

PHYSICS3D: LEARNING PHYSICAL PROPERTIES OF 3D GAUSSIANS VIA VIDEO DIFFUSION

Anonymous authors

Paper under double-blind review

ABSTRACT

In recent years, there has been rapid development in 3D generation models, opening up new possibilities for applications such as simulating the dynamic movements of 3D objects and customizing their behaviors. However, current 3D generative models tend to focus only on surface features such as color and shape, neglecting the inherent physical properties that govern the behavior of objects in the real world. To accurately simulate physics-aligned dynamics, it is essential to predict the physical properties of materials and incorporate them into the behavior prediction process. Nonetheless, predicting the diverse materials of real-world objects is still challenging due to the complex nature of their physical attributes. In this paper, we propose **Physics3D**, a novel method for learning various physical properties of 3D objects through a video diffusion model. Our approach involves designing a highly generalizable physical simulation system based on a viscoelastic material model, which enables us to simulate a wide range of materials with high-fidelity capabilities. Moreover, we distill the physical priors from a video diffusion model that contains more understanding of realistic object materials. Extensive experiments demonstrate the effectiveness of our method with both elastic and plastic materials. Physics3D shows great potential for bridging the gap between the physical world and virtual neural space, providing a better integration and application of realistic physical principles in virtual environments. Project page: <https://physics3d-3dgs.github.io>

1 INTRODUCTION

In recent years, 3D computer vision has witnessed significant advancements, with researchers focusing on reconstructing or generating 3D assets Mildenhall et al. (2021); Kerbl et al. (2023); Tang et al. (2023); Poole et al. (2022); Hong et al. (2023); Li et al. (2024), and even delving into the realm of 4D dynamics Ren et al. (2023); Ling et al. (2023). However, a common feature in these works is the emphasis on color space, which can be difficult in modeling realistic interactive dynamics without any physical priors, especially for applications in areas such as virtual/augmented reality and animation. Physics simulation is one of the most crucial methods to achieve a deeper understanding of the real world and enhance the effectiveness of interactive dynamics. Although conventional methods Stewart (2000); Felippa (2004); Kilian & Ochsendorf (2005) describe behavior using continuous physical dynamic equations based on body-fixed mesh, they are usually difficult and time-consuming to generate complex 3D objects and suffer from highly nonlinear issues, large deformations or fracture-prone physical phenomena Jiang et al. (2016).

Powered by recent advances in implicit and explicit 3D representation techniques (e.g., NeRF Mildenhall et al. (2021) and 3D Gaussian Splatting Kerbl et al. (2023)), some researchers Xie et al. (2023); Li et al. (2022) have attempted to bridge the gap between rendering and simulation using the differentiable Material Point Method (MPM) Hu et al. (2018b), which enables efficient physical simulation driven by 3D particles (i.e., 3D Gaussian kernels). PhysGaussian Xie et al. (2023) extends the capabilities of 3D Gaussian kernels by incorporating physics-based attributes such as velocity, strain, elastic energy, and stress. This unified representation of material substance facilitates both simulation and rendering tasks. However, manual pre-design of physical parameters in PhysGaussian Xie et al. (2023) remains a laborious and imprecise process, where objects are categorized into six material types: jelly, metal, sand, foam, snow, and plasticine. Each category has distinct physical models and parameters, making it inconvenient to manually classify objects and initialize parameters based on

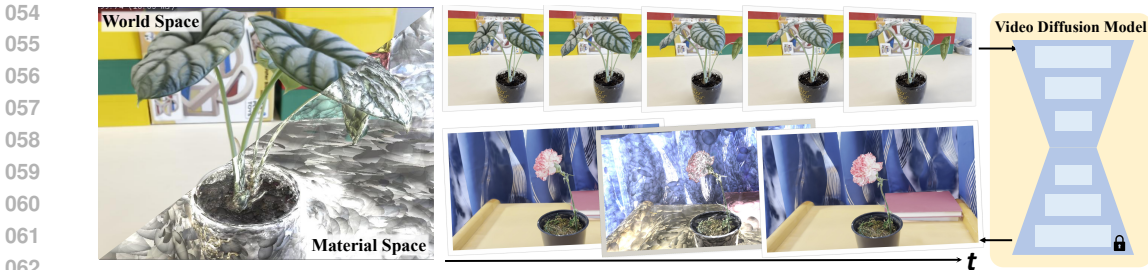


Figure 1: Physics3D is a unified simulation-rendering pipeline based on 3D Gaussians, which learn physics dynamics from video diffusion model. **The top row** shows the dynamic effects in the 2D RGB video space, which represents the starting of the SDS loss in the optimization process. **The bottom row** shows the dynamic effects in the 3D Gaussian space, representing the final optimization goal of SDS. Each 3D Gaussian, when rendered from a specific viewpoint, corresponds to a frame in the 2D video. The different objects in the two rows are intended for visual contrast and clarity.

expert knowledge. To avoid manually setting parameters, PhysDreamer Zhang et al. (2024) leverages object dynamics learned from video generation models Blattmann et al. (2023b); Wang et al. (2023) to estimate a physical material parameter (*i.e.*, Young’s modulus). However, in practical applications, real-world objects often exhibit a complex composite nature, making it challenging for a simulation approach that relies solely on a single physical parameter to fully capture their dynamic behavior. This limits PhysDreamer to be primarily tailored for the simulation of hyper-elastic materials. Specifically, it encounters significant challenges when dealing with materials such as plastics, metals, and non-Newtonian fluids due to its heavy reliance on optimizing Young’s modulus alone. The inherent complexities in these materials surpass the capabilities of PhysDreamer, highlighting the need for a more comprehensive and robust approach that considers a broader range of physical properties for accurate and effective simulation.

In this paper, we propose **Physics3D**, a generalizable physical simulation system to learn various physical properties of 3D objects. Given a 3D Gaussian representation, we first expand the dimension of the physical parameters to capture both elasticity and viscosity. Then we design a viscoelastic Material Point Method (MPM) to simulate 3D dynamics. Through the simulation process, we decompose the deformation gradient into two separate components and calculate them independently to contribute to the overall force. Finally, leveraging the capabilities of the differentiable MPM, we iteratively optimize both 3D Gaussian parameters and physical parameters via the Score Distillation Sampling (SDS) Poole et al. (2022) strategy to distill physical priors from the video diffusion model. Iterating the MPM process and SDS optimization, Physics3D achieves high-fidelity and realistic performance in a wide range of materials. Extensive experiments demonstrate the efficacy and superiority of our proposed Physics3D over existing methods. In summary, our key contributions are as follows.

- We propose a novel generalizable physical simulation system called Physics3D, which is capable of learning physical properties of diverse materials. We model physical properties with both elastoplastic and viscoelastic parts and design a parallel simulation framework.
- We design a physics-driven distillation strategy to iteratively optimize both filling 3D Gaussians and physical parameters, realising to generate realistic, physics-driven and controllable 3D dynamics while maintaining the model’s generalization capability with limited 3D data.
- Experiments show Physics3D is effective in creating high-fidelity and realistic 3D dynamics, ready for various interactions across users and objects in the future.

2 RELATED WORK

Dynamic 3D representations. Rapid advancements in static 3D representations have sparked interest in incorporating temporal dynamics into the 3D modeling of dynamic objects and scenes. Various explicit or hybrid representation techniques have demonstrated impressive outcomes, including planar decomposition for 4D space-time grids Cao & Johnson (2023); Shao et al. (2023); Fridovich-Keil et al. (2023), the utilization of NeRF representation Li et al. (2022); Pumarola et al. (2021); Gao et al.

(2021), and alternative structural approaches Turki et al. (2023); Abou-Chakra et al. (2024); Fang et al. (2022). Recently, 3D Gaussian Splatting Kerbl et al. (2023) has revolutionized the representation via its outstanding rendering efficiency and high-quality results. Efforts have been made to extend static 3D Gaussians into dynamic versions, yielding promising results. Dynamic 3D Gaussians Luiten et al. (2023) refine per-frame Gaussian Splatting through dynamic regularizations and shared properties such as size, color, and opacity. Similarly, the concept of 4D Gaussian Splatting Wu et al. (2023); Yang et al. (2023) employs a deformation network to anticipate time-dependent positional, scaling, and rotational deformations. In addition, DreamGaussian4D Ren et al. (2023) learns motion from image-conditioned generated videos Blattmann et al. (2023a), enabling more controllable and diverse 3D motion representations.

Viscoelastic materials. In the realm of computer graphics and animation, there has been significant interest in accurately simulating the behavior of nonrigid objects and their interactions with physical environments. Conventional elastic models Terzopoulos et al. (1987); Zong et al. (2023) predicated on Hooke’s law Rychlewski (1984) are the cornerstone for simulating the deformation of objects. These models are effective in representing materials that exhibit perfectly elastic behavior, returning to their original shape after the applied force is removed. However, real-world materials often exhibit more complex behaviors that cannot be captured by simple elastic models. The introduction of viscoelastic materials in computer graphics Terzopoulos & Fleischer (1988) has expanded the range of simulated material behaviors. Viscoelastic materials combine the characteristics of both viscous fluids and elastic solids, leading to time-dependent deformations under constant stress, a phenomenon known as creep Christensen (2003). Compared with elastic models, viscoelastic models offer a more versatile framework for animating nonrigid objects. They can simulate the slow restoration of a material’s shape after the cessation of force, as well as the permanent deformation that occurs due to prolonged stress.

Video generation models. With the emergence of models like Sora Brooks et al. (2024), the field of video generation Villegas et al. (2022); Wu et al. (2022); Bar-Tal et al. (2024); Blattmann et al. (2023c) has drawn significant attention. These powerful video models Kondratyuk et al. (2023); Singer et al. (2022); Ho et al. (2022); Hong et al. (2022) are typically trained on extensive datasets of high-quality video content. Sora Brooks et al. (2024) is capable of producing minute-long videos with realistic motions and consistent viewpoints. Furthermore, some large-scale video models, such as Sora Brooks et al. (2024), can even support physically plausible effects. Inspired by these capabilities, we aim to distill the physical principles observed in videos and apply them to our static 3D objects, thereby achieving more realistic and physically accurate results. In our framework, we choose Stable Video Diffusion Blattmann et al. (2023b) to optimize our physical properties.

3 PROBLEM FORMULATION

Given a static representation through 3D Gaussians, our goal is to estimate the physical attributes of each Gaussian particle and generate physics-plausible motions by organizing the interaction of force and velocity among these particles. For pure-elastic models, these physical properties include mass (m), Young’s modulus (E), and Poisson’s ratio (ν). Young’s modulus and Poisson’s ratio control the dynamics of elastic objects. For example, with a fixed amount of external force applied, the system’s higher Young’s modulus will have smaller deformation.

However, only modeling the property of elastic objects is inadequate for recovering diverse physics with heterogeneous materials in real-world applications, which significantly limits the recent work like PhysDreamer Zhang et al. (2024). For example, they usually suffer from complex mixed materials, especially in scenarios of rapid deformation where viscosity emerges as a significant factor in dynamics. Therefore, our key insight is to build a more comprehensive physics model that includes additional parameters, notably viscosity, to enrich the descriptive capacity for real-world objects, especially in inelastic scenarios.

To model viscoelastic stresses with physical fidelity, we explore continuum mechanics, where Lamé constants (also referred to as Lamé coefficients or parameters), denoted by λ and μ , emerge as pertinent material-related quantities within the strain-stress relationship, while viscosity coefficient ν_v and ν_d governs the viscous dynamics. Consequently, our framework pivots towards the estimation of the viscosity coefficient (ν_v and ν_d) and the two Lamé parameters (λ , and μ). As for other physical attributes, we align with the conventional methods Zhang et al. (2024); Xie et al. (2023), where the

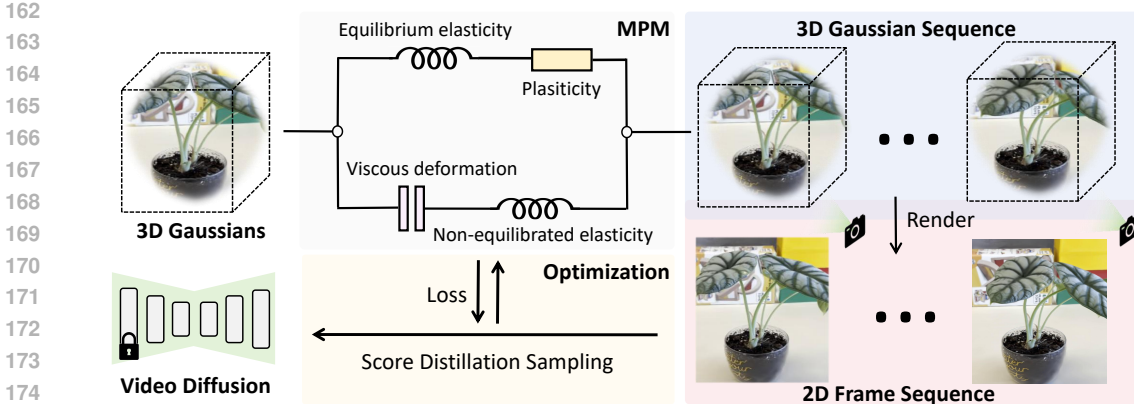


Figure 2: **Pipeline of Physics3D**. Given an object represented as 3D Gaussians, We first simulate it using the Material Point Method (MPM). The simulation comprises two distinct components: an elastoplastic part and a viscoelastic part. These components operate independently to calculate the individual stresses within the object and are combined to determine the overall stress. After the simulation, a series of Gaussians with varying orientations are generated, reflecting the dynamic evolution of the scene. Then, we render these Gaussians from a fixed viewpoint to produce a sequence of video frames. Finally, we utilize a pretrained video diffusion model with Score Distillation Sampling (SDS) strategy to iteratively optimize physical parameters.

particle mass (m_p) is pre-calculated as the product of a constant density (ρ) and the particle volume (V_p), and Poisson’s ratio (ν) is constant across the object. For detailed explanation and clarified summary of notations, please refer to Appendix C.1.

4 METHOD

In this section, we introduce our method, *i.e.*, Physics3D, for learning the dynamics of multi-material 3D systems with physical alignment. Our goal is to estimate the various physical properties of 3D objects. Building upon this, we first review the theory of three foundational techniques (Sec. 4.1) that form the backbone of our algorithm. Then we introduce our physical modeling framework and describe our particle-based simulation process (Sec. 4.2). Finally, we present the physical-based distillation strategy (Sec. 4.3) to iteratively optimize both filling 3D Gaussians and physical parameters with the video diffusion model. An overview of our framework is depicted in Figure 2.

4.1 PRELIMINARY

Continuum Mechanics describes motions by a deformation map $\mathbf{x} = \phi(\mathbf{X}, t)$ from the material space Ω^0 (with coordinate \mathbf{X}) to the world space Ω^n (with coordinate \mathbf{x}). The deformation gradient $\mathbf{F} = \frac{\partial \phi}{\partial \mathbf{X}}(\mathbf{X}, t)$ measures local rotation and strain Bonet & Wood (1997). We consider viscoelastic materials where we have two components Govindjee & Reese (1997), the elastoplastic component $\mathbf{F}_E \mathbf{F}_P$ and the viscoelastic component $\mathbf{F}_N \mathbf{F}_V$. They are in parallel combination shown in Figure 3, and formulated as:

$$\mathbf{F} = \mathbf{F}_E \mathbf{F}_P = \mathbf{F}_N \mathbf{F}_V. \quad (1)$$

Intuitively, we model our material with two different compounds in parallel connection as they deform in the same way thus having the same total strain. However, only the elastic components \mathbf{F}_E and \mathbf{F}_N contributes to the internal stress σ_E and σ_N . We can now evolve the system with dynamical equations. Denoting the velocity field with $\mathbf{v}(\mathbf{x}, t)$ and density field with $\rho(\mathbf{x}, t)$, the conservation of momentum and conservation of mass Germain (1998) is given by:

$$\rho \frac{D\mathbf{v}}{Dt} = \nabla \cdot \boldsymbol{\sigma} + \mathbf{f}, \quad \frac{D\rho}{Dt} + \rho \nabla \cdot \mathbf{v} = 0, \quad (2)$$

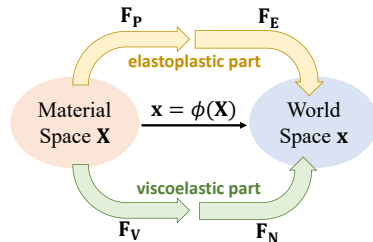


Figure 3: Elastoplastic and viscoelastic decomposition.

where \mathbf{f} denotes an external force, $\boldsymbol{\sigma} = \boldsymbol{\sigma}_E + \boldsymbol{\sigma}_N$ is the total internal stress. We will also need to update the strain tensor after updating the material point.

Material Point Method (MPM) Sulsky et al. (1995); Jiang et al. (2016); Hu et al. (2018b) discretizes the material into deformable particles and employs particles to monitor the complete history of strain and stress states while relying on a background grid for precisely evaluating derivatives during force computations. This methodology has demonstrated its efficacy in simulating diverse materials Jiang et al. (2015); Klár et al. (2016); Stomakhin et al. (2013) and is proved to be capable of simulating some viscoelastic and viscoplastic materials Ram et al. (2015); Yue et al. (2015). MPM operates in a particle-to-grid (P2G) and grid-to-particle (G2P) transfer loop. In P2G process, MPM transfers mass and momentum from particles to grids:

$$m_i^n = \sum_p w_{ip}^n m_p, \quad m_i^n \mathbf{v}_i^n = \sum_p w_{ip}^n m_p (\mathbf{v}_p^n + C_p^n (\mathbf{x}_i - \mathbf{x}_p^n)), \quad (3)$$

Here i and p represent the fields on the Eulerian grid and the Lagrangian particles respectively. Each particle p carries a set of properties including volume V_p , mass m_p , position \mathbf{x}_p^n , velocity \mathbf{v}_p^n , deformation gradient \mathbf{F}_p^n and affine momentum C_p^n at time t_n . w_{ip}^n is the B-spline kernel defined on i -th grid evaluated at \mathbf{x}_p^n . After P2G, the transferred grids can be updated as:

$$\mathbf{v}_i^{n+1} = \mathbf{v}_i^n - \frac{\Delta t}{m_i} \sum_p \tau_p^n \nabla w_{ip}^n V_p^0 + \Delta t \mathbf{g}, \quad (4)$$

here \mathbf{g} represents the acceleration due to gravity. Then G2P transfers velocities back to particles and updates particle states τ (*i.e.*, Kirchhoff tensor).

$$\tau_p^{n+1} = \tau(\mathbf{F}_E^{n+1}, \mathbf{F}_N^{n+1}), \quad (5)$$

where \mathbf{F}_E^{n+1} and \mathbf{F}_N^{n+1} are two parts of strain tensor. We will provide a detailed introduction in Sec. 4.2. In this work, we integrate the physical properties of viscoelastic materials into the Material Point Method, thereby enhancing the generalization capabilities of MPM. This implementation enables the simulation for a wide range of materials commonly found in the real world, including various inelastic materials. We refer to the Appendix for more details.

Score Distillation Sampling (SDS) is first introduced by DreamFusion Poole et al. (2022), which distills the 3D knowledge from large 2D pretrain model Saharia et al. (2022). For a set of particles parameterized by physical properties θ , its rendering \mathbf{x} can be obtained by $\mathbf{x} = g(\theta)$ where g is a differentiable renderer. SDS calculates the gradients of physical parameters θ by,

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}}(\phi, \mathbf{x} = g(\theta)) = \mathbb{E}_{t, \epsilon} \left[w(t) (\epsilon_{\phi}(\mathbf{x}_t; y, t) - \epsilon) \frac{\partial \mathbf{x}_t}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial \theta} \right], \quad (6)$$

where $w(t)$ is a weighting function that depends on the timestep t and y denotes the given condition. $\epsilon_{\theta}(\mathbf{x}_t, t)$ is autoencoder in diffusion model to estimate the origin distribution of the given noise \mathbf{x}_t .

4.2 PHYSICAL MODELING

Our goal is to generate 3D dynamics for diverse materials, especially inelastic materials and composite materials with more complex motion modeling. An intuitive solution is to capture these viscous effects by directly adding a viscosity part following the current elastic pipeline as PhysDreamer Zhang et al. (2024). However, such a simple series connection of components is prone to coupling of the two dissipative behaviors (elastoplastic and viscoelastic), which increases the complexity of the model simulation and disrupts the original elastic properties of the model.

In our paper, we instead design a novel parallel framework to avoid the coupling issue. We decompose the material space into elastoplastic and an additional viscoelastic part, aiming to unify the simulation of a wide range of materials. It showcases the flexibility of this framework, and by adding more blocks into both the static and dynamical modeling, it can be expected to be widely useful for more complicated scenes. To simulate the physical process with Gaussian particles, we employ a particle-based method MLS-MPM Hu et al. (2018a) as our simulator, and we formalize the simulation process for a single sub-step as follows:

$$\mathbf{x}^{n+1}, \mathbf{v}^{n+1}, \mathbf{F}^{n+1}, C^{n+1} = S(\mathbf{x}^n, \mathbf{v}^n, \mathbf{F}^n, C^n, \theta, \Delta t), \quad (7)$$

Here, θ contains the physical properties of all particles: mass m_i , Young’s modulus E_i , Poisson’s ratio ν_i , Lamé coefficients λ, μ , viscosity coefficient ν_{N_i} and volume V_i . Δt denotes the simulation step size. Within the MPM simulation, stress, as depicted in Equation 1, can be divided into two components: one representing elastoplasticity, denoted as \mathbf{F}_E , and the other referring to viscoelasticity, denoted as \mathbf{F}_N . The two parallel parts are also illustrated in Figure 2, consisting of a parallel combination: (a) an elastic part with a frictional element for plasticity and; (b) a viscous part which is assumed to capture the dissipation with an elastic part. Now, we elaborate on the computation for each component individually.

Model for elastoplastic part. We first explain how to compute the internal stress σ_E through \mathbf{F}_E , which is essential in updating kinematical variables like velocities v through Eq. 2, for the elastic part. We will speak out the rule here but postpone explanations and intuitions in Appendix C.4. As we demonstrate in Sec. 4.1 and Appendix C.3, one can compute the Cauchy stress tensor given the energy function. In this work, we choose the fixed corotated constitutive model for the elastic part, whose energy function is

$$\psi(\mathbf{F}_E) = \psi(\Sigma_E) = \mu_E \sum_i (\sigma_{E,i} - 1)^2 + \frac{\lambda_E}{2} (\det(\mathbf{F}_E) - 1)^2, \quad (8)$$

where Σ_E is the diagonal singular value matrix $\mathbf{F}_E = \mathbf{U}\Sigma_E\mathbf{V}^T$. And $\sigma_{E,i}$ are singular values of \mathbf{F}_E . From this energy function, we can compute the Cauchy stress tensor as:

$$\sigma_E = \frac{2\mu}{\det(\mathbf{F}_E)} (\mathbf{F}_E - \mathbf{R}) \mathbf{F}_E^T + \lambda_E (\det(\mathbf{F}_E) - 1), \quad (9)$$

where $\mathbf{R} = \mathbf{U}\mathbf{V}^T$. We now explain how to update the \mathbf{F}_E^n with the velocity field \mathbf{v}^n at n th step. For purely elastic case, this can be simply done via:

$$\mathbf{F}_E^{n+1} = (\mathbf{I} + \Delta t \nabla \mathbf{v}^n) \mathbf{F}_E^n, \quad (10)$$

where Δt is the length of the time segment in the MPM method. This formula represents the fact that the internal velocity field causes a change in strains.

For the plastic part of the branch, we constraint the singular value of \mathbf{F}_E to sit between $[1 - \theta_c, 1 + \theta_s]$, where θ_c and θ_s are learnable parameters quantifying the plasticity. More precisely, in the simulation algorithm, at n -th step, $\mathbf{F}^n = \mathbf{F}_E^n \mathbf{F}_P^n$. We then compute the internal stress $\sigma(\mathbf{F}_E^n)$ and update \mathbf{F}^n to \mathbf{F}^{n+1} with MPM method. Please refer to the Appendix C.4 for more details.

Model for viscoelastic part. We now explain the other branch of our model: the viscoelastic part. The total strain now consist of two parts $\mathbf{F} = \mathbf{F}_N \mathbf{F}_V$. There are two key features of this brunch. First, only the \mathbf{F}_N contributes to the internal stress. Secondly, \mathbf{F}_V only plays a role in the update rule of \mathbf{F}_N . For the first step, we know have a relation between internal stress σ_N (or equivalently Kirchoff tensor $\tau_N = \det(\mathbf{F}_N) \sigma_N$):

$$\tau_N = 2\mu_N \epsilon_N + \lambda_N \text{tr}(\epsilon_N) \mathbf{I} \quad (11)$$

where we denote τ_N as a vector of singular value of τ_N , in another word $\tau_N = U \text{diag}(\tau_N) V^T$. And ϵ_N , called log principle Kirchoff tensor, denotes a vector takes the diagonal element of $\log \Sigma_N$, where Σ_N as the diagonal singular value matrix $\mathbf{F}_N = U \Sigma_N V^T$. We can again update kinematical variables afterward. The second step is we need to modify the update rule for \mathbf{F}_N tensor. This boils down to a trial-and-correction procedure where we first update

$$\mathbf{F}_{N,\text{tr}}^n = (\mathbf{I} + \Delta t \nabla \mathbf{v}^n) \mathbf{F}_N^n, \quad (12)$$

then we modify the trial strain tensor $\mathbf{F}_{N,\text{tr}}^n$ by

$$\epsilon_N^{n+1} = A(\epsilon_{N,\text{tr}}^n - B \text{tr}(\epsilon_{N,\text{tr}}^n) \cdot \mathbf{1}) \quad (13)$$

where ϵ_N^{n+1} denotes the log principle Kirchoff tensor of \mathbf{F}_N^{n+1} . A and B are functions of viscosity parameters ν_d and ν_e . Please refer to Appendix C.4 for more detailed intuitions and explanations.

4.3 PHYSICS-DRIVEN DISTILLATION

Through iterations of the Material Point Method (MPM), we obtain a set of Gaussians evolving over time t . The general practice is that in order to optimize physical properties for each 3D Gaussian, we need high-quality 3D Gaussian models and corresponding motion ground truth. This ground truth could come from various sources, such as single-view or multi-view real-captured videos, 4D videos, or ideally, detailed 3D models with comprehensive physical property labels.

However, acquiring such datasets is currently challenging. Inspired by former works like Zhang et al. (2024); Ren et al. (2023) which learn 3D dynamics from 2D videos, we intend to create 3D dynamics that appear as realistic videos when given different initial external forces and rendered from random angles. Such 3D dynamics can be specified as a differentiable video parameterization, where a differentiable generator g (here indicates MPM processor and following Gaussian rasterizer) transforms both physical parameters and Gaussian parameters θ to create a video $x = g(\theta)$. Therefore, we extend the score distillation sampling (SDS) to iteratively optimize both 3D Gaussian parameters and physical parameters from the video diffusion model. To ensure consistency, the camera remains stationary as we capture and render each Gaussian to produce a reference image I_t^r . To supervise and optimize this process, we employ two optional models: image-to-video and text-to-video diffusion models. Then, we use SDS loss to guide our optimization process:

$$\nabla_{\theta} \mathcal{L}_{\text{SDS}} = \mathbb{E}_{t,p,\epsilon} \left[w(t) (\epsilon_{\phi}(I_t^p; t, I_t^r, \Delta p, y) - \epsilon) \frac{\partial I_t^p}{\partial \theta} \right]. \quad (14)$$

Here, $w(t)$ represents a time-dependent weighting function, $\epsilon_{\phi}(\cdot)$ denotes the predicted noise generated by the 2D diffusion prior ϕ , I_t^p represents a rendered video frame of diffusion timestep t from a camera pose p , Δp signifies the relative change in camera pose from the reference camera r , and y denotes the given condition (*i.e.*, image or text). Furthermore, we adopt a partial filling strategy like Xie et al. (2023), where internal volumes of select solid objects are optionally filled to augment simulation realism. For better rendering quality, we optimize the filled Gaussian along with the SDS process. The optimization objective for filled Gaussians is defined as:

$$\mathcal{L}_{\text{Fill}} = \frac{1}{T} \sum_{t=1}^T \lambda \|I_t^p - I_t^r\|_2^2. \quad (15)$$

For more details about the learnable internal filling strategy, please refer to Appendix C.5.

5 EXPERIMENTS

In this section, we conduct extensive experiments to evaluate our Physics3D and show the comparison results against other methods Xie et al. (2023); Zhang et al. (2024); Ren et al. (2023). We first present our qualitative results and comparisons with baselines (Sec. 5.2). Then we report the quantitative results along with a user study (Sec. 5.3). Finally, we carry out more open settings and ablation studies to further verify the efficacy of our framework design (Sec. 5.4). Please refer to the Appendix for more visualizations and detailed analysis.

5.1 EXPERIMENT SETUP

Datasets. We evaluate our method for generating diverse dynamics using several sources of input. We choose four real-world static scenes from PhysDreamer Zhang et al. (2024) for fair comparison. Each scene includes an object and a background. The objects include a carnation, an alocasia plant, a telephone cord, and a beanie hat. We also employed NeRF datasets Mildenhall et al. (2021) to evaluate the efficacy of our proposed model. Additionally, we utilize BlenderNeRF Raafat (2023) to synthesize several scenes. For more dataset details, please refer to Append A.2.

Implementation Details. In our implementation, we initiate the process by reconstructing 3D Gaussians from multi-view images, establishing a foundational representation of the scene. In complex realistic cases, we undertake a segmentation step to differentiate between the background and foreground elements, focusing solely on the latter for subsequent simulation tasks. Prior to

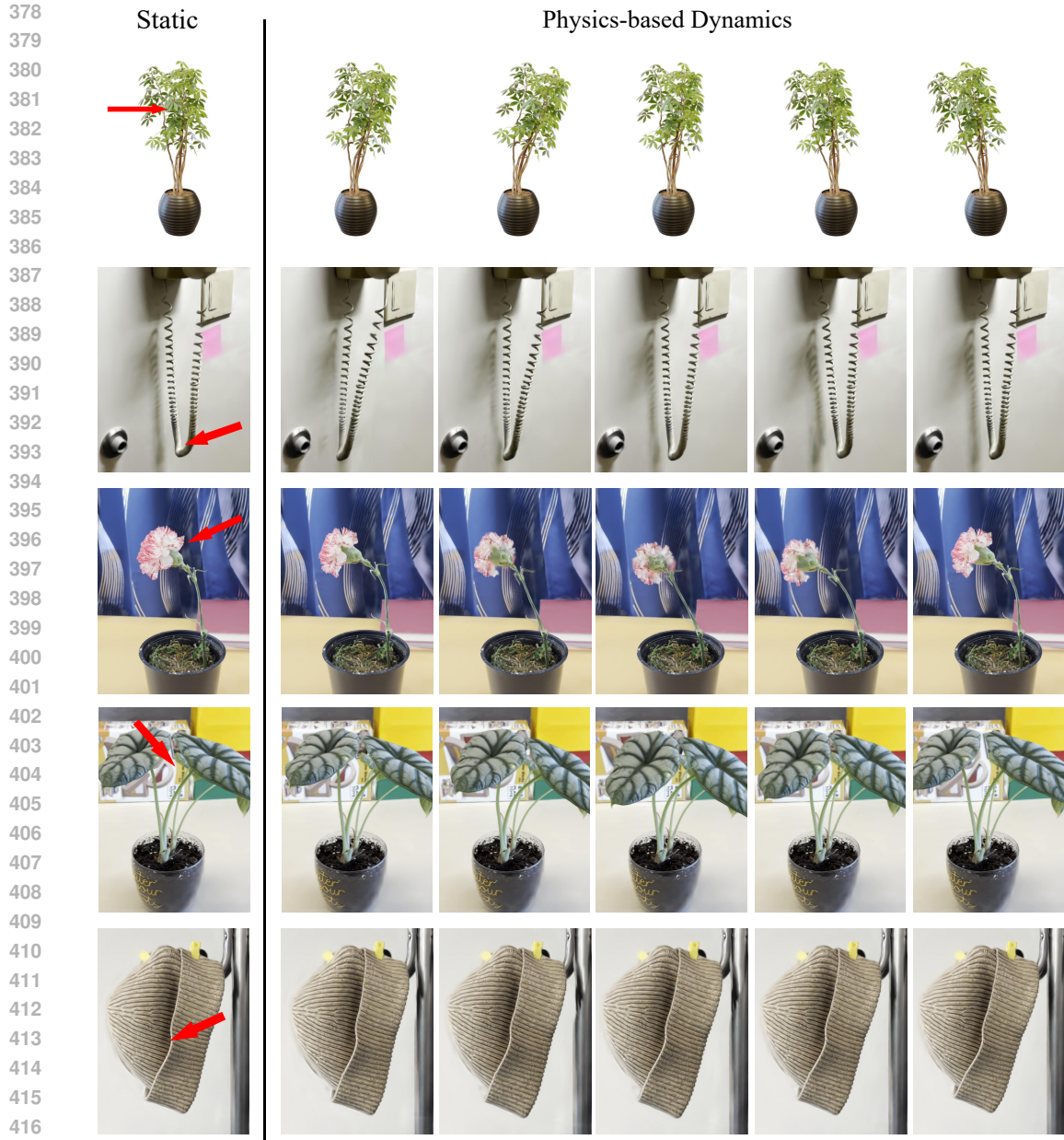


Figure 4: **Visual results of Physics3D** on different subjects with an external force (red arrow). Physics3D is able to generate realistic scene movement while maintaining good motion consistency.

simulation, we execute internal particle filling operations to refine the representation further. Each Gaussian kernel is then associated with a distinct set of physical properties targeted for optimization, facilitating the fine-tuning process. Subsequently, we discretize the foreground region into a grid structure, typically set at dimensions of 50^3 . As for the MPM simulation, we employ 400 sub-steps within each temporal interval spanning successive video frames. This temporal granularity translates to a sub-step duration of $1e - 4$ second, ensuring precision and accuracy in the simulation dynamics. Notably, the optimization process for each object requires approximately 5 minutes to complete on a single NVIDIA A6000 (48GB) GPU.

Baselines and Metrics. We extensively compare our method with three baselines: PhysDreamer Zhang et al. (2024), PhysGaussian Xie et al. (2023) and DreamGaussian4D Ren et al. (2023). Following Zhang et al. (2024), we show our results with notable comparisons with space-time slices. For metrics, we evaluate our approach with the video-quality metrics : PSNR, SSIM, MS-

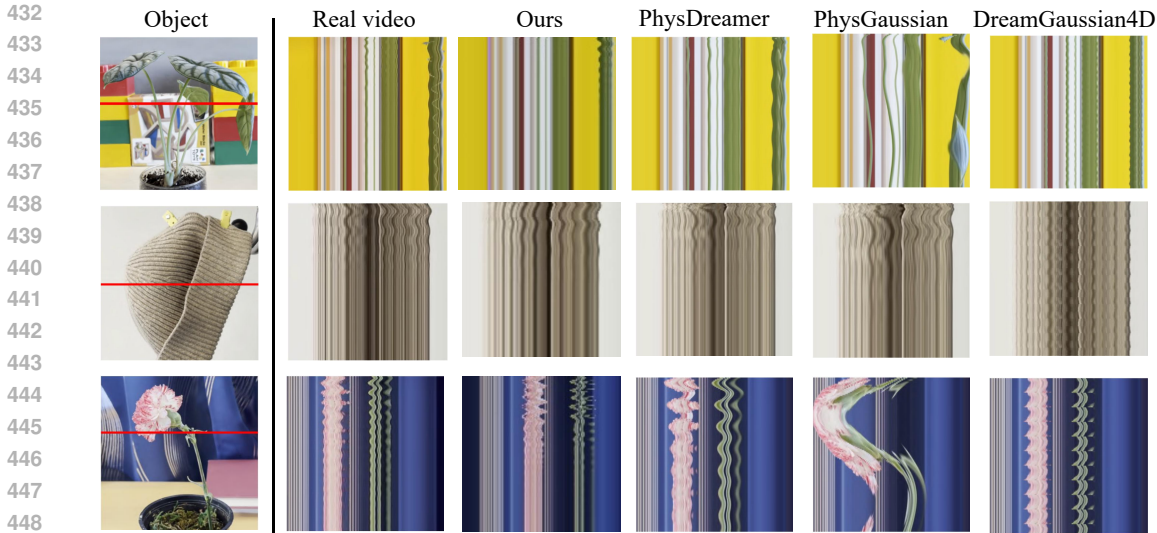


Figure 5: **Comparison results** on different subjects. From the visual results, we observe that PhysDreamer Zhang et al. (2024) only estimates the elastic properties of objects, resulting in the lack of damping. DreamGaussian4D Ren et al. (2023) and PhysGaussian Xie et al. (2023) are respectively limited in unrealistically constant, low-magnitude periodic motion and low-frequency movements. In contrast, our model successfully balances high and low-frequency oscillations with more realistic damping.

SSIM, and VMAF on the rendered videos. Currently, there is a lack of comprehensive 3D dynamic evaluation metrics. PSNR, SSIM and MS-SSIM Free Software Foundation, Inc. (1991) are used to evaluate the quality of generated frames, while video-specific metrics (VMAF Netflix, Inc. (2020)) are employed to assess the dynamic effects of the generated sequences.

5.2 QUALITATIVE RESULTS

Figure 4 shows qualitative results of simulated interactive motion. For each case, we visualize one example with its initial scene and deformation sequence. The results demonstrate our model’s capability of simulating the movement of complex textured objects, presenting a realistic and physically plausible outcome. Additional experiments are included in the Appendix B.1. Please visit our project page: <https://physics3d-3dgs.github.io> for more dynamic visual results.

Following PhysDreamer Zhang et al. (2024), we compare our results with real captured videos and simulations from other methods Zhang et al. (2024); Xie et al. (2023); Ren et al. (2023) in Figure 5. We utilize space-time slices to present our comparisons. These slices depict time along the vertical axis and spatial slices of the object along the horizontal axis, as indicated by the red lines in the "object" column. Through these visualizations, we aim to elucidate the magnitude and frequencies of the oscillating motions under scrutiny. PhysDreamer Zhang et al. (2024) closely approximates the elastic properties of objects, resulting in periodic oscillations with subtle damping, contrasting the unrealistic aspects. PhysGaussian Xie et al. (2023) showcases unrealistically low-frequency movements due to inaccurate parameter settings. DreamGaussian4D Ren et al. (2023) lacks physical prior and generates unrealistically constant, low-magnitude periodic motion. In contrast, our model successfully balances high and low-frequency oscillations with more realistic damping, aligning more closely with the behavior of objects in the real world.

5.3 QUANTITATIVE RESULTS

Table 1 shows the average video-quality metrics over rendered videos from the same fixed perspective. Results clearly demonstrate higher scores for Physics3D, indicating better video

Table 1: **Quantitative comparisons** on rendered videos using different video-quality metrics.

	PSNR↑	SSIM↑	MS-SSIM↑	VMAF↑
PhysDreamer Zhang et al. (2024)	13.89	0.55	0.37	0.52
PhysGaussian Xie et al. (2023)	13.86	0.57	0.39	0.59
Ours	14.72	0.59	0.49	0.59
- w/o viscoelastic part	14.13	0.53	0.36	0.52
- w/o elastoplastic part	13.53	0.55	0.41	0.50

486 quality and motion consistency of our
 487 results. We also conduct numerous user experiments in Appendix B.3.

489 5.4 ABLATION STUDY

491 We conduct ablation study in Figure 6 to evaluate the efficacy of our physics modeling. Specifically,
 492 we investigate the importance of elastoplastic and viscoelastic components from the model architec-
 493 ture. We observe that removing either of the modules leads to a notable degradation in the realism of
 494 the physical simulations. Particularly, the absence of the elastoplastic component results in a lack
 495 of elasticity, making objects more susceptible to shape deformation and fluid-like behavior. On the
 496 other hand, the absence of the viscoelastic component leads to a deficiency in sustained damping and
 497 rebound effects, especially in scenarios with minimal external disturbances where energy dissipation
 498 occurs rapidly. These show the significance of both components of our model in capturing the
 499 dynamics of physical objects. Please refer to our Appendix B.2 for more ablations.



511 Figure 6: Ablation study on elastoplastic and viscoelastic parts of our model.

513 6 CONCLUSION

515 In this paper, we present Physics3D, a novel framework to learn various physical properties of
 516 3D objects from video diffusion model. Our method tackles the challenge of estimating diverse
 517 material properties by incorporating two key components: elastoplastic and viscoelastic modules.
 518 The elastoplastic component facilitates simulations of pure elasticity, while the viscoelastic module
 519 introduces damping effects, crucial for capturing the behavior of materials exhibiting both elasticity
 520 and viscosity. Furthermore, we leverage a video generation model to distill inherent physical priors,
 521 improving our understanding of realistic material characteristics. Extensive experiments show
 522 Physics3D is effective in creating high-fidelity and realistic 3D dynamics.

523 **Limitations and Future Work.** In complex environments with a lot of entangled objects, our method
 524 requires manual intervention to assign the scope of movable objects and define the filling ranges for
 525 objects, which is not efficient for more real applications. In the future, we aim to utilize the prior of
 526 large segmentation models to solve the problem with more comprehensive physics system modeling.
 527 We believe that Physics3D takes a significant step to open up a wide range of applications from
 528 realistic simulations to interactive virtual experiences and will inspire more works in the future.

530 REFERENCES

- 531 Jad Abou-Chakra, Feras Dayoub, and Niko Sünderhauf. Particlenerf: A particle-based encoding for
 532 online neural radiance fields. In *Proceedings of the IEEE/CVF Winter Conference on Applications
 533 of Computer Vision*, pp. 5975–5984, 2024.
- 534 Omer Bar-Tal, Hila Chefer, Omer Tov, Charles Herrmann, Roni Paiss, Shiran Zada, Ariel Ephrat,
 535 Junhwa Hur, Yuanzhen Li, Tomer Michaeli, et al. Lumiere: A space-time diffusion model for
 536 video generation. *arXiv preprint arXiv:2401.12945*, 2024.
- 537 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik
 538 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling
 539 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023a.

- 540 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik
541 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling
542 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023b.
- 543
544 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and
545 Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models.
546 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
547 22563–22575, 2023c.
- 548 Javier Bonet and Richard D Wood. *Nonlinear continuum mechanics for finite element analysis*.
549 Cambridge University Press, 1997.
- 550
551 Tim Brooks, Bill Peebles, Connor Homes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
552 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video
553 generation models as world simulators. 2024. URL [https://openai.com/research/
554 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 555 Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of
556 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 130–141, 2023.
- 557
558 Eric R Chan, Connor Z Lin, Matthew A Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio
559 Gallo, Leonidas J Guibas, Jonathan Tremblay, Sameh Khamis, et al. Efficient geometry-aware 3d
560 generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision
561 and pattern recognition*, pp. 16123–16133, 2022.
- 562 Richard M Christensen. *Theory of viscoelasticity*. Courier Corporation, 2003.
- 563
564 Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner,
565 and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia
566 2022 Conference Papers*, pp. 1–9, 2022.
- 567 Yu Fang, Minchen Li, Ming Gao, and Chenfanfu Jiang. Silly rubber: an implicit material point
568 method for simulating non-equilibrated viscoelastic and elastoplastic solids. *ACM Transactions on
569 Graphics (TOG)*, 38(4):1–13, 2019.
- 570
571 Carlos A Felippa. Introduction to finite element methods. *University of Colorado*, 885, 2004.
- 572 Free Software Foundation, Inc. Video Quality Metrics, 1991. URL [https://github.com/
573 aizvorski/video-quality](https://github.com/aizvorski/video-quality).
- 574
575 Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo
576 Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of
577 the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12479–12488, 2023.
- 578
579 Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic
580 monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
581 pp. 5712–5721, 2021.
- 582 Paul Germain. Functional concepts in continuum mechanics. *Meccanica*, 33:433–444, 1998.
- 583
584 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
585 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the
586 ACM*, 63(11):139–144, 2020.
- 587 Sanjay Govindjee and Stefanie Reese. A presentation and comparison of two large deformation
588 viscoelasticity models. 1997.
- 589
590 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models. *Advances in
591 Neural Information Processing Systems*, 33:6840–6851, 2020.
- 592
593 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P
Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition
video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

- 594 Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale
595 pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
596
- 597 Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan Sunkavalli,
598 Trung Bui, and Hao Tan. Lrm: Large reconstruction model for single image to 3d. In *The Twelfth
599 International Conference on Learning Representations*, 2023.
- 600 Yuanming Hu, Yu Fang, Ziheng Ge, Ziyin Qu, Yixin Zhu, Andre Pradhana, and Chenfanfu Jiang. A
601 moving least squares material point method with displacement discontinuity and two-way rigid
602 body coupling. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018a.
- 603 Yuanming Hu, Yu Fang, Ziheng Ge, Ziyin Qu, Yixin Zhu, Andre Pradhana, and Chenfanfu Jiang. A
604 moving least squares material point method with displacement discontinuity and two-way rigid
605 body coupling. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018b.
606
- 607 Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative
608 and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*,
609 2023.
- 610 Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching.
611 *Journal of Machine Learning Research*, 6(4), 2005.
612
- 613 Chenfanfu Jiang, Craig Schroeder, Andrew Selle, Joseph Teran, and Alexey Stomakhin. The affine
614 particle-in-cell method. *ACM Transactions on Graphics (TOG)*, 34(4):1–10, 2015.
- 615 Chenfanfu Jiang, Craig Schroeder, Joseph Teran, Alexey Stomakhin, and Andrew Selle. The material
616 point method for simulating continuum materials. In *ACM SIGGRAPH 2016 courses*, pp. 1–52,
617 2016.
- 618 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting
619 for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), 2023.
620
- 621 Axel Kilian and John Ochsendorf. Particle-spring systems for structural form finding. *Journal of the
622 international association for shell and spatial structures*, 46(2):77–84, 2005.
- 623 Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint
624 arXiv:1312.6114*, 2013.
625
- 626 Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models.
627 *NeurIPS*, 2021.
- 628 Gergely Klár, Theodore Gast, Andre Pradhana, Chuyuan Fu, Craig Schroeder, Chenfanfu Jiang, and
629 Joseph Teran. Drucker-prager elastoplasticity for sand animation. *ACM Transactions on Graphics
630 (TOG)*, 35(4):1–12, 2016.
- 631 Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig
632 Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model
633 for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023.
634
- 635 Agelos Kratimenos, Jiahui Lei, and Kostas Daniilidis. Dynmf: Neural motion factorization for
636 real-time dynamic view synthesis with 3d gaussian splatting. *arXiv preprint arXiv:2312.00112*,
637 2023.
- 638 Xiaoyu Li, Qi Zhang, Di Kang, Weihao Cheng, Yiming Gao, Jingbo Zhang, Zhihao Liang, Jing
639 Liao, Yan-Pei Cao, and Ying Shan. Advances in 3d generation: A survey. *arXiv preprint
640 arXiv:2401.17807*, 2024.
- 641 Xuan Li, Yi-Ling Qiao, Peter Yichen Chen, Krishna Murthy Jatavallabhula, Ming Lin, Chenfanfu
642 Jiang, and Chuang Gan. Pac-nerf: Physics augmented continuum neural radiance fields for
643 geometry-agnostic system identification. In *The Eleventh International Conference on Learning
644 Representations*, 2022.
645
- 646 Huan Ling, Seung Wook Kim, Antonio Torralba, Sanja Fidler, and Karsten Kreis. Align your
647 gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. *arXiv preprint
arXiv:2312.13763*, 2023.

- 648 Fangfu Liu, Diankun Wu, Yi Wei, Yongming Rao, and Yueqi Duan. Sherpa3D: Boosting high-fidelity
649 text-to-3d generation via coarse 3d prior. *arXiv preprint arXiv:2312.06655*, 2023.
- 650
- 651 Fangfu Liu, Hanyang Wang, Weiliang Chen, Haowen Sun, and Yueqi Duan. Make-Your-3D: Fast
652 and consistent subject-driven 3d content generation. *arXiv preprint arXiv:2403.09625*, 2024.
- 653 Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians:
654 Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023.
- 655
- 656 Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and
657 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications
658 of the ACM*, 65(1):99–106, 2021.
- 659 Netflix, Inc. VMAF: Perceptual video quality assessment based on multi-method fusion, 2020. URL
660 <https://github.com/Netflix/vmaf>.
- 661
- 662 Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking
663 of non-rigid scenes in real-time. In *IEEE Conference on Computer Vision and Pattern Recognition
664 (CVPR)*, 2015.
- 665 Alexander Quinn Nichol and Prafulla Dhariwal. Improved Denoising Diffusion Probabilistic Models.
666 In *International Conference on Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 667
- 668 Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d
669 diffusion. In *The Eleventh International Conference on Learning Representations*, 2022.
- 670 Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural
671 radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer
672 Vision and Pattern Recognition*, pp. 10318–10327, 2021.
- 673
- 674 Maxime Raafat. BlenderNeRF, May 2023. URL [https://github.com/maximeraafat/
675 BlenderNeRF](https://github.com/maximeraafat/BlenderNeRF).
- 676 Daniel Ram, Theodore Gast, Chenfanfu Jiang, Craig Schroeder, Alexey Stomakhin, Joseph Teran,
677 and Pirouz Kavehpour. A material point method for viscoelastic fluids, foams and sponges. In
678 *Proceedings of the 14th ACM SIGGRAPH/Eurographics Symposium on Computer Animation*, pp.
679 157–163, 2015.
- 680 Jiawei Ren, Liang Pan, Jiaxiang Tang, Chi Zhang, Ang Cao, Gang Zeng, and Ziwei Liu. Dreamgaus-
681 sian4d: Generative 4d gaussian splatting. *arXiv preprint arXiv:2312.17142*, 2023.
- 682
- 683 Herbert E. Robbins. *An Empirical Bayes Approach to Statistics*. Springer New York, 1992.
- 684
- 685 Jan Rychlewski. On hooke’s law. *Journal of Applied Mathematics and Mechanics*, 48(3):303–314,
686 1984.
- 687 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
688 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
689 text-to-image diffusion models with deep language understanding. *Advances in neural information
690 processing systems*, 35:36479–36494, 2022.
- 691
- 692 Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. Tensor4d:
693 Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. In
694 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
695 16632–16642, 2023.
- 696
- 697 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry
698 Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video
699 data. *arXiv preprint arXiv:2209.14792*, 2022.
- 700
- 701 David E Stewart. Rigid-body dynamics with friction and impact. *SIAM review*, 42(1):3–39, 2000.
- Alexey Stomakhin, Craig Schroeder, Lawrence Chai, Joseph Teran, and Andrew Selle. A material
point method for snow simulation. *ACM Transactions on Graphics (TOG)*, 32(4):1–10, 2013.

- 702 Deborah Sulsky, Shi-Jian Zhou, and Howard L Schreyer. Application of a particle-in-cell method to
703 solid mechanics. *Computer physics communications*, 87(1-2):236–252, 1995.
704
- 705 Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. Dreamgaussian: Generative
706 gaussian splatting for efficient 3d content creation. *arXiv preprint arXiv:2309.16653*, 2023.
707
- 708 Demetri Terzopoulos and Kurt Fleischer. Modeling inelastic deformation: viscoelasticity, plasticity,
709 fracture. In *Proceedings of the 15th annual conference on Computer graphics and interactive
710 techniques*, pp. 269–278, 1988.
711
- 712 Demetri Terzopoulos, John Platt, Alan Barr, and Kurt Fleischer. Elastically deformable models. In
713 *Proceedings of the 14th annual conference on Computer graphics and interactive techniques*, pp.
205–214, 1987.
- 714 Haithem Turki, Jason Y Zhang, Francesco Ferroni, and Deva Ramanan. Suds: Scalable urban
715 dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
716 Recognition*, pp. 12375–12385, 2023.
717
- 718 Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang,
719 Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable
720 length video generation from open domain textual descriptions. In *International Conference on
721 Learning Representations*, 2022.
- 722 Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Mod-
723 elscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
724
- 725 Zhengyi Wang, Yikai Wang, Yifei Chen, Chendong Xiang, Shuo Chen, Dajiang Yu, Chongxuan Li,
726 Hang Su, and Jun Zhu. Crm: Single image to 3d textured mesh with convolutional reconstruction
727 model. *arXiv preprint arXiv:2403.05034*, 2024.
- 728 Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual
729 synthesis pre-training for neural visual world creation. In *European conference on computer vision*,
730 pp. 720–736. Springer, 2022.
- 731 Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian,
732 and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint
733 arXiv:2310.08528*, 2023.
734
- 735 Tianyi Xie, Zeshun Zong, Yuxin Qiu, Xuan Li, Yutao Feng, Yin Yang, and Chenfanfu Jiang. Physgaus-
736 sian: Physics-integrated 3d gaussians for generative dynamics. *arXiv preprint arXiv:2311.12198*,
737 2023.
- 738 Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. De-
739 formable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint
740 arXiv:2309.13101*, 2023.
741
- 742 Yonghao Yue, Breannan Smith, Christopher Batty, Changxi Zheng, and Eitan Grinspun. Continuum
743 foam: A material point method for shear-dependent flows. *ACM Transactions on Graphics (TOG)*,
744 34(5):1–20, 2015.
- 745 Tianyuan Zhang, Hong-Xing Yu, Rundu Wu, Brandon Y Feng, Changxi Zheng, Noah Snavely, Jiajun
746 Wu, and William T Freeman. Physdreamer: Physics-based interaction with 3d objects via video
747 generation. *arXiv preprint arXiv:2404.13026*, 2024.
- 748 Zeshun Zong, Xuan Li, Minchen Li, Maurizio M Chieramonte, Wojciech Matusik, Eitan Grinspun,
749 Kevin Carlberg, Chenfanfu Jiang, and Peter Yichen Chen. Neural stress fields for reduced-order
750 elastoplasticity and fracture. In *SIGGRAPH Asia 2023 Conference Papers*, pp. 1–11, 2023.
751
752
753
754
755

APPENDIX

A ADDITIONAL IMPLEMENTATION DETAILS

A.1 TRAINING DETAILS.

During the iterative optimization process, we use a specific video diffusion model in our experiment. For the image-to-video model, we utilize Stable Video Diffusion Blattmann et al. (2023b) that learns rich 2D diffusion prior from large real-world data. For a text-to-video diffusion model, we use a diffusion-based model text-to-video-ms-1.7b Wang et al. (2023), which augments the UNet structure with a cross-attention mechanism, which is an effective approach to condition the visual content on texts Huang et al. (2023). Following the text-to-video model, we use the text embedding of the prompt in the spatial attention block as the key and value in the multi-head attention layer. This enables the intermediate UNet features to aggregate text features seamlessly, thereby facilitating an alignment of language and vision embeddings. To ensure a great alignment between language and vision, the text encoder from pre-trained CLIP ViT-H/14 is used to convert the prompt into text embedding. In our experiment, we use text descriptions of object motion as conditional inputs instead of reference images. For example, in Figure 4, the top row uses a text prompt "a ficus swaying in the wind" to generate the video sequence.

In our experiment, we give external force to different objects. For synthetic data, we use a uniform, randomly initialized force for each method. For real-world data in PhysDreamer Zhang et al. (2024), due to the absence of explicit initial external forces and for fair comparisons with PhysDreamer, we choose to estimate an initial force based on the initial velocity field predicted by PhysDreamer. This force is manually tested to perform equally to the motion with the initial velocity field in PhysDreamer. Then we apply the estimated initial force across all baselines in real-world data.

For the representation of the physical parameters to be learned, We do not utilize the tri-plane Chan et al. (2022) approach which is employed in PhysDreamer Zhang et al. (2024). The tri-plane method, while useful, involves a degree of data compression that can compromise the quality of the material representation. Instead, we assign specific physical parameters to each Gaussian kernel directly because of our pursuit of higher fidelity in modeling.

A.2 DATASET DETAILS

For real-world data, we are leveraging the data in Figure 4 from PhysDreamer Zhang et al. (2024) (carnation, alopecia, hat, telephone cord), which is indeed sourced from real-world scenarios. Due to the limited availability of 3D data and its complexities, following the approach of PhysDreamer allows us to verify the effectiveness of Physics3D on real-world data. Our work in this area continues to need more comprehensive 3D datasets, whether synthetic or real.

For a fair comparison with PhysDreamer Zhang et al. (2024) and others, we use the ground truth videos captured from Physdreamer to compute the reconstruction metrics. Since synthetic 3D data lacks dynamic ground truth, we use it simply for qualitative visual results. We also conduct a user study to further demonstrate the overall quality and the alignment with real-world physics in our simulation, which can be found in Figure 9.

