

# Improving Physics-Augmented Continuum Neural Radiance Field-Based Geometry-Agnostic System Identification with Lagrangian Particle Optimization

Takuhiko Kaneko  
NTT Corporation

## Abstract

Geometry-agnostic system identification is a technique for identifying the geometry and physical properties of an object from video sequences without any geometric assumptions. Recently, physics-augmented continuum neural radiance fields (PAC-NeRF) has demonstrated promising results for this technique by utilizing a hybrid Eulerian-Lagrangian representation, in which the geometry is represented by the Eulerian grid representations of NeRF, the physics is described by a material point method (MPM), and they are connected via Lagrangian particles. However, a notable limitation of PAC-NeRF is that its performance is sensitive to the learning of the geometry from the first frames owing to its two-step optimization. First, the grid representations are optimized with the first frames of video sequences, and then the physical properties are optimized through video sequences utilizing the fixed first-frame grid representations. This limitation can be critical when learning of the geometric structure is difficult, for example, in a few-shot (sparse view) setting. To overcome this limitation, we propose Lagrangian particle optimization (LPO), in which the positions and features of particles are optimized through video sequences in Lagrangian space. This method allows for the optimization of the geometric structure across the entire video sequence within the physical constraints imposed by the MPM. The experimental results demonstrate that the LPO is useful for geometric correction and physical identification in sparse-view settings.<sup>1</sup>

## 1. Introduction

System identification is a technique used to identify the geometry and physical properties of an object based on visual observations. It has been actively studied in computer vision and physics because of its wide range of applications, including 3D reproduction, environmental understanding, and content creation. To solve this problem in a realistic environment, it is important to adequately address both the *extrinsic geometric structure* and the *intrinsic physical prop-*

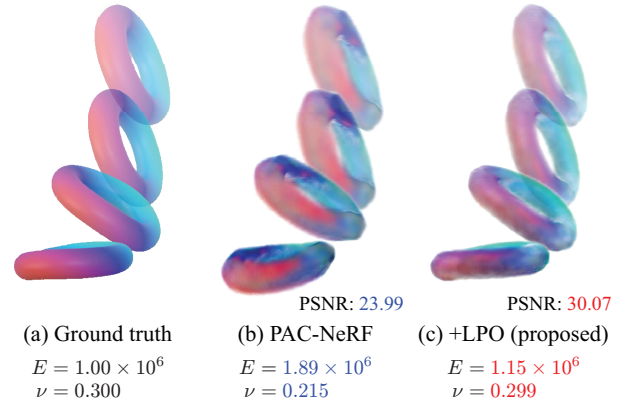


Figure 1. Impact of the proposed Lagrangian particle optimization (LPO) in sparse-view geometric-agnostic system identification. We aim to identify the geometry and physical properties of an object from visual observations *without any geometric assumptions* in severe (e.g., sparse-view) settings. As shown in (b), a standard PAC-NeRF [42] has difficulty learning the geometry in a sparse-view setting (particularly, when the number of views is three). Consequently, it also fails to estimate the physical properties (Young's modulus  $E$  and Poisson's ratio  $\nu$ ). As shown in (c), LPO is useful for correcting this failure and succeeds in improving the identification of both geometry and physical properties.

*erties*. For example, to identify the geometry and physical properties of an object from the observation that it collides with the ground, it is necessary to understand not only the appearance and shape change in the *geometry* but also other internal factors in *physics* which are necessary to explain this phenomenon without any discrepancies over time.

Generally, this problem is ill-posed and challenging to resolve. However, the emergence of powerful geometric representations by neural fields [83], 2D-3D connections by neural rendering [75], and trainable physical simulations using differentiable physical simulators [1] provide clues for solving this problem. For instance, previous studies [34, 35, 48] succeeded in estimating physical properties such as velocity, mass, friction, inertia, and elasticity directly from video sequences by combining differentiable physical simulators with differentiable renderers.

Although these studies have yielded promising results in terms of physical property estimation, their applicability is restricted by the assumption that the geometric struc-

<sup>1</sup>The project page is available at <https://www.kecl.ntt.co.jp/people/kaneko.takuhiko/projects/lpo/>.

ture of a scene is *completely known*, which makes it difficult to apply them in practical scenarios. To eliminate this assumption, *geometry-agnostic system identification*, which is a technique for identifying the geometry and physical properties of an object from visual observations *without any geometric assumptions*, has been studied. Specifically, a pioneering model involves *physics-augmented continuum neural radiance fields (PAC-NeRF)* [42], which is an extension of NeRF [51] that is enforced to follow the conservation laws of continuum mechanics. PAC-NeRF obtains this functionality utilizing a hybrid Eulerian–Lagrangian representation, that is, the geometry (volume density and color fields) is represented with the Eulerian grid representations of NeRF [73], which are transformed into Lagrangian particles in material space, and physical simulation is conducted on the particles using a material point method (MPM) (particularly, differentiable MPM (DiffMPM) [32]). Because all modules are differentiable, PAC-NeRF can be trained directly with multi-view video sequences and can optimize the geometry and physics without explicit supervision.

Although PAC-NeRF has enabled the tackling of new tasks, its notable limitation is that its performance is sensitive to the learning of the geometry from the first frames of the video sequences owing to its two-step optimization. First, the Eulerian grid representations were trained with the first frames of the video sequences, and then the physical properties were optimized across the video sequences by utilizing the frozen first-frame grid representations. This limitation makes it difficult to apply it to cases in which geometry learning is difficult, for example, in a few-shot (sparse-view) setting, as shown in Figure 1(b).

One possible solution is to train the Eulerian grid representations using all video sequences and not just the first frame. However, a critical problem with this approach is that Eulerian grid representations cannot be trained with explicit physical constraints on the particles because they are optimized in Eulerian (i.e., world or grid) space and cannot reflect all physical phenomena occurring in Lagrangian (i.e., material or particle) space; for example, they cannot propagate gradients calculated for particle positions.

To solve this problem, we propose *Lagrangian particle optimization (LPO)*, in which the positions and features (i.e., volume density and color) of the particles are optimized *in Lagrangian space and not in Eulerian space*. Contrary to the abovementioned possible solution, this method allows for the optimization of the geometric structure across video sequences with *explicit* physical constraints on the particles imposed by the MPM. Thus, the gradients calculated for the particle positions are reflected.

Furthermore, LPO is useful not only for geometry correction, but also for corrections to the physical identification because this task is closely related, that is, an accurate geometry is useful for accurately identifying the physical properties, and vice versa. This motivated us to utilize the geometry corrected by LPO to reidentify the physical prop-

erties and propose *iterative geometry–physics optimization* for gradually seeking the optimal states. Figure 1(c) shows the effectiveness of this method.

We evaluated the effectiveness of LPO in sparse-view settings. In particular, we first applied LPO to a pretrained PAC-NeRF and demonstrated that LPO is useful for correcting geometric structures through video sequences. Subsequently, we employed LPO in iterative geometry–physics optimization and demonstrated that LPO contributes to improving the performance of physical identification. Ablation and comparative studies were conducted to determine the importance of each component.

Our contributions can be summarized as follows:

- To improve the performance of *geometry-agnostic system identification*, we propose *LPO*, which optimizes the position and features of the particles *in Lagrangian space* to optimize the geometric structures across video sequences within the physical constraints of an MPM.
- We propose *iterative geometry–physics optimization* to utilize the geometry corrected by LPO for reidentifying physical properties in an iterative manner.
- We experimentally demonstrated that LPO is useful for geometric correction and physical identification in sparse-view settings. We also provide detailed analyses and extended results in the Appendices. Video samples are available at the [project page](#).<sup>1</sup>

## 2. Related work

**Neural radiance fields (NeRF).** Novel view synthesis is a fundamental problem in computer vision and graphics and has been addressed in a large body of research. Recently, a substantial breakthrough has been achieved with the emergence of neural fields that represent a scene utilizing a coordinate-based network (e.g., a survey paper [83]). NeRF [51] is a representative variant of such neural fields and has attracted significant attention because of its geometric consistency and high-fidelity novel view synthesis. Various studies have been conducted since the emergence of NeRF. These can be categorized into three topics. (1) Improvements to the image quality and expanding of the applicable settings, such as wild or few-shot settings (e.g., [2–4, 8, 11, 12, 14, 30, 33, 49, 52, 57, 66, 77, 80, 86, 88, 90, 92]), (2) speeding up the training or inference speed (e.g., [4, 6, 9, 18, 20, 26, 29, 30, 38, 41, 46, 47, 53, 54, 59, 64, 65, 73, 81, 82, 87]), and (3) incorporating other models or functionalities, such as generative models [22, 28, 72] (e.g., [5–7, 15, 23, 37, 55, 56, 60, 68, 71, 85]) and dynamics/physics (e.g., [13, 19, 24, 25, 42, 44, 58, 61, 69, 76, 91]). Among these, this study relates to the first in terms of its application to few-shot settings. This study relates to the third in terms of seeking physics-informed models. As studies on NeRF are benefiting from developments in adjacent categories, applying our ideas to other categories and models will be an interesting future research direction.

**NeRF with dynamics/physics.** As mentioned above,

NeRFs with dynamics/physics have been actively studied (e.g., [13, 19, 24, 25, 42, 44, 58, 61, 69, 76, 91]). These studies have one characteristic in common: they learned time-varying neural fields from video sequences. These can be roughly categorized into two models. (1) Non- (or weak) physics-informed models (e.g., [19, 25, 44, 58, 61, 69, 76, 91]) and (2) physics-informed models (e.g., [13, 24, 42]). The advantage of the first is that it can be applied to scenes or objects that are difficult to describe physically; however, its disadvantage is that it requires a large amount of training data to learn dynamics from scratch, and the learned representations are not necessarily interpretable because they are not based on physics. The advantage of the second is that it can obtain interpretable representations based on physics; however, the disadvantage is that its application is limited to physically describable objects or scenes. In particular, physics informed NeRF [13] targets smoke scenes and does not address the boundary conditions; therefore, it cannot handle solid or contact materials. NeuroFluid [24] focuses on fluid dynamics grounding and solves it using an intuitive physics-based approach in which formal instruction in physics is not explicitly defined. In contrast, PAC-NeRF [42] is based on a principled and interpretable physical simulator and can be applied to various materials, including Newtonian/non-Newtonian fluids, elastic materials, plasticine, and granular media. This study focused on the PAC-NeRF and attempted to widen its applicability while prioritizing its interpretability. However, a Lagrangian particle representation is commonly utilized in physics (e.g., NeuroFluid [24]<sup>2</sup>); therefore, applying our ideas, that is, LPO, to other models is an interesting research topic.

**NeRF for few-shot (sparse-view) settings.** In practice, it is often laborious or impractical to collect multi-view images. Recent studies [8, 12–14, 33, 57, 66, 69, 80, 86, 88, 90] addressed this problem by refining NeRF such that it could be applied to few-shot (sparse-view) settings. These methods can be roughly divided into three approaches. (1) Regularization using external models. For example, normalizing flow [16]-based [57], VGG [70]-based [90], depth-based [14, 66, 80], and CLIP [63]-based [33] regularizations have been proposed. (2) Utilization of general and transferable models [8, 12, 69, 88]. They are trained using external datasets, and in several studies, they are finetuned for each scene. (3) Introduction of advanced training methods, such as gradual training to prevent overfitting to sparse views, including frequency regularization [86] and layer-by-layer growing strategies [13]. Among these categories, the proposed LPO is categorized as the third one, that is, optimization is conducted for each scene, and external models and datasets are not required. The main difference between the LPO and previous methods is that LPO attempts to opti-

<sup>2</sup>Note that NeuroFluid optimizes the *transition probability of particles* under the assumption that the initial particle positions are known, that is, fixed. In contrast, LPO optimizes the *initial particle positions* to correct the geometry estimation of the first frames. Therefore, NeuroFluid and LPO are complementary.

mize geometric structures through video sequences with the *physical constraints* of an MPM. However, previous studies have not sufficiently considered these constraints. Owing to this difference in applications, the proposed method is complementary to, rather than competitive with, other methods.

**System identification of soft bodies.** Soft bodies are not only high-dimensional, but also allow large deformations; therefore, it is challenging to identify their geometry and physical properties. Typical methods can be categorized into three types. (1) Gradient-free methods [74, 78], (2) neural network-based methods [43, 67, 84], and (3) differentiable physics-simulation-based methods [10, 17, 21, 27, 35, 42, 48, 62]. The strength of the methods in the first and second categories is that they are flexible; however, they have difficulty achieving a high accuracy owing to their black-box modeling. Methods in the third category require sophisticated modeling to fill the gap between simulations and the real world; however, they have recently demonstrated promising results owing to advancements in differentiable physical simulators. The third category can be further divided into two subcategories. (i) Methods that assume that watertight geometric mesh representations of objects are available [17, 21, 27, 35, 48, 62] and (2) methods that do not assume this [10, 42]. Among these methods, we focused on the final category, PAC-NeRF [42], which prioritizes general and fast characteristics. Specifically, PAC-NeRF adopts an MPM that can be applied to various materials, including fluids [36], sand [40], foam [89], and elastic objects [10]. In particular, PAC-NeRF utilizes DiffMPM [32] to construct a fully differentiable simulation rendering system. This study is based on this advancement. In particular, sensitivity to geometry learning from the first frames is one of the bottlenecks of PAC-NeRF. Therefore, this study focused on alleviating this bottleneck and attempted to widen its applicability to severe (e.g., sparse-view) settings.

### 3. Method

In this section, we first define the problem statement (Section 3.1) and then briefly review PAC-NeRF [42], upon which our method is built (Section 3.2). Subsequently, we explain the main idea of the proposed method, that is, geometric correction using LPO (Section 3.3), and then introduce its application to physical identification (Section 3.4).

#### 3.1. Problem statement

Given a set of multi-view video sequences, our objectives are as follows: (1) to recover geometric representations<sup>3</sup> of the target dynamic object (particularly continuum materials) through video sequences and (2) identify its physical properties. More formally, the training data comprise ground truth color observations, that is,  $\hat{C}(\mathbf{r}, t)$ . Here,

<sup>3</sup>Here, both appearance and shape are treated as part of the geometric representations. We attempted to recover both.

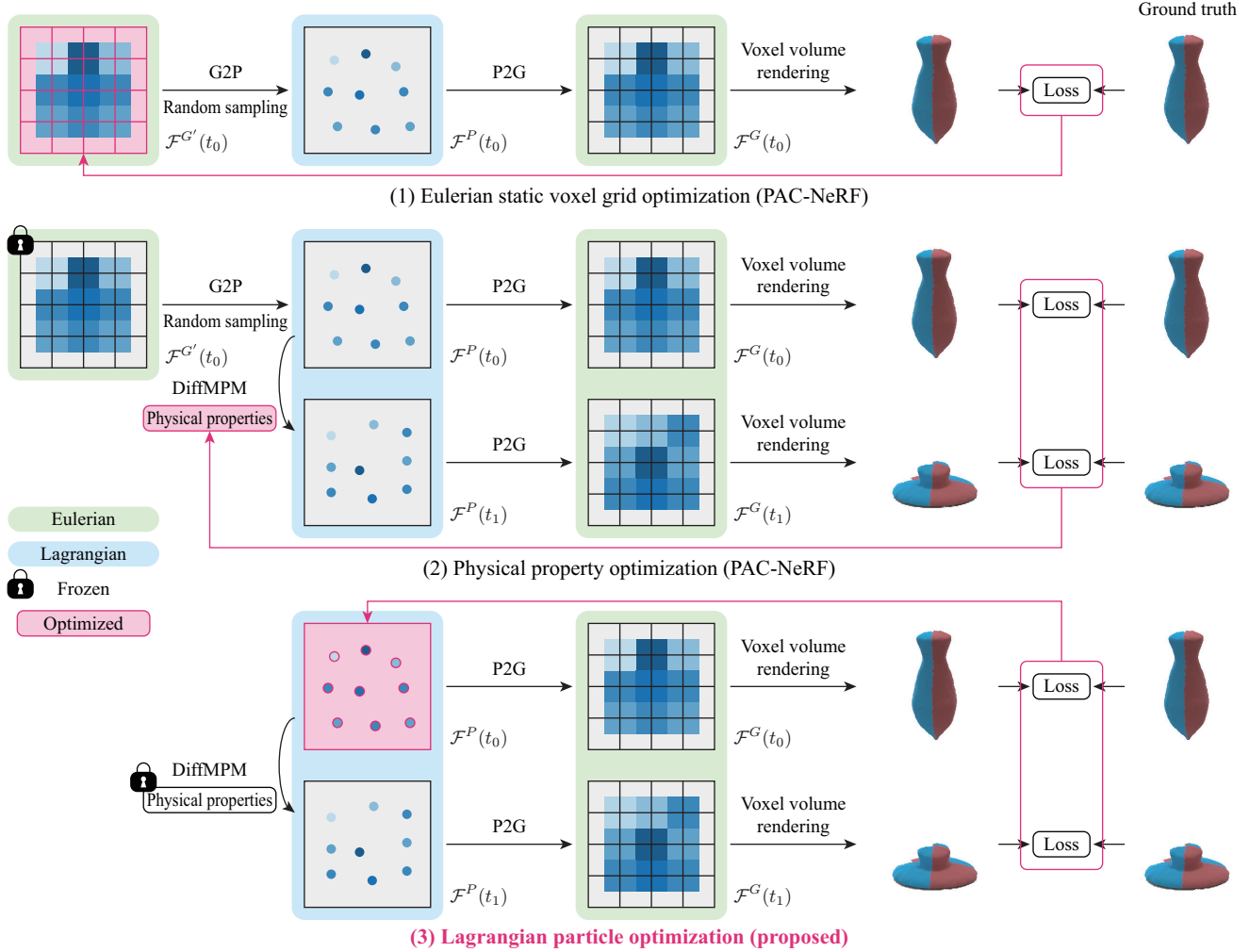


Figure 2. Optimization pipelines of PAC-NeRF (1)(2) and the proposed LPO (3).

$\mathbf{r}(s) \in \mathbb{R}^3$  is a camera ray defined as  $\mathbf{r}(s) = \mathbf{o} + s\mathbf{d}$ , where  $\mathbf{o} \in \mathbb{R}^3$  is the camera origin,  $\mathbf{d} \in \mathbb{S}^2$  is the view direction, and  $s \in [s_n, s_f]$  is a distance from the camera origin. During training,  $\mathbf{r}$  is sampled from a set of ray collections in the training data, that is,  $\hat{\mathcal{R}}$ .  $t \in \mathbb{R}^+$  is a time variable, and during training, it is sampled from  $\{t_0, \dots, t_{N-1}\}$ , where  $N$  is the number of frames. Given these training data, we attempt to (1) predict the color observation  $\mathbf{C}(\mathbf{r}, t)$  that can recover  $\hat{\mathbf{C}}(\mathbf{r}, t)$  from the training data and those for novel views, and (2) identify the physical properties. In particular, to widen its applicability, we addressed a scenario in which input views are sparse, that is,  $\hat{\mathcal{R}}$  is limited.

### 3.2. Preliminary: PAC-NeRF

To solve the abovementioned problem, PAC-NeRF employs three core components: a continuum NeRF, particle-grid interconverters, and a Lagrangian field.

**Continuum NeRF.** PAC-NeRF extends a standard (static) NeRF to a continuum NeRF to address the dynamics of

continuum materials. To achieve this, a standard NeRF is first extended to a dynamic NeRF [61] for dynamic scenarios. In a dynamic NeRF, the volume density and color fields for point  $\mathbf{x}$  and view direction  $\mathbf{d}$  are defined as  $\sigma(\mathbf{x}, t)$  and  $\mathbf{c}(\mathbf{x}, \mathbf{d}, t)$ , respectively, where the time variable  $t$  is introduced to represent the dynamics. The color of each pixel  $\mathbf{C}(\mathbf{r}, t)$  is calculated using volume rendering [50]

$$\mathbf{C}(\mathbf{r}, t) = \int_{s_n}^{s_f} T(s, t) \sigma(\mathbf{r}(s), t) \mathbf{c}(\mathbf{r}(s), \mathbf{d}, t) ds + \mathbf{c}_{bg} T(s_f, t), \quad (1)$$

$$T(s, t) = \exp\left(-\int_{s_n}^s \sigma(\mathbf{r}(u), t) du\right), \quad (2)$$

where  $\mathbf{c}_{bg}$  denotes the background color. The model is trained with pixel-wise loss.

$$\mathcal{L}_{pixel} = \frac{1}{N} \sum_{i=0}^{N-1} \frac{1}{|\hat{\mathcal{R}}|} \sum_{\mathbf{r} \in \hat{\mathcal{R}}} \|\mathbf{C}(\mathbf{r}, t_i) - \hat{\mathbf{C}}(\mathbf{r}, t_i)\|_2^2. \quad (3)$$

Furthermore, a dynamic NeRF is extended to a continuum NeRF to represent continuum materials. To achieve this, conservation laws for the velocity field are imposed on the volume density and color fields.

$$\frac{D\sigma}{Dt} = 0, \quad \frac{D\mathbf{c}}{Dt} = \mathbf{0}, \quad (4)$$

whereas the material derivative of an arbitrary time-dependent field  $\phi(\mathbf{x}, t)$  is defined as  $\frac{D\phi}{Dt} = \frac{\partial\phi}{\partial t} + \mathbf{v} \cdot \nabla\phi$ . Here,  $\mathbf{v}$  is the velocity field and follows momentum conservation for continuum materials:

$$\rho \frac{D\mathbf{v}}{Dt} = \nabla \cdot \mathbf{T} + \rho\mathbf{g}, \quad (5)$$

where  $\rho$  is the physical density field,  $\mathbf{T}$  is the internal Cauchy stress tensor, and  $\mathbf{g}$  is the acceleration of gravity. This equation is solved using DiffMPM [32].

**Particle-grid interconverters.** As shown in Figure 2(1)(2), in PAC-NeRF, a physical simulation is conducted in Lagrangian particle space using DiffMPM [32], whereas neural rendering is performed in Eulerian grid space with a discretized voxel-based NeRF [73]. To bridge these two spaces, grid-to-particle (G2P) and particle-to-grid (P2G) conversions are conducted.

$$\mathcal{F}_p^P \approx \sum_i w_{ip} \mathcal{F}_i^G, \quad \mathcal{F}_i^G \approx \frac{\sum_p w_{ip} \mathcal{F}_p^P}{\sum_p w_{ip}}, \quad (6)$$

where  $\mathcal{F}_*^*$  is a field (e.g., volume density or color field);  $P$  and  $G$  denote particle and grid views, respectively;  $p$  and  $i$  indicate the index of the particle and grid node, respectively; and  $w_{ip}$  is the weight of the trilinear shape function defined on  $i$  and evaluated at  $p$ .

**Lagrangian field.** As shown in Figure 2(1), during training, the Eulerian voxel field  $\mathcal{F}^{G'}(t_0)$  is optimized for the first frames of the video sequences. Subsequently, this field is converted into a Lagrangian particle field  $\mathcal{F}^P(t_0)$  using G2P, in which the positions of the particles are determined by random sampling around voxel grids. The geometric shape is represented in Lagrangian space by removing the particles whose alpha values are small, that is, the remaining particles represent the shape of the object. As shown in Figure 2(2), the particle field in the next step, that is,  $\mathcal{F}^P(t_1)$ , where  $t_1 = t_0 + \delta t$  and  $\delta t$  is the duration of the simulation time step, is calculated based on  $\mathcal{F}^P(t_0)$  using DiffMPM [32]. After this process,  $\mathcal{F}^P(t_1)$  is converted to the Eulerian voxel field  $\mathcal{F}^G(t_1)$  using P2G, and pixels are rendered based on this field using voxel-based volume rendering [73], in which the volume density and color of a point at position  $\mathbf{x}$  are rendered as follows:

$$\sigma(\mathbf{x}, t) = \text{softplus}(\text{interp}(\mathbf{x}, \sigma^G)), \quad (7)$$

$$\mathbf{c}(\mathbf{x}, \mathbf{d}, t) = \text{MLP}(\text{interp}(\mathbf{x}, \mathbf{c}^G), \mathbf{d}), \quad (8)$$

where  $\sigma^G$  and  $\mathbf{c}^G$  denote the volume density and color fields, respectively, discretized on fixed voxel grids, and  $\text{interp}$  denotes trilinear interpolation.

### 3.3. Geometric correction with LPO

As explained above, in PAC-NeRF, the Eulerian voxel fields  $\mathcal{F}^{G'}(t_0)$  are optimized for the first frames of the video sequences (Figure 2(1)). These are fixed while optimizing the physical properties using video sequences (Figure 2(2)). This approach is problematic when the learning of the Eulerian voxel field from the first frames is difficult, for example, in sparse-view settings, because failure of this learning cannot be corrected after the process. This also means that in dynamic scenes, information over time is useful for covering a lack of information owing to sparse views; however, this approach cannot utilize this information to optimize the volume density and color.

One possible solution is to train  $\mathcal{F}^{G'}(t_0)$  with the entire video sequence and not just the first frame. However, in this approach, the optimization targets are limited to the volume density and color of the grids; therefore, it is difficult to reflect the physics-based optimization that occurs in Lagrangian particle space.<sup>4</sup> This can cause excessive modification of the geometry beyond physical constraints.

Alternatively, we developed *Lagrangian particle optimization (LPO)*, in which the geometric structure is optimized *not in Eulerian grid space but in Lagrangian particle space* (Figure 2(3)). More formally, in the Eulerian voxel grid optimization, volume density and color fields of the grids are optimized for the *fixed* voxel grid. In contrast, in LPO, not only the volume density and color fields of the particles, that is,  $\sigma^P$  and  $\mathbf{c}^P$ , but also the particle position, that is  $\mathbf{x}^P$ , are *optimized*.<sup>5</sup> Because  $\mathbf{x}^P$  is defined in Lagrangian space, it can reflect the physical constraints of the MPM. Consequently, we can optimize the geometric structures within the physical allowance.

### 3.4. Physical identification with LPO

Identifying both the geometry and physical properties from limited observations is ill-posed and challenging because there is a chicken-and-egg relationship between the geometric structure and the physical properties, that is, accurate geometry estimation is necessary for accurate physical identification and vice versa. Considering that particle-based geometric correction is highly challenging when there is a large gap between the ground truth and the predicted pixels owing to the high-dimensional nature (e.g., in extreme cases, particles can diverge), in practice, we apply LPO after physical property optimization, as shown in Figure 2.

However, an interesting question is *whether the corrected geometric structures can be utilized to improve the*

<sup>4</sup>More precisely, as explained in Section 3.2, particles are randomly sampled around equally spaced voxel grids, and the Eulerian voxel fields (precisely, alpha values calculated from them) are utilized to mask them with a non-differentiable way. Consequently, gradients calculated for the particle positions cannot propagate to  $\mathcal{F}^{G'}(t_0)$ .

<sup>5</sup>Note that the values of  $\sigma^P$  and  $\mathbf{c}^P$  are changed in training but they are time-invariant, i.e., have the same values across sequences. Therefore, the conservation laws (Equation (4)) are preserved.

---

**Algorithm 1** Iterative geometry–physics optimization

---

**Input:** Ground-truth color observations  $\hat{\mathbf{C}}(\hat{\mathbf{r}}, t)$   
 $\hat{\mathbf{r}} \in \hat{\mathcal{R}}$ : Rays from observed viewpoints  
 $t \in \{t_0, \dots, t_{N-1}\}$ : Time  
**Output:** Physical properties

- 1: **for**  $i = 1$  to  $R$  **do**
- 2:   // (1) Eulerian static voxel grid optimization
- 3:   **if**  $i = 1$  **then**
- 4:     Optimize voxel grids with  $\hat{\mathbf{C}}(\hat{\mathbf{r}}, t_0)$
- 5:   **else**
- 6:     Optimize voxel grids with  $\hat{\mathbf{C}}(\hat{\mathbf{r}}, t_0)$  and  $\check{\mathbf{C}}(\check{\mathbf{r}}, t_0)$
- 7:   **end if**
- 8:   // (2) Physical property optimization
- 9:   Optimize physical properties with  $\hat{\mathbf{C}}(\hat{\mathbf{r}}, t)$
- 10:   // (3) Lagrangian particle optimization
- 11:   Optimize particles with  $\hat{\mathbf{C}}(\hat{\mathbf{r}}, t)$
- 12:   // (4) Color prediction for unobserved views  $\check{\mathbf{r}} \in \check{\mathcal{R}}$
- 13:   Predict  $\check{\mathbf{C}}(\check{\mathbf{r}}, t_0)$
- 14: **end for**

---

*identification of the physical properties.* We developed *iterative geometry–physics optimization* (Algorithm 1) to answer this question. As described above, we first conduct PAC-NeRF optimization ((1) and (2)) and then apply LPO (3). Subsequently, we render the images in the first frame based on the updated  $\sigma^P$ ,  $\mathbf{c}^P$ , and  $\mathbf{x}^P$  (4). We denote them as  $\check{\mathbf{C}}(\check{\mathbf{r}}, t_0)$ , where  $\check{\mathbf{r}} \in \check{\mathcal{R}}$  and  $\check{\mathcal{R}}$  is designed such that it covers the missing area in  $\hat{\mathcal{R}}$  in the training data. We retrain PAC-NeRF from scratch using the combination of the original  $\hat{\mathbf{C}}(\hat{\mathbf{r}}, t_0)$  and synthesized  $\check{\mathbf{C}}(\check{\mathbf{r}}, t_0)$  as training data. Specifically, to alleviate the negative effect caused by incomplete geometry correction, which can occur at the beginning of the iterative calculation, we only utilize  $\check{\mathbf{C}}(\check{\mathbf{r}}, t_0)$  for the Eulerian static voxel grid optimization and do not utilize it to optimize the physical properties and particles.

Another possible solution is to update the voxel grids by applying P2G to the updated  $\sigma^P$ ,  $\mathbf{c}^P$ , and  $\mathbf{x}^P$  instead of (4) and (1). However, we found that the repeated utilization of G2P and P2G tends to produce artifacts (such as those typically occurring in the erosion process) owing to their approximate nature (Equation (6)).<sup>6</sup> Therefore, the full recalculation approach described above was adopted.

## 4. Experiments

### 4.1. Experimental setup

Two experiments were conducted to investigate the effectiveness of the proposed LPO. First, we evaluated the effec-

---

<sup>6</sup>A similar phenomenon was observed in the original study on PAC-NeRF [42]. Based on this observation, they rendered an image in the first frame by utilizing the voxel grids obtained after the P2G and G2P conversions, that is,  $\mathcal{F}^G(t_0)$ , instead of the original voxel grids  $\mathcal{F}^{G'}(t_0)$  to make the rendering conditions the same as those in other frames.

tiveness of LPO for geometric correction (Section 4.2) and then we investigated the utility of LPO for physical identification (Section 4.3). The main results of these experiments are presented here, and the detailed analyses, extended results, and implementation details are provided in the Appendices. Video samples are available at the [project page](#).<sup>1</sup> In this section, we provide the experimental setup common to both experiments and other setups in subsequent sections.

**Dataset.** We investigated the benchmark performance on the dataset provided by the original study on PAC-NeRF [42].<sup>7</sup> This dataset comprised nine scenes and various continuum materials, including Newtonian fluids (*Droplet* and *Letter* with fluid viscosity  $\mu$  and bulk modulus  $\kappa$ ), non-Newtonian fluids (*Cream* and *Toothpaste* with shear modulus  $\mu$ , bulk modulus  $\kappa$ , yield stress  $\tau_Y$ , and plastic viscosity  $\eta$ ), elastic materials (*Torus* and *Bird* with Young’s modulus  $E$  and Poisson’s ratio  $\nu$ ), plasticine (*Playdoh* and *Cat* with  $E$ ,  $\nu$ , and  $\tau_Y$ ), and granular media (*Trophy* with friction angle  $\theta_{fric}$ ). In each scene, the objects fall freely under the influence of gravity and undergo collisions. The ground-truth simulation data were generated using the MLS-MPM framework [31]. A photorealistic simulation engine rendered objects under diverse environmental lighting conditions and ground textures. Each scene was captured from 11 viewpoints with cameras evenly spaced on the upper hemisphere, including the object. To evaluate our method under sparse-view settings, three views were used for training and the remaining eight views were used for testing. To show the robustness to the view settings, we provide the results for the other view settings in Appendix C.

**Assumptions and preprocessing.** For a fair comparison, we follow the assumptions and preprocessing used in PAC-NeRF [42]. It was assumed that the intrinsic and extrinsic parameters of the set of static cameras were previously known. Moreover, it was assumed that collision objects, such as the ground plane, were previously known. As mentioned in [42], this is not difficult to estimate from the observed images. For preprocessing, a video matting framework [45] was applied to remove static background objects and focus the rendering on the target object.

### 4.2. Evaluation of geometric correction

**Compared models.** In the first experiment, we investigated the effectiveness of LPO on the geometric correction (i.e., the method described in Section 3.3). Three models were used as the baseline. (I) *PAC-NeRF* was the same as that in the original [42] and was trained with all views, including the eight views that were used for testing. We examined this model to determine the upper bound of its performance. (II) *PAC-NeRF-3v* was trained using three views. The training settings were the same as those for (I) except that the

---

<sup>7</sup>Note that although only one dataset was used, this dataset is useful for assessing the versatility because it includes a wide variety of materials, surpassing the scope of previous studies (e.g., elastic objects only in [10]).

|                    | PAC-NeRF | PAC-NeRF-3v | +LPO  | +LPO-F | +LPO-P | +GO   | PAC-NeRF-3v <sup>†</sup> | +LPO  | +LPO-F | +LPO-P | +GO   |
|--------------------|----------|-------------|-------|--------|--------|-------|--------------------------|-------|--------|--------|-------|
| PSNR <sup>↑</sup>  | 35.99    | 27.39       | 29.22 | 28.22  | 28.89  | 28.33 | 28.47                    | 30.11 | 29.31  | 29.87  | 29.39 |
| SSIM <sup>↑</sup>  | 0.989    | 0.978       | 0.980 | 0.979  | 0.980  | 0.979 | 0.980                    | 0.982 | 0.981  | 0.981  | 0.981 |
| LPIPS <sup>↓</sup> | 0.020    | 0.034       | 0.032 | 0.033  | 0.032  | 0.033 | 0.033                    | 0.031 | 0.032  | 0.031  | 0.031 |

Table 1. Comparison of PSNR<sup>↑</sup>, SSIM<sup>↑</sup>, and LPIPS<sup>↓</sup> on **geometric correction**. The scores for each scene are provided in Table 5 in Appendix B. The qualitative comparisons are provided in Figures 3–7 in Appendix D.

number of views varied. (III) *PAC-NeRF-3v<sup>†</sup>* was an improved variant of (II) for sparse-view settings. An interesting question is whether LPO, which is a few-shot learning method for *dynamic* scenes, can be used with other few-shot learning methods, such as those for *static* scenes. To answer this question, we adjusted the Eulerian static voxel grid optimization (Figure 2(1)) of PAC-NeRF-3v, such that it became robust to sparse views. Specifically, we scheduled a surface regularizer to reduce unexpected artifacts, introduced a view-invariant pixel-wise loss to compensate for the lack of views, and adjusted the training length to prevent overfitting. Details are provided in Appendix A.<sup>8</sup> We applied LPO to both (II) and (III) to investigate the effect of the initial Eulerian static voxel grid optimization. Hereafter, we denote these variants by *+LPO*. Furthermore, three comparative and ablation models were evaluated to determine the importance of each component. (i) *+LPO-F* and (ii) *+LPO-P* are variants of LPO, in which only the features and positions of the particles are optimized. (iii) *+GO* optimized Eulerian voxel grids through the entire video sequence instead of using Lagrangian particles. (i)–(iii) are used as alternatives to *+LPO*. These variants were applied to both (II) and (III).

**Evaluation metrics.** We evaluated the performance of the geometric correction using metrics commonly used to assess the performance of novel view synthesis in NeRF studies: the peak signal-to-noise ratio (*PSNR*), structural similarity index measure (*SSIM*) [79], and learned perceptual image patch similarity (*LPIPS*) [93]. In particular, we calculated the scores using all test data over time.

**Results.** The results are summarized in Table 1. Our findings are threefold. (1) *PAC-NeRF-3v/3v<sup>†</sup>* vs. *+LPO*. *+LPO* improved both baselines in terms of all metrics. In particular, the utility of *+LPO* for PAC-NeRF-3v<sup>†</sup> indicates an improvement in the Eulerian static voxel grid optimization and that of LPO are complementary. (2) *+LPO* vs. *+LPO-F/P*. We found that *+LPO* was superior or comparable to *+LPO-F/P*. This is because the feature and position optimizations are complementary. Intuitively, feature optimization is useful for correcting features of the *appearance* that are invisible in the first frame but visible after the object

has moved. Similarly, position optimization is effective for correcting features of the *shape* that are invisible in the first frame but visible after the object has moved. In severe (e.g., sparse-view) settings, the utilization of both is important. (3) *+LPO* vs. *+GO*. *+LPO* outperforms *+GO*, possibly because *+LPO* can optimize the geometry adequately within the physical constraints of the MPM, whereas *+GO* cannot because of the absence of an explicit physical constraint.

### 4.3. Evaluation of physical identification

**Compared models.** In the second experiment, we verified the usefulness of LPO for physical identification (i.e., the method described in Section 3.4). In addition to the models described in Section 4.2, we examined *+None*, for which iterative optimization (Algorithm 1) was conducted without geometric correction (i.e., Step (3) was skipped). This model was used to investigate the importance of geometric correction. For a fair comparison, we set the number of iterations (i.e., *R* in Algorithm 1) to four for all models. We use the superscript 4 (e.g., *+LPO<sup>4</sup>*) to specify this.

**Evaluation metrics.** To evaluate the performance of the physical identification, we measured the absolute distance between the ground truth and the estimated physical properties. For an easy comparison, we calculated these distances after adjusting the scale (i.e., either a logarithmic scale or a linear scale) following the study of PAC-NeRF [42]. The smaller the values, the better was the performance.

**Results.** The results are summarized in Table 2. Our findings are fourfold. (1) *PAC-NeRF-3v/3v<sup>†</sup>* vs. *+LPO<sup>4</sup>*. *+LPO<sup>4</sup>* improved the physical identification of PAC-NeRF-3v/3v<sup>†</sup> in most cases. In particular, the effectiveness of *+LPO<sup>4</sup>* for PAC-NeRF-3v<sup>†</sup> indicates that the improvement in the Eulerian static voxel grid optimization and that of LPO are complementary for physical identification. (2) *+LPO<sup>4</sup>* vs. *+LPO-F<sup>4</sup>/P<sup>4</sup>*. In some cases, superiority depends on the physical properties because the physical properties interact with each other, and finding the optimal balance is difficult. However, we found that *+LPO-F<sup>4</sup>/P<sup>4</sup>* sometimes had obvious difficulties (e.g., PAC-NeRF-3v+LPO-F<sup>4</sup> on Bird and PAC-NeRF-3v<sup>†</sup>+LPO-P<sup>4</sup> on Playdoh), whereas *+LPO<sup>4</sup>* exhibited good stability. We consider that the joint optimization of the features and positions is useful for tackling difficult situations. (3) *+LPO<sup>4</sup>* vs. *+GO<sup>4</sup>*. *+GO<sup>4</sup>* experienced difficulties in some cases (e.g., on Bird). *+LPO<sup>4</sup>* exhibited a more stable performance. (4) *+LPO<sup>4</sup>* vs. *+None<sup>4</sup>*. *+None<sup>4</sup>* sometimes worsened the performance. The results indicated that simple iterations were insufficient and that geometric correction was essential.

<sup>8</sup>In the preliminary experiments, we examined the previous representative few-shot learning methods (e.g., DietNeRF [33] and FreeNeRF [86]). However, we found that they were less stable than PAC-NeRF-3v, possibly because in our experimental settings, the number of views was small (three) despite the wide range of the views (upper hemisphere), and explicit voxel representations were more useful than the fully implicit representation in [33, 86]. However, this study and previous studies are complementary. Further investigations will be important in future work.

|            |                       | PAC-NeRF | PAC-NeRF-3v | +LPO <sup>4</sup> | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup> | +None <sup>4</sup> | PAC-NeRF-3v <sup>†</sup> | +LPO <sup>4</sup> | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup> | +None <sup>4</sup> |
|------------|-----------------------|----------|-------------|-------------------|---------------------|---------------------|------------------|--------------------|--------------------------|-------------------|---------------------|---------------------|------------------|--------------------|
| Droplet    | $\log_{10}(\mu)$      | 0.039    | 0.140       | 0.112             | 0.112               | 0.119               | 0.111            | 0.131              | 0.136                    | 0.082             | 0.068               | 0.067               | 0.082            | 0.090              |
|            | $\log_{10}(\kappa)$   | 0.017    | 1.285       | 1.628             | 1.638               | 1.466               | 1.675            | 1.684              | 1.263                    | 0.106             | 1.139               | 0.520               | 1.520            | 1.656              |
| Letter     | $\log_{10}(\mu)$      | 0.041    | 0.674       | 0.015             | 0.026               | 0.079               | 0.048            | 0.530              | 0.379                    | 0.010             | 0.087               | 0.007               | 0.118            | 0.436              |
|            | $\log_{10}(\kappa)$   | 0.039    | 6.772       | 0.174             | 0.411               | 1.040               | 0.865            | 6.083              | 5.229                    | 0.060             | 0.027               | 0.127               | 0.138            | 5.060              |
| Cream      | $\log_{10}(\mu)$      | 0.090    | 0.311       | 0.178             | 0.163               | 0.234               | 0.183            | 0.251              | 0.179                    | 0.100             | 0.073               | 0.065               | 0.092            | 0.076              |
|            | $\log_{10}(\kappa)$   | 0.132    | 0.215       | 0.158             | 0.249               | 0.027               | 0.028            | 0.244              | 0.336                    | 0.121             | 0.197               | 0.142               | 0.193            | 0.339              |
|            | $\log_{10}(\tau_Y)$   | 0.007    | 0.014       | 0.004             | 0.030               | 0.005               | 0.006            | 0.025              | 0.009                    | 0.006             | 0.001               | 0.004               | 0.010            | 0.008              |
|            | $\log_{10}(\eta)$     | 0.015    | 0.281       | 0.183             | 0.256               | 0.083               | 0.061            | 0.252              | 0.195                    | 0.033             | 0.059               | 0.079               | 0.096            | 0.051              |
| Toothpaste | $\log_{10}(\mu)$      | 0.026    | 1.891       | 0.109             | 0.156               | 0.119               | 0.164            | 0.215              | 0.252                    | 0.031             | 0.201               | 0.005               | 0.252            | 0.138              |
|            | $\log_{10}(\kappa)$   | 0.247    | 1.580       | 0.601             | 1.356               | 0.597               | 0.648            | 0.729              | 1.436                    | 0.673             | 0.630               | 0.899               | 0.590            | 0.596              |
|            | $\log_{10}(\tau_Y)$   | 0.066    | 0.201       | 0.114             | 0.191               | 0.075               | 0.250            | 0.061              | 0.199                    | 0.093             | 0.064               | 0.117               | 0.149            | 0.054              |
|            | $\log_{10}(\eta)$     | 0.013    | 0.373       | 0.003             | 0.267               | 0.013               | 0.007            | 0.047              | 0.212                    | 0.009             | 0.016               | 0.005               | 0.007            | 0.042              |
| Torus      | $\log_{10}(E)$        | 0.019    | 0.277       | 0.061             | 0.008               | 0.083               | 0.010            | 0.054              | 0.074                    | 0.036             | 0.026               | 0.039               | 0.015            | 0.074              |
|            | $\nu$                 | 0.023    | 0.085       | 0.001             | 0.120               | 0.031               | 0.074            | 0.316              | 0.131                    | 0.007             | 0.040               | 0.033               | 0.050            | 0.129              |
| Bird       | $\log_{10}(E)$        | 0.013    | 0.449       | 0.067             | 0.240               | 0.074               | 0.313            | 0.246              | 0.123                    | 0.027             | 0.186               | 0.039               | 0.231            | 0.296              |
|            | $\nu$                 | 0.029    | 0.102       | 0.001             | 0.372               | 0.075               | 0.365            | 0.581              | 0.141                    | 0.047             | 0.072               | 0.008               | 0.132            | 0.145              |
| Playdoh    | $\log_{10}(E)$        | 0.286    | 0.290       | 0.116             | 0.580               | 0.120               | 0.304            | 0.170              | 0.521                    | 0.133             | 0.474               | 1.211               | 0.232            | 2.148              |
|            | $\log_{10}(\tau_Y)$   | 0.038    | 0.283       | 0.165             | 0.027               | 0.214               | 0.196            | 0.237              | 0.110                    | 0.173             | 0.079               | 1.179               | 0.197            | 0.967              |
|            | $\nu$                 | 0.076    | 0.495       | 0.111             | 0.173               | 0.109               | 0.248            | 0.128              | 0.212                    | 0.063             | 0.133               | 0.110               | 0.127            | 0.106              |
| Cat        | $\log_{10}(E)$        | 0.855    | 1.301       | 0.973             | 1.068               | 1.071               | 1.127            | 1.066              | 1.192                    | 0.706             | 0.534               | 0.712               | 0.783            | 0.793              |
|            | $\log_{10}(\tau_Y)$   | 0.026    | 0.120       | 0.107             | 0.117               | 0.113               | 0.086            | 0.112              | 0.084                    | 0.067             | 0.014               | 0.076               | 0.069            | 0.072              |
|            | $\nu$                 | 0.027    | 0.044       | 0.004             | 0.036               | 0.069               | 0.190            | 0.079              | 0.118                    | 0.003             | 0.027               | 0.007               | 0.028            | 0.063              |
| Trophy     | $\theta_{fric}$ [rad] | 0.048    | 0.053       | 0.055             | 0.051               | 0.056               | 0.057            | 0.054              | 0.030                    | 0.039             | 0.041               | 0.043               | 0.044            | 0.039              |

Table 2. Comparison of the absolute differences between the ground-truth and the estimated physical properties on **physical identification**. The smaller the values, the better was the performance. The values of the physical properties (i.e., not the absolute differences) are provided in Tables 6 and 7 in Appendix B. The qualitative comparisons are presented in Figures 3–9 in Appendix D.

|                    | PAC-NeRF-3v | +LPO  | +LPO <sup>4</sup> | PAC-NeRF-3v <sup>†</sup> | +LPO  | +LPO <sup>4</sup> |
|--------------------|-------------|-------|-------------------|--------------------------|-------|-------------------|
| PSNR $\uparrow$    | 27.39       | 29.22 | 29.65             | 28.47                    | 30.11 | 30.34             |
| SSIM $\uparrow$    | 0.978       | 0.980 | 0.982             | 0.980                    | 0.982 | 0.983             |
| LPIPS $\downarrow$ | 0.034       | 0.032 | 0.030             | 0.033                    | 0.031 | 0.029             |

Table 3. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  on **geometric rectification**. The scores for each scene are provided in Table 8 in Appendix B. The qualitative comparisons are provided in Figures 3–9 in Appendix D.

**Applications.** As mentioned above, there is a chicken-and-egg relationship between the geometric structure and physical properties. An interesting question is *whether the corrected physical properties can be used to reestimate the geometry* (the opposite problem). To answer this question, we investigated the geometry identification performance after physical identification. The results are summarized in Table 3. These results indicate that accurate physical identification is useful for improving the geometry identification performance.

## 5. Discussion

Based on the above experiments, we demonstrated promising results for geometric correction and physical identification. However, our methods have some limitations. (1) Our methods require a longer training time due to the introduction of LPO (Section 3.3)<sup>9</sup> and the iterative algorithm (Sec-

<sup>9</sup>The calculation time of LPO (Figure 2(3)) is almost identical to that of the main process of physical property optimization (Figure 2(2)) because the forward and backward processes are identical with different opt-

tion 3.4).<sup>10</sup> However, it is not trivial to obtain robustness to sparse views only by increasing the training time because overfitting is a typical factor that causes learning to fail. (2) LPO is sensitive to the state before applying LPO (e.g., either PAC-NeRF-3v or -3v<sup>†</sup>) because solving geometry-agnostic system identification in sparse-view settings is ill-posed and challenging. We believe that the advancements in previous few-shot learning methods (Section 2) and the newly introduced few-shot learning method with physical constraints (LPO) will solve this problem.

## 6. Conclusion

We proposed LPO to improve PAC-NeRF-based geometric-agnostic system identification. Our core idea is to optimize the geometry *not in Eulerian space but in Lagrangian space*, utilizing the particles to directly reflect the physical constraints of an MPM. The results demonstrate that LPO is useful for both geometric correction and physical identification. Although we focused on PAC-NeRF while prioritizing its high generality, Lagrangian particles are commonly employed in physics-informed models (e.g., several studies discussed in Section 2). We expect that our ideas can be utilized in other models or tasks.

Similarly, the calculation times for +LPO-F, +LPO-P, and +GO were almost identical to those for +LPO, indicating that the improvement was attributable to the ingenuity of the algorithm and not to an increase in the calculation cost. The computation times are discussed in detail in Appendix B.2.

<sup>10</sup>The total computation time increases linearly when running Algorithm 1 repeatedly but is adjustable under a quality-and-time trade-off as discussed in Appendix B.3.



## References

- [1] Chayan Banerjee, Kien Nguyen, Clinton Fookes, and George Karniadakis. Physics-informed computer vision: A review and perspectives. *arXiv preprint arXiv:2305.18035*, 2023. 1
- [2] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. In *CVPR*, 2022.
- [4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-NeRF: Anti-aliased grid-based neural radiance fields. In *ICCV*, 2023. 2
- [5] Eric R. Chan, Marco Monteiro, Petr Kellnhöfer, Jiajun Wu, and Gordon Wetzstein. pi-GAN: Periodic implicit generative adversarial networks for 3D-aware image synthesis. In *CVPR*, 2021. 2
- [6] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022. 2
- [7] Eric R. Chan, Koki Nagano, Matthew A. Chan, Alexander W. Bergman, Jeong Joon Park, Axel Levy, Miika Aittala, Shalini De Mello, Tero Karras, and Gordon Wetzstein. Generative novel view synthesis with 3D-aware diffusion models. In *ICCV*, 2023. 2
- [8] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. MVSNeRF: Fast generalizable radiance field reconstruction from multi-view stereo. In *ICCV*, 2021. 2, 3
- [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. TensorRF: Tensorial radiance fields. In *ECCV*, 2022. 2
- [10] Hsiao-yu Chen, Edith Tretschk, Tuur Stuyck, Petr Kadlec, Ladislav Kavan, Etienne Vouga, and Christoph Lassner. Virtual elastic objects. In *CVPR*, 2022. 3, 6
- [11] Xingyu Chen, Qi Zhang, Xiaoyu Li, Yue Chen, Ying Feng, Xuan Wang, and Jue Wang. Hallucinated neural radiance fields in the wild. In *CVPR*, 2022. 2
- [12] Julian Chibane, Aayush Bansal, Verica Lazova, and Gerard Pons-Moll. Stereo radiance fields (SRF): Learning view synthesis for sparse views of novel scenes. In *CVPR*, 2021. 2, 3
- [13] 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 Trans. Graph.*, 41(4), 2022. 2, 3, 12
- [14] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *CVPR*, 2022. 2, 3
- [15] Yu Deng, Jiaolong Yang, Jianfeng Xiang, and Xin Tong. GRAM: Generative radiance manifolds for 3D-aware image generation. In *CVPR*, 2022. 2
- [16] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. In *ICLR*, 2017. 3
- [17] Tao Du, Kui Wu, Pingchuan Ma, Sebastien Wah, Andrew Spielberg, Daniela Rus, and Wojciech Matusik. DiffPD: Differentiable projective dynamics. *ACM Trans. Graph.*, 41(2), 2021. 3
- [18] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 2
- [19] Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4D facial avatar reconstruction. In *CVPR*, 2021. 2, 3
- [20] Stephan J. Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. FastNeRF: High-fidelity neural rendering at 200FPS. In *ICCV*, 2021. 2
- [21] 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 Trans. Graph.*, 39(6), 2020. 3
- [22] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, 2014. 2
- [23] Jiatao Gu, Lingjie Liu, Peng Wang, and Christian Theobalt. StyleNeRF: A style-based 3D-aware generator for high-resolution image synthesis. In *ICLR*, 2022. 2
- [24] Shanyan Guan, Huayu Deng, Yunbo Wang, and Xiaokang Yang. NeuroFluid: Fluid dynamics grounding with particle-driven neural radiance fields. In *ICML*, 2022. 2, 3
- [25] Xiang Guo, Jiadai Sun, Yuchao Dai, Guanying Chen, Xiaoqing Ye, Xiao Tan, Errui Ding, Yumeng Zhang, and Jingdong Wang. Forward flow for novel view synthesis of dynamic scenes. In *ICCV*, 2023. 2, 3
- [26] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021. 2
- [27] Eric Heiden, Miles Macklin, Yashraj Narang, Dieter Fox, Animesh Garg, and Fabio Ramos. DiSECT: A differentiable simulation engine for autonomous robotic cutting. In *RSS*, 2021. 3
- [28] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 2
- [29] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. EfficientNeRF: Efficient neural radiance fields. In *CVPR*, 2022. 2
- [30] Wenbo Hu, Yuling Wang, Lin Ma, Bangbang Yang, Lin Gao, Xiao Liu, and Yuewen Ma. Tri-MipRF: Tri-Mip representation for efficient anti-aliasing neural radiance fields. In *ICCV*, 2023. 2
- [31] Yuanming Hu, Yu Fang, Ziheng Ge, Ziyin Qu, Yixin Zhu, Andre Pradhana, and Chenfanfu Jiang. A moving least squares material point method with displacement discontinuity and two-way rigid body coupling. *ACM Trans. Graph.*, 37(4), 2018. 6, 31
- [32] Yuanming Hu, Luke Anderson, Tzu-Mao Li, Qi Sun, Nathan Carr, Jonathan Ragan-Kelley, and Frédo Durand. Diff-

- Taichi: Differentiable programming for physical simulation. In *ICLR*, 2020. 2, 3, 5, 31
- [33] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting NeRF on a diet: Semantically consistent few-shot view synthesis. In *ICCV*, 2021. 2, 3, 7, 12
- [34] Miguel Jaques, Michael Burke, and Timothy Hospedales. Physics-as-inverse-graphics: Unsupervised physical parameter estimation from video. In *ICLR*, 2020. 1
- [35] Krishna Murthy Jatavallabhula, Miles Macklin, Florian Golemo, Vikram Voleti, Linda Pettrini, Martin Weiss, Brendan Considine, Jérôme Parent-Lévesque, Kevin Xie, Kenny Erleben, Liam Paull, Florian Shkurti, Derek Nowrouzezahrai, and Sanja Fidler. gradSim: Differentiable simulation for system identification and visuomotor control. In *ICLR*, 2021. 1, 3
- [36] Chenfanfu Jiang, Craig Schroeder, Andrew Selle, Joseph Teran, and Alexey Stomakhin. The affine particle-in-cell method. *ACM Trans. Graph.*, 34(4), 2015. 3
- [37] Takuhiro Kaneko. AR-NeRF: Unsupervised learning of depth and defocus effects from natural images with aperture rendering neural radiance fields. In *CVPR*, 2022. 2
- [38] Takuhiro Kaneko. MIMO-NeRF: Fast neural rendering with multi-input multi-output neural radiance fields. In *ICCV*, 2023. 2
- [39] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 31
- [40] Gergely Klár, Theodore Gast, Andre Pradhana, Chuyuan Fu, Craig Schroeder, Chenfanfu Jiang, and Joseph Teran. Drucker-prager elastoplasticity for sand animation. *ACM Trans. Graph.*, 35(4), 2016. 3
- [41] Andreas Kurz, Thomas Neff, Zhaoyang Lv, Michael Zollhöfer, and Markus Steinberger. AdaNeRF: Adaptive sampling for real-time rendering of neural radiance fields. In *ECCV*, 2022. 2
- [42] 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. In *ICLR*, 2023. 1, 2, 3, 6, 7, 12, 20, 31
- [43] Yunzhu Li, Jiajun Wu, Russ Tedrake, Joshua B. Tenenbaum, and Antonio Torralba. Learning particle dynamics for manipulating rigid bodies, deformable objects, and fluids. In *ICLR*, 2019. 3
- [44] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 2, 3
- [45] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L. Curless, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *CVPR*, 2021. 6
- [46] David B. Lindell, Julien N. P. Martel, and Gordon Wetzstein. AutoInt: Automatic integration for fast neural volume rendering. In *CVPR*, 2021. 2
- [47] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *NeurIPS*, 2020. 2
- [48] 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. In *ICLR*, 2022. 1, 3
- [49] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the wild: Neural radiance fields for unconstrained photo collections. In *CVPR*, 2021. 2
- [50] Nelson Max. Optical models for direct volume rendering. *IEEE Trans. Vis. Comput. Graph.*, 1(2), 1995. 4
- [51] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 12
- [52] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul P. Srinivasan, and Jonathan T. Barron. NeRF in the dark: High dynamic range view synthesis from noisy raw images. In *CVPR*, 2022. 2
- [53] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4), 2022. 2
- [54] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. DONeRF: Towards real-time rendering of compact neural radiance fields using depth oracle networks. *Comput. Graph. Forum*, 40(4), 2021. 2
- [55] Michael Niemeyer and Andreas Geiger. CAMPARI: Camera-aware decomposed generative neural radiance fields. In *3DV*, 2021. 2
- [56] Michael Niemeyer and Andreas Geiger. GIRAFFE: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 2
- [57] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. RegNeRF: Regularizing neural radiance fields for view synthesis from sparse inputs. In *CVPR*, 2022. 2, 3
- [58] 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 *ICCV*, 2021. 2, 3
- [59] Martin Píala and Ronald Clark. TerminiNeRF: Ray termination prediction for efficient neural rendering. In *3DV*, 2021. 2
- [60] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. DreamFusion: Text-to-3D using 2D diffusion. In *ICLR*, 2023. 2
- [61] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2, 3, 4
- [62] Yiling Qiao, Junbang Liang, Vladlen Koltun, and Ming Lin. Differentiable simulation of soft multi-body systems. In *NeurIPS*, 2021. 3
- [63] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual

- models from natural language supervision. In *ICML*, 2021. 3
- [64] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. DeRF: Decomposed radiance fields. In *CVPR*, 2021. 2
- [65] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. KiloNeRF: Speeding up neural radiance fields with thousands of tiny MLPs. In *ICCV*, 2021. 2
- [66] Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. Dense depth priors for neural radiance fields from sparse input views. In *CVPR*, 2022. 2, 3
- [67] Alvaro Sanchez-Gonzalez, Jonathan Godwin, Tobias Pfaff, Rex Ying, Jure Leskovec, and Peter Battaglia. Learning to simulate complex physics with graph networks. In *ICML*, 2020. 3
- [68] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. GRAF: Generative radiance fields for 3D-aware image synthesis. In *NeurIPS*, 2020. 2
- [69] Shuai Shen, Wanhua Li, Zheng Zhu, Yueqi Duan, Jie Zhou, and Jiwen Lu. Learning dynamic facial radiance fields for few-shot talking head synthesis. In *ECCV*, 2022. 2, 3
- [70] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 3
- [71] Ivan Skorokhodov, Sergey Tulyakov, Yiqun Wang, and Peter Wonka. EpiGRAF: Rethinking training of 3D GANs. In *NeurIPS*, 2022. 2
- [72] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019. 2
- [73] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 2, 5, 12, 31
- [74] Tetsuya Takahashi and Ming C. Lin. Video-guided real-to-virtual parameter transfer for viscous fluids. *ACM Trans. Graph.*, 38(6), 2019. 3
- [75] Ayush Tewari, Justus Thies, Ben Mildenhall, Pratul P. Srinivasan, Edgar Tretschk, Yifan Wang, Christoph Lassner, Vincent Sitzmann, Ricardo Martin-Brualla, Stephen Lombardi, Tomas Simon, Christian Theobalt, Matthias Nießner, Jonathan T. Barron, Gordon Wetzstein, Michael Zollhöfer, and Vladislav Golyanik. Advances in neural rendering. *Comput. Graph. Forum*, 41, 2022. 1
- [76] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 2, 3
- [77] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T. Barron, and Pratul P. Srinivasan. Ref-NeRF: Structured view-dependent appearance for neural radiance fields. In *CVPR*, 2022. 2
- [78] Bin Wang, Longhua Wu, KangKang Yin, Uri Ascher, Libin Liu, and Hui Huang. Deformation capture and modeling of soft objects. *ACM Trans. Graph.*, 94, 2015. 3
- [79] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: From error visibility to structural similarity. *IEEE Trans. Image Process.*, 13(4), 2004. 7, 31
- [80] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. NerfingMVS: Guided optimization of neural radiance fields for indoor multi-view stereo. In *ICCV*, 2021. 2, 3
- [81] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. NeX: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2
- [82] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yu-Xiong Wang, and David Forsyth. DIVER: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. In *CVPR*, 2022. 2
- [83] Yiheng Xie, Towaki Takikawa, Shunsuke Saito, Or Litany, Shiqin Yan, Numair Khan, Federico Tombari, James Tompkin, Vincent Sitzmann, and Srinath Sridhar. Neural fields in visual computing and beyond. *Comput. Graph. Forum*, 41, 2022. 1, 2
- [84] Zhenjia Xu, Jiajun Wu, Andy Zeng, Joshua B. Tenenbaum, and Shuran Song. DensePhysNet: Learning dense physical object representations via multi-step dynamic interactions. In *RSS*, 2019. 3
- [85] Yang Xue, Yuheng Li, Krishna Kumar Singh, and Yong Jae Lee. GIRAFFE HD: A high-resolution 3D-aware generative model. In *CVPR*, 2022. 2
- [86] Jiawei Yang, Marco Pavone, and Yue Wang. FreeNeRF: Improving few-shot neural rendering with free frequency regularization. In *CVPR*, 2023. 2, 3, 7, 12
- [87] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. PlenOctrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 2
- [88] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2, 3
- [89] Yonghao Yue, Breannan Smith, Christopher Batty, Changxi Zheng, and Eitan Grinspun. Continuum foam: A material point method for shear-dependent flows. *ACM Trans. Graph.*, 34(5), 2015. 3
- [90] Jason Zhang, Gengshan Yang, Shubham Tulsiani, and Deva Ramanan. NeRS: Neural reflectance surfaces for sparse-view 3D reconstruction in the wild. In *NIPS*, 2021. 2, 3
- [91] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. FdNeRF: Few-shot dynamic neural radiance fields for face reconstruction and expression editing. In *SIGGRAPH Asia*, 2022. 2, 3
- [92] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. NeRF++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2
- [93] 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, 31

## A. Details of PAC-NeRF-3v<sup>†</sup>

### A.1. Method

In the experiments (Section 4), we used PAC-NeRF-3v<sup>†</sup> as the stronger baseline. PAC-NeRF-3v<sup>†</sup> is an improved variant of PAC-NeRF-3v for sparse-view settings in which Eulerian static voxel grid optimization (Figure 2(a)) is improved. We examined this model to determine whether the proposed LPO, a few-shot learning method for *dynamic* scenes, can be combined with other few-shot learning methods, such as those for *static* scenes.<sup>11</sup> This appendix explains the details of the model. PAC-NeRF-3v<sup>†</sup> adopts three modifications: scheduling a surface regularizer, introducing view-invariant pixel-wise loss, and adjusting the training length.

**Scheduling of surface regularizer.** In the original PAC-NeRF [42], a surface regularizer  $\mathcal{L}_{surf}$  is applied to regularize the volume density field

$$\mathcal{L}_{surf} = \sum_p \text{clamp}(\alpha_p, 10^{-4}, 10^{-1}) \left( \frac{\Delta x}{2} \right)^2, \quad (9)$$

where  $\alpha_p$  indicates the alpha value of a particle with length 1 and is calculated as  $\alpha_p = 1 - \exp(-\text{softplus}(\sigma_p))$ , where  $\sigma_p$  denotes the volume density of the particle. This regularizer minimizes the total surface area, making the spread of the particles more compact and tightening their shapes. Consequently, the quality of the reconstructed geometries improved [42].

In a preliminary experiment, we found that this regularizer was also effective in eliminating unexpected masses that tend to appear in places where there are few clues owing to the lack of views. However, we also found that this regularizer has a side effect: it removes necessary components when it is too strong. Based on these observations, we scheduled this regularizer.

1. At the beginning of the training, the weight of  $\mathcal{L}_{surf}$  is initialized to a default value of PAC-NeRF [42].
2. From the beginning of the training, the weight of  $\mathcal{L}_{surf}$  is gradually increased for a certain period of time.
3. After a certain period, the weight of  $\mathcal{L}_{surf}$  is gradually decreased until it reaches the default value.

We do not impose a large weight on  $\mathcal{L}_{surf}$  from the beginning of the training (Step 1) because it can eliminate all particles, leading to learning “none” object. We decreased the

<sup>11</sup>As described in the footnote of the main text,<sup>8</sup> in preliminary experiments, we found that previous representative few-shot learning methods (for example, DietNeRF [33] and FreeNeRF [86]) were less stable than the standard PAC-NeRF-3v. This is possible because, in our experimental settings, the number of views was small (three) despite the wide range of views (upper hemisphere), and explicit voxel representations were more effective than the fully implicit representation in [33, 86]. Therefore, we used an improved variant that we developed. However, this study and previous studies are not competitive but complementary. Therefore, further investigation is important.

weight of  $\mathcal{L}_{surf}$  in Step 3 to alleviate the negative effects caused by the introduction of strong regularization.

**Implementation details.** In the experiments, we doubled the weight of  $\mathcal{L}_{surf}$  in Step 2 every 100 iterations until the weight reached eight times its default value. In Step 3, we halved the weight of  $\mathcal{L}_{surf}$  each time the resolution of voxel grids was scaled. This halving process was conducted three times; therefore, the weight of  $\mathcal{L}_{surf}$  returned to the default value when all halving processes were finished.

**View-invariant pixel-wise loss.** In sparse view settings, it is challenging to distinguish view-dependent from view-independent factors because there are few clues. Specifically, in NeRF [51] (mainly, voxel-based NeRF [73]), the view-dependent and view-independent factors (i.e., colors) are represented by a multilayer perceptron (MLP), which additionally receives a view direction  $\mathbf{d}$ , and color fields  $\mathbf{c}^G$  that do not receive  $\mathbf{d}$ , respectively (Equation (8)). In sparse-view settings, dividing the roles between the MLP and  $\mathbf{c}^G$  is not trivial. In extreme cases, the MLP can overfit specific views in the training data (in such cases,  $\mathbf{c}^G$  no longer plays an essential role in representing colors), making it difficult to represent colors in novel views. To alleviate this difficulty, we introduce a view-invariant (VI) pixel-wise loss  $\mathcal{L}_{pixel}^{VI}$ , which is a variant of the pixel-wise loss  $\mathcal{L}_{pixel}$  (Equation (3)) where the color of a sample on a ray, i.e.,  $\mathbf{c}$ , is calculated by the following equation instead of Equation (8)

$$\tilde{\mathbf{c}}(\mathbf{x}, \tilde{\mathbf{d}}, t) = \text{MLP}(\text{interp}(\mathbf{x}, \mathbf{c}^G), \tilde{\mathbf{d}}), \quad (10)$$

where  $\tilde{\mathbf{d}} \in \mathbb{S}^2$  denotes the view direction randomly sampled from a unit sphere. Here, we use  $\tilde{\mathbf{c}}$  to denote  $\mathbf{c}$  and distinguish it from the original  $\mathbf{c}$ . This loss encourages MLP and  $\mathbf{c}^G$  to capture the colors of the training images independent of the viewing direction. Consequently, the model makes it possible to avoid extreme cases (i.e., it mitigates the MLP to overfit specific views in the training data and prevents  $\mathbf{c}^G$  from losing a role) and provides some colors even for novel views. The total pixel-wise loss  $\mathcal{L}_{pixel}^\dagger$  is given by

$$\mathcal{L}_{pixel}^\dagger = \mathcal{L}_{pixel} + \lambda \mathcal{L}_{pixel}^{VI}, \quad (11)$$

where  $\lambda$  is a hyperparameter balancing the two losses. During training,  $\mathcal{L}_{pixel}^\dagger$  was used instead of  $\mathcal{L}_{pixel}$ .

**Implementation details.** We set  $\lambda = 0.1$  in the experiments.

**Adjustment of training length.** As discussed in previous studies [13, 86], one of the factors that make learning difficult in few-shot settings is overfitting to sparse views in the training data, causing a loss of generalization ability. In preliminary experiments, we observed a similar tendency in our settings. Based on this observation, we adjusted the training length (TL). In particular, the number of iterations was reduced. Despite its simplicity, we empirically found that this solution works well in severe view settings (e.g. when the number of views is three).

|            |                    | (a) PAC-NeRF-3v | (b) PAC-NeRF-3v <sup>†</sup> | (c) w/o scheduling of $\mathcal{L}_{surf}$ | (d) w/o $\mathcal{L}_{pixel}^{VI}$ | (e) w/o adjustment of TL |
|------------|--------------------|-----------------|------------------------------|--|------------------------------------|--------------------------|
| Droplet    | PSNR $\uparrow$    | 24.42           | 26.83                        | 26.51                                      | 26.41                              | 24.88                    |
|            | SSIM $\uparrow$    | 0.975           | 0.981                        | 0.981                                      | 0.979                              | 0.977                    |
|            | LPIPS $\downarrow$ | 0.048           | 0.045                        | 0.050                                      | 0.048                              | 0.047                    |
| Letter     | PSNR $\uparrow$    | 28.56           | 29.37                        | 28.52                                      | 27.75                              | 27.62                    |
|            | SSIM $\uparrow$    | 0.979           | 0.981                        | 0.978                                      | 0.977                              | 0.977                    |
|            | LPIPS $\downarrow$ | 0.032           | 0.030                        | 0.036                                      | 0.032                              | 0.034                    |
| Cream      | PSNR $\uparrow$    | 26.10           | 27.09                        | 26.51                                      | 25.73                              | 25.67                    |
|            | SSIM $\uparrow$    | 0.982           | 0.983                        | 0.981                                      | 0.981                              | 0.980                    |
|            | LPIPS $\downarrow$ | 0.027           | 0.025                        | 0.027                                      | 0.029                              | 0.024                    |
| Toothpaste | PSNR $\uparrow$    | 28.70           | 31.77                        | 31.70                                      | 31.60                              | 30.76                    |
|            | SSIM $\uparrow$    | 0.988           | 0.991                        | 0.991                                      | 0.991                              | 0.990                    |
|            | LPIPS $\downarrow$ | 0.013           | 0.012                        | 0.011                                      | 0.011                              | 0.012                    |
| Torus      | PSNR $\uparrow$    | 25.17           | 26.77                        | 25.33                                      | 26.15                              | 26.55                    |
|            | SSIM $\uparrow$    | 0.972           | 0.974                        | 0.972                                      | 0.973                              | 0.972                    |
|            | LPIPS $\downarrow$ | 0.047           | 0.043                        | 0.047                                      | 0.043                              | 0.044                    |
| Bird       | PSNR $\uparrow$    | 26.29           | 27.64                        | 28.82                                      | 27.45                              | 26.26                    |
|            | SSIM $\uparrow$    | 0.979           | 0.982                        | 0.982                                      | 0.982                              | 0.981                    |
|            | LPIPS $\downarrow$ | 0.034           | 0.029                        | 0.033                                      | 0.031                              | 0.030                    |
| Playdoh    | PSNR $\uparrow$    | 27.87           | 30.05                        | 29.03                                      | 29.77                              | 28.32                    |
|            | SSIM $\uparrow$    | 0.976           | 0.981                        | 0.981                                      | 0.981                              | 0.978                    |
|            | LPIPS $\downarrow$ | 0.046           | 0.042                        | 0.046                                      | 0.043                              | 0.043                    |
| Cat        | PSNR $\uparrow$    | 30.82           | 31.93                        | 31.90                                      | 30.36                              | 31.75                    |
|            | SSIM $\uparrow$    | 0.987           | 0.989                        | 0.989                                      | 0.987                              | 0.988                    |
|            | LPIPS $\downarrow$ | 0.035           | 0.035                        | 0.035                                      | 0.038                              | 0.036                    |
| Trophy     | PSNR $\uparrow$    | 28.93           | 29.05                        | 29.23                                      | 29.17                              | 28.85                    |
|            | SSIM $\uparrow$    | 0.963           | 0.963                        | 0.964                                      | 0.964                              | 0.964                    |
|            | LPIPS $\downarrow$ | 0.038           | 0.039                        | 0.039                                      | 0.038                              | 0.036                    |
| Average    | PSNR $\uparrow$    | 27.43           | 28.94                        | 28.62                                      | 28.27                              | 27.85                    |
|            | SSIM $\uparrow$    | 0.978           | 0.980                        | 0.980                                      | 0.979                              | 0.979                    |
|            | LPIPS $\downarrow$ | 0.036           | 0.033                        | 0.036                                      | 0.035                              | 0.034                    |

Table 4. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  on **static voxel grid optimization**. The scores are calculated for the first frames of the video sequences in the test set. Here, PAC-NeRF-3v<sup>†</sup> (b), an improved variant of PAC-NeRF-3v, is compared with the original PAC-NeRF-3v (a) and the ablated models, including PAC-NeRF-3v<sup>†</sup> without scheduling of  $\mathcal{L}_{surf}$  (c), that without  $\mathcal{L}_{pixel}^{VI}$  (d), and that without adjustment of TL (e). PAC-NeRF-3v<sup>†</sup> (b) achieved the 20 best and five second-best scores among 27 evaluation items.

*Implementation details.* In the experiments, we reduced the number of iterations to one-third; the default number (6000) was reduced to 2000. Based on this change, we reduced the timing of scaling the resolution of voxel grids by one-third.

## A.2. Experiments

Ablation studies were conducted to confirm the importance of each modification. Specifically, we compared *PAC-NeRF-3v<sup>†</sup>* with three ablated models: *PAC-NeRF-3v<sup>†</sup> without scheduling of  $\mathcal{L}_{surf}$* , *PAC-NeRF-3v<sup>†</sup> without  $\mathcal{L}_{pixel}^{VI}$* , and *PAC-NeRF-3v<sup>†</sup> without adjustment of TL*. We also examined the original *PAC-NeRF-3v*, that is, PAC-NeRF-3v<sup>†</sup> without all three modifications.

**Results.** Table 4 summarizes the results. We found that PAC-NeRF-3v<sup>†</sup> achieved the best or second-best scores in most cases (20 best scores and five second-best scores among the 27 evaluation items). Consequently, PAC-NeRF-3v<sup>†</sup> achieved the best average scores in terms of all metrics. These results indicate the importance of each modification.

## B. Detailed analyses on main experiments

The main text focuses on representative results because of space limitations. This appendix provides the extended results that further clarify the effectiveness of the proposed method. In particular, we provide the scores for each scene (Appendix B.1), discuss the computation times (Appendix B.2), and discuss the impact of the number of iterations in Algorithm 1 (Appendix B.3).

### B.1. Scores for each scene

In this appendix, we provide the scores for each scene. To evaluate the effectiveness of the proposed method from various perspectives, we describe three experiments in the main text. (I) Evaluation of the geometric correction (Table 1), where the image quality after LPO was applied once, was evaluated. (II) Evaluation of the physical identification (Table 2), in which the accuracy of the physical property estimation after Algorithm 1 was applied, was evaluated. (III) Evaluation of the geometric recorection (Table 3), in which the image quality after Algorithm 1 was applied, was evaluated. In this appendix, we discuss these issues in detail.

**I. Evaluation of geometric correction.** First, we discuss

|            |                    | PAC-NeRF | PAC-NeRF-3v | +LPO  | +LPO-F | +LPO-P | +GO   | PAC-NeRF-3v <sup>†</sup> | +LPO  | +LPO-F | +LPO-P | +GO   |
|------------|--------------------|----------|-------------|-------|--------|--------|-------|--------------------------|-------|--------|--------|-------|
| Droplet    | PSNR $\uparrow$    | 35.30    | 25.42       | 27.56 | 27.05  | 27.26  | 27.55 | 26.40                    | 28.18 | 27.36  | 28.00  | 27.36 |
|            | SSIM $\uparrow$    | 0.990    | 0.975       | 0.978 | 0.977  | 0.977  | 0.978 | 0.978                    | 0.980 | 0.979  | 0.980  | 0.979 |
|            | LPIPS $\downarrow$ | 0.029    | 0.047       | 0.043 | 0.045  | 0.044  | 0.044 | 0.046                    | 0.044 | 0.045  | 0.044  | 0.046 |
| Letter     | PSNR $\uparrow$    | 36.01    | 28.94       | 29.99 | 29.53  | 29.90  | 29.56 | 29.59                    | 30.44 | 30.23  | 30.29  | 30.09 |
|            | SSIM $\uparrow$    | 0.991    | 0.981       | 0.982 | 0.982  | 0.982  | 0.982 | 0.983                    | 0.982 | 0.983  | 0.982  | 0.983 |
|            | LPIPS $\downarrow$ | 0.012    | 0.028       | 0.029 | 0.028  | 0.029  | 0.027 | 0.028                    | 0.029 | 0.027  | 0.030  | 0.026 |
| Cream      | PSNR $\uparrow$    | 36.70    | 26.61       | 28.48 | 27.79  | 27.80  | 28.39 | 27.43                    | 28.82 | 28.34  | 28.50  | 28.73 |
|            | SSIM $\uparrow$    | 0.993    | 0.983       | 0.983 | 0.984  | 0.983  | 0.984 | 0.983                    | 0.984 | 0.984  | 0.984  | 0.984 |
|            | LPIPS $\downarrow$ | 0.014    | 0.026       | 0.023 | 0.024  | 0.025  | 0.023 | 0.024                    | 0.023 | 0.023  | 0.024  | 0.023 |
| Toothpaste | PSNR $\uparrow$    | 39.46    | 29.23       | 31.27 | 30.90  | 30.66  | 31.63 | 31.72                    | 33.54 | 33.15  | 33.15  | 33.77 |
|            | SSIM $\uparrow$    | 0.996    | 0.991       | 0.991 | 0.991  | 0.991  | 0.991 | 0.993                    | 0.993 | 0.993  | 0.993  | 0.993 |
|            | LPIPS $\downarrow$ | 0.006    | 0.010       | 0.010 | 0.010  | 0.011  | 0.010 | 0.010                    | 0.010 | 0.010  | 0.010  | 0.009 |
| Torus      | PSNR $\uparrow$    | 34.54    | 23.99       | 28.91 | 25.22  | 27.82  | 24.92 | 26.65                    | 30.05 | 29.18  | 29.03  | 29.27 |
|            | SSIM $\uparrow$    | 0.988    | 0.970       | 0.978 | 0.971  | 0.977  | 0.970 | 0.974                    | 0.980 | 0.978  | 0.978  | 0.978 |
|            | LPIPS $\downarrow$ | 0.026    | 0.048       | 0.038 | 0.046  | 0.040  | 0.045 | 0.040                    | 0.034 | 0.036  | 0.037  | 0.036 |
| Bird       | PSNR $\uparrow$    | 35.55    | 24.82       | 27.20 | 25.47  | 27.18  | 25.30 | 24.91                    | 28.51 | 25.73  | 28.43  | 26.18 |
|            | SSIM $\uparrow$    | 0.991    | 0.978       | 0.980 | 0.978  | 0.980  | 0.978 | 0.979                    | 0.983 | 0.980  | 0.983  | 0.980 |
|            | LPIPS $\downarrow$ | 0.019    | 0.036       | 0.032 | 0.035  | 0.032  | 0.036 | 0.036                    | 0.028 | 0.033  | 0.028  | 0.033 |
| Playdoh    | PSNR $\uparrow$    | 36.43    | 27.70       | 28.39 | 27.85  | 28.41  | 27.80 | 29.66                    | 30.07 | 29.71  | 30.13  | 29.63 |
|            | SSIM $\uparrow$    | 0.991    | 0.978       | 0.978 | 0.977  | 0.978  | 0.977 | 0.982                    | 0.982 | 0.982  | 0.982  | 0.982 |
|            | LPIPS $\downarrow$ | 0.026    | 0.042       | 0.041 | 0.042  | 0.041  | 0.042 | 0.038                    | 0.038 | 0.038  | 0.038  | 0.037 |
| Cat        | PSNR $\uparrow$    | 37.53    | 30.45       | 31.00 | 30.53  | 30.99  | 30.19 | 31.11                    | 31.41 | 30.98  | 31.48  | 30.55 |
|            | SSIM $\uparrow$    | 0.993    | 0.987       | 0.987 | 0.987  | 0.987  | 0.986 | 0.988                    | 0.988 | 0.988  | 0.988  | 0.988 |
|            | LPIPS $\downarrow$ | 0.016    | 0.033       | 0.032 | 0.031  | 0.031  | 0.031 | 0.034                    | 0.032 | 0.032  | 0.032  | 0.032 |
| Trophy     | PSNR $\uparrow$    | 32.43    | 29.33       | 30.16 | 29.69  | 30.01  | 29.66 | 28.80                    | 29.92 | 29.10  | 29.78  | 28.96 |
|            | SSIM $\uparrow$    | 0.967    | 0.963       | 0.964 | 0.963  | 0.964  | 0.963 | 0.962                    | 0.964 | 0.963  | 0.964  | 0.963 |
|            | LPIPS $\downarrow$ | 0.034    | 0.039       | 0.037 | 0.039  | 0.038  | 0.039 | 0.039                    | 0.037 | 0.039  | 0.038  | 0.039 |
| Average    | PSNR $\uparrow$    | 35.99    | 27.39       | 29.22 | 28.22  | 28.89  | 28.33 | 28.47                    | 30.11 | 29.31  | 29.87  | 29.39 |
|            | SSIM $\uparrow$    | 0.989    | 0.978       | 0.980 | 0.979  | 0.980  | 0.979 | 0.980                    | 0.982 | 0.981  | 0.981  | 0.981 |
|            | LPIPS $\downarrow$ | 0.020    | 0.034       | 0.032 | 0.033  | 0.032  | 0.033 | 0.033                    | 0.031 | 0.032  | 0.031  | 0.031 |

Table 5. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  for each scene on **geometric correction**. This table is an extended version of Table 1. The qualitative comparisons are provided in Figures 3–7.

the scores for each scene on geometric correction (image quality after LPO was applied once) in detail. Table 5 summarizes the results. This is an extension of the list in Table 1. We discuss the results from three perspectives.

(1) *PAC-NeRF-3v/3v<sup>†</sup> vs. +LPO*. When PAC-NeRF-3v was used as the baseline, +LPO outperformed the baseline in most cases (21 wins, 5 draws, and 1 loss). Similarly, when PAC-NeRF-3v<sup>†</sup> was used as the baseline, +LPO outperformed the baseline in most cases (20 wins, 5 draws, and 2 losses). Consequently, in both cases, +LPO yielded better average scores for all metrics. These results demonstrate that +LPO is effective for geometric correction independent of the baseline models. The qualitative comparisons are shown in Figures 3–7.

(2) *+LPO vs. +LPO-F/P*. +LPO-F and +LPO-P are ablated variants of +LPO. In +LPO-F, position (shape) optimization is ablated, and only feature (appearance) optimization is conducted. In +LPO-P, feature (i.e., appearance) optimization is ablated, and only position (i.e., shape) optimization is conducted. +LPO outperformed +LPO-F/P in most cases when both PAC-NeRF-3v and PAC-NeRF-3v<sup>†</sup> were used as the baselines. Specifically, when PAC-NeRF-3v was used as the baseline, +LPO outperformed +LPO-F with 20 wins, 4 draws, and 3 losses and outperformed +LPO-P

with 15 wins, 10 draws, and 2 losses. When PAC-NeRF-3v<sup>†</sup> was used as the baseline, +LPO outperformed +LPO-F with 17 wins, 8 draws, and 2 losses, and outperformed +LPO-P with 12 wins, 13 draws, and 2 losses, respectively. Between +LPO-F and +LPO-P, +LPO-P tends to outperform +LPO-F. Specifically, when PAC-NeRF-3v was used as the baseline, +LPO-P outperformed +LPO-F with 17 wins, 5 draws, and 5 losses. When PAC-NeRF-3v<sup>†</sup> was used as the baseline, +LPO-P outperformed +LPO-F with 13 wins, 9 draws, and 5 losses. We consider this because shape correction by +LPO-P can correct the failure estimation of the geometry within the physical constraints of MPM. In contrast, appearance correction by +LPO-F cannot do so and can cause overcorrection beyond the physical constraints, as shown in Figure 8. These results indicate that the feature and position optimizations are complementary rather than competitive.

(3) *+LPO vs. +GO*. The difference between these two models is that +LPO conducts optimization in Lagrangian particle space and can optimize not only the features but also the positions of the particles, whereas +GO performs optimization in Eulerian grid space and can optimize the features of the grids, but cannot optimize their positions. Owing to these characteristics, the performance of +GO is close to that of +LPO-F, which also optimizes the features but not

|            |                 | Ground truth       | PAC-NeRF           | PAC-NeRF-3v           | +LPO <sup>4</sup>  | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup>   | +None <sup>4</sup>    |
|------------|-----------------|--------------------|--------------------|-----------------------|--------------------|---------------------|---------------------|--------------------|-----------------------|
| Droplet    | $\mu$           | $2.00 \times 10^2$ | $2.19 \times 10^2$ | $2.76 \times 10^2$    | $2.59 \times 10^2$ | $2.59 \times 10^2$  | $2.63 \times 10^2$  | $2.58 \times 10^2$ | $2.70 \times 10^2$    |
|            | $\kappa$        | $1.00 \times 10^5$ | $9.62 \times 10^4$ | $5.19 \times 10^3$    | $2.36 \times 10^3$ | $2.30 \times 10^3$  | $3.42 \times 10^3$  | $2.11 \times 10^3$ | $2.07 \times 10^3$    |
| Letter     | $\mu$           | $1.00 \times 10^2$ | $9.10 \times 10^1$ | $2.12 \times 10^1$    | $9.67 \times 10^1$ | $9.42 \times 10^1$  | $1.20 \times 10^2$  | $1.12 \times 10^2$ | $2.95 \times 10^1$    |
|            | $\kappa$        | $1.00 \times 10^5$ | $9.14 \times 10^4$ | $1.69 \times 10^{-2}$ | $6.69 \times 10^4$ | $3.88 \times 10^4$  | $9.12 \times 10^3$  | $1.37 \times 10^4$ | $8.26 \times 10^{-2}$ |
| Cream      | $\mu$           | $1.00 \times 10^4$ | $1.23 \times 10^4$ | $2.05 \times 10^4$    | $1.51 \times 10^4$ | $1.46 \times 10^4$  | $1.72 \times 10^4$  | $1.52 \times 10^4$ | $1.78 \times 10^4$    |
|            | $\kappa$        | $1.00 \times 10^6$ | $1.35 \times 10^6$ | $1.64 \times 10^6$    | $1.44 \times 10^6$ | $1.78 \times 10^6$  | $9.39 \times 10^5$  | $9.37 \times 10^5$ | $1.75 \times 10^6$    |
|            | $\tau_Y$        | $3.00 \times 10^3$ | $3.05 \times 10^3$ | $2.90 \times 10^3$    | $2.97 \times 10^3$ | $2.80 \times 10^3$  | $2.96 \times 10^3$  | $3.04 \times 10^3$ | $2.83 \times 10^3$    |
|            | $\eta$          | 10.00              | 10.36              | 19.10                 | 15.22              | 18.02               | 12.11               | 8.69               | 17.85                 |
| Toothpaste | $\mu$           | $5.00 \times 10^3$ | $5.31 \times 10^3$ | $3.89 \times 10^5$    | $6.42 \times 10^3$ | $7.16 \times 10^3$  | $3.80 \times 10^3$  | $3.43 \times 10^3$ | $3.05 \times 10^3$    |
|            | $\kappa$        | $1.00 \times 10^5$ | $5.66 \times 10^4$ | $2.63 \times 10^3$    | $2.51 \times 10^4$ | $4.41 \times 10^3$  | $2.53 \times 10^4$  | $2.25 \times 10^4$ | $1.87 \times 10^4$    |
|            | $\tau_Y$        | $2.00 \times 10^2$ | $2.33 \times 10^2$ | $3.18 \times 10^2$    | $1.54 \times 10^2$ | $3.11 \times 10^2$  | $1.68 \times 10^2$  | $1.12 \times 10^2$ | $1.74 \times 10^2$    |
|            | $\eta$          | 10.00              | 9.71               | 4.23                  | 9.93               | 5.41                | 10.31               | 10.16              | 8.97                  |
| Torus      | $E$             | $1.00 \times 10^6$ | $1.05 \times 10^6$ | $1.89 \times 10^6$    | $1.15 \times 10^6$ | $1.02 \times 10^6$  | $1.21 \times 10^6$  | $1.02 \times 10^6$ | $8.83 \times 10^5$    |
|            | $\nu$           | 0.300              | 0.323              | 0.215                 | 0.299              | 0.420               | 0.331               | 0.374              | -0.016                |
| Bird       | $E$             | $3.00 \times 10^5$ | $2.91 \times 10^5$ | $8.43 \times 10^5$    | $3.50 \times 10^5$ | $1.73 \times 10^5$  | $3.56 \times 10^5$  | $1.46 \times 10^5$ | $1.70 \times 10^5$    |
|            | $\nu$           | 0.300              | 0.329              | 0.402                 | 0.301              | -0.072              | 0.225               | -0.065             | -0.281                |
| Playdoh    | $E$             | $2.00 \times 10^6$ | $3.87 \times 10^6$ | $3.90 \times 10^6$    | $2.61 \times 10^6$ | $7.61 \times 10^6$  | $2.64 \times 10^6$  | $4.03 \times 10^6$ | $2.96 \times 10^6$    |
|            | $\tau_Y$        | $1.54 \times 10^4$ | $1.68 \times 10^4$ | $2.96 \times 10^4$    | $2.25 \times 10^4$ | $1.64 \times 10^4$  | $2.52 \times 10^4$  | $2.42 \times 10^4$ | $2.66 \times 10^4$    |
|            | $\nu$           | 0.300              | 0.224              | -0.195                | 0.189              | 0.127               | 0.191               | 0.052              | 0.172                 |
| Cat        | $E$             | $1.00 \times 10^6$ | $1.39 \times 10^5$ | $5.00 \times 10^4$    | $1.06 \times 10^5$ | $8.55 \times 10^4$  | $8.49 \times 10^4$  | $7.46 \times 10^4$ | $8.58 \times 10^4$    |
|            | $\tau_Y$        | $3.85 \times 10^3$ | $3.62 \times 10^3$ | $5.08 \times 10^3$    | $4.93 \times 10^3$ | $5.04 \times 10^3$  | $4.99 \times 10^3$  | $4.69 \times 10^3$ | $4.99 \times 10^3$    |
|            | $\nu$           | 0.300              | 0.327              | 0.344                 | 0.304              | 0.336               | 0.231               | 0.110              | 0.379                 |
| Trophy     | $\theta_{fric}$ | 40.00°             | 37.28°             | 36.97°                | 36.84°             | 37.07°              | 36.79°              | 36.73°             | 36.90°                |

Table 6. Comparison of the values of the physical properties for each scene on **physical identification** when PAC-NeRF-3v was used as a baseline. The absolute differences between the ground truth and the estimated physical properties are provided in Table 2. The qualitative comparisons are provided in Figures 3–9.

|            |                 | Ground truth       | PAC-NeRF           | PAC-NeRF-3v <sup>†</sup> | +LPO <sup>4</sup>  | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup>   | +None <sup>4</sup>    |
|------------|-----------------|--------------------|--------------------|--------------------------|--------------------|---------------------|---------------------|--------------------|-----------------------|
| Droplet    | $\mu$           | $2.00 \times 10^2$ | $2.19 \times 10^2$ | $2.73 \times 10^2$       | $2.41 \times 10^2$ | $2.34 \times 10^2$  | $2.33 \times 10^2$  | $2.42 \times 10^2$ | $2.46 \times 10^2$    |
|            | $\kappa$        | $1.00 \times 10^5$ | $9.62 \times 10^4$ | $5.45 \times 10^3$       | $1.28 \times 10^5$ | $7.26 \times 10^3$  | $3.02 \times 10^4$  | $3.02 \times 10^3$ | $2.21 \times 10^3$    |
| Letter     | $\mu$           | $1.00 \times 10^2$ | $9.10 \times 10^1$ | $4.18 \times 10^1$       | $1.02 \times 10^2$ | $8.18 \times 10^1$  | $9.84 \times 10^1$  | $7.62 \times 10^1$ | $3.66 \times 10^1$    |
|            | $\kappa$        | $1.00 \times 10^5$ | $9.14 \times 10^4$ | $5.91 \times 10^{-1}$    | $8.70 \times 10^4$ | $1.06 \times 10^5$  | $1.34 \times 10^5$  | $7.28 \times 10^4$ | $8.72 \times 10^{-1}$ |
| Cream      | $\mu$           | $1.00 \times 10^4$ | $1.23 \times 10^4$ | $1.51 \times 10^4$       | $1.26 \times 10^4$ | $1.18 \times 10^4$  | $1.16 \times 10^4$  | $1.24 \times 10^4$ | $1.19 \times 10^4$    |
|            | $\kappa$        | $1.00 \times 10^6$ | $1.35 \times 10^6$ | $2.17 \times 10^6$       | $1.32 \times 10^6$ | $1.57 \times 10^6$  | $1.39 \times 10^6$  | $1.56 \times 10^6$ | $2.18 \times 10^6$    |
|            | $\tau_Y$        | $3.00 \times 10^3$ | $3.05 \times 10^3$ | $2.94 \times 10^3$       | $3.04 \times 10^3$ | $2.99 \times 10^3$  | $2.97 \times 10^3$  | $2.93 \times 10^3$ | $2.95 \times 10^3$    |
|            | $\eta$          | 10.00              | 10.36              | 15.67                    | 10.80              | 11.44               | 12.01               | 12.47              | 11.25                 |
| Toothpaste | $\mu$           | $5.00 \times 10^3$ | $5.31 \times 10^3$ | $2.80 \times 10^3$       | $4.66 \times 10^3$ | $3.15 \times 10^3$  | $4.94 \times 10^3$  | $2.80 \times 10^3$ | $3.64 \times 10^3$    |
|            | $\kappa$        | $1.00 \times 10^5$ | $5.66 \times 10^4$ | $3.67 \times 10^3$       | $2.12 \times 10^4$ | $2.34 \times 10^4$  | $1.26 \times 10^4$  | $2.57 \times 10^4$ | $2.53 \times 10^4$    |
|            | $\tau_Y$        | $2.00 \times 10^2$ | $2.33 \times 10^2$ | $3.16 \times 10^2$       | $1.62 \times 10^2$ | $1.72 \times 10^2$  | $1.53 \times 10^2$  | $1.42 \times 10^2$ | $1.77 \times 10^2$    |
|            | $\eta$          | 10.00              | 9.71               | 6.13                     | 10.20              | 9.65                | 10.12               | 10.16              | 9.07                  |
| Torus      | $E$             | $1.00 \times 10^6$ | $1.05 \times 10^6$ | $8.43 \times 10^5$       | $1.09 \times 10^6$ | $1.06 \times 10^6$  | $1.09 \times 10^6$  | $9.66 \times 10^5$ | $1.18 \times 10^6$    |
|            | $\nu$           | 0.300              | 0.323              | 0.431                    | 0.307              | 0.340               | 0.267               | 0.350              | 0.429                 |
| Bird       | $E$             | $3.00 \times 10^5$ | $2.91 \times 10^5$ | $3.99 \times 10^5$       | $3.19 \times 10^5$ | $4.61 \times 10^5$  | $3.28 \times 10^5$  | $5.10 \times 10^5$ | $5.93 \times 10^5$    |
|            | $\nu$           | 0.300              | 0.329              | 0.441                    | 0.347              | 0.372               | 0.308               | 0.432              | 0.445                 |
| Playdoh    | $E$             | $2.00 \times 10^6$ | $3.87 \times 10^6$ | $6.63 \times 10^6$       | $2.72 \times 10^6$ | $5.95 \times 10^6$  | $1.51 \times 10^4$  | $3.41 \times 10^6$ | $1.42 \times 10^4$    |
|            | $\tau_Y$        | $1.54 \times 10^4$ | $1.68 \times 10^4$ | $1.99 \times 10^4$       | $2.30 \times 10^4$ | $1.85 \times 10^4$  | $2.33 \times 10^5$  | $2.42 \times 10^4$ | $1.43 \times 10^5$    |
|            | $\nu$           | 0.300              | 0.224              | 0.088                    | 0.237              | 0.167               | 0.410               | 0.173              | 0.406                 |
| Cat        | $E$             | $1.00 \times 10^6$ | $1.39 \times 10^5$ | $6.42 \times 10^4$       | $1.97 \times 10^5$ | $2.93 \times 10^5$  | $1.94 \times 10^5$  | $1.65 \times 10^5$ | $1.61 \times 10^5$    |
|            | $\tau_Y$        | $3.85 \times 10^3$ | $3.62 \times 10^3$ | $4.67 \times 10^3$       | $4.50 \times 10^3$ | $3.97 \times 10^3$  | $4.59 \times 10^3$  | $4.52 \times 10^3$ | $4.54 \times 10^3$    |
|            | $\nu$           | 0.300              | 0.327              | 0.418                    | 0.303              | 0.327               | 0.293               | 0.272              | 0.363                 |
| Trophy     | $\theta_{fric}$ | 40.00°             | 37.28°             | 38.31°                   | 37.75°             | 37.66°              | 37.51°              | 37.46°             | 37.79°                |

Table 7. Comparison of the values of the physical properties for each scene on **physical identification** when PAC-NeRF-3v<sup>†</sup> was used as a baseline. The absolute differences between the ground truth and the estimated physical properties are provided in Table 2. The qualitative comparisons are provided in Figures 3–9.

|            |                    | PAC-NeRF | PAC-NeRF-3v | +LPO <sup>4</sup> | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup> | +None <sup>4</sup> | PAC-NeRF-3v <sup>†</sup> | +LPO <sup>4</sup> | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup> | +None <sup>4</sup> |
|------------|--------------------|----------|-------------|-------------------|---------------------|---------------------|------------------|--------------------|--------------------------|-------------------|---------------------|---------------------|------------------|--------------------|
| Droplet    | PSNR <sup>↑</sup>  | 35.30    | 25.42       | 28.21             | 26.52               | 28.21               | 26.03            | 25.42              | 26.40                    | 28.91             | 27.85               | 29.03               | 26.23            | 27.52              |
|            | SSIM <sup>↑</sup>  | 0.990    | 0.975       | 0.981             | 0.978               | 0.980               | 0.978            | 0.977              | 0.978                    | 0.983             | 0.982               | 0.983               | 0.980            | 0.981              |
|            | LPIPS <sup>↓</sup> | 0.029    | 0.047       | 0.042             | 0.046               | 0.045               | 0.046            | 0.047              | 0.046                    | 0.041             | 0.042               | 0.043               | 0.045            | 0.043              |
| Letter     | PSNR <sup>↑</sup>  | 36.01    | 28.94       | 30.07             | 30.07               | 29.83               | 30.00            | 29.25              | 29.59                    | 30.26             | 31.09               | 30.50               | 29.70            | 29.80              |
|            | SSIM <sup>↑</sup>  | 0.991    | 0.981       | 0.983             | 0.983               | 0.982               | 0.983            | 0.983              | 0.983                    | 0.983             | 0.984               | 0.982               | 0.982            | 0.984              |
|            | LPIPS <sup>↓</sup> | 0.012    | 0.028       | 0.026             | 0.026               | 0.029               | 0.026            | 0.026              | 0.028                    | 0.026             | 0.024               | 0.029               | 0.025            | 0.025              |
| Cream      | PSNR <sup>↑</sup>  | 36.70    | 26.61       | 29.65             | 29.02               | 28.83               | 28.95            | 26.57              | 27.43                    | 30.23             | 29.87               | 29.65               | 29.79            | 27.88              |
|            | SSIM <sup>↑</sup>  | 0.993    | 0.983       | 0.985             | 0.985               | 0.984               | 0.985            | 0.983              | 0.983                    | 0.985             | 0.985               | 0.985               | 0.985            | 0.984              |
|            | LPIPS <sup>↓</sup> | 0.014    | 0.026       | 0.022             | 0.023               | 0.024               | 0.022            | 0.026              | 0.024                    | 0.021             | 0.022               | 0.022               | 0.021            | 0.024              |
| Toothpaste | PSNR <sup>↑</sup>  | 39.46    | 29.23       | 32.18             | 32.41               | 31.87               | 32.87            | 28.96              | 31.72                    | 34.00             | 34.51               | 33.94               | 34.40            | 31.49              |
|            | SSIM <sup>↑</sup>  | 0.996    | 0.991       | 0.991             | 0.992               | 0.991               | 0.992            | 0.991              | 0.993                    | 0.992             | 0.993               | 0.993               | 0.993            | 0.993              |
|            | LPIPS <sup>↓</sup> | 0.006    | 0.010       | 0.011             | 0.010               | 0.011               | 0.010            | 0.010              | 0.010                    | 0.010             | 0.009               | 0.010               | 0.009            | 0.009              |
| Torus      | PSNR <sup>↑</sup>  | 34.54    | 23.99       | 30.07             | 27.06               | 29.37               | 25.36            | 24.82              | 26.65                    | 30.48             | 30.01               | 30.23               | 28.93            | 26.74              |
|            | SSIM <sup>↑</sup>  | 0.988    | 0.970       | 0.982             | 0.975               | 0.981               | 0.970            | 0.973              | 0.974                    | 0.983             | 0.982               | 0.982               | 0.979            | 0.975              |
|            | LPIPS <sup>↓</sup> | 0.026    | 0.048       | 0.032             | 0.039               | 0.035               | 0.041            | 0.044              | 0.040                    | 0.031             | 0.032               | 0.033               | 0.033            | 0.040              |
| Bird       | PSNR <sup>↑</sup>  | 35.55    | 24.82       | 27.97             | 24.86               | 27.71               | 24.20            | 23.98              | 24.91                    | 28.98             | 26.51               | 28.74               | 26.08            | 24.95              |
|            | SSIM <sup>↑</sup>  | 0.991    | 0.978       | 0.980             | 0.979               | 0.980               | 0.977            | 0.978              | 0.979                    | 0.982             | 0.980               | 0.982               | 0.979            | 0.980              |
|            | LPIPS <sup>↓</sup> | 0.019    | 0.036       | 0.034             | 0.039               | 0.034               | 0.041            | 0.040              | 0.036                    | 0.031             | 0.033               | 0.032               | 0.034            | 0.036              |
| Playdoh    | PSNR <sup>↑</sup>  | 36.43    | 27.70       | 29.32             | 28.05               | 29.01               | 28.02            | 27.94              | 29.66                    | 30.01             | 29.45               | 29.25               | 28.94            | 28.02              |
|            | SSIM <sup>↑</sup>  | 0.991    | 0.978       | 0.981             | 0.978               | 0.980               | 0.979            | 0.979              | 0.982                    | 0.983             | 0.982               | 0.982               | 0.982            | 0.981              |
|            | LPIPS <sup>↓</sup> | 0.026    | 0.042       | 0.040             | 0.043               | 0.043               | 0.042            | 0.042              | 0.038                    | 0.038             | 0.039               | 0.045               | 0.041            | 0.047              |
| Cat        | PSNR <sup>↑</sup>  | 37.53    | 30.45       | 29.81             | 29.71               | 29.95               | 28.35            | 30.41              | 31.11                    | 30.61             | 30.47               | 31.47               | 29.21            | 31.78              |
|            | SSIM <sup>↑</sup>  | 0.993    | 0.987       | 0.987             | 0.987               | 0.987               | 0.985            | 0.987              | 0.988                    | 0.988             | 0.988               | 0.989               | 0.987            | 0.989              |
|            | LPIPS <sup>↓</sup> | 0.016    | 0.033       | 0.028             | 0.028               | 0.028               | 0.028            | 0.028              | 0.034                    | 0.025             | 0.026               | 0.025               | 0.026            | 0.025              |
| Trophy     | PSNR <sup>↑</sup>  | 32.43    | 29.33       | 29.57             | 29.13               | 29.89               | 28.31            | 29.44              | 28.80                    | 29.58             | 29.23               | 30.01               | 27.87            | 29.18              |
|            | SSIM <sup>↑</sup>  | 0.967    | 0.963       | 0.964             | 0.963               | 0.964               | 0.963            | 0.963              | 0.962                    | 0.964             | 0.963               | 0.964               | 0.963            | 0.963              |
|            | LPIPS <sup>↓</sup> | 0.034    | 0.039       | 0.037             | 0.039               | 0.039               | 0.039            | 0.039              | 0.039                    | 0.036             | 0.039               | 0.038               | 0.039            | 0.039              |
| Average    | PSNR <sup>↑</sup>  | 35.99    | 27.39       | 29.65             | 28.54               | 29.41               | 28.01            | 27.42              | 28.47                    | 30.34             | 29.89               | 30.31               | 29.02            | 28.60              |
|            | SSIM <sup>↑</sup>  | 0.989    | 0.978       | 0.982             | 0.980               | 0.981               | 0.979            | 0.979              | 0.980                    | 0.983             | 0.982               | 0.982               | 0.981            | 0.981              |
|            | LPIPS <sup>↓</sup> | 0.020    | 0.034       | 0.030             | 0.032               | 0.032               | 0.033            | 0.034              | 0.033                    | 0.029             | 0.030               | 0.031               | 0.030            | 0.032              |

Table 8. Comparison of PSNR<sup>↑</sup>, SSIM<sup>↑</sup>, and LPIPS<sup>↓</sup> for each scene on **geometric recorection**. This table is an extended version of Table 3. The qualitative comparisons are provided in Figures 3–9.

the positions.<sup>12</sup> Specifically, when PAC-NeRF-3v was used as the baseline, +LPO outperformed +GO with 18 wins, 5 draws, and 4 losses. When PAC-NeRF-3v<sup>†</sup> was used as the baseline, +LPO outperformed +GO with 16 wins, 6 draws, and 5 losses. These results confirm the importance of position optimization in the +LPO.

**II. Evaluation of physical identification.** Next, we discuss the scores for each scene for physical identification (the accuracy of the physical property estimation after Algorithm 1 was applied) in detail. Table 2 in the main text summarizes the absolute differences between the ground truth and the estimated physical properties. Tables 6 and 7 summarize the physical property values when PAC-NeRF-3v and PAC-NeRF-3v<sup>†</sup> were used as the baselines, respectively. These results are discussed from four perspectives.

(1) *PAC-NeRF-3v/3v<sup>†</sup> vs. +LPO<sup>4</sup>.* +LPO<sup>4</sup> improved the physical identification of PAC-NeRF-3v/3v<sup>†</sup> in most cases. Specifically, when PAC-NeRF-3v was used as the baseline, +LPO<sup>4</sup> outperformed the baseline for 21 of the 23 evaluation items. When PAC-NeRF-3v<sup>†</sup> was used as the base-

<sup>12</sup>The main differences between these two models are that, in +LPO-F, the positions of particles are fixed during training, while in +GO, the positions of particles are changed in each iteration by random sampling (note that they are not trainable).

line, +LPO<sup>4</sup> outperformed the baseline for 21 of 23 evaluation items. These results indicate that +LPO is effective for physical identification independent of the baseline models. The qualitative comparisons are shown in Figures 3–7.

(2) *+LPO<sup>4</sup> vs. +LPO-F<sup>4</sup>/P<sup>4</sup>.* When comparing these models, we found that superiority depends on the physical properties. This is because the physical properties interact, and finding the optimal balance for performance is challenging. However, +LPO-F<sup>4</sup>/P<sup>4</sup> sometimes encountered apparent difficulties. For example, the difference in  $\log_{10}(E)$  and that in  $\nu$  on Bird was large when +LPO-F<sup>4</sup> was used with PAC-NeRF-3v and the difference in  $\log_{10}(E)$  on Playdoh was large when +LPO-P<sup>4</sup> was used with PAC-NeRF-3v<sup>†</sup>. In contrast, +LPO<sup>4</sup> exhibited stable performance. Owing to this stability, +LPO<sup>4</sup> outperformed +LPO-F<sup>4</sup> and +LPO-P<sup>4</sup> in most cases. Specifically, when PAC-NeRF-3v was used as the baseline, +LPO<sup>4</sup> outperformed +LPO-F<sup>4</sup> for 18 items and outperformed +LPO-P<sup>4</sup> for 17 of the 23 evaluation items. When PAC-NeRF-3v<sup>†</sup> was used as the baseline, +LPO<sup>4</sup> outperformed +LPO-F<sup>4</sup> for 13 items and outperformed +LPO-P<sup>4</sup> for 16 of the 23 evaluation items. These results indicate that the joint optimization of features and positions in Lagrangian space is useful for obtaining stability and tackling difficult situations.



(3)  $+LPO^4$  vs.  $+GO^4$ .  $+GO^4$  also sometimes suffers from critical difficulties. For example, the difference in  $\log_{10}(E)$  and that in  $\nu$  on Bird was large when  $+GO^4$  was used with PAC-NeRF-3v. In contrast,  $+LPO^4$  exhibited stable performance and outperformed  $+GO^4$  in most cases. Specifically, when PAC-NeRF-3v was used as the baseline,  $+LPO^4$  outperformed  $+GO^4$  for 18 of the 23 evaluation items. When PAC-NeRF-3v<sup>†</sup> was used as the baseline,  $+LPO^4$  outperformed  $+GO^4$  for 18 of the 23 evaluation items.

(4)  $+LPO^4$  vs.  $+None^4$ .  $+None^4$  is outperformed by the dynamic optimization methods (i.e.,  $+LPO^4$ ,  $+LPO-F^4$ ,  $+LPO-P^4$ , and  $+GO^4$ ) in most cases. In particular, when PAC-NeRF-3v was used as the baseline,  $+LPO^4$  outperformed  $+None^4$  for 20 of the 23 evaluation items. When PAC-NeRF-3v<sup>†</sup> was used as the baseline,  $+LPO^4$  outperformed  $+None^4$  for 19 of the 23 evaluation items. These results indicate that simple iterative updates without geometric correction in Algorithm 1 are insufficient to improve physical identification and that it is crucial to correct geometric structures using a dynamic optimization method.

**III. Evaluation of geometric recorection.** Finally, we discuss the scores for each scene in the geometry recorection (image quality after Algorithm 1 was applied) in detail. Table 8 summarizes these results. This is an extension of the results in Table 3. We observed tendencies similar to those for geometry correction. Specifically,  $+LPO^4$  outperformed not only the baselines (PAC-NeRF-3v and PAC-NeRF-3v<sup>†</sup>) but also the ablated and comparative models, including  $+LPO-F^4$ ,  $+LPO-P^4$ ,  $+GO^4$ , and  $+None^4$ , in most cases. These results indicate that joint optimization of the features and positions of particles in Lagrangian space is essential not only when geometric correction is conducted once but also when geometric correction is repeatedly conducted along with physical reidentification through Algorithm 1. Qualitative comparisons of PAC-NeRF-3v/3v<sup>†</sup>,  $+LPO$ , and  $+LPO^4$  are shown in Figure 3–7. Qualitative comparisons of  $+LPO^4$ ,  $+LPO-F^4$ ,  $+LPO-P^4$ ,  $+GO^4$ , and  $+None^4$  are presented in Figures 8 and 9.

## B.2. Computation times

Table 9 lists the computation times for executing Algorithm 1 for one iteration on Droplet with PAC-NeRF-3v+ $LPO$ .<sup>13</sup> The total computation time increases linearly when repeatedly running Algorithm 1. However, it is adjustable under a quality-and-time trade-off, as discussed in Appendix B.3. The computation time of  $LPO$  (3) is almost identical to that of the main process of physical property optimization (2-c) because the forward and backward processes are identical with different optimization targets, as shown in Figures 2(2) and (3). Similarly, the calculation times for  $+LPO-F$ ,  $+LPO-P$ , and  $+GO$  were almost identical

<sup>13</sup>Computation times vary with scenes because of the difference in number of frames and object size (which affects the number of particles); however, we observed similar tendencies.

| Process                                      | # Frames | Time (s) |
|--|----------|----------|
| (1) Eulerian static voxel grid optimization  | 1        | 284.6    |
| (2) Physical property optimization           |          |          |
| (a) Velocity optimization                    | 4        | 101.9    |
| (b) Warm-up optimization with partial frames | 7        | 765.2    |
| (c) Main optimization with entire frames     | 13       | 3153.3   |
| (3) Lagrangian particle optimization         | 13       | 3225.4   |
| (4) Color prediction                         | 1        | 1.0      |

Table 9. Computation times of Algorithm 1 on NVIDIA A100 GPU. (1) and (2) are processes driven from the original PAC-NeRF. (3) and (4) are processes newly introduced.

to those for  $+LPO$ , indicating that the performance improvement was attributable to the ingenuity of the algorithm and not to an increase in calculation cost.

## B.3. Impact of the number of iterations

In the main experiments, we fixed the number of iterations in Algorithm 1, that is,  $R$ , to four for simplicity and fair comparison. However, it is interesting and important to investigate how the number of iterations affects performance. To answer this question, we analyze the impact of the number of iterations in this appendix.

**Results.** Table 10 summarizes the performance changes with geometry (re)correction when the number of iterations was changed. This table calculates the scores using the images in the *test* set. Table 11 lists the scores calculated using images in the *training* set. We investigated these two cases to examine the trade-off between the feasible reconstruction of the training data and the generalization ability for the test data. Table 12 presents the performance changes with physical identification when the number of iterations is changed.

When PAC-NeRF-3v was used as the baseline, the performances for geometry (re)correction (for both the test and training sets) and physical identification were the best when the number of iterations was four in most cases. Specifically, as listed in Table 10, 6, 8, 13, and 18 times achieved the best scores among the 27 evaluation items when the number of iterations was one, two, three, and four, respectively, for geometry (re)correction for the test set. When scores were the same, they were counted multiple times. As listed in Table 11, 1, 3, 11, and 25 items achieved the best scores among the 27 evaluation items when the number of iterations was one, two, three, and four, respectively, for geometry (re)correction for the training set. As listed in Table 12, 2, 3, 3, and 16 items achieved the best scores among the 23 evaluation items when the numbers of iterations were one, two, three, and four for physical identification.

In contrast, when PAC-NeRF-3v<sup>†</sup> was used as the baseline, the performance of the geometry (re)correction for the test set was the best when the number of iterations was three, whereas that for the training set was the best when the number of iterations was four. The performance of the physical identification was the best when the number of iterations was three or four. Specifically, as listed in Ta-

|            |                    | PAC-NeRF-3v | +LPO <sup>1</sup> | +LPO <sup>2</sup> | +LPO <sup>3</sup> | +LPO <sup>4</sup> | PAC-NeRF-3v <sup>†</sup> | +LPO <sup>1</sup> | +LPO <sup>2</sup> | +LPO <sup>3</sup> | +LPO <sup>4</sup> |
|------------|--------------------|-------------|-------------------|-------------------|-------------------|-------------------|--------------------------|-------------------|-------------------|-------------------|-------------------|
| Droplet    | PSNR $\uparrow$    | 25.42       | 27.56             | 28.19             | 28.06             | 28.21             | 26.40                    | 28.18             | 28.95             | 29.14             | 28.91             |
|            | SSIM $\uparrow$    | 0.975       | 0.978             | 0.980             | 0.981             | 0.981             | 0.978                    | 0.980             | 0.982             | 0.983             | 0.983             |
|            | LPIPS $\downarrow$ | 0.047       | 0.043             | 0.043             | 0.042             | 0.042             | 0.046                    | 0.044             | 0.042             | 0.041             | 0.041             |
| Letter     | PSNR $\uparrow$    | 28.94       | 29.99             | 29.90             | 29.96             | 30.07             | 29.59                    | 30.44             | 30.76             | 30.70             | 30.26             |
|            | SSIM $\uparrow$    | 0.981       | 0.982             | 0.983             | 0.982             | 0.983             | 0.983                    | 0.982             | 0.983             | 0.984             | 0.983             |
|            | LPIPS $\downarrow$ | 0.028       | 0.029             | 0.027             | 0.031             | 0.026             | 0.028                    | 0.029             | 0.026             | 0.026             | 0.026             |
| Cream      | PSNR $\uparrow$    | 26.61       | 28.48             | 29.18             | 29.38             | 29.65             | 27.43                    | 28.82             | 29.69             | 29.93             | 30.23             |
|            | SSIM $\uparrow$    | 0.983       | 0.983             | 0.984             | 0.985             | 0.985             | 0.983                    | 0.984             | 0.985             | 0.985             | 0.985             |
|            | LPIPS $\downarrow$ | 0.026       | 0.023             | 0.022             | 0.023             | 0.022             | 0.024                    | 0.023             | 0.022             | 0.022             | 0.021             |
| Toothpaste | PSNR $\uparrow$    | 29.23       | 31.27             | 31.78             | 31.93             | 32.18             | 31.72                    | 33.54             | 34.16             | 34.04             | 34.00             |
|            | SSIM $\uparrow$    | 0.991       | 0.991             | 0.991             | 0.991             | 0.991             | 0.993                    | 0.993             | 0.993             | 0.992             | 0.992             |
|            | LPIPS $\downarrow$ | 0.010       | 0.010             | 0.011             | 0.011             | 0.011             | 0.010                    | 0.010             | 0.010             | 0.010             | 0.010             |
| Torus      | PSNR $\uparrow$    | 23.99       | 28.91             | 29.87             | 30.20             | 30.07             | 26.65                    | 30.05             | 30.90             | 30.97             | 30.48             |
|            | SSIM $\uparrow$    | 0.970       | 0.978             | 0.981             | 0.982             | 0.982             | 0.974                    | 0.980             | 0.983             | 0.984             | 0.983             |
|            | LPIPS $\downarrow$ | 0.048       | 0.038             | 0.034             | 0.032             | 0.032             | 0.040                    | 0.034             | 0.031             | 0.030             | 0.031             |
| Bird       | PSNR $\uparrow$    | 24.82       | 27.20             | 26.55             | 27.89             | 27.97             | 24.91                    | 28.51             | 28.88             | 28.85             | 28.98             |
|            | SSIM $\uparrow$    | 0.978       | 0.980             | 0.979             | 0.981             | 0.980             | 0.979                    | 0.983             | 0.983             | 0.982             | 0.982             |
|            | LPIPS $\downarrow$ | 0.036       | 0.032             | 0.035             | 0.033             | 0.034             | 0.036                    | 0.028             | 0.030             | 0.031             | 0.031             |
| Playdoh    | PSNR $\uparrow$    | 27.70       | 28.39             | 29.14             | 29.32             | 29.32             | 29.66                    | 30.07             | 30.16             | 30.26             | 30.01             |
|            | SSIM $\uparrow$    | 0.978       | 0.978             | 0.981             | 0.981             | 0.981             | 0.982                    | 0.982             | 0.983             | 0.983             | 0.983             |
|            | LPIPS $\downarrow$ | 0.042       | 0.041             | 0.040             | 0.040             | 0.040             | 0.038                    | 0.038             | 0.038             | 0.038             | 0.038             |
| Cat        | PSNR $\uparrow$    | 30.45       | 31.00             | 30.64             | 30.21             | 29.81             | 31.11                    | 31.41             | 31.25             | 31.29             | 30.61             |
|            | SSIM $\uparrow$    | 0.987       | 0.987             | 0.988             | 0.987             | 0.987             | 0.988                    | 0.988             | 0.989             | 0.989             | 0.988             |
|            | LPIPS $\downarrow$ | 0.033       | 0.032             | 0.027             | 0.027             | 0.028             | 0.034                    | 0.032             | 0.027             | 0.024             | 0.025             |
| Trophy     | PSNR $\uparrow$    | 29.33       | 30.16             | 29.99             | 29.78             | 29.57             | 28.80                    | 29.92             | 29.97             | 29.82             | 29.58             |
|            | SSIM $\uparrow$    | 0.963       | 0.964             | 0.964             | 0.964             | 0.964             | 0.962                    | 0.964             | 0.964             | 0.964             | 0.964             |
|            | LPIPS $\downarrow$ | 0.039       | 0.037             | 0.037             | 0.037             | 0.037             | 0.039                    | 0.037             | 0.037             | 0.036             | 0.036             |
| Average    | PSNR $\uparrow$    | 27.39       | 29.22             | 29.47             | 29.64             | 29.65             | 28.47                    | 30.11             | 30.52             | 30.56             | 30.34             |
|            | SSIM $\uparrow$    | 0.978       | 0.980             | 0.981             | 0.981             | 0.982             | 0.980                    | 0.982             | 0.983             | 0.983             | 0.983             |
|            | LPIPS $\downarrow$ | 0.034       | 0.032             | 0.031             | 0.031             | 0.030             | 0.033                    | 0.031             | 0.029             | 0.029             | 0.029             |

Table 10. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  on **geometric (re)correction** when the number of iterations in Algorithm 1 is changed. The scores were calculated using the images in the *test* set.

ble 10, 5, 10, 15, and 10 items achieved the best scores among the 27 evaluation items when the number of iterations was one, two, three, and four, respectively, for geometry (re)correction for the test set. As listed in Table 11, 0, 4, 14, and 26 items achieved the best scores among the 27 evaluation items when the number of iterations was one, two, three, and four, respectively, for geometry (re)correction for the training set. As listed in Table 12, 1, 4, 7, and 11 items achieved the best scores among the 23 evaluation items when the numbers of iterations were one, two, three, and four for physical identification.

We consider that these differences arise from differences in the initial static voxel grids. PAC-NeRF-3v<sup>†</sup> had better static voxel grids; therefore, fewer updates were required to obtain the best performance (i.e., optimal geometric structures). Note that it is generally not trivial to determine when to stop iterative updates because there is an intractable trade-off between faithful reproduction of training data and overfitting. When PAC-NeRF-3v was used as the baseline, the performance improvement continued until four iterations because the geometric reconstruction was more influential than overfitting. In contrast, when PAC-NeRF-3v<sup>†</sup> was used as the baseline, the performance improvement (particularly for geometric (re)correction for the test set) was saturated at three iterations because the perfor-

mance achieved the upper bound in an earlier phase. From this time on, the overfitting problem became non-negligible. Conducting further studies on this topic and exploring an improved method (for example, determining the number of iterations adaptively during training) will be interesting future research topics. Furthermore, the effect of the above-mentioned trade-off on the physical identification performance is an open issue. A detailed investigation of this will be an interesting topic for future research.

### C. Experiments on other view settings

In the main text (Section 4), we investigate the performance when three specific views are selected as training data. In this appendix, experiments were conducted using other settings to investigate the versatility of the proposed method. In particular, we investigated the robustness of the selection of views (Appendix C.1) and the number of views (Appendix C.2).

#### C.1. Robustness of the selection of views

In the main text (Section 4), we investigate the performance when three specific views are selected as training data. We examined the performance when the training set was changed to investigate the robustness of the selection of

|            |                    | PAC-NeRF-3v | +LPO <sup>1</sup> | +LPO <sup>2</sup> | +LPO <sup>3</sup> | +LPO <sup>4</sup> | PAC-NeRF-3v <sup>†</sup> | +LPO <sup>1</sup> | +LPO <sup>2</sup> | +LPO <sup>3</sup> | +LPO <sup>4</sup> |
|------------|--------------------|-------------|-------------------|-------------------|-------------------|-------------------|--------------------------|-------------------|-------------------|-------------------|-------------------|
| Droplet    | PSNR $\uparrow$    | 32.48       | 36.61             | 36.97             | 37.12             | 37.41             | 32.31                    | 36.09             | 37.21             | 37.83             | 37.96             |
|            | SSIM $\uparrow$    | 0.985       | 0.990             | 0.991             | 0.991             | 0.992             | 0.985                    | 0.990             | 0.992             | 0.992             | 0.992             |
|            | LPIPS $\downarrow$ | 0.034       | 0.029             | 0.029             | 0.029             | 0.028             | 0.036                    | 0.033             | 0.029             | 0.027             | 0.027             |
| Letter     | PSNR $\uparrow$    | 32.57       | 35.89             | 38.74             | 37.62             | 39.59             | 32.92                    | 36.84             | 38.81             | 39.09             | 39.90             |
|            | SSIM $\uparrow$    | 0.986       | 0.990             | 0.994             | 0.992             | 0.995             | 0.986                    | 0.992             | 0.994             | 0.995             | 0.995             |
|            | LPIPS $\downarrow$ | 0.020       | 0.020             | 0.015             | 0.019             | 0.013             | 0.021                    | 0.019             | 0.015             | 0.013             | 0.012             |
| Cream      | PSNR $\uparrow$    | 34.91       | 38.45             | 39.50             | 37.90             | 40.11             | 35.18                    | 38.50             | 38.70             | 39.60             | 40.39             |
|            | SSIM $\uparrow$    | 0.991       | 0.990             | 0.994             | 0.992             | 0.995             | 0.990                    | 0.993             | 0.994             | 0.995             | 0.995             |
|            | LPIPS $\downarrow$ | 0.016       | 0.014             | 0.013             | 0.014             | 0.013             | 0.017                    | 0.014             | 0.014             | 0.013             | 0.012             |
| Toothpaste | PSNR $\uparrow$    | 37.41       | 41.49             | 42.10             | 42.54             | 42.91             | 37.36                    | 41.43             | 42.56             | 43.07             | 43.41             |
|            | SSIM $\uparrow$    | 0.996       | 0.997             | 0.997             | 0.997             | 0.997             | 0.996                    | 0.997             | 0.997             | 0.997             | 0.998             |
|            | LPIPS $\downarrow$ | 0.006       | 0.005             | 0.005             | 0.005             | 0.004             | 0.006                    | 0.005             | 0.005             | 0.004             | 0.004             |
| Torus      | PSNR $\uparrow$    | 28.45       | 34.29             | 35.96             | 36.75             | 37.20             | 31.70                    | 35.30             | 36.66             | 37.22             | 37.56             |
|            | SSIM $\uparrow$    | 0.977       | 0.986             | 0.989             | 0.990             | 0.990             | 0.981                    | 0.987             | 0.990             | 0.990             | 0.990             |
|            | LPIPS $\downarrow$ | 0.041       | 0.031             | 0.028             | 0.027             | 0.027             | 0.034                    | 0.029             | 0.027             | 0.026             | 0.026             |
| Bird       | PSNR $\uparrow$    | 30.72       | 33.22             | 32.06             | 34.91             | 34.88             | 30.97                    | 33.84             | 34.73             | 35.47             | 35.91             |
|            | SSIM $\uparrow$    | 0.985       | 0.987             | 0.986             | 0.990             | 0.990             | 0.987                    | 0.989             | 0.990             | 0.990             | 0.991             |
|            | LPIPS $\downarrow$ | 0.029       | 0.027             | 0.030             | 0.025             | 0.025             | 0.028                    | 0.026             | 0.025             | 0.024             | 0.023             |
| Playdoh    | PSNR $\uparrow$    | 34.50       | 40.38             | 41.47             | 41.23             | 41.80             | 35.51                    | 40.92             | 41.22             | 41.93             | 41.74             |
|            | SSIM $\uparrow$    | 0.988       | 0.993             | 0.994             | 0.994             | 0.994             | 0.989                    | 0.993             | 0.994             | 0.994             | 0.994             |
|            | LPIPS $\downarrow$ | 0.027       | 0.022             | 0.022             | 0.021             | 0.021             | 0.027                    | 0.023             | 0.022             | 0.021             | 0.021             |
| Cat        | PSNR $\uparrow$    | 37.25       | 41.34             | 41.77             | 42.37             | 42.94             | 36.66                    | 40.43             | 42.17             | 43.32             | 43.49             |
|            | SSIM $\uparrow$    | 0.993       | 0.996             | 0.996             | 0.997             | 0.997             | 0.993                    | 0.995             | 0.996             | 0.997             | 0.997             |
|            | LPIPS $\downarrow$ | 0.020       | 0.018             | 0.017             | 0.016             | 0.016             | 0.026                    | 0.024             | 0.018             | 0.014             | 0.014             |
| Trophy     | PSNR $\uparrow$    | 32.35       | 34.76             | 34.82             | 35.01             | 35.13             | 31.86                    | 34.54             | 34.78             | 35.00             | 35.19             |
|            | SSIM $\uparrow$    | 0.966       | 0.969             | 0.969             | 0.970             | 0.970             | 0.966                    | 0.969             | 0.969             | 0.970             | 0.970             |
|            | LPIPS $\downarrow$ | 0.033       | 0.031             | 0.031             | 0.031             | 0.031             | 0.035                    | 0.032             | 0.032             | 0.031             | 0.031             |
| Average    | PSNR $\uparrow$    | 33.41       | 37.38             | 38.16             | 38.38             | 39.11             | 33.83                    | 37.54             | 38.54             | 39.17             | 39.51             |
|            | SSIM $\uparrow$    | 0.985       | 0.989             | 0.990             | 0.990             | 0.991             | 0.986                    | 0.990             | 0.991             | 0.991             | 0.991             |
|            | LPIPS $\downarrow$ | 0.025       | 0.022             | 0.021             | 0.021             | 0.020             | 0.026                    | 0.023             | 0.021             | 0.019             | 0.019             |

Table 11. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  on **geometric (re)correction** when the number of iterations in Algorithm 1 is changed. The scores were calculated using the images in the *training* set.

|            |                       | PAC-NeRF-3v | +LPO <sup>2</sup> | +LPO <sup>3</sup> | +LPO <sup>4</sup> | PAC-NeRF-3v <sup>†</sup> | +LPO <sup>2</sup> | +LPO <sup>3</sup> | +LPO <sup>4</sup> |
|------------|-----------------------|-------------|-------------------|-------------------|-------------------|--------------------------|-------------------|-------------------|-------------------|
| Droplet    | $\log_{10}(\mu)$      | 0.140       | 0.136             | 0.111             | 0.112             | 0.136                    | 0.129             | 0.084             | 0.082             |
|            | $\log_{10}(\kappa)$   | 1.285       | 1.447             | 1.784             | 1.628             | 1.263                    | 1.053             | 0.097             | 0.106             |
| Letter     | $\log_{10}(\mu)$      | 0.674       | 0.023             | 1.264             | 0.015             | 0.379                    | 0.013             | 0.053             | 0.010             |
|            | $\log_{10}(\kappa)$   | 6.772       | 0.325             | 1.424             | 0.174             | 5.229                    | 0.507             | 0.589             | 0.060             |
| Cream      | $\log_{10}(\mu)$      | 0.311       | 0.200             | 0.115             | 0.178             | 0.179                    | 0.031             | 0.066             | 0.100             |
|            | $\log_{10}(\kappa)$   | 0.215       | 0.384             | 0.392             | 0.158             | 0.336                    | 0.157             | 0.060             | 0.121             |
|            | $\log_{10}(\tau_Y)$   | 0.014       | 0.032             | 0.031             | 0.004             | 0.009                    | 0.007             | 0.005             | 0.006             |
|            | $\log_{10}(\eta)$     | 0.281       | 0.209             | 0.198             | 0.183             | 0.195                    | 0.105             | 0.026             | 0.033             |
| Toothpaste | $\log_{10}(\mu)$      | 1.891       | 0.264             | 0.246             | 0.109             | 0.252                    | 0.259             | 0.026             | 0.031             |
|            | $\log_{10}(\kappa)$   | 1.580       | 1.439             | 1.434             | 0.601             | 1.436                    | 1.382             | 0.888             | 0.673             |
|            | $\log_{10}(\tau_Y)$   | 0.201       | 0.118             | 0.116             | 0.114             | 0.199                    | 0.168             | 0.137             | 0.093             |
|            | $\log_{10}(\eta)$     | 0.373       | 0.200             | 0.180             | 0.003             | 0.212                    | 0.187             | 0.005             | 0.009             |
| Torus      | $\log_{10}(E)$        | 0.277       | 0.055             | 0.053             | 0.061             | 0.074                    | 0.049             | 0.040             | 0.036             |
|            | $\nu$                 | 0.085       | 0.129             | 0.050             | 0.001             | 0.131                    | 0.032             | 0.060             | 0.007             |
| Bird       | $\log_{10}(E)$        | 0.449       | 0.136             | 0.146             | 0.067             | 0.123                    | 0.158             | 0.084             | 0.027             |
|            | $\nu$                 | 0.102       | 0.466             | 0.037             | 0.001             | 0.141                    | 0.009             | 0.056             | 0.047             |
| Playdoh    | $\log_{10}(E)$        | 0.290       | 0.403             | 0.190             | 0.116             | 0.521                    | 0.303             | 0.260             | 0.133             |
|            | $\log_{10}(\tau_Y)$   | 0.283       | 0.075             | 0.215             | 0.165             | 0.110                    | 0.119             | 0.105             | 0.173             |
|            | $\nu$                 | 0.495       | 0.092             | 0.236             | 0.111             | 0.212                    | 0.218             | 0.082             | 0.063             |
| Cat        | $\log_{10}(E)$        | 1.301       | 1.120             | 1.051             | 0.973             | 1.192                    | 0.584             | 0.687             | 0.706             |
|            | $\log_{10}(\tau_Y)$   | 0.120       | 0.119             | 0.111             | 0.107             | 0.084                    | 0.062             | 0.070             | 0.067             |
|            | $\nu$                 | 0.044       | 0.079             | 0.104             | 0.004             | 0.118                    | 0.054             | 0.024             | 0.003             |
| Trophy     | $\theta_{fric}$ [rad] | 0.053       | 0.058             | 0.053             | 0.055             | 0.030                    | 0.036             | 0.046             | 0.039             |

Table 12. Comparison of the absolute differences between the ground-truth and estimated physical properties on **physical identification** when the number of iterations in Algorithm 1 is changed. The smaller the values, the better the performance.

|                    | PAC-NeRF-3v <sup>†</sup> | +LPO       | +LPO-F     | +LPO-P     | +GO        | +LPO <sup>4</sup> | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup> | +None <sup>4</sup> |
|--------------------|--------------------------|------------|------------|------------|------------|-------------------|---------------------|---------------------|------------------|--------------------|
| PSNR <sup>↑</sup>  | 28.18±.43                | 29.60±.59  | 28.89±.45  | 29.40±.55  | 28.91±.45  | 29.88±.53         | 29.34±.51           | 29.81±.54           | 28.79±.44        | 28.49±.45          |
| SSIM <sup>↑</sup>  | 0.979±.001               | 0.981±.001 | 0.980±.001 | 0.981±.001 | 0.980±.001 | 0.982±.001        | 0.981±.001          | 0.982±.001          | 0.981±.001       | 0.981±.001         |
| LPIPS <sup>↓</sup> | 0.036±.002               | 0.034±.003 | 0.035±.002 | 0.035±.003 | 0.035±.002 | 0.031±.002        | 0.032±.002          | 0.033±.003          | 0.033±.002       | 0.034±.002         |

Table 13. Comparison of PSNR<sup>↑</sup>, SSIM<sup>↑</sup>, and LPIPS<sup>↓</sup> averaged over five view settings on **geometric correction** and **recorection**.

views. Specifically, we investigated performance when five different view sets were used for training. The number of views used for training was fixed at three, and the remaining eight were used for testing.

**Compared models.** We used *PAC-NeRF-3v<sup>†</sup>* as the baseline and applied +LPO and +LPO<sup>4</sup>. Furthermore, we examined the performances of the ablated and comparative models. Specifically, in the evaluation of the geometry correction, we examined the performance when +LPO-F, +LPO-P, and +GO were applied to the baseline. In the evaluation of the physical identification and geometric recorection, we examined the performance when +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup> were applied to the baseline.

**Results.** Table 13 summarizes the geometric correction and recorection results. We observed tendencies similar to those for the results when three specific views were used for training (see Appendix B.1). Specifically, with respect to geometry correction, +LPO outperformed not only the baseline (PAC-NeRF-3v<sup>†</sup>) but was also superior to or comparable to the ablated and comparative models, including +LPO-F, +LPO-P, and +GO. Similarly, with respect to geometry recorection, +LPO<sup>4</sup> outperformed not only the baseline (PAC-NeRF-3v<sup>†</sup>) but was also superior to or comparable to the ablated and comparative models, including +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup>. In terms of physical property identification, +LPO<sup>4</sup> outperformed PAC-NeRF-3v<sup>†</sup>, +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup> in the 20.8±1.2, 15.2±2.0, 16.6±2.0, 17.8±1.2, and 18.8±1.9 cases, respectively, across the 23 properties for each view setting. These results indicate that the joint optimization of features and positions in Lagrangian space is effective for geometry-agnostic system identification, regardless of the selection of views used for training.

## C.2. Robustness of the number of views

In the main experiment (Section 4) and the above experiments (Appendix C.1), we investigated the performance when the number of views in the training set was three. To investigate the robustness of the number of views, we examined the performance when the number of views in the training set was changed to six. Specifically, among the 11 views in the dataset, six were used for training, and the remaining five were used for testing.

**Compared models.** We used *PAC-NeRF-6v*, that is, PAC-NeRF [42] trained with six views, as the baseline. In preliminary experiments, we found that PAC-NeRF-6v demonstrated reasonable performance without the advanced techniques described in Appendix A owing to the increase in

the number of views. Therefore, we only used this model as the baseline. We applied +LPO and +LPO<sup>4</sup> to the baseline. Furthermore, we investigated the performances of the ablated and comparative models. Specifically, in the evaluation of the geometry correction, we investigated the performance when +LPO-F, +LPO-P, and +GO were applied to the baseline. In the evaluation of the physical identification and geometric recorection, we investigated the performance when +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup> were applied to the baseline.

**Results.** Tables 14, 15, and 16 summarize the results of geometry correction, physical identification, and geometry reconstruction, respectively. Similar tendencies were observed when the number of views in the training set was three (as discussed in Section 4 and Appendix C.1). Specifically, with respect to geometry correction, +LPO outperformed not only the baseline (PAC-NeRF-6v) but also the ablated and comparative models, including +LPO-F, +LPO-P, and +GO, in most cases. On physical identification and geometry recorection, +LPO<sup>4</sup> outperformed not only the baseline (PAC-NeRF-6v) but also the ablated and comparative models, including +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup>, in most cases. These results indicate that the joint optimization of features and positions in Lagrangian space is effective for geometry-agnostic system identification, regardless of the number of views used for training.

|            |                    | PAC-NeRF | PAC-NeRF-6v | +LPO  | +LPO-F | +LPO-P | +GO   |
|------------|--------------------|----------|-------------|-------|--------|--------|-------|
| Droplet    | PSNR $\uparrow$    | 35.29    | 28.92       | 30.58 | 29.82  | 30.47  | 29.13 |
|            | SSIM $\uparrow$    | 0.989    | 0.982       | 0.984 | 0.983  | 0.984  | 0.983 |
|            | LPIPS $\downarrow$ | 0.030    | 0.042       | 0.040 | 0.041  | 0.040  | 0.042 |
| Letter     | PSNR $\uparrow$    | 36.99    | 31.42       | 32.85 | 32.36  | 32.62  | 32.32 |
|            | SSIM $\uparrow$    | 0.992    | 0.985       | 0.987 | 0.987  | 0.986  | 0.986 |
|            | LPIPS $\downarrow$ | 0.011    | 0.020       | 0.019 | 0.019  | 0.019  | 0.020 |
| Cream      | PSNR $\uparrow$    | 36.46    | 30.36       | 31.31 | 31.13  | 30.89  | 30.80 |
|            | SSIM $\uparrow$    | 0.993    | 0.986       | 0.987 | 0.987  | 0.986  | 0.987 |
|            | LPIPS $\downarrow$ | 0.014    | 0.021       | 0.019 | 0.019  | 0.020  | 0.019 |
| Toothpaste | PSNR $\uparrow$    | 38.84    | 34.74       | 35.82 | 35.27  | 35.58  | 34.79 |
|            | SSIM $\uparrow$    | 0.996    | 0.993       | 0.994 | 0.993  | 0.993  | 0.993 |
|            | LPIPS $\downarrow$ | 0.006    | 0.009       | 0.009 | 0.009  | 0.009  | 0.009 |
| Torus      | PSNR $\uparrow$    | 34.60    | 29.28       | 32.31 | 30.82  | 31.80  | 30.80 |
|            | SSIM $\uparrow$    | 0.988    | 0.979       | 0.984 | 0.981  | 0.984  | 0.982 |
|            | LPIPS $\downarrow$ | 0.026    | 0.035       | 0.030 | 0.033  | 0.030  | 0.034 |
| Bird       | PSNR $\uparrow$    | 35.70    | 27.38       | 30.48 | 28.67  | 30.39  | 28.09 |
|            | SSIM $\uparrow$    | 0.992    | 0.981       | 0.984 | 0.982  | 0.984  | 0.981 |
|            | LPIPS $\downarrow$ | 0.019    | 0.029       | 0.025 | 0.027  | 0.025  | 0.028 |
| Playdoh    | PSNR $\uparrow$    | 36.55    | 29.74       | 31.01 | 29.93  | 30.99  | 29.84 |
|            | SSIM $\uparrow$    | 0.991    | 0.982       | 0.983 | 0.982  | 0.983  | 0.982 |
|            | LPIPS $\downarrow$ | 0.026    | 0.039       | 0.037 | 0.039  | 0.037  | 0.039 |
| Cat        | PSNR $\uparrow$    | 37.10    | 30.77       | 31.40 | 30.85  | 31.33  | 30.72 |
|            | SSIM $\uparrow$    | 0.993    | 0.989       | 0.987 | 0.989  | 0.988  | 0.988 |
|            | LPIPS $\downarrow$ | 0.016    | 0.024       | 0.026 | 0.025  | 0.026  | 0.025 |
| Trophy     | PSNR $\uparrow$    | 32.23    | 30.22       | 30.97 | 30.38  | 30.85  | 30.23 |
|            | SSIM $\uparrow$    | 0.965    | 0.962       | 0.964 | 0.963  | 0.963  | 0.962 |
|            | LPIPS $\downarrow$ | 0.036    | 0.039       | 0.037 | 0.039  | 0.038  | 0.039 |
| Average    | PSNR $\uparrow$    | 35.97    | 30.31       | 31.86 | 31.03  | 31.66  | 30.75 |
|            | SSIM $\uparrow$    | 0.989    | 0.982       | 0.984 | 0.983  | 0.984  | 0.983 |
|            | LPIPS $\downarrow$ | 0.020    | 0.029       | 0.027 | 0.028  | 0.027  | 0.028 |

Table 14. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  for each scene on **geometric correction** when the number of views in a training set was six.

|            |                       | PAC-NeRF | PAC-NeRF-6v | +LPO <sup>4</sup> | +LPO-F <sup>4</sup> | +LPO-P <sup>4</sup> | +GO <sup>4</sup> | +None <sup>4</sup> |
|------------|-----------------------|----------|-------------|-------------------|---------------------|---------------------|------------------|--------------------|
| Droplet    | $\log_{10}(\mu)$      | 0.039    | 0.026       | 0.016             | 0.039               | 0.025               | 0.042            | 0.039              |
|            | $\log_{10}(\kappa)$   | 0.017    | 0.386       | 0.045             | 0.128               | 0.076               | 0.404            | 0.144              |
| Letter     | $\log_{10}(\mu)$      | 0.041    | 0.229       | 0.175             | 0.220               | 0.199               | 0.222            | 0.224              |
|            | $\log_{10}(\kappa)$   | 0.039    | 0.166       | 0.006             | 0.021               | 0.060               | 0.064            | 0.089              |
| Cream      | $\log_{10}(\mu)$      | 0.090    | 0.131       | 0.056             | 0.047               | 0.065               | 0.057            | 0.065              |
|            | $\log_{10}(\kappa)$   | 0.132    | 0.835       | 0.146             | 0.020               | 0.143               | 0.072            | 0.132              |
|            | $\log_{10}(\tau_Y)$   | 0.007    | 0.011       | 0.001             | 0.003               | 0.001               | 0.003            | 0.012              |
|            | $\log_{10}(\eta)$     | 0.015    | 0.287       | 0.049             | 0.060               | 0.075               | 0.088            | 0.134              |
| Toothpaste | $\log_{10}(\mu)$      | 0.026    | 0.084       | 0.045             | 0.179               | 0.144               | 0.208            | 0.166              |
|            | $\log_{10}(\kappa)$   | 0.247    | 0.629       | 0.580             | 0.370               | 0.655               | 0.654            | 0.573              |
|            | $\log_{10}(\tau_Y)$   | 0.066    | 0.088       | 0.036             | 0.090               | 0.062               | 0.101            | 0.097              |
|            | $\log_{10}(\eta)$     | 0.013    | 0.015       | 0.006             | 0.022               | 0.006               | 0.018            | 0.018              |
| Torus      | $\log_{10}(E)$        | 0.019    | 0.020       | 0.023             | 0.019               | 0.025               | 0.032            | 0.025              |
|            | $\nu$                 | 0.023    | 0.024       | 0.001             | 0.057               | 0.035               | 0.139            | 0.058              |
| Bird       | $\log_{10}(E)$        | 0.013    | 0.106       | 0.047             | 0.105               | 0.058               | 0.119            | 1.769              |
|            | $\nu$                 | 0.029    | 0.133       | 0.013             | 0.134               | 0.147               | 0.146            | 0.741              |
| Playdoh    | $\log_{10}(E)$        | 0.286    | 0.483       | 0.252             | 0.311               | 0.282               | 0.310            | 0.289              |
|            | $\log_{10}(\tau_Y)$   | 0.038    | 0.209       | 0.126             | 0.203               | 0.171               | 0.184            | 0.180              |
|            | $\nu$                 | 0.076    | 0.317       | 0.110             | 0.135               | 0.141               | 0.130            | 0.193              |
| Cat        | $\log_{10}(E)$        | 0.855    | 2.821       | 0.787             | 0.949               | 0.843               | 2.423            | 0.745              |
|            | $\log_{10}(\tau_Y)$   | 0.026    | 0.978       | 0.008             | 0.005               | 0.024               | 0.393            | 0.020              |
|            | $\nu$                 | 0.027    | 0.071       | 0.017             | 0.065               | 0.026               | 0.218            | 0.026              |
| Trophy     | $\theta_{fric}$ [rad] | 0.048    | 0.052       | 0.035             | 0.040               | 0.041               | 0.044            | 0.044              |

Table 15. Comparison of the absolute differences between the ground-truth and estimated physical properties for each scene on **physical identification** when the number of views in a training set was six. The smaller the values, the better the performance.

|            |                    | PAC-NeRF | PAC-NeRF-6v | +LPO <sup>d</sup> | +LPO-F <sup>d</sup> | +LPO-P <sup>d</sup> | +GO <sup>d</sup> | +None <sup>d</sup> |
|------------|--------------------|----------|-------------|-------------------|---------------------|---------------------|------------------|--------------------|
| Droplet    | PSNR $\uparrow$    | 35.29    | 28.92       | 30.72             | 29.64               | 31.26               | 27.66            | 28.73              |
|            | SSIM $\uparrow$    | 0.989    | 0.982       | 0.986             | 0.984               | 0.986               | 0.981            | 0.983              |
|            | LPIPS $\downarrow$ | 0.030    | 0.042       | 0.037             | 0.041               | 0.039               | 0.044            | 0.042              |
| Letter     | PSNR $\uparrow$    | 36.99    | 31.42       | 32.70             | 32.69               | 32.66               | 32.13            | 32.25              |
|            | SSIM $\uparrow$    | 0.992    | 0.985       | 0.989             | 0.989               | 0.988               | 0.987            | 0.987              |
|            | LPIPS $\downarrow$ | 0.011    | 0.020       | 0.017             | 0.016               | 0.018               | 0.018            | 0.018              |
| Cream      | PSNR $\uparrow$    | 36.46    | 30.36       | 32.11             | 31.78               | 31.48               | 30.65            | 30.34              |
|            | SSIM $\uparrow$    | 0.993    | 0.986       | 0.988             | 0.988               | 0.987               | 0.987            | 0.987              |
|            | LPIPS $\downarrow$ | 0.014    | 0.021       | 0.018             | 0.019               | 0.020               | 0.019            | 0.020              |
| Toothpaste | PSNR $\uparrow$    | 38.84    | 34.74       | 34.91             | 34.44               | 35.35               | 33.52            | 34.84              |
|            | SSIM $\uparrow$    | 0.996    | 0.993       | 0.994             | 0.994               | 0.994               | 0.993            | 0.994              |
|            | LPIPS $\downarrow$ | 0.006    | 0.009       | 0.009             | 0.009               | 0.009               | 0.009            | 0.009              |
| Torus      | PSNR $\uparrow$    | 34.60    | 29.28       | 32.33             | 31.57               | 31.94               | 30.64            | 29.58              |
|            | SSIM $\uparrow$    | 0.988    | 0.979       | 0.986             | 0.984               | 0.986               | 0.982            | 0.982              |
|            | LPIPS $\downarrow$ | 0.026    | 0.035       | 0.028             | 0.031               | 0.030               | 0.032            | 0.034              |
| Bird       | PSNR $\uparrow$    | 35.70    | 27.38       | 30.57             | 28.27               | 29.67               | 27.84            | 25.96              |
|            | SSIM $\uparrow$    | 0.992    | 0.981       | 0.985             | 0.981               | 0.983               | 0.981            | 0.983              |
|            | LPIPS $\downarrow$ | 0.019    | 0.029       | 0.028             | 0.029               | 0.028               | 0.030            | 0.035              |
| Playdoh    | PSNR $\uparrow$    | 36.55    | 29.74       | 30.26             | 28.83               | 30.90               | 29.03            | 29.65              |
|            | SSIM $\uparrow$    | 0.991    | 0.982       | 0.984             | 0.981               | 0.984               | 0.982            | 0.984              |
|            | LPIPS $\downarrow$ | 0.026    | 0.039       | 0.037             | 0.041               | 0.037               | 0.040            | 0.039              |
| Cat        | PSNR $\uparrow$    | 37.10    | 30.77       | 31.45             | 31.14               | 31.11               | 29.95            | 31.79              |
|            | SSIM $\uparrow$    | 0.993    | 0.989       | 0.988             | 0.989               | 0.988               | 0.987            | 0.989              |
|            | LPIPS $\downarrow$ | 0.016    | 0.024       | 0.024             | 0.024               | 0.026               | 0.025            | 0.023              |
| Trophy     | PSNR $\uparrow$    | 32.23    | 30.22       | 30.44             | 29.69               | 30.64               | 29.20            | 30.23              |
|            | SSIM $\uparrow$    | 0.965    | 0.962       | 0.963             | 0.962               | 0.963               | 0.962            | 0.963              |
|            | LPIPS $\downarrow$ | 0.036    | 0.039       | 0.037             | 0.039               | 0.038               | 0.040            | 0.039              |
| Average    | PSNR $\uparrow$    | 35.97    | 30.31       | 31.72             | 30.89               | 31.67               | 30.07            | 30.37              |
|            | SSIM $\uparrow$    | 0.989    | 0.982       | 0.985             | 0.983               | 0.984               | 0.982            | 0.983              |
|            | LPIPS $\downarrow$ | 0.020    | 0.029       | 0.026             | 0.028               | 0.027               | 0.029            | 0.029              |

Table 16. Comparison of PSNR $\uparrow$ , SSIM $\uparrow$ , and LPIPS $\downarrow$  for each scene on **geometric recorection** when the number of views in a training set was six.

## D. Qualitative results

This appendix discusses the qualitative results obtained using several frames selected from the video sequences. We provide video samples at the [project page](#).<sup>1</sup>

### D.1. Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup>

- **Figure 3:**  
Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on Newtonian fluids (Droplet and Letter).
- **Figure 4:**  
Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on non-Newtonian fluids (Cream and Toothpaste).
- **Figure 5:**  
Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on elastic materials (Torus and Bird).
- **Figure 6:**  
Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on plasticine (Playdoh and Cat).
- **Figure 7:**  
Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on granular media (Trophy).

### D.2. Qualitative comparisons among +LPO<sup>4</sup>, +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup>

- **Figure 8:**  
Qualitative comparisons among +LPO<sup>4</sup>, +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup> in the scenes where +LPO-P<sup>4</sup> outperformed +LPO-F<sup>4</sup>.
- **Figure 9:**  
Qualitative comparisons among +LPO<sup>4</sup>, +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup> in the scenes where +LPO-F<sup>4</sup> outperformed +LPO-P<sup>4</sup>.

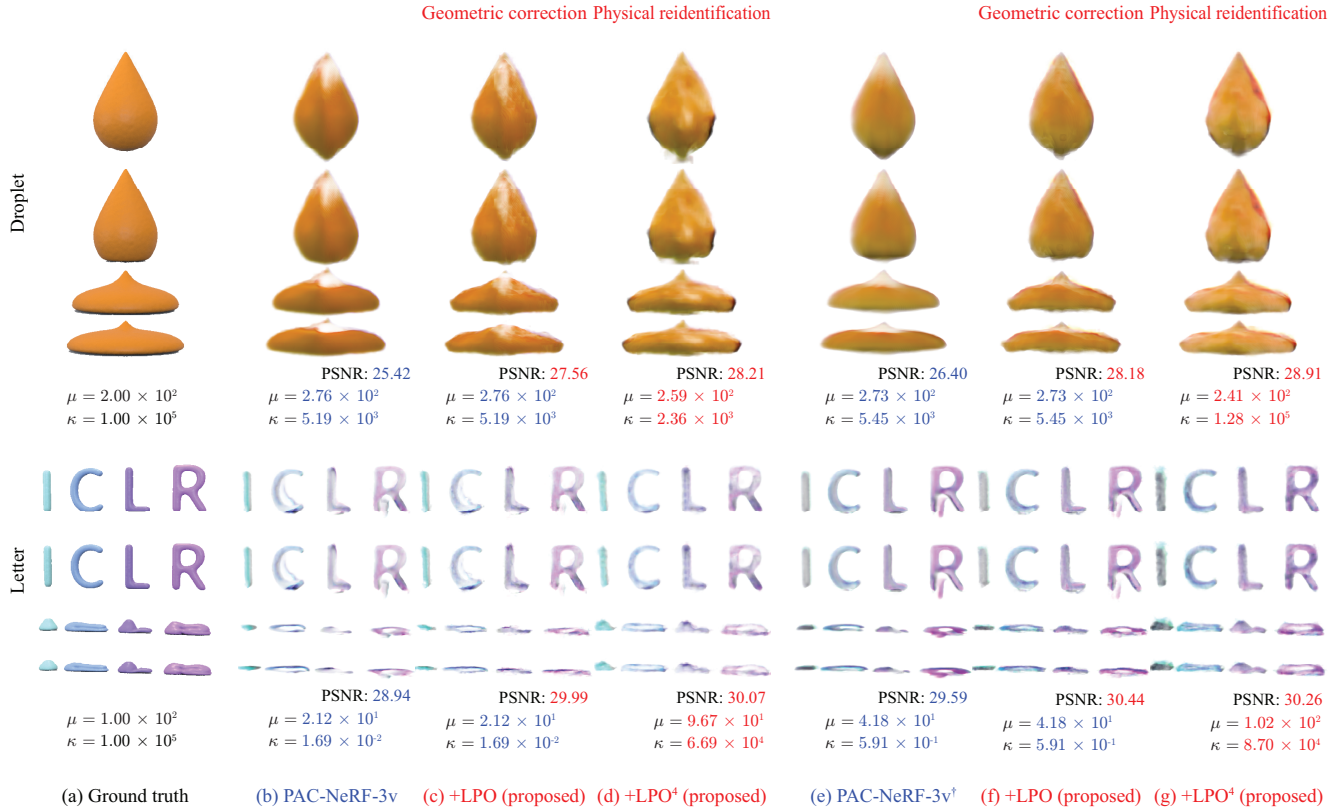


Figure 3. Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on Newtonian fluids (Droplet and Letter). Blue fonts indicate the scores obtained by the baselines (PAC-NeRF-3v/3v<sup>†</sup>). Red fonts indicate the scores obtained by the proposed methods (+LPO and +LPO<sup>4</sup>). Given the initial estimation by the baseline (b)(e), +LPO first corrects the geometric structures (including appearance and shape) (c)(f). By repeatedly conducting physical identification and geometric correction via Algorithm 1, +LPO<sup>4</sup> reidentifies physical properties and recorrects geometric structures (d)(g). In the Droplet scene, the bottom of the droplet is sharply pointed, and its tip is whitened in the baseline (b)(e). They are gradually mitigated by applying +LPO (c)(f) and +LPO<sup>4</sup> (d)(g). In the Letter scene, +LPO (c)(f) and +LPO<sup>4</sup> (d)(g) succeed in gradually eliminating artifacts existing in the vicinity of the left line of the letter “R.”, which arise in the baselines (b)(e).



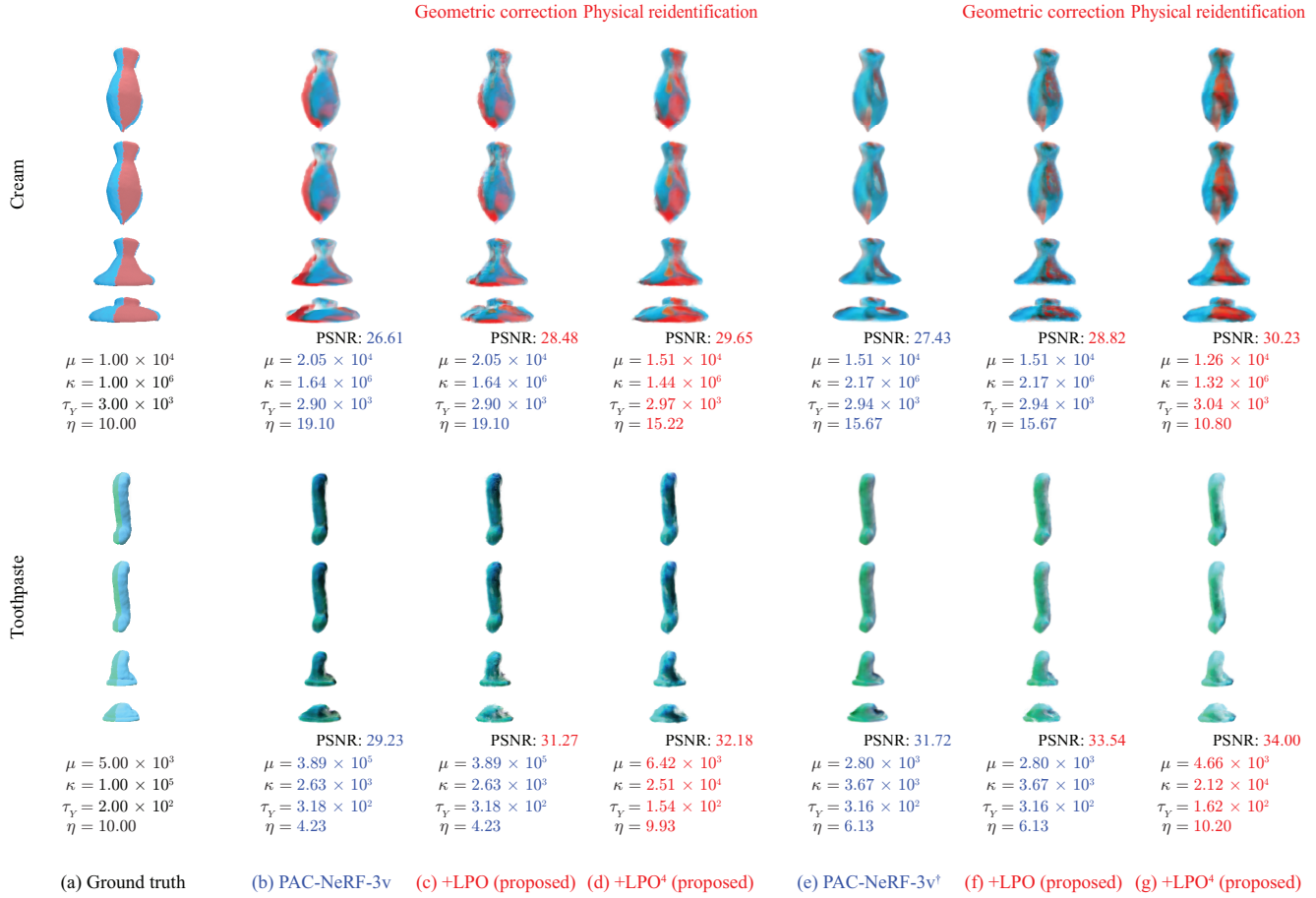


Figure 4. Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on non-Newtonian fluids (Cream and Toothpaste). Blue fonts indicate the scores obtained by the baselines (PAC-NeRF-3v/3v<sup>†</sup>). Red fonts indicate the scores obtained by the proposed methods (+LPO and +LPO<sup>4</sup>). Given the initial estimation by the baseline (b)(e), +LPO first corrects the geometric structures (including appearance and shape) (c)(f). By repeatedly conducting physical identification and geometric correction via Algorithm 1, +LPO<sup>4</sup> reidentifies physical properties and recorrects geometric structures (d)(g). In the Cream scene, the baselines (b)(e) fail to color the materials correctly. This failure is alleviated by +LPO (c)(f) and further mitigated by +LPO<sup>4</sup> (d)(g). In the Toothpaste scene, the baselines (b)(e) make the material darker color than that in the ground truth (a). +LPO (c)(f) makes the material brighter, and +LPO<sup>4</sup> (d)(g) obtains the color closer to the ground truth (a). This effect is pronounced when PAC-NeRF-3v<sup>†</sup> (e) is used as a baseline.

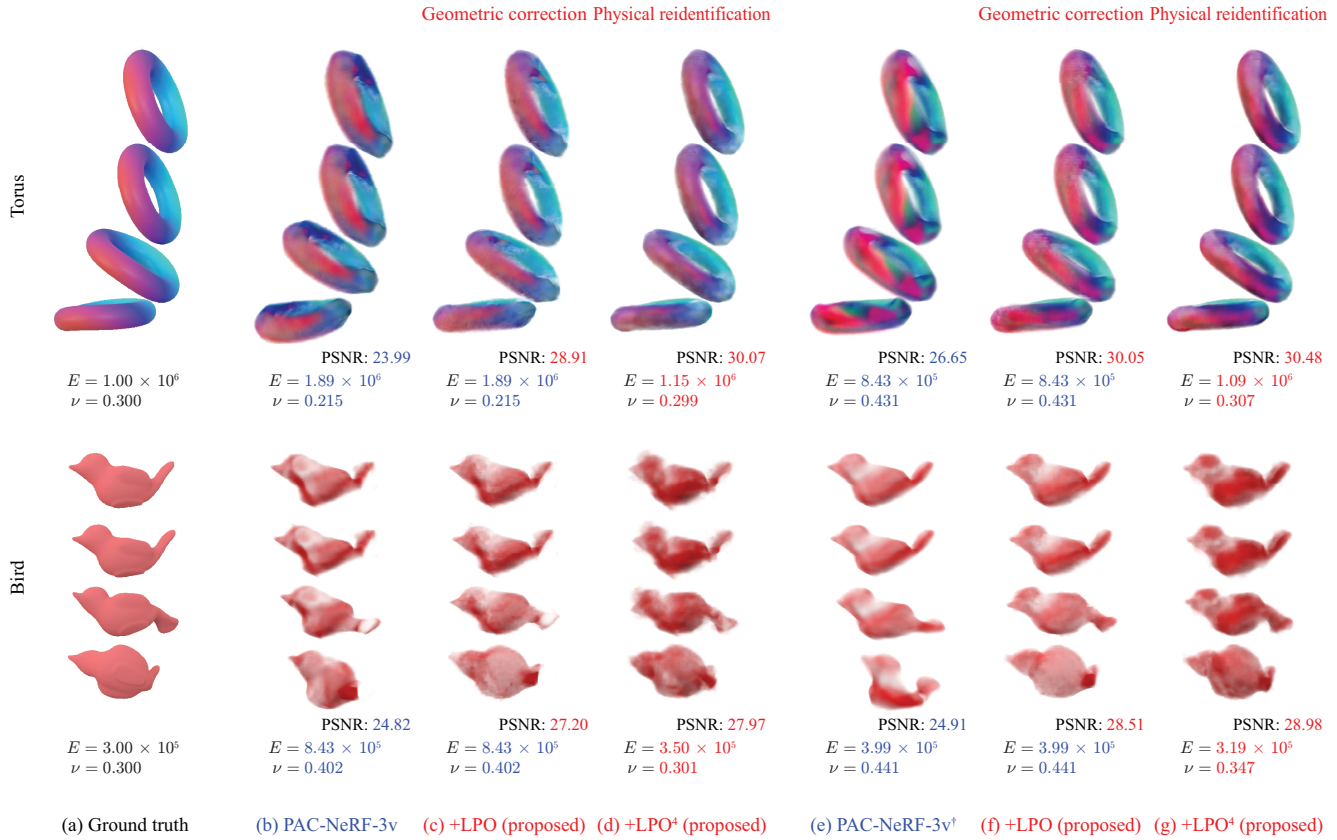


Figure 5. Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on elastic materials (Torus and Bird). