

Gaussian-Flow: 4D Reconstruction with Dynamic 3D Gaussian Particle

Youtian Lin¹ Zuozyuo Dai² Siyu Zhu³ Yao Yao¹✉
¹Nanjing University ²Alibaba Group ³Fudan University



Figure 1. Dynamic reconstruction results of the proposed Gaussian-Flow on the monocular HyperNeRF Dataset [28] (left) and the multi-view Plenoptic Dataset [18] (right). Our method achieves a $5\times$ faster training and rendering speed compared with the per-frame 3DGS modeling and significantly outperforms previous methods in novel view rendering quality.

Abstract

We introduce *Gaussian-Flow*, a novel point-based approach for fast dynamic scene reconstruction and real-time rendering from both multi-view and monocular videos. In contrast to the prevalent NeRF-based approaches hampered by slow training and rendering speeds, our approach harnesses recent advancements in point-based 3D Gaussian Splatting (3DGS). Specifically, a novel *Dual-Domain Deformation Model (DDDM)* is proposed to explicitly model attribute deformations of each Gaussian point, where the time-dependent residual of each attribute is captured by a polynomial fitting in the time domain, and a Fourier series fitting in the frequency domain. The proposed DDDM is capable of modeling complex scene deformations across long video footage, eliminating the need for training separate 3DGS for each frame or introducing an additional implicit neural field to model 3D dynamics. Moreover, the explicit deformation modeling for discretized Gaussian points ensures ultra-fast training and rendering of a 4D scene, which is comparable to the original 3DGS designed for static 3D reconstruction. Our proposed approach showcases a substantial efficiency improvement, achieving a $5\times$ faster training speed compared to the per-frame 3DGS modeling.

In addition, quantitative results demonstrate that the proposed *Gaussian-Flow* significantly outperforms previous leading methods in novel view rendering quality. Project page: <https://nju-3dv.github.io/projects/Gaussian-Flow>.

1. Introduction

In the realm of digital scene synthesis, achieving a balance between high-quality reconstructions and real-time rendering is paramount, especially for applications like virtual reality (VR) playback, where immediate feedback and immersive experiences are essential. Neural Radiance Fields (NeRFs) [25] has risen as a promising method for synthesizing intricate scenes. However, despite their ability to produce visually stunning results, NeRFs require costly sampling and evaluation of the neural radiance field at multiple points along each ray. Consequently, the substantial computational demands impede the real-time rendering capabilities of NeRFs. These are attempts to accelerate the rendering process of NeRFs, such as direct volume representation [21, 33, 40], neural hasing [26], and tri-plane structures [2, 4, 9]. However, it still remains a challenge for

high-fidelity real-time rendering. Nevertheless, the issue becomes even more severe when turning to dynamic scene reconstruction and rendering.

Recent progress on 3D Gaussian Splatting (3DGS) [14] has drawn attention from the 3D computer vision community. With tile-based rasterization instead of plain volume rendering, 3DGS can render images two orders of magnitude faster than the vanilla NeRF. The technique has also been quickly applied to 4D scene reconstruction by extending to the separate per-frame 3DGS optimization [23]. However, such direct extension is storage-intensive and is not applicable to monocular video input. Some other concurrent works [35, 38] try to mix the explicit point-based 3DGS and an implicit neural field for dynamic information modeling, however, require computationally expensive forward passes of the neural network, which significantly lowers the rendering speed of the original 3DGS.

In this work, we propose Gaussian-Flow, an explicit particle-based deformation model designed specifically for 3DGS to model the dynamic scene without using any neural network. Gaussian-Flow can recover a high-fidelity 4D scene from captured videos, while still preserving the ultra-fast training and rendering speed of the original 3DGS. In particular, we formulate a 4D scene as a set of deformable 3D Gaussian points. A novel Dual-Domain Deformation Model (DDDM) is proposed to explicitly model deformations of each Gaussian point’s attributes, including position, rotation, and radiance. The time-dependent deformation residual is modeled simultaneously in time and frequency domains: we apply joint polynomial and Fourier series fitting for each deformable attribute. This compact dynamic representation greatly reduces the computation cost of the deformation model, which is a key factor in preserving the rendering speed of 3DGS. Moreover, an adaptive timestamp scaling technique is introduced to avoid over-fitting the scene to only frames with violent motions. For robust estimation, we also regularize the motion trajectory by the KNN-based rigid and the time smooth constraints. It is also noteworthy that our discretized point-based 4D representation naturally supports the edition of both static and dynamic 3D scenes, showing the potential for unlocking a variety of downstream applications related to dynamic 3D reconstruction and rendering.

We have conducted extensive experiments to demonstrate the effectiveness of the proposed method on several multi-view and monocular datasets. The proposed Gaussian-Flow achieves a $5\times$ faster training speed compared with the separate per-frame 3DGS modeling, and significantly outperforms prior leading methods in novel view rendering quality. Our major contributions can be summarized as follows:

- We introduce Gaussian-Flow, a novel point-based differentiable rendering approach for dynamic 3D scene recon-

struction, setting a new state-of-the-art for training speed, rendering FPS, and novel view synthesis quality for 4D scene reconstruction.

- We propose a Dual-Domain Deformation Model for efficient 4D scene training and rendering, eliminating the need for per-frame 3DGS optimization and sampling on implicit neural fields. This preserves a running speed on par with the original 3DGS with minimum overhead.
- We demonstrate that our discretized point-based representation supports the segmentation, edition, and composition of both static and dynamic 3D scenes.

2. Related Works

2.1. Dynamic Neural Radiance Field

Dynamic NeRF modeling has become a heated research topic in recent years due to the development of the neural radiance field and differentiable rendering. By treating time as an extended input dimension to NeRF, researchers successfully achieve qualified image-based 4D scene rendering [5, 10, 19, 22, 37]. To further improve the reconstruction quality and incorporate prior knowledge of motions and structures, dynamic neural scene flow methods have been proposed [27, 30], where a canonical space is constructed and then transferred to each frame with scene flow or motion fields. HyperNeRF [28] models the deformation of the object topologies by using higher-dimensional inputs, while DyNeRF [18] utilizes time-conditioned NeRF to represent a 4D scene. However, the aforementioned approaches are all based on the vanilla NeRF, which requires a long training time and does not meet the requirement of real-time rendering.

2.2. Accelerated Neural Radiance Field

To expedite NeRF training and rendering, numerous approaches have been suggested, employing more streamlined strategies [6, 11, 12, 16, 20, 29, 36]. Other methods propose the integration of neural implicit functions with explicit 3D structures, forming a hybrid representation for faster radiance field sampling [4, 26, 33, 34]. These approaches establish strong foundations for enhancing dynamic NeRF, concurrently decreasing required training and inference time. Apart from implicit neural representations, explicit NeRF modeling has shown promising results for real-time rendering: NSVF [21] employs a neural sparse voxel field for efficient NeRF sampling, which stands for the earliest attempt on the explicit NeRF modeling; PlenOctrees [40] utilizes the explicit octree structure for rendering acceleration.

Recent efforts have also emerged to accelerate the intricate dynamic neural radiance field. TensorRF [4] employs multiple planes as explicit representations for direct dynamic scene modeling. More recent approaches of this kind include K-Planes [8], Tensor4D [31], and Hex-

Plane [3]. Alternatively, NeRFPlayer[32] introduces a unified streaming representation for both grid-based [26] and plane-based methods, utilizing separate models to distinguish static and dynamic scene components, however, leading to slow rendering times. HyperReel [1] further suggests a flexible sampling network coupled with two planes for dynamic scene representation. While these methods improve the rendering speed of a dynamic scene to some extent, it is still hard to achieve real-time rendering, let alone a good balance between the running speed and rendering quality. In contrast, we resort to the recent 3D Gaussians splatting, which applies an explicit soft point cloud representation for real-time image-based rendering.

2.3. Differentiable Point-based Rendering

The original idea of using 3D points as rendering primitives was first introduced in [17]. By incorporating differentiable rendering, recent approaches have made remarkable progress in image-based rendering, representative methods include PointRF [41], DSS [39], and 3D Gaussians splatting (3DGS) [13]. Specifically, 3DGS [13] has demonstrated extraordinary performance in novel-view synthesis, achieving real-time rendering speed and state-of-the-art rendering quality. The method adopts a soft point representation with attributes of position, rotation, density, and radiance, and applies differentiable point-based rendering for scene optimization. 3DGS has quickly been extended to dynamic scene modeling by direct separate per-frame optimization [23], however, requires a long optimization time and a large amount of storage for long video footage. Other works [35, 38] apply an implicit motion field to model scene dynamics, but the introduction of the implicit neural network significantly slows down the sampling and rendering speed. In this work, we adopt the approach of representing 4D scenes through a purely discretized point cloud model, ensuring a fast training and rendering speed comparable with the original 3DGS.

3. Gaussian-Flow

In this section, we introduce the proposed Gaussian-Flow for dynamic scene modeling. We first review the 3DGS in Sec. 3.1. Then, we introduce our explicit motion modeling of each Gaussian point by using a novel Dual-Domain Deformation Model (DDDM), as outlined in Sec. 3.2. An adaptive timestamp scaling technique is described in Sec. 3.3 for balanced training of each frame. To ensure the continuity of the motion in both spatial and temporal dimensions, we incorporate appropriate regularizations on each point during the optimization, as detailed in Sec. 3.4.

3.1. Recap on 3D Gaussian Splatting

3D Gaussian Splatting [14] is designed to efficiently optimize a 3D scene for real-time and high-quality novel view

synthesis. The 3DGS framework has garnered significant attention within the community due to its remarkable enhancements in both training and rendering times, concurrently achieving state-of-the-art rendering quality. In contrast to the volume rendering in the vanilla NeRF which relies on ray marching, 3DGS adopts a tile-based rasterization on a distinctive soft point cloud representation to achieve fast rendering. Specifically, 3DGS models a 3D scene as a large amount of 3D Gaussian points in the world space, where each point is represented by:

$$G(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \quad (1)$$

where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are the mean position and covariance matrix of a 3D Gaussian particle. The 3DGS takes sparse Structure-from-Motion (SfM) points or even random points as input, and initializes each point to a 3D Gaussian based on its neighbors. Besides, each 3D Gaussian is associated with a learnable view-dependent radiance \mathbf{c} and a learnable opacity α for rendering. Subsequently, an efficient 3D to 2D Gaussian mapping [42] is employed to project the point onto the image plane:

$$\boldsymbol{\mu}' = \mathbf{P}\mathbf{W}\boldsymbol{\mu}, \quad (2)$$

$$\boldsymbol{\Sigma}' = \mathbf{J}\mathbf{W}\boldsymbol{\Sigma}\mathbf{W}^T\mathbf{J}^T, \quad (3)$$

where $\boldsymbol{\mu}'$ and $\boldsymbol{\Sigma}'$ represent the 2D mean position and 2D covariance of the projected 3D Gaussian. \mathbf{P} , \mathbf{W} and \mathbf{J} denote the projective transformation, viewing transformation, and Jacobian of the affine approximation of \mathbf{P} . After that, α -blending is executed to merge the overlapping Gaussians for each pixel, yielding a final color:

$$\mathbf{C} = \sum_{i=1}^n \mathbf{c}_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j), \quad (4)$$

where \mathbf{c}_i and α_i are the color and opacity of the i -th Gaussian, and n is the number of overlapping Gaussians.

The attributes of the 3D Gaussian, including the mean $\boldsymbol{\mu}$, covariance matrix $\boldsymbol{\Sigma}$, opacity α , and color \mathbf{c} , are optimized through backward propagation of the gradient flow. In particular, to ensure positive definiteness during the optimization process, the covariance matrix is parameterized as a scaling vector \mathbf{s} and a rotation matrix \mathbf{R} , i.e., $\boldsymbol{\Sigma} = \mathbf{R}\boldsymbol{\Lambda}(\mathbf{s})\boldsymbol{\Lambda}(\mathbf{s})^T\mathbf{R}^T$, where $\boldsymbol{\Lambda}(\mathbf{s})$ is the diagonal matrix of \mathbf{s} . To facilitate optimization, the rotation matrix is further parameterized using a quaternion \mathbf{q} . Leveraging the inherent flexibility of discrete points, 3DGS incorporates adaptive density control for points. This mechanism utilizes the gradient flow to identify where geometric reconstruction is suboptimal, and employs cloning and splitting to augment the density of points for a higher rendering quality.

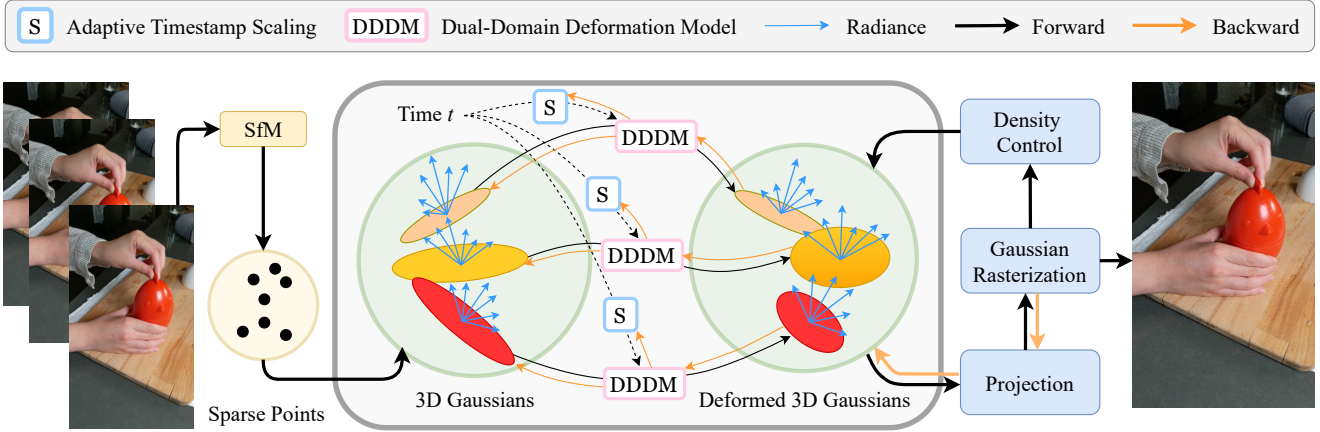


Figure 2. **Overview of the Gaussian-Flow pipeline.** We model the deformation of attributes of each 3D Gaussian point independently by using the Dual-Domain Deformation Model (DDDM), which preserves the discretized nature of the 3D Gaussian points, and thus achieves ultra-fast training and rendering speed comparable with the original 3DGS.

3.2. Dual-Domain Deformation Model

We target to directly model the dynamics of each 3D Gaussian point by fitting each of its attributes into a time-dependent curve. Among different approaches, Polynomials fitting in the time domain and Fourier series fitting in the frequency domain are the two most widely used approaches [5, 10, 19, 22, 37], due to their simplicity and effectiveness. However, each method comes with its own advantages and drawbacks: describing the motion of a Gaussian particle in terms of polynomials yields a good fit with smooth motion with a small order of polynomials, however, can easily overfit to a violent motion if assuming a larger order of polynomials, resulting in unreasonable oscillations in the fitted trajectory. Whereas, the Fourier series excels at capturing the variations associated with violent motion, however, requires a manually reduced order when dealing with smooth motion.

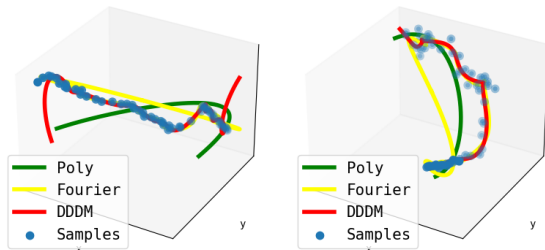


Figure 3. Two exemplar motion fittings using polynomial, Fourier Series, and the joint DDDM functions. Our DDDM is able to accurately fit complex trajectories denoted by the sampling points.

In this work, our key insight is to use a Dual-Domain Deformation Model (DDDM) for fitting the scene dynamics,

which integrates both the time domain polynomials and the frequency domain Fourier series into a unified fitting model. We assume that only the rotation \mathbf{q} , radiance c , and position $\boldsymbol{\mu}$ of a 3D Gaussian particle change over time, while the scaling s and opacity α remain constant. Specifically, we conceptualize the change in each particle’s attributes as its base attributes $\mathbf{S}_0 \in \{\boldsymbol{\mu}_0, c_0, \mathbf{q}_0\}$ at the reference time frame t_0 (usually set to the first frame), superimposed on a time-dependent attribute residual $\mathbf{D}(t)$. For simplicity, we use lowercase characters to represent a single attribute in \mathbf{S} . The time-dependent residual of each attribute is modeled through polynomial fitting in the time domain and Fourier series fitting in the frequency domain, expressed as:

$$S(t) = S_0 + D(t), \quad (5)$$

where $D(t) = P_N(t) + F_L(t)$ is combined by a polynomial $P_N(t)$ with coefficients $\mathbf{a} = \{a_n\}_{n=0}^N$ and a Fourier series $F_L(t)$ with coefficients $\mathbf{f} = \{f_{sin}^l, f_{cos}^l\}_{l=0}^L$. These are respectively denoted as:

$$P_N(t) = \sum_{n=0}^N a_n t^n, \quad (6)$$

$$F_L(t) = \sum_{l=1}^L (f_{sin}^l \cos(lt) + f_{cos}^l \sin(lt)). \quad (7)$$

It is important to note that we assume different dimensions of an attribute are independently changed over time. Therefore, we assign a different $D(t)$ for each dimension of an attribute. For instance, we utilize $\{D_{\mu_i}(t)\}_{i=0}^3$ to describe the motion of a 3D position $\boldsymbol{\mu}$.

Figure 3 illustrates a comparative analysis of trajectory fitting using polynomial, Fourier series, and the proposed

joint DDDM functions. The figure highlights the superior fitting capabilities of the DDDM approach in capturing complex motion trajectories as represented by the sampled data points.

3.3. Adaptive Timestep Scaling

In a typical scenario, a normalized frame index t ranging from 0 to 1 will be used as the temporal input of $D(t)$. However, this poses a challenge when endeavoring to model substantial motions within a very short time using polynomials and Fourier series. Adhering to the standard temporal division would necessitate an exceedingly large coefficient to accommodate highly intense movements within a very short time frame. This circumstance has the potential to induce instability or even a breakdown in the optimization process. To address this issue, we introduce a time dilation factor λ to scale the temporal input for each Gaussian point, which is formulated as:

$$t_s = \lambda_s \cdot t + \lambda_b \quad (8)$$

where t_s represents the scaled time input, serving as the input of $D(t)$, $t \in [0, 1]$ denotes the normalized frame index, and λ and λ_b stand for the dilation factor and base factor of a Gaussian, respectively. In all our experiments, λ_s and λ_b are initialized to 1 and 0, respectively.

To summarize, in our dynamic scene setting, a Gaussian particle contains multiple attributes to be optimized, including base attributes $\{\mu_0, q_0, s_0, c_0, \alpha_0\}$ at reference frame t_0 , polynomial coefficients and Fourier coefficients in $\{D_\mu(t), D_q(t), D_c(t)\}$. Since the proposed DDDM optimizes the time-dependent residuals without the need for an intricate neural field structure, our Gaussian-Flow inherits the benefits of extreme-fast training and rendering speed from the vanilla 3DGS.

3.4. Regularizations

While the utilization of discrete points as scene representations accelerates rendering, several challenges remain. First, points are optimized individually, losing connections with their spatial neighbors, which do not align with the real-world scenario. Optimizing these Gaussian points without considering continuity will inevitably lead to a degradation in reconstruction quality and spatial coherence. Additionally, motions should be smoothed over time. Based on these observations, we employ two regularizations, namely time smoothness and a KNN rigid regularization, for robust optimization of Gaussian points and their motions.

Time Smoothness Loss To ensure temporal smoothness over time, we apply a perturbing ϵ on the input timestamp t and encourage the time-dependent attributes (i.e., position

μ , rotation q and radiance c) at time $t + \epsilon$ to be consistent with those at time t . The time smoothness term is defined as:

$$\mathcal{L}_t = \|D(t) - D(t + \epsilon)\|_2. \quad (9)$$

It is noteworthy that the magnitude of the perturbing value is adaptively set according to the number of total frames, i.e., $\epsilon = 0.1/\text{frames}$.

KNN Rigid Loss During the optimization of 3D Gaussians with adaptive density control, points will be dynamically added or removed. This dynamic nature implies that the neighbors of points within a local space are subject to constant changes, posing a challenge for directly enforcing a spatial local consistency constraint. To sidestep this problem, we propose to divide the optimization part into two alternating stages: in the former stage, we optimize all variables with adaptive density control; while in the latter stage, we optimize the attributes without adding or removing points. The local rigid constraint is incorporated in every latter stage, and it is defined as:

$$\mathcal{L}_s = \sum_{j \in \mathcal{N}_i} \|D(t)_i - D(t)_j\|_2 \quad (10)$$

where, \mathcal{N}_i represents the K nearest neighbor (KNN) of i -th Gaussian.

4. Experiments

4.1. Implementation Details

We train our model using the Adam [15] optimizer with separate learning rates for different attributes of the Gaussian point. We set a learning rate of 4×10^{-4} for the point position with an exponential decay of 8×10^{-7} . The learning rates for point rotation and all DDDM parameters are set to 0.002 and 4×10^{-4} respectively. We apply a weight decay of 8×10^{-7} to all parameters. The rest of the learning rate follows the 3DGS setting. We train the model for 30K steps and 60K steps for all scenes. All experiments are conducted on a single NVIDIA RTX 4090 GPU with 24GB memory. In addition, we use Taichi to implement our DDDM model, which can parallelize the computation of the DDDM (i.e. polynomial and Fourier series computation) of each Gaussian point.

4.2. Datasets

We evaluate our method on both multi-view and monocular datasets, to demonstrate the effectiveness of our method in both settings.

Plenoptic Video dataset [18] The dataset was captured using 21 cameras at a resolution of 2704×2028 , with each camera recording a 10-second video. Six scenes from this

dataset are publicly available. For a fair comparison, we downsampled the images to 1352×1014 resolution in our experiments to keep the same setting from the concurrent work of 4D Gaussian [35].

HyperNeRF dataset [28] This dataset uses a monocular camera (e.g., iPhone) to record real-world motions, which includes real rigid and non-rigidly deforming scenes, such as a person splitting a cookie. The dataset is rather challenging due to large motions, complex lighting conditions, and thin object structures. To ensure a fair comparison, we downsampled images to 540×960 in our experiments and followed the training and validation camera split provided by [28]. We conducted experiments on the four "vrig" scenes, and we also provided results of the "interp" scenes in the supplementary material.

4.3. Ablation Study

Deformation Models The deformation model is the core component of the proposed Gaussian-Flow. We conduct ablation studies to validate the effectiveness of the particular choice of DDDM. As DDDM consists of a Fourier series and a polynomial function, hence we first study the Fourier series and the polynomial functions separately, in which we only use the Fourier series or the polynomial function as the deformation model. As shown in Figure. 4, the Fourier series contains more high-frequency components than the polynomial function, thus the Fourier series has sharper image details but results in more artifacts. The polynomial function is smoother than the Fourier series, leading to fewer artifacts but blurry scene renderings. Finally, the hybrid DDDM function is able to generate sharper details with fewer artifacts.

Furthermore, we study the orders of the polynomial and Fourier series functions in our DDDM, which are related to the complexity of the scene and are crucial to the final performance. As shown in Figure. 5, the performance of our method increases with the number of Fourier series orders but starts to drop after the order over 32, which should be related to the over-parameterization of the deformation model.

Regularizations Next, we study the effectiveness of the two proposed regularizations in Gaussian-Flow optimization. As shown in Table. 1, adding separate KNN rigid regularization or the time smooth regularization can both improve the novel view rendering quality, and the full model that contains both regularizations can achieve the best performance.

4.4. Quantitative Comparisons

We compare our method against previous SOTA NeRF-based methods, including NeRF [25], Nerfies [27], Hyper-

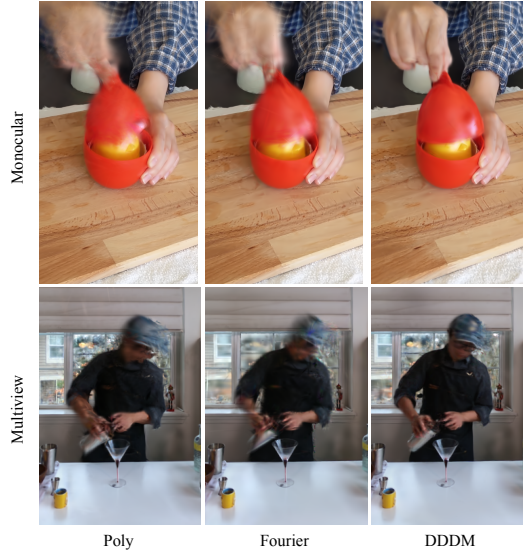


Figure 4. Ablation study on different deformation models. From left to right are deformation fitting with polynomial function only, Fourier series only, and our dual-domain deformation fitting. The proposed DDDM achieves the best rendering quality qualitatively.

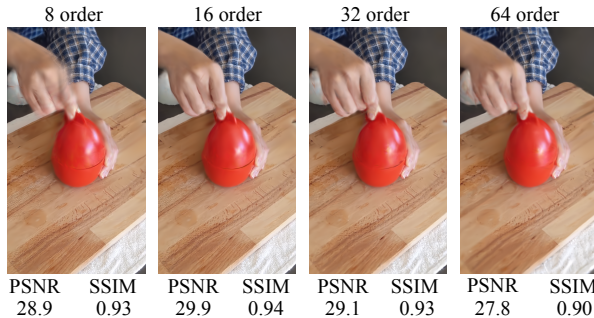


Figure 5. Ablation study on different orders of the DDDM. We find an order number of 16 leads to the best novel view rendering results in the HyperNeRF dataset.

reg,	w/o KNN rigid	w/o Time smooth	full model
PSNR	28.48	29.12	29.92
SSIM	0.92	0.93	0.94

Table 1. Ablation study on the proposed KNN rigid and time smooth regularizations. The quantitative results demonstrate the effectiveness of both regularizations.

NeRF [28], NeRFPlayer [32], and TiNeuVox [7]. We also provide comparisons with other 3DGS-based approaches concurrently proposed with our Gaussian-Flow. The training time, rendering FPS, and novel view synthesis PSNR of different methods can be found in Table. 2. Previous NeRF-based methods require at least 30 minutes to train the scene, and fail to achieve real-time rendering of the dynamic scene.

Table 2. Per-scene quantitative comparisons on HyperNeRF [28] dataset. Results are gathered from papers of the corresponding methods. Our method achieves the fastest training time, the highest rendering FPS, and the highest PSNR score for novel view synthesis, setting a new state-of-the-art for image-based dynamic scene rendering.

Method	Train Time↓	Render FPS↑	Broom		3D Printer		Chicken		Peel Banana		Mean	
			PSNR↑	SSIM ↑	PSNR↑	SSIM ↑	PSNR↑	SSIM ↑	PSNR↑	SSIM ↑	PSNR↑	SSIM↑
NeRF [25]	16 hours	0.013	19.9	0.653	20.7	0.780	19.9	0.777	20.0	0.769	20.1	0.745
Nerfies [27]	16 hours	0.011	19.2	0.567	20.6	0.830	26.7	0.943	22.4	0.872	22.2	0.803
HyperNeRF [28]	32 hours	0.011	19.3	0.591	20.0	0.821	26.9	0.948	23.3	0.896	22.4	0.814
NeRFPlayer [32]	6 hours	0.208	21.7	0.635	22.9	0.810	26.3	0.905	24.0	0.863	23.7	0.803
TiNeuVox [7]	30 min	0.5	21.5	0.686	22.8	0.841	28.3	0.947	24.4	0.873	24.3	0.837
Ours (30K)	7 min	125	22.5	0.690	24.3	0.857	29.4	0.934	26.3	0.906	25.6	0.847
Ours (60K)	12 min	125	22.8	0.709	25.0	0.877	30.4	0.945	27.0	0.917	26.3	0.862

Our method only requires **7 minutes** of training time and can achieve real-time rendering speed, which is much faster than previous methods. Moreover, our method can achieve better performance than previous SOTA methods in terms of PSNR.

Table 3. Quantitative comparison on Plenoptic Video dataset [18]. Our training speed is 5× faster of magnitude faster than previous leading approaches. Also, we achieved the highest PSNR score among all methods.

Method	Train Time ↓	PSNR ↑	SSIM ↑
LLFF [24]	-	23.2	-
DyNeRF [18]	1344 hours	29.6	0.96
NeRFPlayer [32]	5.5 hours	30.7	-
K-Planes [8]	1.8 hours	31.6	0.96
4D-GS [35]	2 hours	31.0	0.94
Ours (30K)	22.5 min	30.5	0.97
Ours (60K)	41.8 min	32.0	0.97

We evaluated various methods on the Plenoptic Video dataset, as summarized in Table. 3. The comparison focuses on training time efficiency and image quality, assessed through PSNR and SSIM. Our approach demonstrated a significant advancement in training efficiency, requiring only 22.5 minutes, a drastic reduction compared to the hours needed by methods like DyNeRF [18] and K-Planes [8]. This efficiency is paramount for practical applications, where reduced training time can be a critical factor. In terms of image quality, our method achieved a PSNR of 30.5 with 30K steps, which, while not the highest, is competitive with the leading methods. However, our method scored 0.97 in SSIM, higher than K-Planes’ leading score of 0.96. This indicates a potential trade-off between training efficiency and the ability to preserve structural details in images.