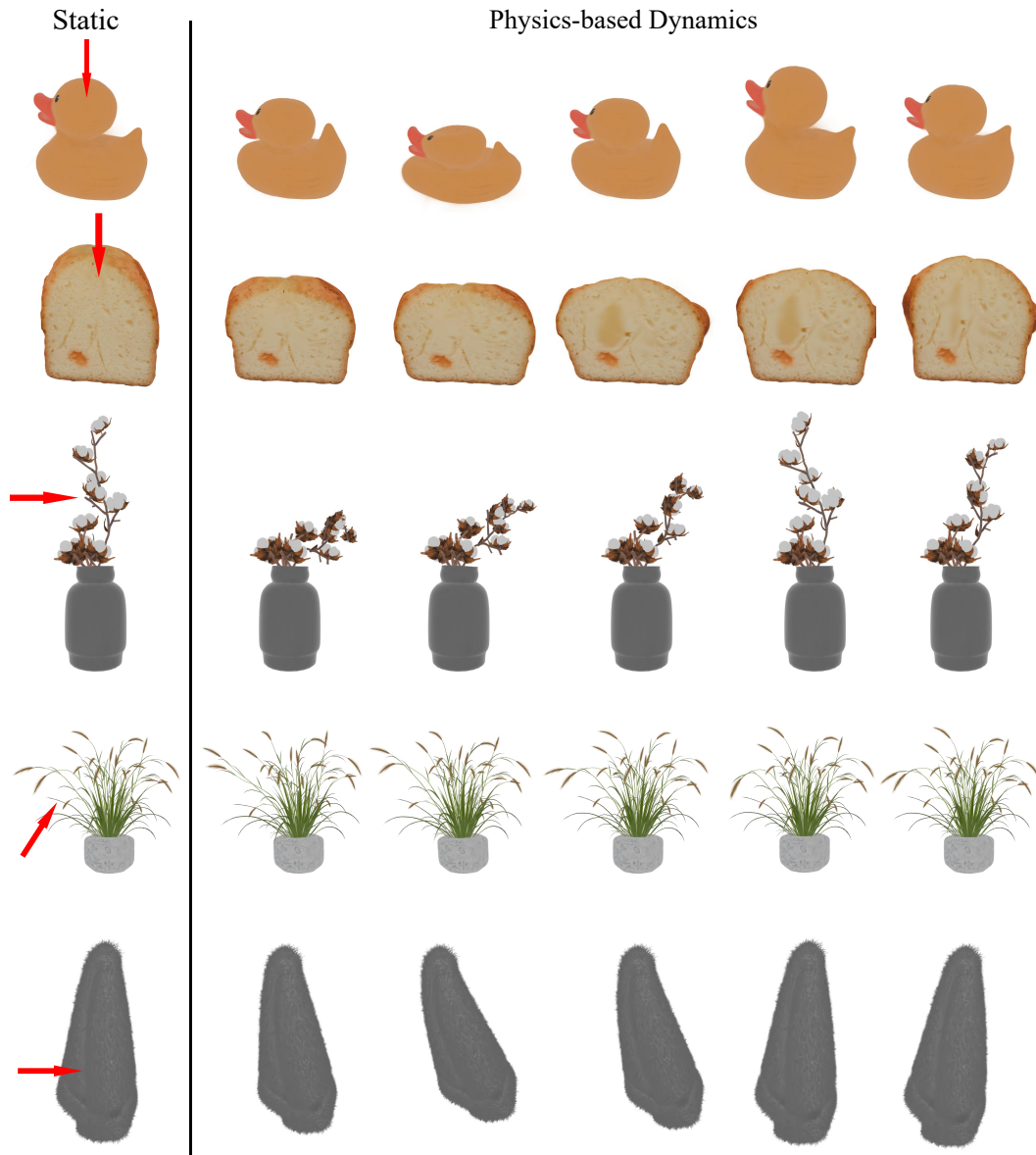
B MORE RESULTS

To further demonstrate the effectiveness and impressive visualization results of our Physics3D, we conducted more experiments including additional visual results and user study.

B.1 ADDITIONAL VISUAL RESULTS

Figure 7 shows additional visual results of interactive video sequences. We use BlenderNeRF Raafat (2023) to synthesize the cases below. For each case, We apply external forces of different directions and magnitudes to static objects and render video frames to show their motion states.

810 For more complex scenes, we have done more experiments on complex synthetic data such as
 811 collisions or fluid dynamics. Please refer to Figure 8, 9 for visual results of complex synthetic data.
 812 We are also planning to construct more complex real-world dynamic scenes using the segmentation
 813 method to separate objects and their movements, thereby enabling a more detailed analysis of dynamic
 814 behaviors. While the scarcity of diverse and high-quality 3D data poses remains challenging, we are
 815 making efforts to expand our dataset and refine our methodology. We are still committed to advancing
 816 our research by incorporating more complex synthetic and real-world data.



855

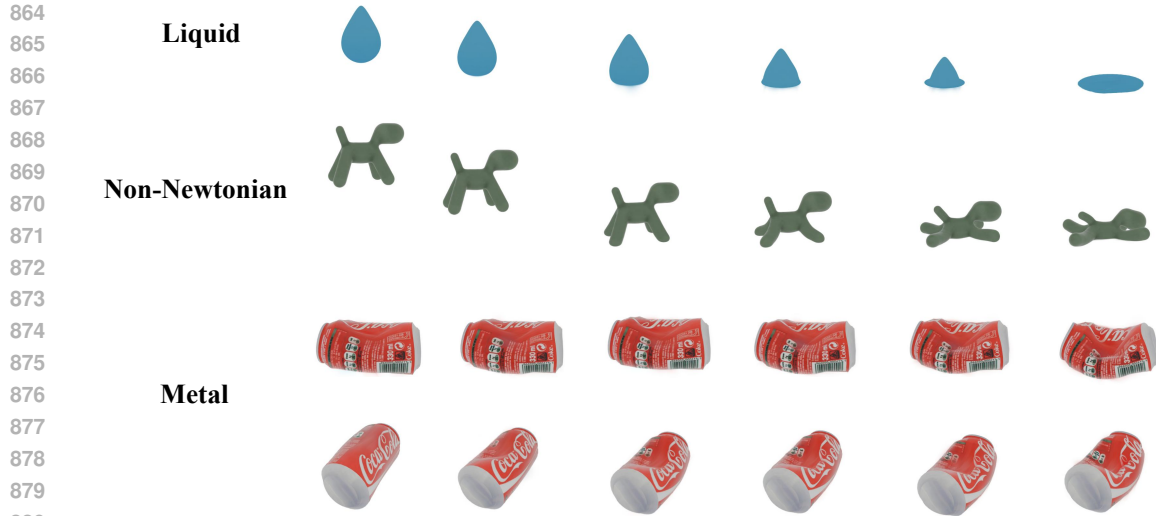
856 Figure 7: More visual results on different objects. Notice ours can simulate a variety of non-elastic
 857 and composite materials, including *rubber*, *fluffy bread*, *fabric* etc.

858

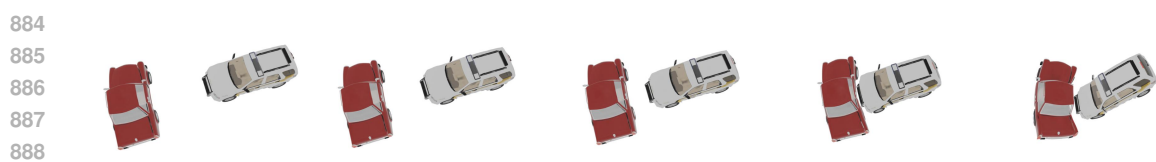
859 B.2 ADDITIONAL ABLATION STUDY

860

861 We carry out additional ablation analyses on the Physics3D design in Figure 6 to assess the efficacy
 862 of our physics modeling process. Specifically, we investigate the impact of removing the elastoplastic
 863 and viscoelastic components from the model architecture. The findings highlight that the removal
 of either of these modules leads to a notable decrease in the realism of the physical simulations.



881 Figure 8: **More visual results of various materials** such as *liquid*, *non-Newtonian fluids* and *metal*
882 (*a coke can is pressed*).



889 Figure 9: **Variety of dynamic interactions** such as *car collision*.

891
892
893
894
895
896
897
898

Particularly, the absence of the elastoplastic component results in a lack of elasticity, thereby making objects more susceptible to shape deformation and fluid-like behavior. On the other hand, the absence of the viscoelastic component leads to a deficiency in sustained damping and rebound effects, especially in scenarios with minimal external disturbances where energy dissipation occurs rapidly. These outcomes underscore the significance of both components of our model in capturing the dynamics of physical objects.

899
900
901
902
903
904

We also conduct an ablation study on the internal filling strategy shown in Figure 10. Notice that without internal filling, the hollow part of the cake is prone to collapse in the simulation process, but with internal filling, the object's deformation under external force is more realistic and smooth. For static and learnable internal filling, we can find that learnable internal filling can significantly reduce noise and artifacts introduced by initialization errors of internal Gaussian particles compared to static internal filling in Figure 11.

905
906
907
908
909
910
911
912
913
914
915
916
917



Figure 10: Ablation study on internal filling.

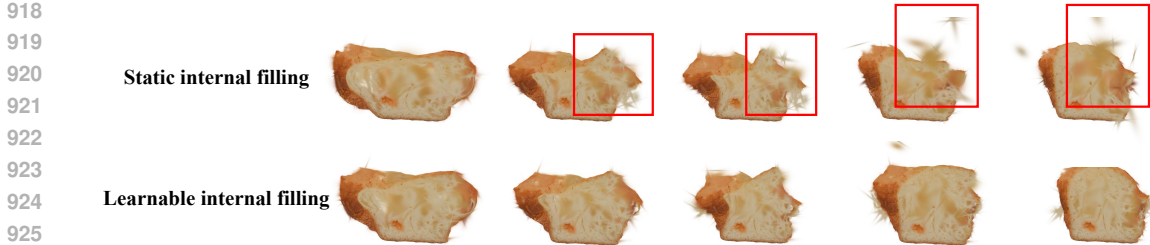


Figure 11: Comparison on static and learning internal filling. **Top row:** static filling strategy encounters challenges in dynamic situations, such as *tearing a bread*, where blur internal regions are exposed to the outside. **Bottom row:** learnable internal filling with iterative optimization strategy of SDS ensures thorough optimization and stability.

B.3 USER STUDY

For user study, we show each volunteer with five samples of rendered video from a random method (PhysDreamer Zhang et al. (2024), PhysGaussian Xie et al. (2023), DreamGaussian4D Ren et al. (2023), and ours) and real-captured videos. The study engaged 30 volunteers to assess the generated results in 20 rounds. In each round, they were asked to select the video they preferred the most, based on quality, realism, alignment with input 3D object, and fluency. We find our method is significantly preferred by users over these aspects.

Table 2: **Quantitative comparison results** of PhysGaussian Xie et al. (2023), PhysDreamer Zhang et al. (2024) and our Physics3D on the action coherence, motion realism, and overall quality score in a user study, rated on a range of 1-10, with higher scores indicating better performance.

Method	Action Coherence	Motion Realism	Overall Quality
PhysGaussian Xie et al. (2023)	7.82	6.89	6.93
PhysDreamer Zhang et al. (2024)	8.76	7.73	7.89
Physics3D (Ours)	8.95	8.57	9.05

C MORE ANALYSIS

C.1 PROBLEM FORMULATION AND NOTATION

In modeling viscoelastic stresses with physical fidelity, λ and μ are respectively denoted as Lamé’s first and second parameters. Their significance varies in different contexts. For example, in fluid dynamics, μ is related to the dynamic viscosity of fluids while in elastic environments, Lamé parameters λ and μ intertwine with Young’s modulus (E) and Poisson’s ratio (ν) via a specific relationship. On the other hand, the viscosity coefficient corrects the elastic strain-stress relation by dynamically controlling the relation to account for viscosity. With a closer look at the viscosity coefficient ν_v and ν_d , it serves as a viscous damper to exert resistance against rapid deformation, thus enhancing the model’s fidelity in capturing viscoelasticity. We also summarize the main physical parameters needed in the simulation process in Table 3.

