002

003

004

005

006

007

800

009

010

011

012

013 014

015

016

017

018

019

020

021

022

023

024

025

026

027

028

029

030

031

032

033

034

035

036

037

038

039

040

041

# PhysTwin: Physics-Informed Reconstruction and Simulation of Deformable Objects from Videos

# Anonymous ICCV submission

# Paper ID 11564

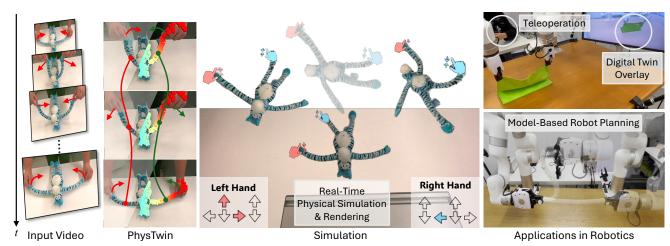


Figure 1. **PhysTwin** takes sparse videos of deformable objects under interaction as input and automatically reconstructs a simulatable digital twin with complete geometry, high-fidelity appearance, and accurate physical parameters. This enables multiple applications, such as real-time interactive simulation using keyboards and robotic teleoperation devices, as well as model-based robot planning.

#### **Abstract**

Creating a physical digital twin of a real-world object has immense potential in robotics, content creation, and XR. In this paper, we present PhysTwin, a novel framework that uses sparse videos of dynamic objects in interaction to produce a photo- and physically realistic, real-time interactive virtual replica. Our approach centers on two key components: (1) a physics-informed representation that combines springmass models for realistic physical simulation, generative shape models for geometry, and Gaussian splats for rendering, and (2) a novel multi-stage optimization-based inverse modeling framework that reconstructs complete geometry, infers dense physical properties, and replicates realistic appearance from videos. Our method integrates an inverse physics framework with visual perception cues, enabling high-fidelity reconstruction even from partial, occluded, and limited viewpoints. PhysTwin supports modeling various deformable objects, including ropes, stuffed animals, cloth, and delivery packages. Experiments show that PhysTwin outperforms competing methods in reconstruction, rendering, future prediction, and simulation under novel interactions. We further demonstrate its applications in interactive realtime simulation and model-based robotic motion planning. (See our supplement webpage for all videos and demos.)

## 1. Introduction

The construction of interactive digital twins is essential for modeling the world and simulating future states, with applications in virtual reality, augmented reality, and robotic manipulation. A physically realistic digital twin (PhysTwin) should accurately capture the geometry, appearance, and physical properties of an object, allowing simulations that closely match observations in the real world. However, constructing such a representation from sparse observations remains a significant challenge.

The creation of digital twins for deformable objects has long been a challenging topic in the vision community. While dynamic 3D methods (e.g., dynamic NeRFs [2, 5, 8, 13, 14, 17, 27, 29–31, 39–41, 43, 55, 56, 58, 61], dynamic 3D Gaussians [10, 20, 24, 33, 34, 59, 65, 66, 68]) capture observed motion, appearance, and geometry from videos, they omit the underlying physics and are thus unsuitable for simulating outcomes in unseen interactions. While recent neural-

043

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

based models [4, 11, 28, 32, 36, 42, 49, 51, 52, 60, 64, 69] learn intuitive physics models from videos, they require large amounts of data and remain limited to specific objects or motions, whereas physics-driven approaches [9, 12, 27, 44, 63, 71, 72] often rely on pre-scanned shapes or dense observations to mitigate ill-posedness. Additionally, it requires dense viewpoint coverage and supports only limited motion types, making it unsuitable for general dynamics modeling.

In this work, we aim to build an interactive PhysTwin from sparse-viewpoint RGB-D video sequences, capturing object geometry, non-rigid dynamic physics, and appearance for realistic physical simulation and rendering. We model deformable object dynamics with a spring-mass-based representation, enabling efficient physical simulation and handling a wide range of common objects, such as ropes, stuffed animals, cloth, and delivery packages. To address challenges posed by sparse observations, we leverage shape priors and motion evidence from advanced 3D generative models [62] and vision foundation models [23, 46, 48] to estimate the topology, geometry, and physical parameters of our physical representation. Since some physical parameters (such as topology-related properties) are non-differentiable and optimizing them efficiently is non-trivial, we design a hierarchical sparse-to-dense optimization strategy. This strategy integrates zero-order optimization [18] for non-differentiable topology and sparse physical parameters (e.g., collision parameters and homogeneous spring stiffness), while employing first-order gradient-based optimization to refine dense spring stiffness and further optimize collision parameters. For appearance modeling, we adopt a Gaussian blending strategy, initializing static Gaussians from sparse observations in the first frame using shape priors and deforming them with a linear blending algorithm to generate realistic dynamic appearances.

Our inverse modeling framework effectively constructs interactive PhysTwin from videos of objects under interaction. We create a real-world deformable object interaction dataset and evaluate our method on three key tasks: reconstruction and resimulation, future prediction, and generalization to unseen interactions. Both quantitative and qualitative results demonstrate that our reconstructed PhysTwin aligns accurately with real-world observations, achieves precise future predictions, and generates realistic simulations under diverse unseen interactions. Furthermore, the high computational efficiency of our physics simulator enables real-time dynamics and rendering of our constructed PhysTwin, facilitating multiple applications, including real-time interactive simulation and model-based robotic motion planning.

# 2. Related Works

**Dynamic Scene Reconstruction.** Dynamic scene reconstruction aims to recover the underlying representation of dynamic scenes from inputs like depth scans [6, 26], RGBD

videos [38], or monocular or multi-view videos [1, 5, 24, 31, 34, 39, 40, 43, 56, 58, 61, 67, 68]. Recent advancements in dynamic scene modeling have involved the adaptation of novel scene representations, including Neural Radiance Fields (NeRF) [2, 5, 8, 13, 14, 16, 17, 27, 29, 30, 30, 31, 39– 41, 43, 55, 56, 58, 61] and 3D Gaussian splats [10, 20, 24, 33, 34, 59, 65, 66, 68]. D-NeRF [43] extends a canonical NeRF on dynamic scenes by optimizing a deformable field. Similarly, Deformable 3D-GS [66] optimizes a deformation field of each Gaussian kernel. Dynamic 3D-GS [34] optimizes the motion of Gaussian kernels for each frame to capture scene dynamics. 4D-GS [59] modulates 3D Gaussians with 4D neural voxels for dynamic multi-view synthesis. Although these methods achieve high-fidelity results in dynamic multiview synthesis, they primarily focus on reconstructing scene appearance and geometry without capturing real-world dynamics, limiting their ability to support action-conditioned future predictions and interactive simulations.

Physics-Based Simulation of Deformable Objects. Another line of work incorporates physical simulators to perform system identification of physical parameters during reconstruction. Earlier methods relied on pre-scanned static objects and required clean point cloud observations [9, 15, 19, 21, 35, 44, 47, 57]. Most recent approaches build upon SDF [45], NeRF [3, 12, 27] or Gaussian Splatting [22, 63, 71, 72] to support more flexible physical digital twin reconstruction. Several works [12, 22, 63] manually specify physics parameters, resulting in a mismatch between the simulation and real-world video observations. Other works [3, 27, 45, 71, 72] attempt to estimate physical parameters from videos. However, they are often constrained to synthetic data, limited motion, or the need for dense viewpoints to accurately reconstruct static geometry, limiting their practical applicability. The closest related work to ours is Spring-Gaus [72], which also utilizes a 3D Spring-Mass model for learning from videos. However, their physical model is overly regularized and violates real-world physics, lacking momentum conservation and realistic gravity. Moreover, Spring-Gaus requires dense viewpoint coverage to reconstruct the full geometry at the initial state, which is impractical in many real-world settings. The motions are also only limited to tabletop collisions and lack action inputs, making it unsuitable as a general dynamics model for downstream applications.

Learning-Based Simulation of Deformable Objects. Analytically modeling the dynamics of deformable objects is extremely challenging due to the high complexity of state space and the variance of physical properties. Recent works [4, 11, 36, 60, 64] have chosen to use neural network-based simulators to model object dynamics. Specifically, graph-based networks effectively learn the dynamics of various types of objects such as plasticine [51, 52], cloth [32, 42], fluid [28, 49], and stuffed animals [69]. GS-Dynamics [69]

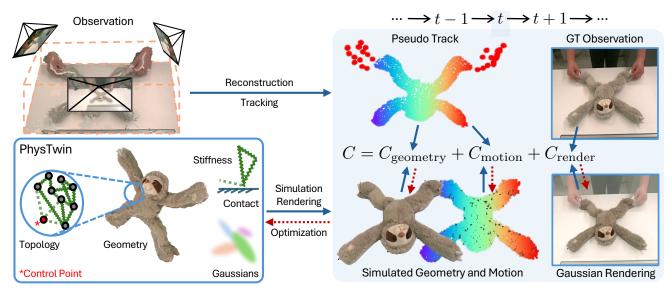


Figure 2. **Overview of Our PhysTwin Framework.** We present an overview of our PhysTwin framework, where the core representation includes geometry, topology, physical parameters (associated with springs and contacts), and Gaussian kernels. To optimize PhysTwin, we minimize the rendering loss and the discrepancy between simulated and observed geometry/motion. The rendering loss optimizes the Gaussian kernels, while the geometry and motion losses refine the overall geometry, topology, and physical parameters in PhysTwin.

attempted to learn object dynamics directly from real-world videos using tracking and appearance priors from Dynamic Gaussians [34], and generalized well to unseen actions. However, these learned models need extensive training samples and are often limited to specific environments with limited motion ranges. In contrast, our method requires only one interaction trial while achieving a broader range of motions.

## 3. Preliminary: Spring-Mass Model

Spring-mass models are widely used for simulating deformable objects due to their simplicity and computational efficiency. A deformable object is represented as a set of spring-connected mass nodes, forming a graph structure  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is the set of mass points and  $\mathcal{E}$  is the set of springs. Each mass node i has a position  $\mathbf{x}_i \in \mathbb{R}^3$  and velocity  $\mathbf{v}_i \in \mathbb{R}^3$ , which evolve over time according to Newtonian dynamics. Springs are constructed between neighboring nodes based on a predefined topology, defining the elastic structure of the object.

The force on node i is the result of the combined effects of adjacent nodes connected by springs:

$$\mathbf{F}_{i} = \sum_{(i,j)\in\mathcal{E}} \mathbf{F}_{i,j}^{\text{spring}} + \mathbf{F}_{i,j}^{\text{dashpot}} + \mathbf{F}_{i}^{\text{ext}}, \tag{1}$$

where the spring force and dashpot damping force between nodes i and j are given by  $\mathbf{F}_{i,j}^{\text{spring}} = k_{ij}(\|\mathbf{x}_j - \mathbf{x}_i\| - l_{ij})\frac{\mathbf{x}_j - \mathbf{x}_i}{\|\mathbf{x}_j - \mathbf{x}_i\|}$  and  $\mathbf{F}_{i,j}^{\text{dashpot}} = -\gamma(\mathbf{v}_i - \mathbf{v}_j)$ , respectively. Here,  $k_{ij}$  is the spring stiffness,  $l_{ij}$  is the rest length, and  $\gamma$  is the dashpot damping coefficient. The external force  $\mathbf{F}_i^{\text{ext}}$  accounts for factors such as gravity, collisions, and user interactions. The spring force restores the system to its rest

shape, while the dashpot damping dissipates energy, preventing oscillations. For collisions, we use impulse-based collision handling when two mass points are very close, including collisions between the object and the collider, as well as between two object points.

The spring-mass model updates the system state with a dynamic model  $\mathbf{X}_{t+1} = f_{\alpha,\mathcal{G}_0}(\mathbf{X}_t,a_t)$  by applying explicit Euler integration to both velocity and position. More formally, for all  $i, \mathbf{v}_i^{t+1} = \delta\left(\mathbf{v}_i^t + \Delta t \, \frac{\mathbf{F}_i}{m_i}\right), \quad \mathbf{x}_i^{t+1} = \mathbf{x}_i^t + \Delta t \, \mathbf{v}_i^{t+1},$  where  $\mathbf{X}_t$  represents the system state at time t. In this formulation,  $\alpha$  denotes all physical parameters of the spring-mass model, including spring stiffness, collision parameters, and damping. It also encompasses the parameters related to the control interaction.  $\mathcal{G}_0$  represents the "canonical" geometry and topology for the spring-mass system  $^1$ , and  $a_t$  represents the actions at time t.

# 4. Method

In this section, we formulate the construction of PhysTwin as an optimization problem. We then present our two-stage strategy, where the first stage addresses the physics-related optimization, followed by the appearance-based optimization in the second stage. Finally, we demonstrate the capability of our framework to perform real-time simulation using the constructed PhysTwin.

#### 4.1. Problem Formulation

Given three RGBD videos of a deformable object under interaction, our objective is to construct a PhysTwin model that captures the geometry, appearance, and physical parameters

<sup>&</sup>lt;sup>1</sup>In practice, we use the first-frame object state as the canonical state.

of the object over time. At each time frame t, we denote the RGBD observations from each video as  $\mathbf{O}_{t,i}$  for the i-th observation at time t, where  $\mathbf{O} = (\mathbf{I}, \mathbf{D})$  represents the RGB image  $\mathbf{I}$  and depth map  $\mathbf{D}$ .

The goal of our optimization problem is to minimize the discrepancy between the predicted observation  $\hat{\mathbf{O}}_{t,i}$  and the actual observation  $\mathbf{O}_{t,i}$ . The predicted observation is derived by projecting and rendering the predicted state  $\hat{\mathbf{X}}_t$  onto images through a function  $g_{\theta}$ , where  $\theta$  encodes the appearance of the objects represented by Gaussian splats. The 3D state  $\hat{\mathbf{X}}_t$  evolves over time according to the Spring-Mass model, which captures the deformable object's dynamics and updates the state using the explicit Euler integration method. The optimization problem is formulated as:

$$\min_{\alpha, \mathcal{G}_0, \theta} \sum_{t, i} C(\hat{\mathbf{O}}_{t, i}, \mathbf{O}_{t, i})$$
s.t.  $\hat{\mathbf{O}}_{t, i} = q_{\theta}(\hat{\mathbf{X}}_t, i), \qquad \hat{\mathbf{X}}_{t+1} = f_{\alpha, \mathcal{G}}(\hat{\mathbf{X}}_t, a_t),$  (2)

where  $\alpha$ ,  $\mathcal{G}_0$ ,  $\theta$  captures the physics, geometry, topology and appearance parameters (Sec. 3); the cost function quantifies the difference between the predicted observation  $\hat{\mathbf{O}}_{t,i}$  and the actual observation  $\mathbf{O}_{t,i}$ . This cost function is decomposed into three components:  $C = C_{\text{geometry}} + C_{\text{motion}} + C_{\text{render}}$ , each capturing the discrepancy between the inferred system states and the corresponding observations from 3D geometry, 3D motion tracking, and 2D color, respectively (we defer the details of each cost component to Sec.4.2.1 and Sec.4.2.2). The function  $g_{\theta}$  is the observation model, describing the projection from the predicted state to the image plane and render the i-th image-space sensory observation, and  $f_{\alpha,\mathcal{G}}$  models the dynamic evolution of the object's state under the Spring-Mass model, as detailed in Sec. 3.

#### 4.2. PhysTwin Framework

Given the complexity of the overall optimization defined in Eq. 2, our PhysTwin framework decomposes it into two stages. The first stage focuses on optimizing the geometry and physical parameters, while the second stage is dedicated to optimizing the appearance-related parameters.

## 4.2.1. Physics and Geometry Optimization.

As outlined in our optimization formulation in Sec. 4.1, the objective is to minimize the discrepancy between the predicted observation  $\hat{\mathbf{O}}_{t,i}$  and the actual observation  $\mathbf{O}_{t,i}$ . First, we convert the depth observations  $\mathbf{D}_t$  at each time frame t into the observed partial 3D point cloud  $\mathbf{X}_t$ . In the first stage, we consider the following formulation for the optimization:

$$\min_{\alpha, \mathcal{G}_0} \sum_{t} \left( C_{\text{geometry}}(\hat{\mathbf{X}}_t, \mathbf{X}_t) + C_{\text{motion}}(\hat{\mathbf{X}}_t, \mathbf{X}_t) \right) 
\text{s.t.} \quad \hat{\mathbf{X}}_{t+1} = f_{\alpha, \mathcal{G}_0}(\hat{\mathbf{X}}_t, g_t).$$
(3)

where the  $C_{\rm geometry}$  function quantifies the single-direction Chamfer distance between the partial observed point cloud

 $\mathbf{X}_t$  and the inferred state  $\hat{\mathbf{X}}_t$ , and  $C_{\mathrm{motion}}$  quantifies the tracking error between the predicted point  $\hat{x}_i^t$  and its corresponding observed tracking  $x_i^t$ . The observed tracking is obtained using the vision foundation model CoTracker3 [23], followed by lifting the result to 3D via depth map unprojection.

There are three main challenges in the first-stage optimization: 1) partial observations from sparse viewpoints, 2) joint discrete topology and physical parameter optimization, and 3) discontinuities in the dynamic model, as well as the long-time horizon and dense properties, making continuous optimization difficult. To address these challenges, we handle the geometry and other parameters separately. Specifically, we first leverage generative shape initialization to obtain the full geometry, then employ our two-stage sparse-to-dense optimization to refine the other parameters.

Generative Shape Prior. Due to partial observations, recovering full geometry is challenging. We use a shape prior from the image-to-3D generative model, TRELLIS [62], to generate the full mesh conditioned on a single RGB observation of the masked object. To improve mesh quality, we apply a super-resolution model [48] to upscale the foreground, segmented using Grounded-SAM2 [46]. While the mesh corresponds reasonably with the observation, inconsistencies in scale, pose, and deformation remain.

To address this, we design a registration module using 2D matching for scale estimation, rigid registration, and non-rigid deformation. A coarse-to-fine strategy estimates initial rotation via 2D correspondences matched using SuperGlue [50], followed by refinement with the Perspectiven-Point [25] algorithm. We resolve scale and translation ambiguities by optimizing the matched point distances in the camera coordinate system. After applying these transformations, the objects align in pose, with some deformations handled by as-rigid-as-possible registration [53]. Finally, ray-casting alignment ensures that observed points match the deformed mesh without occlusions.

These steps yield a shape prior aligned with the first-frame observations, which serves as a crucial initialization for the inverse physics and appearance optimization stages.

**Sparse-to-Dense Optimization** The Spring-Mass model consists of both the topological structure (i.e., the connectivity of the springs) and the physical parameters defined on the springs. As mentioned in Sec. 3, we also include the control parameters to connect the springs between control points and object points through radius and max neighbours. For topology optimization, we employ a heuristic approach to connect the nearest neighbor points, parameterized by a connection radius and a maximum number of neighbors, thereby controlling the density of the springs. The same parameterization for the springs between control points and object points. To extract control points from video data, we utilize Grounded-SAM2 [46] to segment the hand mask and CoTracker3 [23] to track hand movements. After lifting the

points to 3D, we apply farthest-point sampling to obtain the final set of control points.

All the aforementioned parameters constitute the parameter space we aim to optimize. The two main challenges are: i) some parameters are not differentiable (e.g., radius and maximum number of neighbors); ii) to represent a wide range of objects, we model dense spring stiffness, leading to a parameter space that includes thousands of springs.

To address these challenges, we introduce a hierarchical sparse-to-dense optimization strategy. Initially, we employ zero-order sampling-based optimization to optimize the parameters, which naturally overcome the differentiability issue. However, zero-order optimization is inefficient when the parameter space is too large. Therefore, in the first stage, we assume homogeneous stiffness, allowing the topology and other physical parameters to obtain a good initialization. In the second stage, we further refine the parameters using first-order gradient descent, leveraging our built differentiable spring-mass simulator. This stage optimizes the dense spring stiffness and collision parameters simultaneously.

Beyond the optimization strategy, we incorporate additional supervision by utilizing tracking priors from vision foundation models. We lift the 2D tracking prediction into 3D to obtain pseudo-ground-truth tracking data for the 3D points, which forms a crucial component of our cost function as mentioned above.

By integrating our optimization strategy with a cost function that leverages additional tracking priors, our PhysTwin framework can effectively and efficiently model the dynamics of diverse interactable objects from videos.

#### 4.2.2. Appearance Optimization

For the second-stage appearance optimization, to model object appearance, we construct a set of static 3D Gaussian kernels parameterized by  $\theta$ , with each Gaussian defined by a 3D center position  $\mu$ , a rotation matrix represented by a quaternion  $q \in \mathbf{SO}(3)$ , a scaling matrix represented by a 3D vector s, an opacity value  $\alpha$ , and color coefficients c. We optimize  $\theta$  here via

$$\min_{\theta} \sum_{t,i} C_{\text{render}}(\hat{\mathbf{I}}_{i,t}, \mathbf{I}_{i,t}) \text{ s.t. } \hat{\mathbf{I}}_{i,t} = g_{\theta}(\hat{\mathbf{X}}_t, i), \quad (4)$$

where  $\hat{\mathbf{X}}_t$  is the optimized system states at time t,i is the camera index, and  $\mathbf{I}_{i,t}$ ,  $\hat{\mathbf{I}}_{i,t}$  are the ground truth image and rendered image from camera view i at time t, respectively.  $C_{\text{render}}$  computes the  $\mathcal{L}_1$  loss with a D-SSIM term between the rendering and ground truth image. For simplicity, we set t=0 to optimize appearance only at the first frame. We restrict the Gaussian shape to be isotropic to prevent spiky artifacts during deformation.

To ensure realistic rendering under deformation, we need to dynamically adjust each Gaussian at each timestep t based on the transition between states  $\hat{\mathbf{X}}_t$  and  $\hat{\mathbf{X}}_{t+1}$ . To achieve

this, we adopt a Gaussian updating algorithm using Linear Blend Skinning (LBS) [20, 54, 69], which interpolates the motions of 3D Gaussians using the motions of neighboring mass nodes. Please refer to the supplementary for details.

## 4.3. Capabilities of PhysTwin

Our constructed PhysTwin supports real-time simulation of deformable objects under various motions while maintaining realistic appearance. This real-time, photorealistic simulation enables interactive exploration of object dynamics.

By introducing control points and dynamically connecting them to object points via springs, our system can simulate diverse motion patterns and interactions. These capabilities make PhysTwin a powerful representation for real-time interactive simulation and model-based robotic motion planning, which are further described in Sec. 5.3.

# 5. Experiments

In this section, we evaluate the performance of our PhysT-win framework across three distinct tasks involving different types of objects. Our primary objective is to address the following three questions: 1) How accurately does our framework reconstruct and re-simulate deformable objects and predict its future states? 2) How well does the constructed PhysTwin generalize to unseen interactions? 3) What is the utility of PhysTwin in downstream tasks?

# **5.1.** Experiment Settings

**Dataset.** We collect a dataset of RGBD videos capturing human interactions with various deformable objects with different physical properties, such as ropes, stuffed animals, cloth, and delivery packages. Three RealSense-D455 RGBD cameras are used to record the interactions. Each video is 1 to 10 sec long and covers different interactions, including quick lifting, stretching, pushing, and squeezing with either one or both hands. We collect 22 scenarios encompassing various object types, interaction types, and hand configurations. For each scenario, the RGBD videos are split into a training set and a test set following a 7:3 ratio, where only the training set is used to construct PhysTwin. We manually annotate 9 ground-truth tracking points for each video to evaluate tracking performance with the semi-auto tool introduced in [7].

