scene normalization strategies. (a) Normalizing the scene based on the *range of 3D point cloud* often results in an inefficient use of representational power (because the scene scale can be very large). This typically leads to a blurry foreground and an incorrectly rendered background. (b) When normalizing the scene using the *range of camera poses* and applying ∞ -norm contraction, the foreground becomes sharper. However, the background remains blurry due to inaccurately predicted motion. (c) We propose to normalize the scene using the range of camera poses and map points outside this range to the boundary. Our results show that this achieves higher-quality synthesis in both the foreground and background regions.

render depth map:

$$\hat{D} = \sum_{i \in N} z_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j) \quad (4)$$

Similar to the photometric loss, we acquire depth regularization through the L2 loss between the depth map \hat{D} generated by Gaussian splatting and NeRF-depth D in the pre-processing stage. This regularization effectively improves the reconstruction quality, as shown in Fig. 4. The large unbounded scenes often significantly increase the number of Gaussians, making training expensive and less efficient. To address this issue, we employ a mask loss and mask pruning technique [Lee et al. 2023] to reduce the number of Gaussians used in our model effectively. Furthermore, we implement gradient accumulation and develop a multi-pass rendering scheme for memory-efficient training, which we detail in Sec. 3.5.

Dynamic scene reconstruction. After achieving a satisfying static reconstruction in the canonical space, we now focus on bringing back the motion to model the dynamic scene. For each 3D Gaussian, we want to predict a time-dependent translation $\Delta\mu_t$ and rotation quaternion Δq_t to deform it from the canonical space to a time-dependent space at time t . Specifically, the center position μ and rotation q become μ_t and q_t , where:

$$\mu_t = \mu + \Delta\mu_t \quad (5)$$

$$q_t = \frac{(q + \Delta q_t)}{\|q + \Delta q_t\|} \quad (6)$$

However, instead of directly predicting $\Delta\mu_t$ and Δq_t , we propose to first model a motion *trajectory*, and then sample time-varying

$\Delta\mu_t$ and Δq_t by querying the trajectory at time-step t . This novel design enables us to extrapolate motion trajectories beyond the observed data, thereby synthesizing unseen motions.

To model the motion trajectory of a 3D Gaussian centered at μ , we encode the position using a voxel-based encoder [Cao and Johnson 2023]. Unlike the position μ , which is unbounded, the voxel-based encoder has a predefined range, $[-1, 1]$. We have to normalize the position μ to fit this range. We have observed that normalizing the scene based on the range of point clouds can lead to inefficient use of representational power on the static background. This often results in a blurry foreground and an incorrectly moving background, as depicted in Fig. 5(a). We experiment with the ∞ -norm contraction proposed in Mip-NeRF-360 [Barron et al. 2022], specifically designed for unbounded scenes. However, the background remains blurry due to inaccurately predicted motion. Therefore, we propose normalizing the scene based on the range of the input camera positions, ensuring all camera positions fall within the $[-1, 1]$ range. We reduce their resolution for Gaussians outside this range by mapping them to the boundary. This strategy preserves higher resolution for those salient motions in the foreground because the range of camera positions is typically smaller than that of the initial point cloud. Additionally, it serves as a *regularization* for motions further from the camera, where significant motion is less expected. This approach effectively achieves high-quality representation in both foreground and background regions, as demonstrated in Fig 5(c). The encoder output is then passed through an MLP to predict coefficients for the pre-defined basis. We take the linear combination of all these bases to obtain the final trajectory. Please refer to Sec. 3.4 and Fig. 3 for a detailed explanation of how we model the trajectory.

By jointly optimizing the time-varying deformation along with their static counterparts within the Gaussians, as outlined in Eq. 5 and 6, we can determine the status of the Gaussians at different time steps. Then, we follow the same rasterization pipeline in 3D-GS [Kerbl et al. 2023b] to render the image at time t . The encoder and MLP are optimized by minimizing the photometric loss (i.e., L1 loss) between the rendered image and the corresponding input image. After the training of dynamic scene reconstruction, the blurry and ghosting artifacts due to the inability to model ambient motion in static scene reconstruction is greatly reduced (see Fig. 6).

3.4 Modeling periodic motion

Various forms of periodic motion, especially those observed in nature, such as the motion of plants, can be effectively modeled using periodic functions and their related coefficients. Inspired by DCT-NeRF [Wang et al. 2021], we utilize basis functions derived from the discrete cosine transform (DCT) to represent the time-varying deformation (Fig. 3). We pre-define a set of DCT basis functions and predict the coefficients for each basis. Each 3D Gaussian’s time-varying deformation is independently modeled:

$$v(t) = \sqrt{\frac{2}{K+1}} \sum_{k=1}^K \phi_{v,k} \cos\left(\frac{\pi}{2T}(2t+1)k\right) \quad (7)$$

Here, $\phi_{v,k}$ represents the k -th coefficient associated with the time-varying deformation v , while $v(t)$ represents any time-varying deformation scalar at time step t . For example, $v(t)$ could be the translation $\Delta\mu_t$ in the x-axis. Utilizing DCT to represent time-varying deformation in our scenario offers two distinct advantages: storage efficiency and improved generalization capabilities. As highlighted in concurrent research [Katsumata et al. 2023], representing time-varying deformation with basis functions enhances storage efficiency. In all of our experiments, we set $K = \lceil \frac{1}{4}T \rceil$, and significantly reduce the number of deformation parameters to only a quarter of the total. Moreover, thanks to the inherent characteristics of DCT, this approach accommodates motion fitting when motion is discernible in the input images, and it replicates motion in the absence of direct supervision (e.g., when a Gaussian is occluded or out of the camera frustum). Given that the motion in our scenarios is also periodic (repetitive), the natural properties of DCT contribute to enhanced generalization capabilities for time-varying deformation.

We have observed that neighboring Gaussians often exhibit spatially inconsistent deformations, resulting in unnatural motion during dynamic novel view synthesis. To tackle this issue, we leverage rigidity regularization [Luiten et al. 2024]. The movement of each 3D Gaussian G ’s neighbors should adhere to the rigid-body transformation of G ’s coordinate system across time steps. Additionally, the neighbors of each 3D Gaussian G should maintain similar rotations over time to ensure spatial and temporal coherence.

3.5 Memory efficient multi-pass rendering

While we have reduced the number of Gaussians in static scene reconstruction, GPU out-of-memory issues persist in unbounded scenes due to the optimization demands of the time-varying MLP. As we cache features at the MLP’s hidden layers for each input sample (i.e. each Gaussian), we face significant memory usage. We employ a *two-pass* rendering approach to handle it. In the first pass, we avoid intense memory usage by not caching any feature and proceed to obtain the time-varying deformation parameters from the MLP. We then deform the Gaussians according to these parameters, rasterize them, and identify those Gaussians involved in rendering the current image. We store the deformed status of the involved Gaussians, which will be used in the next pass.

In the second pass, we obtain the deformation using the MLP again. Two things are different here: First, we only consider the Gaussians identified in the previous pass to reduce memory usage. Second, we cache features for back-propagation. We then deform, rasterize, and render the involved Gaussians into an image. Finally, we compute the loss, back-propagate to obtain gradients, release the cache, and optimize the MLP.

To further address the memory constraints, in the second pass, we divide the involved Gaussians into chunks, each containing up to 500k samples. Sequentially, we cache features corresponding to each chunk for back-propagation and release them once we have obtained the gradients. This approach, known as gradient accumulation, circumvents the memory bottleneck. One issue remains: each time we only have a certain chunk of deformed Gaussian status and cannot render a complete image. Thus, we simply take the

deformed status of the other chunks from the first pass. With these strategies, we resolve the memory issue.

4 EXPERIMENTAL RESULTS

4.1 Experimental Setup

Hyperparameters. We implement our approach using a single NVIDIA A100 GPU with 40GB of memory. During the pre-processing stage, we start by training the Instant-NGP model for 10k iterations. In the static scene reconstruction stage, we train static 3D Gaussians for 80k iterations. In the first 9k iterations, we use the split and clone operation to enhance the density of Gaussians. Subsequently, from the 9k to 80k iterations, we deactivate the densification operation and switch to the mask pruning method [Lee et al. 2023] to reduce the number of Gaussians while preserving the visual quality of rendered images. The training time of this stage is about 1.5 hours. Moving on to dynamic scene reconstruction, we jointly train the MLP and 3D Gaussians through 30k iterations. The overall training procedure exhibits a variable duration, ranging from 1.5 to 8 hours, contingent on the complexity and scale of the scene.

4.2 Dataset

Due to the scarcity of the unbounded real-world ambient scene dataset, we introduce the Forest dataset. This dataset includes 10 dynamic ambient scenes captured in various environments. Each video lasts between 10 and 30 seconds, recorded at 1080p resolution and 30 fps. We use COLMAP [Schonberger and Frahm 2016] to reconstruct the sparse point cloud and obtain camera parameters.