C.2 MATERIAL POINT METHOD (MPM) ALGORITHM

The Material Point Method (MPM) simulates the behavior of materials by discretizing a continuum body into particles and updating their properties over time. Here’s an additional summary of the MPM algorithm:

Table 3: **Material Parameters.**

Notation	Meaning	Value
m	Mass of the Gaussian particle.	/
E	Young’s Modulus: Controls the dynamics of elastic objects; higher E results in smaller deformation under a fixed external force.	Learnable parameter
ν	Poisson’s Ratio: Determines the relationship between lateral strain and longitudinal strain in elastic materials.	/
λ_E	Lamé’s First Parameter: Relevant in continuum mechanics for stress-strain relationships in elastoplastic part.	$\lambda_E = \frac{E\nu}{(1+\nu)(1-2\nu)}$
μ_E	Lamé’s Second Parameter (Shear Modulus): Related to elastic materials in elastoplastic part.	$\mu_E = \frac{E}{2(1+\nu)}$
ν_d	Dynamic Viscosity Coefficient: Serves as a viscous damper exerting resistance against rapid deformation.	Learnable parameter
ν_v	Kinematic Viscosity Coefficient: Governs the dynamic aspect of viscosity in the model for the viscoelastic part.	Learnable parameter
μ_N	Lamé’s First Parameter: Relevant in continuum mechanics for stress-strain relationships in viscoelastic part. In fluid dynamics, related to dynamic viscosity.	Learnable parameter
λ_N	Lamé’s Second Parameter (Shear Modulus): Related to elastic materials in viscoelastic part.	Learnable parameter

Particle to Grid Transfer. Mass and momentum are transferred from particles to grid nodes. This step involves distributing the particle properties (mass and velocity) to nearby grid points.

$$m_i^n = \sum_p w_{ip}^n m_p,$$

$$m_i^n \mathbf{v}_i^n = \sum_p w_{ip}^n m_p (\mathbf{v}_p^n + \mathbf{C}_p^n (\mathbf{x}_i - \mathbf{x}_p^n)).$$

Grid Update. Grid velocities are updated based on external forces and the forces from neighboring particles. This step moves the grid points according to the applied forces.

$$\mathbf{v}_i^{n+1} = \mathbf{v}_i^n - \frac{\Delta t}{m_i} \sum_p \tau_p^n \nabla w_{ip}^n V_p^0 + \Delta t \mathbf{g}.$$

Grid to Particle Transfer. Velocities are transferred back to particles, and particle states are updated. This step brings the changes in grid velocities back to the particles.

$$\mathbf{v}_p^{n+1} = \sum_i \mathbf{v}_i^{n+1} w_{ip}^n,$$

$$\mathbf{x}_p^{n+1} = \mathbf{x}_p^n + \Delta t \mathbf{v}_p^{n+1},$$

$$\mathbf{C}_p^{n+1} = \frac{12}{\Delta x^2 (b+1)} \sum_i w_{ip}^n \mathbf{v}_i^{n+1} (\mathbf{x}_i^n - \mathbf{x}_p^n)^T,$$

$$\nabla \mathbf{v}_p^{n+1} = \sum_i \mathbf{v}_i^{n+1} \nabla w_{ip}^n{}^T,$$

$$\tau_p^{n+1} = \tau(\mathbf{F}_E^{n+1}, \mathbf{F}_N^{n+1}),$$

Here b is the B-spline degree, and Δx is the Eulerian grid spacing. We extensively demonstrate how to update \mathbf{F}_E , \mathbf{F}_N and τ in Appendix C.3.

1026 C.3 DYNAMIC 3D GENERATION.

1027
1028 Dynamic 3D generation aims to synthesize the dynamic behavior of 3D objects or scenes in the
1029 process of generating three-dimensional representations. Unlike static 3D generation methods Liu
1030 et al. (2023); Hong et al. (2023); Wang et al. (2024); Liu et al. (2024) that solely focus on the spatial
1031 morphology of objects, the field of 3D dynamics imposes a higher requirement, necessitating the
1032 incorporation of information from all three spatial dimensions as well as the temporal dimension.
1033 As advancements in 3D representation techniques continue to emerge, certain parameterizable 3D
1034 representations have empowered us to inject dynamic information into 3D models Li et al. (2022);
1035 Kratimenos et al. (2023). One popular approach Newcombe et al. (2015); Ren et al. (2023) integrates
1036 three essential components:

1037 **Gaussian Splatting** Kerbl et al. (2023) represents 3D information with a set of 3D Gaussian
1038 kernels. Each Gaussian can be described with a center $x \in \mathbb{R}^3$, a scaling factor $s \in \mathbb{R}^3$, and a
1039 rotation quaternion $q \in \mathbb{R}^4$. Additionally, an opacity value $\alpha \in \mathbb{R}$ and a color feature $c \in \mathbb{R}^3$ are for
1040 volumetric rendering. These parameters can be collectively denoted by θ , with $\theta_i = \{x_i, s_i, q_i, \alpha_i, c_i\}$
1041 representing the parameters for the i -th Gaussian. The volume rendering color C of each pixel is
1042 computed by blending N ordered points overlapping the pixel:

$$1043 C = \sum_{i \in N} c_i \alpha_i \prod_{j=1}^{i-1} (1 - \alpha_j). \quad (16)$$

1046 **Diffusion Model** Ho et al. (2020); Nichol & Dhariwal (2021) is usually pre-trained on large 2D
1047 datasets to provide a foundational motion prior for dynamic 3D generation. These models, charac-
1048 terized by their probabilistic nature, are tailored to acquire knowledge of the data distribution $p(x_0)$
1049 through a step-by-step denoising process applied to a normally distributed variable. Throughout
1050 the training phase, the data distribution is perturbed towards an isotropic Gaussian distribution over
1051 T timesteps, guided by a predefined noising schedule $\alpha_t \in (0, 1)$, where $\bar{\alpha}_t = \sum_{s=1}^t \alpha_s$ and t
1052 uniformly sampled from $1, \dots, T$:

$$1053 \mathbf{z}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \text{ where } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (17)$$

1054 The backward denoising process estimates the origin distribution by autoencoder $\epsilon_{\theta}(\mathbf{z}_t, t)$. The final
1055 training loss can be simplified as:

$$1056 \mathcal{L}_{DM} = \mathbb{E}_{\mathbf{x}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\|\epsilon - \epsilon_{\theta}(\mathbf{z}_t, t)\|_2^2 \right]. \quad (18)$$

1059 **Score Distillation Sampling** (SDS) is a technique used in machine learning, particularly in the
1060 context of generative models, to improve sample quality and diversity. It addresses the challenge of
1061 generating high-quality samples from complex probability distributions, such as those learned by
1062 deep neural networks.

1063 In traditional generative models like Generative Adversarial Networks (GANs) Goodfellow et al.
1064 (2020) or Variational Autoencoders (VAEs) Kingma & Welling (2013), generating samples involves
1065 directly sampling from the learned latent space. However, this approach often results in low-quality
1066 samples with poor diversity, especially in regions of low probability density.

1067 SDS introduces a novel sampling strategy that leverages the notion of score matching. Score matching,
1068 introduced by Hyvärinen in 2005 Hyvärinen & Dayan (2005), is a technique for training generative
1069 models by matching the score function (gradient of the log-density) of the model distribution to that
1070 of the true data distribution.

1071 Given a target distribution $p_{\text{data}}(\mathbf{x})$ and a model distribution $p_{\phi}(\mathbf{x})$, where ϕ represents the parameters
1072 of the model, the score function is defined as:

$$1073 \nabla_{\mathbf{x}} \log p_{\phi}(\mathbf{x}) \quad (19)$$

1075 The score function provides valuable information about the local geometry of the probability distribu-
1076 tion, helping to guide the sampling process towards regions of high probability density.

1077
1078 In SDS, the goal is to distill the score function learned by a complex generative model into a simpler,
1079 more tractable form. This distilled score function can then be used to guide the sampling process
efficiently, resulting in higher-quality samples with improved diversity.

Formally, given a set of generated samples $\{\mathbf{x}_i\}_{i=1}^N$ from the model distribution $p_\phi(\mathbf{x})$, the distilled score function $\hat{s}(\mathbf{x})$ is learned to approximate the score function of the model distribution. This is achieved by minimizing the score matching loss:

$$\mathcal{L}_{\text{SM}}(\hat{s}) = \frac{1}{N} \sum_{i=1}^N \|\nabla_{\mathbf{x}} \log p_\theta(\mathbf{x}_i) - \hat{s}(\mathbf{x}_i)\|^2 \quad (20)$$

where $\|\cdot\|$ denotes some norm (e.g., L^2 norm) and $\hat{s}(\mathbf{x}_i)$ is the estimated score at sample \mathbf{x}_i .

As for particular diffusion model, its score function can be related to the predicted noise (shown in Eq. 18) for the smoothed density through Tweedie’s formula Robbins (1992):

$$\epsilon_\phi(\mathbf{z}_t, t) = -\sigma_t s_\phi(\mathbf{z}_t; t). \quad (21)$$

Training the diffusion model with a (weighted) evidence lower bound (ELBO) simplifies to a weighted denoising score matching objective for parameters ϕ (Ho et al., 2020; Kingma et al., 2021):

$$\mathcal{L}_{\text{Diff}}(\phi, \mathbf{x}) = \mathbb{E}_{t \sim \mathcal{U}(0,1), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [w(t) \|\epsilon_\phi(\alpha_t \mathbf{x} + \sigma_t \epsilon; t) - \epsilon\|_2^2], \quad (22)$$

where $w(t)$ is a weighting function that depends on the timestep t . To understand the difficulties of this approach, consider the gradient of $\mathcal{L}_{\text{Diff}}$:

$$\nabla_\theta \mathcal{L}_{\text{Diff}}(\phi, \mathbf{x} = g(\theta)) = \mathbb{E}_{t, \epsilon} \left[w(t) \underbrace{(\hat{\epsilon}_\phi(\mathbf{z}_t; y, t) - \epsilon)}_{\text{Noise Residual}} \underbrace{\frac{\partial \hat{\epsilon}_\phi(\mathbf{z}_t; y, t)}{\partial \mathbf{z}_t}}_{\text{U-Net Jacobian}} \underbrace{\frac{\partial \mathbf{x}}{\partial \theta}}_{\text{Generator Jacobian}} \right] \quad (23)$$

where following DreamFusion Poole et al. (2022), we absorb the constant $\alpha_t \mathbf{I} = \partial \mathbf{z}_t / \partial \mathbf{x}$ into $w(t)$. In practice, the U-Net Jacobian term is expensive to compute (requires backpropagating through the diffusion model U-Net), and poorly conditioned for small noise levels as it is trained to approximate the scaled Hessian of the marginal density. In Poole et al. (2022), It is found that omitting the U-Net Jacobian term leads to an effective gradient:

$$\nabla_\theta \mathcal{L}_{\text{SDS}}(\phi, \mathbf{x} = g(\theta)) \triangleq \mathbb{E}_{t, \epsilon} \left[w(t) (\hat{\epsilon}_\phi(\mathbf{z}_t; y, t) - \epsilon) \frac{\partial \mathbf{x}}{\partial \theta} \right], \quad (24)$$

here, we get the gradient of a weighted probability density distillation loss in Eq. 14.

C.4 ADDITIONAL ANALYSIS IN CONTINUUM MECHANICS

Basic intuitions in continuum mechanics involves explaining the strain tensor \mathbf{F} which describes the internal deformation of the material, and $\boldsymbol{\sigma}$ which describe the internal stress tensor. In an actual material, there might be multiple compounds, each of them might have their own strain and stress tensor and they might interact with each other. For purpose of this subsection, we explain the intuition when there is only one such compounds. However we will generalize it into multiple compounds case later in the paper.

One could describe an equilibrium material as a point cloud \mathbf{X} . In actually simulation, one take a discretization of the system. For our convenience, we will just state the intuition in continuum, and generalization to discrete points should be straightforward. When the material is away from its equilibrium position, the position of point is given by $\mathbf{x} \neq \mathbf{X}$. However, we notice that if, for example, we move \mathbf{x} and the adjacent point $\tilde{\mathbf{x}}$ in the same way, then there is actually no internal deformation of the material. In another word, internal deformation describe how relative position between two points \mathbf{x} and its adjacent $\tilde{\mathbf{x}}$ change. To make this quantity well-defined, we normalize it with the original difference between \mathbf{X} and $\tilde{\mathbf{X}}$

$$\tilde{\mathbf{F}}_{ij} = \frac{\tilde{\mathbf{x}}_i - \mathbf{x}_i}{\tilde{\mathbf{X}}_j - \mathbf{X}_j} \quad (25)$$

Taking continuum limit, the equation becomes differential and reproduce

$$\mathbf{F} = \frac{\partial \phi}{\partial \mathbf{X}}(\mathbf{X}, t) \quad (26)$$

1134 in the main text.