**Tasks.** To assess the effectiveness of our PhysTwin framework and the quality of our constructed PhysTwin, we formulate three tasks: 1) Reconstruction & Resimulation; 2) Future Prediction; 3) Generalization to Unseen Actions.

For the Reconstruction & Resimulation task, the objective is to construct PhysTwin such that it can accurately reconstruct and resimulate the motion of deformable objects given the actions represented by the control point positions.

For the Future Prediction task, we aim to assess whether PhysTwin can perform well on unseen future frames during



Figure 3. Qualitative Results on Reconstruction & Resimulation and Future Prediction. We visualize the rendering results of different methods on two tasks. For the reconstruction & resimulation task, our method achieves a better match with the observations. For the future prediction task, our method accurately predicts the future state of the objects. In contrast, the baselines fail in most cases: GS-Dynamics tends to remain static, while Spring-Gauss frequently causes the physical model to crash.

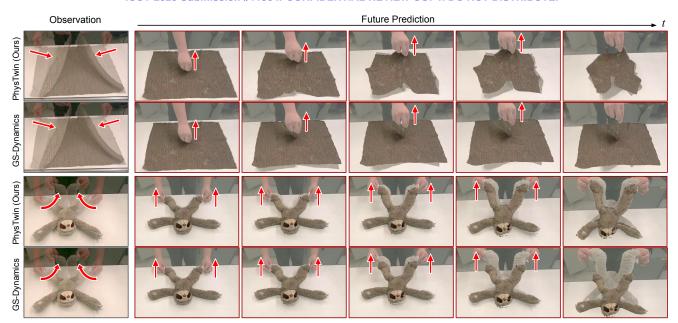


Figure 4. Qualitative Results on Generalization to Unseen Interactions. We visualize the simulation of a deformable object under unseen interactions using our method and GS-Dynamics. The leftmost image illustrates the interaction on which the dynamics models are trained, while the right images demonstrate their generalization ability to unseen interactions. Our PhysTwin significantly outperforms prior work.

Table 1. Quantitative Results on Reconstruction & Resimulation and Future Prediction. We compare the performance of our method with two prior work, GS-Dynamics and Spring-Gaus, on two tasks: reconstruction & resimulation and future prediction. Our PhysTwin framework consistently outperforms the baselines across all metrics.

Task	Reconstruction & Resimulation						Future Prediction					
Method	CD↓	Track Error ↓	IoU %↑	PSNR ↑	SSIM ↑	LPIPS ↓	CD↓	Track Error ↓	IoU %↑	PSNR ↑	SSIM ↑	LPIPS ↓
Spring-Gaus [72]	0.041	0.050	57.6	23.445	0.928	0.102	0.062	0.094	46.4	22.488	0.924	0.113
GS-Dynamics [69]	0.014	0.022	72.1	26.260	0.940	0.052	0.041	0.070	49.8	22.540	0.924	0.097
PhysTwin (Ours)	0.005	0.009	84.4	28.214	0.945	0.034	0.012	0.022	72.5	25.617	0.941	0.055

Table 2. Quantitative Results on Generalization to Unseen Interactions. We compare our method with GS-Dynamics on generalization to unseen interactions. Both methods are trained on the same video with a specific interaction and tested on unseen interactions. Our method achieves significantly better results.

Method	CD↓	Track Error ↓	IoU %↑	PSNR ↑	SSIM ↑	LPIPS ↓
GS-Dynamics [69]	0.029	0.038	63.4	25.053	0.934	0.067
PhysTwin (Ours)	0.013	0.018	72.18	26.199	0.938	0.047

its construction. For the Generalization to Unseen Interactions task, the goal is to assess whether PhysTwin can adapt to different interactions. To evaluate this, we construct a generalization dataset consisting of interaction pairs performed on the same object but with varying motions, including differences in hand configuration and interaction type.

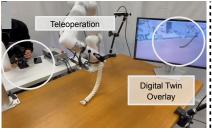
**Baselines.** To the best of our knowledge, there is currently no existing work that demonstrates good performance across all three tasks. Therefore, we select two main research directions as baselines and further augment them to match the tasks in our setting (full details in the supplementary).

The first baseline we consider is a physics-based simulation method for identifying the material properties of deformable objects, Spring-Gaus [72]. Their work has demonstrated the statement of the constant of the statement of

strated strong capabilities in reconstruction, resimulation, and future prediction in its original setting. However, their framework does not support human control, so we augment their method with additional control support.

The second baseline is based on a learning-based simulation approach, GS-Dynamics [69], specifically a GNN-based neural dynamics model. This model directly learns the system's dynamics from two partial states. In their original setting, video preprocessing is required using Dyn3DGS [34] to obtain tracking information. In our case, we leverage our 3D-lifting tracking from CoTracker3 [23] as supervision for the neural dynamics model.

**Evaluation.** To better understand whether our prediction matches the observations, we evaluate predictions in both 3D and 2D. For the 3D evaluation, we use the single-direction Chamfer Distance (partial ground truth with our full-state prediction) and the tracking error (based on our manually annotated ground-truth tracking points). For the 2D evaluation, we assess image quality using PSNR, SSIM, and LPIPS [70], and silhouette alignment using IoU. We perform 2D evaluation only at the center viewpoint due to optimal visibility of objects, with metrics averaged across all frames and scenarios. Specially, for the Spring-Gaus [72] baseline,







(a) Real-Time Interactive Simulation using Keyboard

(b) Real-Time Interactive Simulation using Teleoperation

(c) Model-Based Robot Planning

Figure 5. **Applications of our PhysTwin.** Our constructed PhysTwin supports a variety of tasks, including real-time interactive simulation, which can accept input from either a keyboard or a robot teleoperation setup. Meanwhile, PhysTwin also enables model-based robot planning to accomplish tasks such as lifting a rope into some specific configuration.

its optimization process is unstable due to inaccurate physics modeling. Therefore, we report the above metrics only for its successful cases.