In addition, we extend our method to 60K steps on both datasets, which can further improve the performance on the

HyperNeRF dataset [28], and achieve the highest performance on the Plenoptic Video dataset [18]. However, the training time is also increased by approximately 2× for HyperNeRF dataset [28] and about 1.5× for Plenoptic Video dataset [18], which is still much faster than previous methods.

4.5. Qualitative Comparisons

In this section, we show qualitative comparisons of our method and previous SOTA methods on the HyperNeRF [28] dataset and the Plenoptic Video dataset [18]. Figure. 6 shows the qualitative comparisons of our method and TiNeuVox [7], HyperNeRF [28], Nerfies [27] and NeRFPlayer [8] on HyperNeRF dataset. Notice that our method can produce comparably clear and sharp images than previous SOTA methods, which highlights our method’s superior performance in monocular conditions. Despite its overall efficacy, our method does encounter limitations in extremely thin structure, as shown in the 3D Printer scene, the thread of the 3D Printer is not clear, while other methods can produce a clear thread.

We also compare our method with previous SOTA methods on the Plenoptic Video dataset [18], as shown in Figure. 7. Compared with previous SOTA methods, our method can produce more accurate color and correct structure. Moreover, our method successfully reconstructs the flame in the scene, while NeRFPlayer [8] fails to reconstruct the flame. These results suggest that our method can achieve comparable image quality with previous SOTA methods, which demonstrates the effectiveness of our method in the multiview conditions.

5. Conclusion

In this paper, we introduced Gaussian-Flow, a novel framework for dynamic 3D scene reconstruction using a point-based differentiable rendering approach. The core of our innovation lies in the DDDM, which efficiently models deformations of each 3D Gaussian point in both the time and

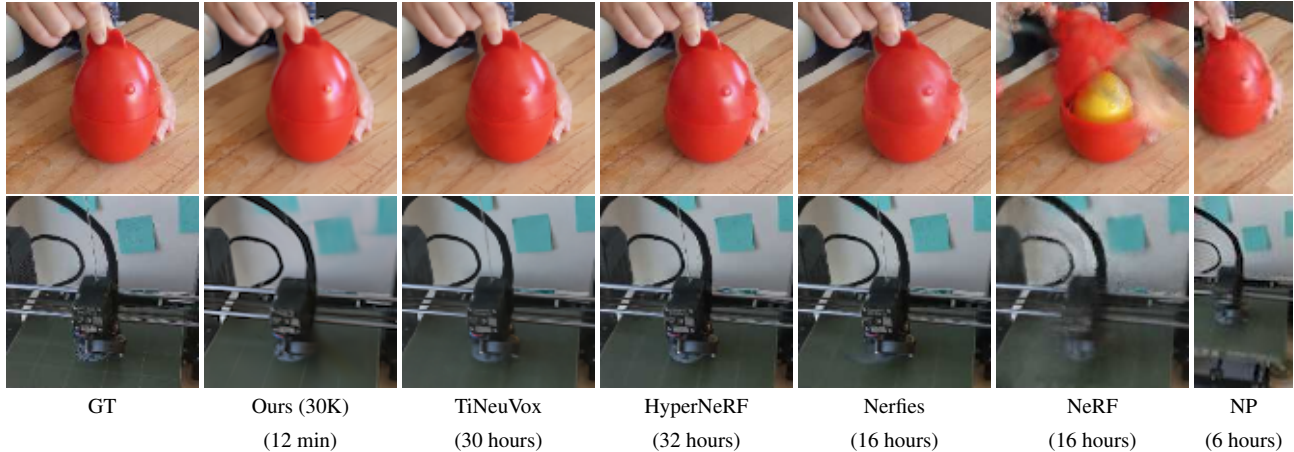


Figure 6. Qualitative comparisons of our method and TiNeuVox [7], HyperNeRF [28], Nerfies [27] and NeRFPlayer (NP) [8] on the HyperNeRF [28] dataset. The training time of each method is shown in the brackets.