4.3 Baseline

We adopt two different baselines that have demonstrated their ability to handle unbounded scenes. The first one is Robust-DyNeRF (RoDynRF), which leverages the scene contraction technique from [Barron et al. 2022] to deal with unbounded dynamic scenes. The second one is 4D Gaussian Splatting (4D-GS) [Wu et al. 2023], which also shows its ability to render high-quality images of dynamic scenes efficiently and outperforms opponents that use NeRF representation. We modify the default settings in 4D-GS by disabling the time-variant scaling parameters Δs_t , which we found to be unstable and prone to producing NaN issues.

4.4 Qualitative Comparison

We show visual comparisons of our results against RoDyNeRF [Liu et al. 2023] and 4D-GS [Wu et al. 2023] in Figs. 10 and 2. RoDyNeRF struggles to produce sharp images due to the limitations of its implicit representation and constraints on optical flow. While 4D-GS manages to capture parts of the static scene effectively, it introduces significant artifacts in distant areas and fails to accurately represent ambient motion. Our method successfully captures detailed appearances and handles ambient motion efficiently. The improvements our approach offers over 4D-GS stem from various design strategies discussed in Sec 3. To further quantify how our method compares with these baselines, we have conducted a user study on the dataset, detailed in Sec. 4.5.

Table 1: Quantitative Evaluation using FR metrics. Comparison to baselines on the forest sequence.

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
4D-GS	19.04	0.472	0.414
RoDyNeRF	18.12	0.393	0.512
Ours	21.94	0.675	0.278

4.5 Quantitative Comparison

User study. We randomly selected 7 examples from our dataset for the user study. For each example, we rendered two types of novel view synthesis results. Given a video with a duration of T and a total of N frames, we first fixed the time at $\frac{T}{2}$ and rendered a video using the camera trajectory from the input video. This first video assesses the quality of the reconstructed geometry and appearance. Secondly, we used the viewpoint from the $\frac{N}{2}$ -th frame of the input video and rendered a video that progresses from time 0 to T . This second video evaluates the quality of the reconstructed motion. With a side-by-side comparison of two baselines, we presented a total of 14 videos for users to judge in each of the two settings. Based on feedback from 24 users, we display the win rate in Fig. 9. It is clear from these results that our method significantly outperforms the others.

Full-reference(FR) metrics. To evaluate the performance in terms of FR metrics (i.e. PSNR/SSIM/LPIPS [Zhang et al. 2018]), we hold out 3 segments, each containing 30 frames (totaling 90 frames), from the video to form the test-val set, using the remaining frames as the training set. We show the results in Tab. 1, where it is evident that our method outperforms the two baselines.

4.6 More diverse scenes

Our method can handle more than just plant motion. We demonstrate our method on a variety of examples in Fig. 11, including the dancing flames of a candle, chimes swaying in the wind, and fluttering notes. We also show the XT-YT slices of the rendered videos. Our method produces realistic, smooth, and periodic motion. For a more detailed visualization of the motion, please visit our project website.

4.7 Ablation Study

We conduct a quantitative analysis to highlight the importance of each component in our method, as shown in Table 2. Specifically, we ablate depth regularization, DCT trajectory, trajectory MLP, rigidity regularization, and scene normalization, following the settings described in Section 4.5 to demonstrate the effectiveness of each component.

We illustrate the effects of depth regularization and rigidity regularization in Fig.4 and Fig.8. For the DCT trajectory prediction, we explore an alternative by replacing it with the direct prediction of Δx_t and Δq_t as proposed in 4D-GS (Fig. 7). Additionally, we ablate the trajectory MLP by optimizing DCT coefficients directly per Gaussian, which we discovered leads to motion flickering. Furthermore, we replace the normalization method for the Gaussian position μ fed into the MLP with the one proposed in 4D-GS (Fig. 5).

Table 2: Ablation study. We report PSNR, SSIM and LPIPS on the forest sequence.

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
Ours w/o depth regularization	20.17	0.588	0.335
Ours w/o DCT transform	21.65	0.663	0.282
Ours w/o trajectory MLP	18.27	0.423	0.452
Ours w/o rigidity regularization	21.74	0.666	0.280
Ours w/o scene normalization	21.67	0.662	0.281
Ours (full model)	21.94	0.675	0.278

4.8 Limitations and Future Work

Our method has certain limitations. One constraint is the inability to handle complex non-periodic motion. Additionally, while we can synthesize views freely at any time step within the range of input viewpoints, the quality of view synthesis deteriorates when attempting to extrapolate beyond this range (see Fig. 12). Therefore, exploring the potential of generative models in future research could provide a promising avenue to enhance our method’s ability to generalize to these complex scenarios.

REFERENCES

- Aayush Bansal, Minh Vo, Yaser Sheikh, Deva Ramanan, and Srinivasa Narasimhan. 2020. 4D Visualization of Dynamic Events from Unconstrained Multi-View Videos. In *CVPR*.
- Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. 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.
- Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. 2022. Mip-NeRF 360: Unbounded Anti-Aliased Neural Radiance Fields. *CVPR* (2022).
- Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. 2023. Zip-NeRF: Anti-aliased grid-based neural radiance fields. *arXiv preprint arXiv:2304.06706* (2023).
- Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew Duvall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. 2020. Immersive light field video with a layered mesh representation. *ACM TOG* 39, 4 (2020), 86–1.
- Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. 2001. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*. 425–432.
- Ang Cao and Justin Johnson. 2023. HexPlane: A Fast Representation for Dynamic Scenes. *CVPR* (2023).
- Rodrigo Ortiz Cayon, Abdelaziz Djelouah, and George Drettakis. 2015. A bayesian approach for selective image-based rendering using superpixels. In *2015 International Conference on 3D Vision*. IEEE, 469–477.
- Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. 2015. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (TOG)* 34, 4 (2015), 1–13.
- Devikalyan Das, Christopher Wewer, Raza Yunus, Eddy Ilg, and Jan Eric Lenssen. 2023. Neural Parametric Gaussians for Monocular Non-Rigid Object Reconstruction. *arXiv preprint arXiv:2312.01196* (2023).
- Abe Davis, Justin G Chen, and Frédo Durand. 2015. Image-space modal bases for plausible manipulation of objects in video. *ACM Transactions on Graphics (TOG)* 34, 6 (2015), 1–7.
- Paul E. Debevec, Camillo J. Taylor, and Jitendra Malik. 1996. *Modeling and Rendering Architecture from Photographs*. Technical Report. USA.
- Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. 2021. Neural radiance flow for 4d view synthesis and video processing. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE Computer Society, 14304–14314.
- Jieming Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. 2022. Fast Dynamic Radiance Fields with Time-Aware Neural Voxels. In *SIGGRAPH Asia 2022 Conference Papers*.
- John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. 2019. Deepview: View synthesis with learned gradient descent. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2367–2376.
- Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. 2023. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12479–12488.
- Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. 2022. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5501–5510.
- Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. 2021. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5712–5721.
- Steven J. Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F. Cohen. 1996. The lumigraph. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH '96)*. Association for Computing Machinery, New York, NY, USA, 43–54. <https://doi.org/10.1145/237170.237200>
- Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. 2022. Neural Deformable Voxel Grid for Fast Optimization of Dynamic View Synthesis. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*.
- MARC LEVOY Pat Hanrahan. 1996. Light field rendering. *SIGGRAPH96, Computer Graphics Proceeding* (1996).
- Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. 2018. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (ToG)* 37, 6 (2018), 1–15.
- B. Heigl, R. Koch, M. Pollefeys, J. Denzler, and L. Van Gool. 1999. Plenoptic Modeling and Rendering from Image Sequences Taken by a Hand-Held Camera. In *Mustererkennung 1999*, Wolfgang Förstner, Joachim M. Buhmann, Annett Faber, and Petko Faber (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 94–101.
- Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. 2023. SC-GS: Sparse-Controlled Gaussian Splatting for Editable Dynamic Scenes. *arXiv preprint arXiv:2312.14937* (2023).
- Mustafa Işık, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. 2023. Humanrf: High-fidelity neural radiance fields for humans in motion. *arXiv preprint arXiv:2305.06356* (2023).
- Ramesh Jain and Koji Wakimoto. 1995. Multiple perspective interactive video. In *Proceedings of the international conference on multimedia computing and systems*. IEEE, 202–211.
- Takeo Kanade, Peter Rander, and PJ Narayanan. 1997. Virtualized reality: Constructing virtual worlds from real scenes. *IEEE multimedia* 4, 1 (1997), 34–47.
- Kai Katsumata, Duc Minh Vo, and Hideki Nakayama. 2023. An Efficient 3D Gaussian Representation for Monocular/Multi-view Dynamic Scenes. *arXiv preprint arXiv:2311.12897* (2023).
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023a. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (July 2023). <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>
- Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 2023b. 3D Gaussian Splatting for Real-Time Radiance Field Rendering. *ACM Transactions on Graphics* 42, 4 (2023).
- Johannes Kopf, Michael F Cohen, and Richard Szeliski. 2014. First-person hyper-lapse videos. *ACM Transactions on Graphics (TOG)* 33, 4 (2014), 1–10.
- Agelos Kratimenos, Jiahui Lei, and Kostas Daniilidis. 2023. DynMF: Neural Motion Factorization for Real-time Dynamic View Synthesis with 3D Gaussian Splatting. *arXiv preprint arXiv:2312.00112* (2023).
- Joo Chan Lee, Daniel Rho, Xiangyu Sun, Jong Hwan Ko, and Eunbyung Park. 2023. Compact 3D Gaussian Representation for Radiance Field. *arXiv preprint arXiv:2311.13681* (2023).
- Lingzhi Li, Zhen Shen, Zhongshu Wang, Li Shen, and Ping Tan. 2022a. Streaming radiance fields for 3d video synthesis. *Advances in Neural Information Processing Systems* 35 (2022), 13485–13498.
- Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. 2022b. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5521–5531.
- Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. 2023a. Spacetime Gaussian Feature Splatting for Real-Time Dynamic View Synthesis. *arXiv preprint arXiv:2312.16812* (2023).
- Zhong Li, Yu Ji, Wei Yang, Jinwei Ye, and Jingyi Yu. 2017. Robust 3D human motion reconstruction via dynamic template construction. In *2017 International Conference on 3D Vision (3DV)*. IEEE, 496–505.
- Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. 2021. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6498–6508.
- Zhengqi Li, Richard Tucker, Noah Snavely, and Aleksander Holynski. 2023b. Generative image dynamics. *arXiv preprint arXiv:2309.07906* (2023).
- Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. 2023c. Dynibar: Neural dynamic image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4273–4284.

- Zhong Li, Minye Wu, Wangyiteng Zhou, and Jingyi Yu. 2018. 4D Human Body Correspondences from Panoramic Depth Maps. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2877–2886.
- Yiqing Liang, Numair Khan, Zhengqin Li, Thu Nguyen-Phuoc, Douglas Lanman, James Tompkin, and Lei Xiao. 2023. GauFR: Gaussian Deformation Fields for Real-time Dynamic Novel View Synthesis. *arXiv preprint arXiv:2312.11458* (2023).
- Haocong Lin, Sida Peng, Zhen Xu, Tao Xie, Xingyi He, Hujun Bao, and Xiaowei Zhou. 2023b. Im4D: High-Fidelity and Real-Time Novel View Synthesis for Dynamic Scenes. *arXiv preprint arXiv:2310.08585* (2023).
- Youtian Lin, Zuozhuo Dai, Siyu Zhu, and Yao Yao. 2023a. Gaussian-Flow: 4D Reconstruction with Dynamic 3D Gaussian Particle. *arXiv preprint arXiv:2312.03431* (2023).
- Yu-Lun Liu, Chen Gao, Andreas Meuleman, Hung-Yu Tseng, Ayush Saraf, Changil Kim, Yung-Yu Chuang, Johannes Kopf, and Jia-Bin Huang. 2023. Robust Dynamic Radiance Fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. 2024. Dynamic 3D Gaussians: Tracking by Persistent Dynamic View Synthesis. In *3DV*.
- Li Ma, Xiaoyu Li, Jing Liao, and Pedro V. Sander. 2023. 3D Video Loops from Asynchronous Input. <https://doi.org/10.48550/ARXIV.2303.05312>
- Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. 2020. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*.
- Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. 2022. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM Trans. Graph.* 41, 4, Article 102 (July 2022), 15 pages. <https://doi.org/10.1145/3528223.3530127>
- Sergio Orts-Escolano, Christoph Rhemann, Sean Fanello, Wayne Chang, Adarsh Kowdle, Yury Degtyarev, David Kim, Philip L Davidson, Sameh Khamis, Mingsong Dou, et al. 2016. Holoportation: Virtual 3d teleportation in real-time. In *Proceedings of the 29th annual symposium on user interface software and technology*. 741–754.
- Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. 2021a. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5865–5874.
- Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. 2021b. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228* (2021).
- Automne Petitjean, Yohan Poirier-Ginter, Ayush Tewari, Guillaume Cordonnier, and George Drettakis. 2023. ModalNeRF: Neural Modal Analysis and Synthesis for Free-Viewpoint Navigation in Dynamically Vibrating Scenes. In *Computer Graphics Forum*, Vol. 42.
- Johannes L Schonberger and Jan-Michael Frahm. 2016. Structure-from-motion revisited. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4104–4113.
- Ruizhi Shao, Jingxiang Sun, Cheng Peng, Zerong Zheng, Boyao Zhou, Hongwen Zhang, and Yebin Liu. 2023a. Control4D: Dynamic Portrait Editing by Learning 4D GAN from 2D Diffusion-based Editor. *arXiv preprint arXiv:2305.20082* (2023).
- Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. 2023b. Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16632–16642.
- Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. 2023. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *IEEE Transactions on Visualization and Computer Graphics* 29, 5 (2023), 2732–2742.
- Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. 2019. Pushing the boundaries of view extrapolation with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 175–184.
- Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5459–5469.
- Théo Thonat, Yagiz Aksoy, Miika Aittala, Sylvain Paris, Frédéric Durand, and George Drettakis. 2021. Video-Based Rendering of Dynamic Stationary Environments from Unsynchronized Inputs. In *Computer Graphics Forum*, Vol. 40. Wiley Online Library, 73–86.
- Chaoyang Wang, Ben Eckart, Simon Lucey, and Orazio Gallo. 2021. Neural trajectory fields for dynamic novel view synthesis. *arXiv preprint arXiv:2105.05994* (2021).
- Feng Wang, Sinan Tan, Xinghang Li, Zeyue Tian, Yafei Song, and Huaping Liu. 2023. Mixed neural voxels for fast multi-view video synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 19706–19716.
- Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. 2022. Fourier plenotrees for dynamic radiance field rendering in real-time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13524–13534.
- Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Wang Xinggang. 2023. 4D Gaussian Splatting for Real-Time Dynamic Scene Rendering. *arXiv preprint arXiv:2310.08528* (2023).
- Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. 2022. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5438–5448.
- Jason C Yang, Matthew Everett, Chris Buehler, and Leonard McMillan. 2002. A real-time distributed light field camera. *Rendering Techniques* 2002, 77–86 (2002), 2.
- Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. 2023a. Deformable 3D Gaussians for High-Fidelity Monocular Dynamic Scene Reconstruction. *arXiv preprint arXiv:2309.13101* (2023).
- Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. 2023b. Real-time Photorealistic Dynamic Scene Representation and Rendering with 4D Gaussian Splatting. *arXiv preprint arXiv:2310.10642* (2023).
- Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. 2023c. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv:2310.10642* (2023).
- Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. 2021. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5752–5761.
- Heng Yu, Joel Julin, Zoltán Á Milacski, Koichiro Niinuma, and László A Jeni. 2023. CoGS: Controllable Gaussian Splatting. *arXiv preprint arXiv:2312.05664* (2023).
- Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *CVPR*.
- Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. 2018. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817* (2018).

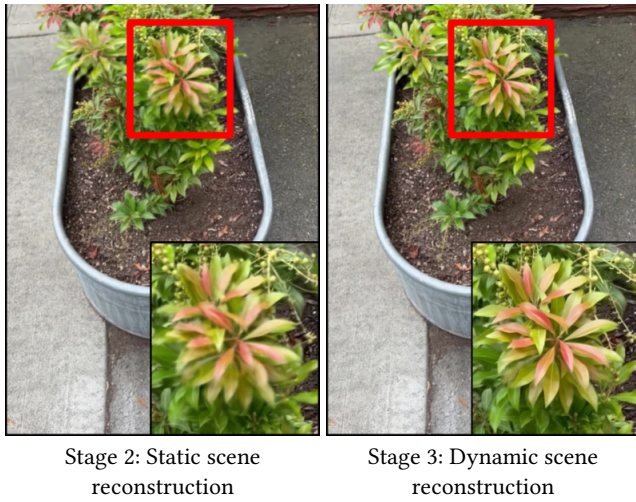


Figure 6: Effect of time-varying parameters. In stage 2, the absence of time-varying parameters leads to blurry and ghosting artifacts. In contrast, stage 3’s joint optimization of time-varying and time-independent parameters allows for accurate reconstruction of ambient motion and 3-D geometry.

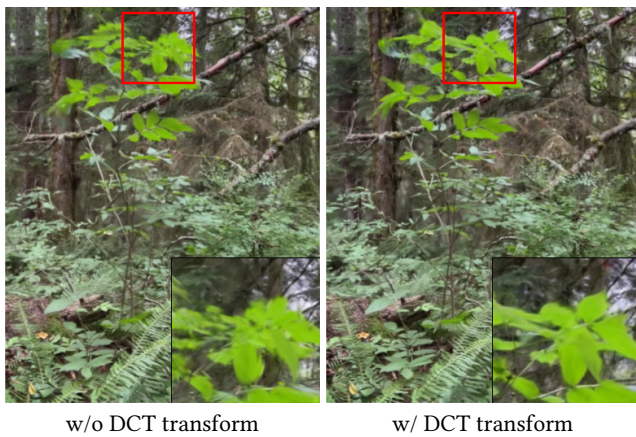


Figure 7: Effect of DCT coefficients. We synthesize an image of concealed leaves from a specific time when they’re not visible in training data, comparing two MLP predictions: DCT coefficients and direct time-varying parameters. Direct predictions lack constraints on invisible motion, causing blurry or erroneous results. Using DCT coefficients, the model generalizes leaf motion from times when leaves are visible.



Figure 8: Effect of rigidity and relative rotation loss. Introducing rigidity and relative rotation loss methods [Luiten et al. 2024] ensures that neighboring Gaussians maintain consistent behavior, reducing artifacts from random rotations. This regularization helps the Gaussians move together, creating a more natural scene with fewer odd rotations.

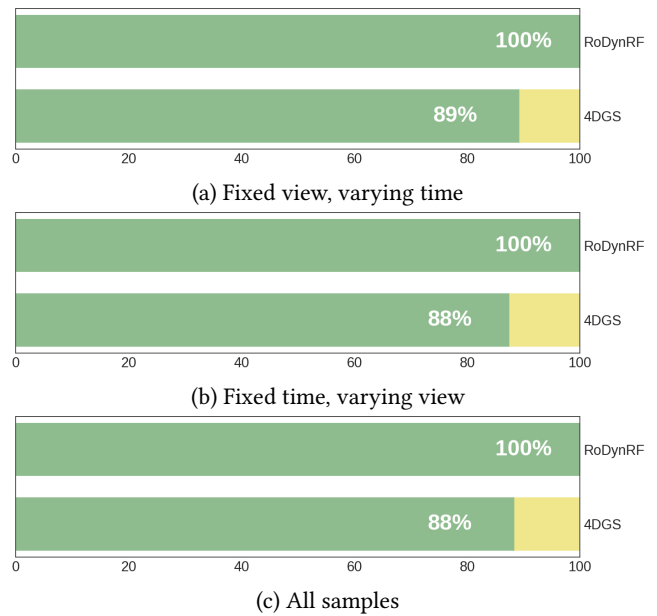


Figure 9: User study. Winrate of our method in comparison to other methods (labeled on the Y-axis). Lengths of the bars indicate the percentage of times users rated visual quality of a method higher than the competing method. The green bars with percentages correspond to our method.

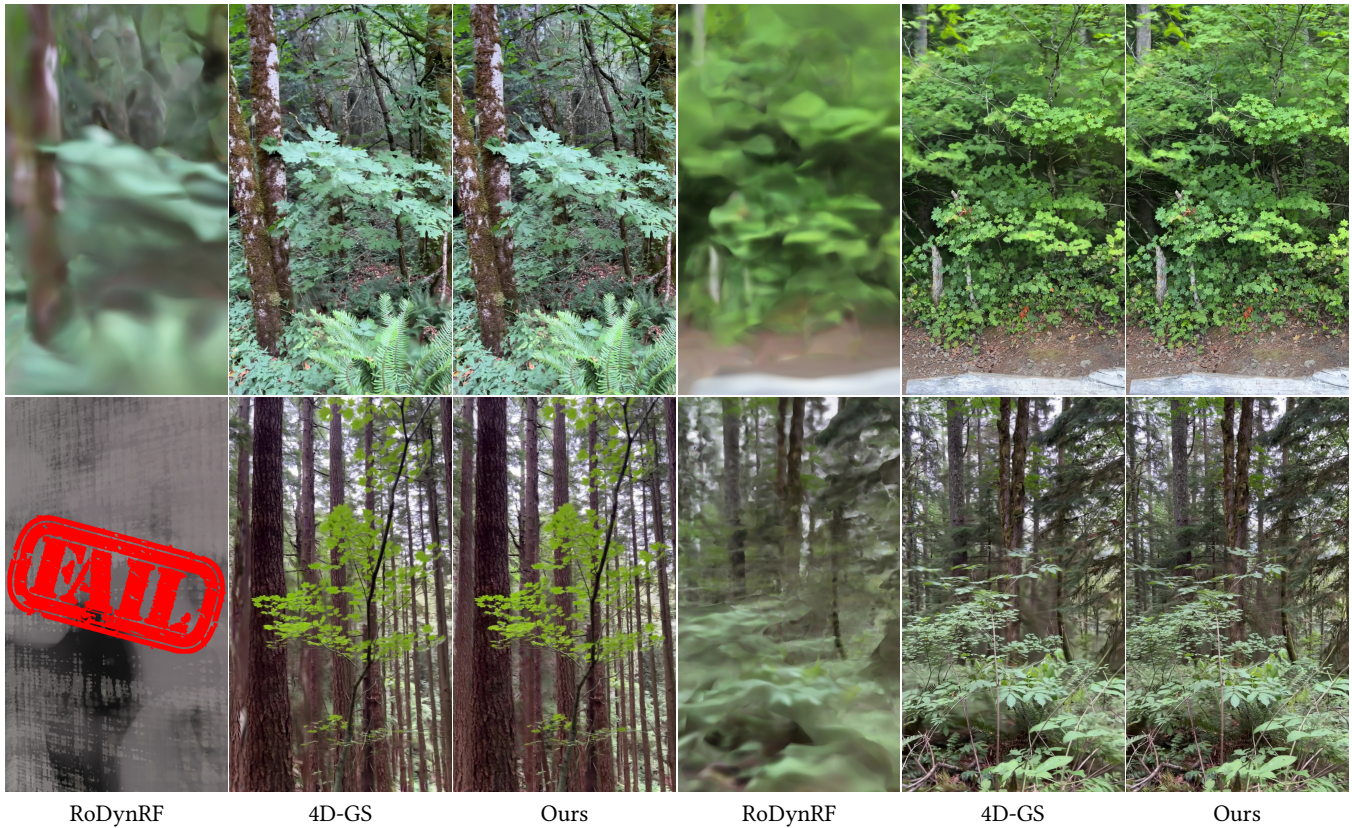


Figure 10: Novel view synthesis comparisons. We compare our method with RoDyNeRF [Liu et al. 2023] and 4D-GS [Wu et al. 2023]. Both show promising results for transient motion like moving people but struggle with complex ambient motion. RoDyNeRF struggles to produce sharp images due to the limitations of its implicit representation and constraints on optical flow. 4D-GS models part of the static scene but introduces artifacts in distant areas and fails to accurately represent ambient motion. In contrast, our method successfully captures detailed appearances and handles ambient motion efficiently.

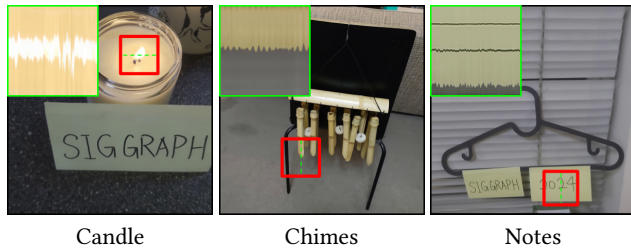


Figure 11: Diverse scenes. We showcase our method’s ability to model candles, chimes, and notes using XT- or YT-slices.

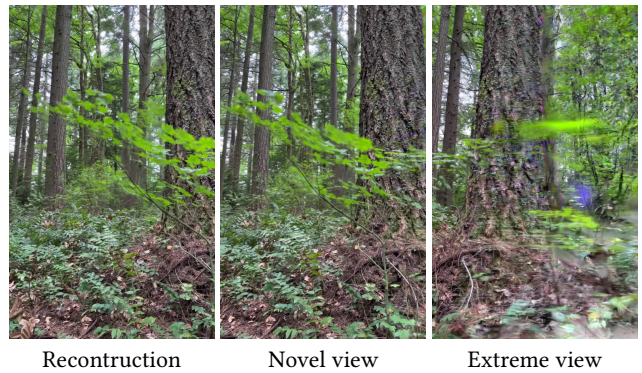


Figure 12: Failure case of extreme view synthesis (view extrapolation). We observe color shift due to the extreme, unregularized viewpoints, unseen during training.