#### 5.2. Results

We compare with two augmented baselines across three task settings. Our quantitative analysis reveals that the PhysT-win framework consistently outperforms the baselines across various tasks. Full video results can be found on our supplementary webpage.

Reconstruction & Resimulation. The quantitative results in Tab. 1 Reconstruction & Resimulation section demonstrate the superior performance of our PhysTwin method over baselines. Our approach significantly improves all evaluated metrics, including Chamfer Distance, tracking error, and 2D IoU, confirming that our reconstruction and resimulation align more closely with the original observations. This highlights the effectiveness of our model in learning a more accurate dynamics model under sparse observations. Additionally, rendering metrics show that our method produces more realistic 2D images, benefiting from the Gaussian blending strategy and enhanced dynamic modeling. Fig. 3 further provides qualitative visualizations across different objects, illustrating precise alignment with original observations. Notably, our physics-based representation inherently improves point tracking. After physicsconstrained optimization, our tracking surpasses the original CoTracker3 [23] predictions used for training, achieving better alignment after global optimization (See supplement for more details).

**Future Prediction.** Table 1, in the Future Prediction section, demonstrates that our method achieves superior performance in predicting unseen frames, excelling in both dynamics alignment and rendering quality. Fig. 3 further provides qualitative results, illustrating the accuracy of our predictions on unseen frames.

Generalization to Unseen Interactions. We also evaluate the generalization performance to unseen interactions. We directly use our constructed PhysTwin and leverage our registration pipeline to align it with the first frame of the target case. Fig. 4 shows that our method closely matches the

ground truth observations in terms of dynamics. Quantitative results further demonstrate the robustness of our method across different actions. In contrast, the neural dynamics model struggles to adapt to environmental changes and diverse interactions as effectively as our approach. Moreover, for unseen interaction scenarios, our method achieves performance comparable to the future prediction task, highlighting the robustness and practicality of our constructed PhysTwin.

## 5.3. Application

The high-speed forward simulation of our Spring-Mass simulator implemented using Warp [37] enables a variety of downstream applications. Fig. 5 illustrates three key applications of PhysTwin. 1) Interactive Simulation: Users can interact with objects in real time using keyboard controls, either with one hand or both hands. 2) Real-Time Future Prediction: Our method enables real-time simulation of an object's future state while a human teleoperates the robotic arms to interact with a real object. This can serve as a crucial tool for predicting object dynamics during manipulation. 3) Model-Based Robotic Planning: Due to the high fidelity of our constructed PhysTwin, it can also be used purely as a dynamic function. By integrating it with model-based planning techniques, we can generate motion plans for the robot to complete different types of tasks effectively.

### 6. Conclusion

We introduced PhysTwin, a novel framework for constructing physical digital twins from sparse videos, enabling effective reconstruction and resimulation of deformable objects. Our approach excels in predicting future states and simulating object interactions that generalize to unseen scenarios. We showed the superior performance of our method across various object types, control configurations, and task settings, significantly outperforming prior work. PhysTwin enables various downstream tasks that demand high-speed simulation and accurate future prediction. Moreover, our approach provides valuable insights for robotic manipulation. By bridging perception and physics-based simulation, PhysTwin serves as a crucial tool for guiding robot interactions, making real-world deployment more efficient and reliable.

## 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 rayconditioned sampling. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 16610–16620, 2023. 2
- [2] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Con*ference on Computer Vision and Pattern Recognition, pages 130–141, 2023. 1, 2
- [3] Hsiao-yu Chen, Edith Tretschk, Tuur Stuyck, Petr Kadle-cek, Ladislav Kavan, Etienne Vouga, and Christoph Lassner. Virtual elastic objects. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 15827–15837, 2022.
- [4] Zhenfang Chen, Kexin Yi, Yunzhu Li, Mingyu Ding, Antonio Torralba, Joshua B Tenenbaum, and Chuang Gan. Comphy: Compositional physical reasoning of objects and events from videos. *arXiv* preprint arXiv:2205.01089, 2022. 2
- [5] Mengyu Chu, Lingjie Liu, Quan Zheng, Erik Franz, Hans-Peter Seidel, Christian Theobalt, and Rhaleb Zayer. Physics informed neural fields for smoke reconstruction with sparse data. ACM Transactions on Graphics, 41(4):119:1–119:14, 2022. 1, 2
- [6] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings* of the 23rd annual conference on Computer graphics and interactive techniques, pages 303–312, 1996. 2
- [7] Carl Doersch, Yi Yang, Mel Vecerik, Dilara Gokay, Ankush Gupta, Yusuf Aytar, Joao Carreira, and Andrew Zisserman. TAPIR: Tracking any point with per-frame initialization and temporal refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10061– 10072, 2023. 5
- [8] Danny Driess, Zhiao Huang, Yunzhu Li, Russ Tedrake, and Marc Toussaint. Learning multi-object dynamics with compositional neural radiance fields. In *Conference on robot learning*, pages 1755–1768. PMLR, 2023. 1, 2
- [9] Tao Du, Kui Wu, Pingchuan Ma, Sebastien Wah, Andrew Spielberg, Daniela Rus, and Wojciech Matusik. Diffpd: Differentiable projective dynamics. ACM Transactions on Graphics (ToG), 41(2):1–21, 2021. 2
- [10] Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen. 4d-rotor gaussian splatting: towards efficient novel view synthesis for dynamic scenes. In ACM SIGGRAPH 2024 Conference Papers, pages 1–11, 2024.