1135
1136 To get the dynamics, we will further need the stress tensor σ . However, the general relation between
1137 σ and \mathbf{F} is usually complicated. In fact only when internal stress is a conservative force, we can
1138 related it to \mathbf{F} through an energy function $\psi(\mathbf{F})$. An analog for this case is the Hooke's force in a
1139 spring, where we can introduce a potential energy for the force. There are two main types of internal
1140 stress people use in the literature. The first Piola-Kirchoff stress

$$1141 \quad \mathbf{P} = \frac{\partial \psi}{\partial \mathbf{F}} \quad (27)$$

1142 and a related Cauchy stress

$$1143 \quad \sigma = \frac{1}{\det(\mathbf{F})} \frac{\partial \psi}{\partial \mathbf{F}} \mathbf{F}^T \quad (28)$$

1144 We will use Cauchy stress for most part of our paper. However, we could sometimes use the other
1145 one interchangeably. Physically, if we want the force acting on a unit surface with normal vector \hat{n} ,
1146 the force will be $\sigma \cdot \hat{n}$. Another version people use is the Kirchoff stress $\tau = \det(\mathbf{F})\sigma$. The only
1147 difference between them are whether it is measured with the deformed volume(σ) or the undeformed
1148 volume(τ). Sometimes using one or the other could be easier for technical reason, as we can see later
1149 in the viscoelastic model.

1150 Beyond this simple case, the relation will be quite complicated and non-universal. In this paper, we
1151 discuss a specific variant of them, which we explain below. One can see both the relation between σ
1152 and \mathbf{F} , and the update rule for \mathbf{F} get modified.

1153 With these background, we can explain the intuition behind Eq.2. The second equation

$$1154 \quad \frac{D\rho}{Dt} + \rho \nabla \cdot \mathbf{v} = 0, \quad (29)$$

1155 is nothing but mass conservation, $\rho \nabla \cdot \mathbf{v}$ is the local density times the local divergence, which
1156 physically correspond to the mass flows out of a local unit volume in unit time.

1157 The first equation is Newton's law

$$1158 \quad \rho \frac{D\mathbf{v}}{Dt} = \nabla \cdot \sigma + \mathbf{f}, \quad (30)$$

1159 the \mathbf{f} term is the external force. The term $\nabla \cdot \sigma$ is the divergence for the stress tensor, thus giving the
1160 total force on a unit volume material.

1161 **Model for elastoplastic part.** We explain how to compute the internal stress σ_E through \mathbf{F}_E , and
1162 how to update \mathbf{F}_E for the elastic part in this appendix. Again for current paper we consider the
1163 elastoplastic branch with only elastic part. As we have reviewed, one can compute the Cauchy stress
1164 tensor given the energy function. In this work, we choose the Fixed corotated constitutive model for
1165 the elastic part, whose energy function is

$$1166 \quad \psi(\mathbf{F}_E) = \psi(\Sigma_E) = \mu \sum_i (\sigma_{E,i} - 1)^2 + \frac{\lambda}{2} (\det(\mathbf{F}_E) - 1)^2 \quad (31)$$

1167 in which Σ_E is the diagonal singular value matrix $\mathbf{F}_E = \mathbf{U}\Sigma_E\mathbf{V}^T$. And $\sigma_{E,i}$ s are singular values
1168 of \mathbf{F}_E . From this energy function, one can compute the Cauchy stress tensor

$$1169 \quad \sigma_E = \frac{2\mu}{\det(\mathbf{F}_E)} (\mathbf{F}_E - \mathbf{R}) \mathbf{F}_E^T + \lambda (\det(\mathbf{F}_E) - 1) \quad (32)$$

1170 in which $\mathbf{R} = \mathbf{U}\mathbf{V}^T$.

1171 We now explain how to update the \mathbf{F}_E^n with the velocity field \mathbf{v}^n at n th step. For purely elastic case,
1172 this can be simply done via

$$1173 \quad \mathbf{F}_E^{n+1} = (\mathbf{I} + \Delta t \nabla \mathbf{v}^n) \mathbf{F}_E^n \quad (33)$$

1174 in which Δt is the length of the time segment in the MPM method. This formula represents the fact
1175 that internal velocity field cause change in strains. However, as one will say for the viscoelastic part,
1176 this rule will change.

1177 For the plastic part of the branch, we have $\mathbf{F} = \mathbf{F}_E \mathbf{F}_P$, where we constraint the singular value of \mathbf{F}_E
1178 to sit between $[1 - \theta_c, 1 + \theta_s]$, where θ_c and θ_s are learnable parameters quantifying the plasticity.

More precisely, in the simulation algorithm, at n -th step, $\mathbf{F}^n = \mathbf{F}_E^n \mathbf{F}_P^n$. We then compute the internal stress $\boldsymbol{\sigma}(\mathbf{F}_E^n)$ and update \mathbf{F}^n to \mathbf{F}^{n+1} with MPM method. At step $n + 1$, we first have a trial

$$\tilde{\mathbf{F}}_E^{n+1} = \mathbf{F}^{n+1} (\mathbf{F}_P^n)^{-1}. \quad (34)$$

This $\tilde{\mathbf{F}}_E^{n+1}$ might have singular value lies out of $[1 - \theta_c, 1 + \theta_s]$, so we truncate the singular value again to get \mathbf{F}_E^{n+1} , and get

$$\mathbf{F}_P^{n+1} = \mathbf{F}^{n+1} (\mathbf{F}_E^{n+1})^{-1}. \quad (35)$$

This procedure updates all strain tensor $\mathbf{F}_E, \mathbf{F}_P$ from step n to step $n + 1$. Noticing a key difference here between \mathbf{F}_E and \mathbf{F}_P is only \mathbf{F}_E contributes to the stress tensor. This is a manifestation of plasticity, where a large enough deformation will not contribute forces that move the object back to its original position.

Model for viscoelastic part. We now explain the other branch of our model: the viscoelastic part. Now we should imagine a elastic strain series connect with a viscous dissipator. In another word, the total strain $\mathbf{F} = \mathbf{F}_N \mathbf{F}_V$. We now explain two key features for this brunch. First, only the \mathbf{F}_N contributes to the internal stress. Secondly, \mathbf{F}_V only plays a role in the update rule of \mathbf{F}_N . Very roughly, one can view the dissipator \mathbf{F}_V as a reservoir for strain and each time during the update, people will be able to relocate part of the strain into the \mathbf{F}_V , such that only part of the strain contributes to the internal stress.

For the \mathbf{F}_N part, we again take it to be a elastic system, so the relation between strain and stress is again given by the energy function. Here we choose the energy function following Fang et al. (2019)

$$\psi_N(\Sigma_N) = \mu_N \text{tr}((\log \Sigma_N)^2) + \frac{1}{2} \lambda_N (\text{tr}(\log \Sigma_N))^2 \quad (36)$$

In this formula, we introduce Σ_N as the diagonal singular value matrix $\mathbf{F}_N = U \Sigma_N V^T$. As one can see, this potential is only a function of the singular values, so the Kirchoff tensor $\boldsymbol{\tau}_N$ will be diagonal in the same basis as the strain \mathbf{F}_N . So we will just write down the relation between \mathbf{F}_N and $\boldsymbol{\tau}_N$ just in the singular value basis

$$\boldsymbol{\tau}_N = \frac{\partial \psi_N(\Sigma_N)}{\partial \epsilon_N} = 2\mu_N \epsilon_N + \lambda_N \text{tr}(\epsilon_N) \mathbf{I} \quad (37)$$

where we denote τ_N as a vector of singular value of $\boldsymbol{\tau}_N$, in another word $\boldsymbol{\tau}_N = U \text{diag}(\tau_N) V^T$. This is sometimes called principle Kirchoff tensor in literature. And ϵ_N , called log principle Kirchoff tensor, denotes a vector takes the diagonal element of $\log \Sigma_N$. So this relation is essentially a vector equation. We can recover the Kirchoff tensor $\boldsymbol{\tau}_N$ thus the Cauchy stress tensor

$$\boldsymbol{\sigma}_N = \frac{1}{\det(\mathbf{F}_N)} \boldsymbol{\tau}_N. \quad (38)$$

Now we can discuss the update rule for the strain tensor \mathbf{F}_N . As we explained before, the rough intuition was part of the strain tensor \mathbf{F}_N can dissipate into the dissipator \mathbf{F}_V . More quantitatively, we can follow the dissipator model in Fang et al. (2019), which results into a trial-and-correction procedure. The idea is we first imagine \mathbf{F}_N evolve elastically

$$\mathbf{F}_{N,\text{tr}}^n = (\mathbf{I} + \Delta t \nabla \mathbf{v}^n) \mathbf{F}_N^n \quad (39)$$

Next step, we modify this trial strain tensor by introducing a dissipation into the dissipator

$$\epsilon_N^{n+1} = \epsilon_{N,\text{tr}}^n - \Delta t \frac{\partial \psi_V}{\partial \tau_N} \quad (40)$$

Here again, we assume the dissipation happen to only the singular value of the strain tensor, so $\epsilon_{N,\text{tr}}^n$ is the log principle Kirchoff tensor of $\mathbf{F}_{N,\text{tr}}^n$ and ψ_V is the dissipation potential. We take a model where

$$\psi_V(\tau_N) = \frac{1}{2\nu_d} |\text{dev}(\tau_N)|^2 + \frac{1}{9\nu_v} (\tau_N \cdot \mathbf{1})^2 \quad (41)$$

where ν_d and ν_v are parameters controlling the dissipation of the deviatoric and dilational parts. Physically, parameter ν_d controls dissipation in deviatoric deformation (deformation that does not change volume, but only the shape), and ν_e controls dissipation in dilational deformation (deformation

that only changes the volume but not the shape). These two parameters will be the same when materials have a certain kind of homogenous property. We will not use them directly instead, we can put this equation back into Eq.40, and get a version of this formula purely in terms of ϵ_N^{n+1} and $\epsilon_{N,tr}^n$

$$\epsilon_N^{n+1} = A(\epsilon_{N,tr}^n - B\text{tr}(\epsilon_{N,tr}^n) \cdot \mathbf{1}) \quad (42)$$

In principle, we can write the parameters A, B in terms of ν_d, ν_v , however since in our algorithm, all these parameters will be learnt from videos, we will directly learn parameter A and B in the update rule without bothering ν_d, ν_v . But it worth remembering our model comes from a dissipation potential.

C.5 LEARNABLE INTERNAL FILLING

In Xie et al. (2023), it employs a static filling strategy where the filled particles inherit the parameters σ_p and C_p from their nearest Gaussian kernels. These parameters correspond to the Gaussian properties of the filled particles. However, this static filling strategy encounters challenges in dynamic situations, such as tearing or compression shown in Figure 10, where large internal regions are exposed to the outside. The issue arises because the internal particles are not well-initialized and remain fixed during optimization.

To address these challenges, we propose an iterative optimization strategy that refines both physical parameters and internal Gaussian parameters. This approach ensures that the distribution of these two components does not interfere with each other, thereby preventing mode collapse. In other words, one set of parameters can be optimized independently, avoiding "self-sacrifice" for the optimization of the other set. Please refer to the visual comparison of the original and learnable strategy in Figure 11. Specifically, for each epoch in the iterative optimization, we first fix all physical parameters and focus on optimizing the internal filling Gaussian parameters (e.g., Gaussian center $x \in \mathbb{R}^3$, an opacity value $\alpha \in \mathbb{R}$, a color feature $c \in \mathbb{R}^3$ and so on) by Eq. 15. Subsequently, we fix the Gaussian parameters and optimize the physical properties (e.g., Young's modulus E , Poisson's ratio ν , and so on) of each Gaussian by Eq. 14. This iterative process is designed to ensure thorough optimization and stability. Moreover, to enhance the generalizability, we choose a random view for each epoch. This ensures that our model is robust and performs well across different perspectives.

D ETHICAL STATEMENT

We confirm that all data used in this paper for research and publication have been obtained and used in a manner compliant with ethical standards. The individuals engaged in all experiments have given consent for their use, or the data are sourced from publicly available datasets and were used in accordance with the terms of use and permissions. Furthermore, the publication and use of these data and models do not pose any societal or ethical harm. We have taken necessary precautions to ensure that the research presented in this paper respects individual rights, including the right to privacy and the fundamental principles of ethical research conduct.