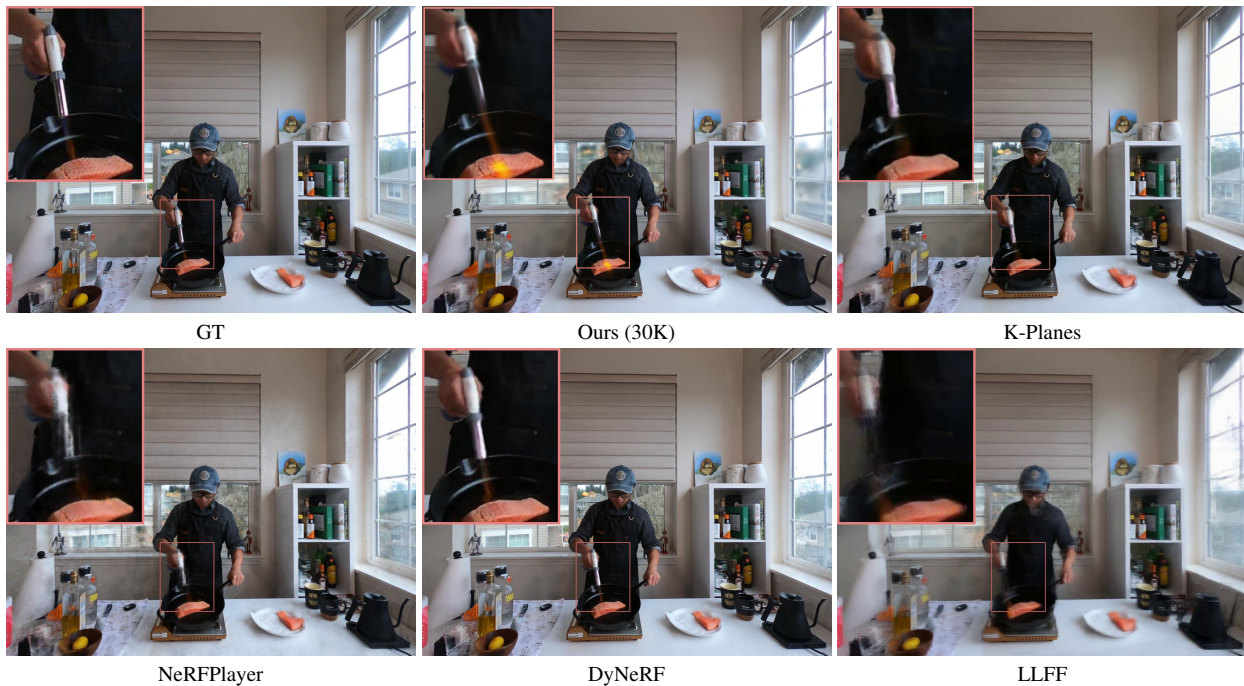


Figure 7. Qualitative comparisons of our method and K-Planes [8], NeRFPlayer [8], DyNeRF [18], and LLFF on Plenoptic Video dataset.

frequency domains. This approach has enabled us to set a new state-of-the-art for 4D scene reconstruction in terms of training speed, rendering frames per second, and novel view synthesis quality. Our extensive experiments and ablation studies have demonstrated the efficacy of proposed Gaussian-Flow across various datasets. We achieved significant improvements over existing methods, particularly in training speed and rendering performance. The ability to efficiently handle dynamic scenes without the computational overhead of neural networks marks a substantial leap

forward in this domain.

6. Limitations

While our method excels in rendering speed and training efficiency, there is room for improvement in maintaining high-fidelity thin structures in the final rendering. Future work could focus on enhancing the balance between speed and image detail preservation, potentially through more refined deformation models or advanced regularization techniques.

References

- [1] Benjamin Attal, Jia-Bin Huang, Christian Richardt, Michael Zollhoefer, Johannes Kopf, Matthew O’Toole, and Changil Kim. HyperReel: High-fidelity 6-DoF video with ray-conditioned sampling. *arXiv preprint arXiv:2301.02238*, 2023. [3](#)
- [2] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. *arXiv:2301.09632*, 2023. [1](#)
- [3] Ang Cao and Justin Johnson. Hexplane: a fast representation for dynamic scenes. *arXiv preprint arXiv:2301.09632*, 2023. [3](#)
- [4] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *Proceedings of the European Conference on Computer Vision*, 2022. [1](#), [2](#)
- [5] Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14304–14314. IEEE Computer Society, 2021. [2](#), [4](#)
- [6] Jiemin Fang, Lingxi Xie, Xinggang Wang, Xiaopeng Zhang, Wenyu Liu, and Qi Tian. Neusample: Neural sample field for efficient view synthesis. *arXiv:2111.15552*, 2021. [2](#)
- [7] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. *arXiv preprint arXiv:2205.15285*, 2022. [6](#), [7](#), [8](#)
- [8] Sara Fridovich-Keil, Giacomo Meanti, Frederik Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. *arXiv preprint arXiv:2301.10241*, 2023. [2](#), [7](#), [8](#)
- [9] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance, 2023. [1](#)
- [10] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5712–5721, 2021. [2](#), [4](#)
- [11] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14346–14355, 2021. [2](#)
- [12] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5875–5884, 2021. [2](#)
- [13] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (ToG)*, 42(4):1–14, 2023. [3](#)
- [14] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023. [2](#), [3](#)
- [15] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015. [5](#)
- [16] Andreas Kurz, Thomas Neff, Zhaoyang Lv, Michael Zollhoefer, and Markus Steinberger. Adanerf: Adaptive sampling for real-time rendering of neural radiance fields. In *European Conference on Computer Vision*, 2022. [2](#)
- [17] Marc Levoy and Turner Whitted. The use of points as a display primitive. 1985. [3](#)
- [18] Tianye Li, Mira Slavcheva, Michael Zollhöfer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5521–5531, 2022. [1](#), [2](#), [5](#), [7](#), [8](#)
- [19] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. 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*, 2021. [2](#), [4](#)
- [20] D. B.* Lindell, J. N. P.* Martel, and G. Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021. [2](#)
- [21] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *Advances in Neural Information Processing Systems*, 2020. [1](#), [2](#)
- [22] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: learning dynamic renderable volumes from images. *ACM Transactions on Graphics*, 38(4):1–14, 2019. [2](#), [4](#)
- [23] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023. [2](#), [3](#)
- [24] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (Proceedings of SIGGRAPH)*, 38(4), 2019. [7](#)
- [25] 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. In *European Conference on Computer Vision*, pages 405–421. Springer, 2020. [1](#), [6](#), [7](#)
- [26] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. [1](#), [2](#), [3](#)
- [27] 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. [2](#), [6](#), [7](#), [8](#)
- [28] 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. *ACM Trans. Graph.*, 40(6), 2021. 1, 2, 6, 7, 8
- [29] Martin Píala and Ronald Clark. Terminerf: Ray termination prediction for efficient neural rendering. *International Conference on 3D Vision*, pages 1106–1114, 2021. 2
- [30] 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. 2
- [31] Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. *arXiv preprint arXiv:2211.11610*, 2022. 2
- [32] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *arXiv preprint arXiv:2210.15947*, 2022. 3, 6, 7
- [33] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5459–5469, 2022. 1, 2
- [34] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Improved direct voxel grid optimization for radiance fields reconstruction. *arXiv preprint arXiv:2206.05085*, 2022. 2
- [35] 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. 2, 3, 6, 7
- [36] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yu-Xiong Wang, and David Forsyth. Diver: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16200–16209, 2022. 2
- [37] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9421–9431, 2021. 2, 4
- [38] 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. 2, 3
- [39] Wang Yifan, Felice Serena, Shihao Wu, Cengiz Öztireli, and Olga Sorkine-Hornung. Differentiable surface splatting for point-based geometry processing. *ACM Transactions on Graphics (TOG)*, 38(6):1–14, 2019. 3
- [40] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 1, 2
- [41] 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*, 2022. 3
- [42] Matthias Zwicker, Hanspeter Pfister, Jeroen Van Baar, and Markus Gross. Ewa volume splatting. In *Proceedings Visualization, 2001. VIS'01.*, pages 29–538. IEEE, 2001. 3

Gaussian-Flow: 4D Reconstruction with Dynamic 3D Gaussian Particle

Supplementary Material

A. Implementation Details

In this paper, we focus exclusively on modeling three key attributes of the 3D Gaussian Splatting (3DGS) using our novel DDDM model. These attributes are: 1) the position of the Gaussian, 2) the rotation represented by a quaternion, and 3) the first three coefficients of the Spherical Harmonics (SHs). The learning approach employed for each of these attributes within the DDDM framework mirrors that of the corresponding 3DGS attribute, ensuring consistency in our modeling strategy.

We first train each scene with no deformation (as a 3DGS) for 2000 iterations, and then train the scene with deformation (with DDDM) for the remaining training phase. We stop the process of adding (through splitting and cloning as delineated in 3DGS) and pruning Gaussian points as 15K iterations. Then, We start using the KNN rigid loss at 5000 iterations, since the number of Gaussian points is fixed, which is more computationally efficient, because we can only compute the KNN index once instead of calculating the KNN for each iteration.

B. More Results

This section presents additional visual results. These include a broader range of viewpoints and scenes, highlighting the capability of our method in rendering novel view variants across both spatial and temporal dimensions. We also showcase the proficiency of our approach in reconstructing depth maps.

As shown in Figure 8, we show more results on scene *americano*, *chickchicken* and *split cookie*. As shown in Figure 9, we show the rendering and depth map results on Plenoptic Video dataset rendered at more viewpoints and time.



Figure 8. View Synthesis Results on HyperNeRF Dataset.

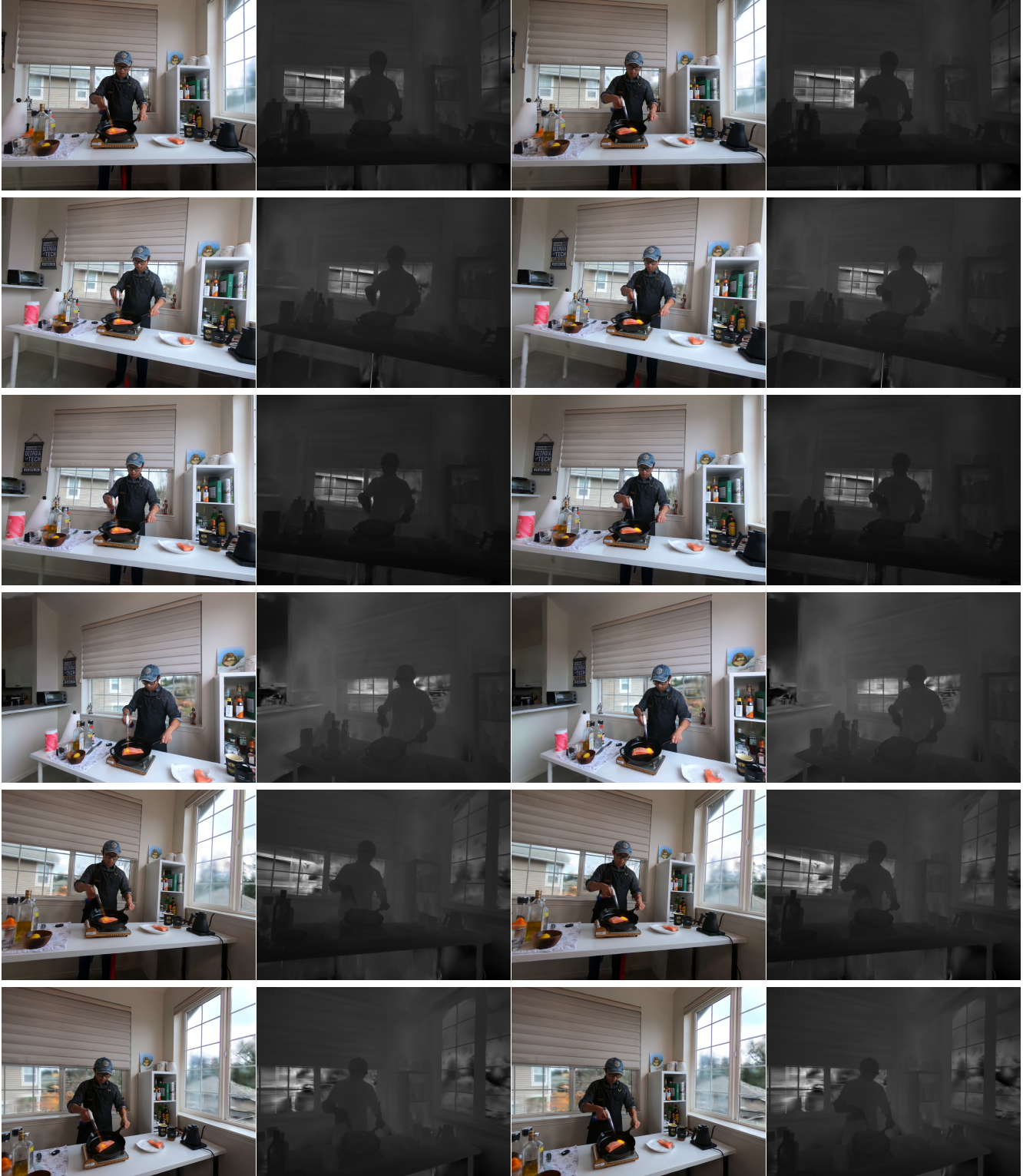


Figure 9. View Synthesis Results and Depths on Plenoptic Video Dataset.