   2
- [11] Ben Evans, Abitha Thankaraj, and Lerrel Pinto. Context is everything: Implicit identification for dynamics adaptation. In 2022 International Conference on Robotics and Automation (ICRA), pages 2642–2648. IEEE, 2022. 2
- [12] Yutao Feng, Yintong Shang, Xuan Li, Tianjia Shao, Chen-fanfu Jiang, and Yin Yang. Pie-nerf: Physics-based interactive elastodynamics with nerf. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4450–4461, 2024. 2

- [13] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12479–12488, 2023. 1, 2
- [14] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular dynamic view synthesis: A reality check. Advances in Neural Information Processing Systems, 35:33768–33780, 2022. 1, 2
- [15] Moritz Geilinger, David Hahn, Jonas Zehnder, Moritz Bächer, Bernhard Thomaszewski, and Stelian Coros. Add: Analytically differentiable dynamics for multi-body systems with frictional contact. ACM Transactions on Graphics (TOG), 39 (6):1–15, 2020.
- [16] Shanyan Guan, Huayu Deng, Yunbo Wang, and Xiaokang Yang. Neurofluid: Fluid dynamics grounding with particledriven neural radiance fields. In *International conference on machine learning*, pages 7919–7929. PMLR, 2022. 2
- [17] Xiang Guo, Jiadai Sun, Yuchao Dai, Guanying Chen, Xiao-qing Ye, Xiao Tan, Errui Ding, Yumeng Zhang, and Jingdong Wang. Forward flow for novel view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16022–16033, 2023. 1,
- [18] Nikolaus Hansen. The cma evolution strategy: a comparing review. *Towards a new evolutionary computation: Advances in the estimation of distribution algorithms*, pages 75–102, 2006. 2
- [19] Eric Heiden, Miles Macklin, Yashraj Narang, Dieter Fox, Animesh Garg, and Fabio Ramos. Disect: A differentiable simulation engine for autonomous robotic cutting. *arXiv* preprint arXiv:2105.12244, 2021. 2
- [20] Yi-Hua Huang, Yang-Tian Sun, Ziyi Yang, Xiaoyang Lyu, Yan-Pei Cao, and Xiaojuan Qi. Sc-gs: Sparse-controlled gaussian splatting for editable dynamic scenes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4220–4230, 2024. 1, 2, 5
- [21] Krishna Murthy Jatavallabhula, Miles Macklin, Florian Golemo, Vikram Voleti, Linda Petrini, Martin Weiss, Breandan Considine, Jérôme Parent-Lévesque, Kevin Xie, Kenny Erleben, et al. gradsim: Differentiable simulation for system identification and visuomotor control. arXiv preprint arXiv:2104.02646, 2021. 2
- [22] Ying Jiang, Chang Yu, Tianyi Xie, Xuan Li, Yutao Feng, Huamin Wang, Minchen Li, Henry Lau, Feng Gao, Yin Yang, et al. Vr-gs: A physical dynamics-aware interactive gaussian splatting system in virtual reality. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–1, 2024. 2
- [23] Nikita Karaev, Iurii Makarov, Jianyuan Wang, Natalia Neverova, Andrea Vedaldi, and Christian Rupprecht. Cotracker3: Simpler and better point tracking by pseudolabelling real videos. arXiv preprint arXiv:2410.11831, 2024. 2, 4, 7, 8
- [24] Agelos Kratimenos, Jiahui Lei, and Kostas Daniilidis. Dynmf: Neural motion factorization for real-time dynamic view synthesis with 3d gaussian splatting. In *European Conference on Computer Vision*, pages 252–269. Springer, 2024. 1, 2

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687 688

689

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709 710

711

712

713

714

715

716

717

718

719

720 721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

- [25] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. Ep n p: An accurate o (n) solution to the p n p problem. International journal of computer vision, 81:155–166, 2009.
- [26] Hao Li, Robert W Sumner, and Mark Pauly. Global correspondence optimization for non-rigid registration of depth scans. In *Computer graphics forum*, pages 1421–1430. Wiley Online Library, 2008.
- [27] Xuan Li, Yi-Ling Qiao, Peter Yichen Chen, Krishna Murthy Jatavallabhula, Ming Lin, Chenfanfu Jiang, and Chuang Gan. Pac-nerf: Physics augmented continuum neural radiance fields for geometry-agnostic system identification. arXiv preprint arXiv:2303.05512, 2023. 1, 2
- [28] Yunzhu Li, Jiajun Wu, Russ Tedrake, Joshua B Tenenbaum, and Antonio Torralba. Learning particle dynamics for manipulating rigid bodies, deformable objects, and fluids. *arXiv* preprint arXiv:1810.01566, 2018. 2
- [29] Yunzhu Li, Shuang Li, Vincent Sitzmann, Pulkit Agrawal, and Antonio Torralba. 3d neural scene representations for visuomotor control. In *Conference on Robot Learning*, pages 112–123. PMLR, 2022. 1, 2
- [30] 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, pages 6498– 6508, 2021. 2
- [31] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4273–4284, 2023. 1, 2
- [32] Xingyu Lin, Yufei Wang, Zixuan Huang, and David Held. Learning visible connectivity dynamics for cloth smoothing. In Conference on Robot Learning, pages 256–266. PMLR, 2022.
- [33] Youtian Lin, Zuozhuo Dai, Siyu Zhu, and Yao Yao. Gaussian-flow: 4d reconstruction with dynamic 3d gaussian particle. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 21136–21145, 2024.

   2
- [34] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. In 2024 International Conference on 3D Vision (3DV), pages 800–809. IEEE, 2024. 1, 2, 3, 7
- [35] Pingchuan Ma, Tao Du, Joshua B Tenenbaum, Wojciech Matusik, and Chuang Gan. Risp: Rendering-invariant state predictor with differentiable simulation and rendering for cross-domain parameter estimation. *arXiv* preprint *arXiv*:2205.05678, 2022. 2
- [36] Pingchuan Ma, Peter Yichen Chen, Bolei Deng, Joshua B Tenenbaum, Tao Du, Chuang Gan, and Wojciech Matusik. Learning neural constitutive laws from motion observations for generalizable pde dynamics. In *International Conference on Machine Learning*, pages 23279–23300. PMLR, 2023. 2
- [37] Miles Macklin. Warp: A high-performance python framework for gpu simulation and graphics. https://github.com/nvidia/warp, 2022. NVIDIA GPU Technology Conference (GTC). 8

- [38] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Proceedings of the IEEE conference on com*puter vision and pattern recognition, pages 343–352, 2015.
- [39] 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. 1, 2
- [40] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higherdimensional representation for topologically varying neural radiance fields. arXiv preprint arXiv:2106.13228, 2021. 2
- [41] Yicong Peng, Yichao Yan, Shenqi Liu, Yuhao Cheng, Shanyan Guan, Bowen Pan, Guangtao Zhai, and Xiaokang Yang. Cagenerf: Cage-based neural radiance fields for genrenlized 3d deformation and animation. In *Thirty-Sixth Con*ference on Neural Information Processing Systems, 2022. 1, 2
- [42] Tobias Pfaff, Meire Fortunato, Alvaro Sanchez-Gonzalez, and Peter Battaglia. Learning mesh-based simulation with graph networks. In *International conference on learning* representations, 2020. 2
- [43] 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 con*ference on computer vision and pattern recognition, pages 10318–10327, 2021. 1, 2
- [44] Yi-Ling Qiao, Junbang Liang, Vladlen Koltun, and Ming C. Lin. Differentiable simulation of soft multi-body systems. In Conference on Neural Information Processing Systems (NeurIPS), 2021. 2
- [45] Yi-Ling Qiao, Alexander Gao, and Ming Lin. Neuphysics: Editable neural geometry and physics from monocular videos. Advances in Neural Information Processing Systems, 35: 12841–12854, 2022. 2
- [46] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kunchang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, et al. Grounded sam: Assembling open-world models for diverse visual tasks. arXiv preprint arXiv:2401.14159, 2024. 2, 4
- [47] Junior Rojas, Eftychios Sifakis, and Ladislav Kavan. Differentiable implicit soft-body physics. arXiv preprint arXiv:2102.05791, 2021. 2
- [48] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2, 4
- [49] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and Peter Battaglia. Learning to simulate complex physics with graph networks. In *Interna*tional conference on machine learning, pages 8459–8468. PMLR, 2020. 2
- [50] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Superglue: Learning feature match-

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791792

793

794

795

796 797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

- ing with graph neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4938–4947, 2020. 4
  - [51] Haochen Shi, Huazhe Xu, Samuel Clarke, Yunzhu Li, and Jiajun Wu. Robocook: Long-horizon elasto-plastic object manipulation with diverse tools. arXiv preprint arXiv:2306.14447, 2023. 2
  - [52] Haochen Shi, Huazhe Xu, Zhiao Huang, Yunzhu Li, and Jiajun Wu. Robocraft: Learning to see, simulate, and shape elasto-plastic objects in 3d with graph networks. *The International Journal of Robotics Research*, 43(4):533–549, 2024.
  - [53] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry processing*, pages 109–116. Citeseer, 2007. 4
  - [54] Robert W Sumner, Johannes Schmid, and Mark Pauly. Embedded deformation for shape manipulation. In SIGGRAPH, pages 80–es. 2007. 5
  - [55] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Nonrigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 12959–12970, 2021. 1, 2
  - [56] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Nonrigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *IEEE International Conference on Computer Vision (ICCV)*. IEEE, 2021. 1, 2
  - [57] Bin Wang, Longhua Wu, KangKang Yin, Uri M Ascher, Libin Liu, and Hui Huang. Deformation capture and modeling of soft objects. ACM Trans. Graph., 34(4):94–1, 2015.
  - [58] Chaoyang Wang, Lachlan Ewen MacDonald, Laszlo A Jeni, and Simon Lucey. Flow supervision for deformable nerf. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 21128–21137, 2023. 1, 2
  - [59] 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. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 20310–20320, 2024. 1,
  - [60] Yilin Wu, Wilson Yan, Thanard Kurutach, Lerrel Pinto, and Pieter Abbeel. Learning to manipulate deformable objects without demonstrations. arXiv preprint arXiv:1910.13439, 2019. 2
  - [61] 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 (CVPR), pages 9421–9431, 2021. 1, 2
  - [62] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong Yang. Structured 3d latents for scalable and versatile 3d generation. *arXiv preprint arXiv:2412.01506*, 2024. 2, 4

- [63] Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. Physgaussian: Physicsintegrated 3d gaussians for generative dynamics. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 4389–4398, 2024. 2
- [64] Zhenjia Xu, Jiajun Wu, Andy Zeng, Joshua B Tenenbaum, and Shuran Song. Densephysnet: Learning dense physical object representations via multi-step dynamic interactions. arXiv preprint arXiv:1906.03853, 2019.
- [65] Zeyu Yang, Hongye Yang, Zijie Pan, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. arXiv preprint arXiv:2310.10642, 2023. 1, 2
- [66] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for highfidelity monocular dynamic scene reconstruction. In *Proceed*ings of the IEEE/CVF conference on computer vision and pattern recognition, pages 20331–20341, 2024. 1, 2
- [67] Heng Yu, Joel Julin, Zoltan A Milacski, Koichiro Niinuma, and Laszlo A Jeni. Dylin: Making light field networks dynamic. *arXiv preprint arXiv:2303.14243*, 2023. 2
- [68] Heng Yu, Joel Julin, Zoltán Á Milacski, Koichiro Niinuma, and László A Jeni. Cogs: Controllable gaussian splatting. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 21624–21633, 2024. 1, 2
- [69] Mingtong Zhang, Kaifeng Zhang, and Yunzhu Li. Dynamic 3d gaussian tracking for graph-based neural dynamics modeling. *arXiv preprint arXiv:2410.18912*, 2024. 2, 5, 7
- [70] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In CVPR, 2018. 7
- [71] Tianyuan Zhang, Hong-Xing Yu, Rundi Wu, Brandon Y Feng, Changxi Zheng, Noah Snavely, Jiajun Wu, and William T Freeman. Physdreamer: Physics-based interaction with 3d objects via video generation. In European Conference on Computer Vision, pages 388–406. Springer, 2024. 2
- [72] Licheng Zhong, Hong-Xing Yu, Jiajun Wu, and Yunzhu Li. Reconstruction and simulation of elastic objects with springmass 3d gaussians. In *European Conference on Computer Vision*, pages 407–423. Springer, 2024. 2, 7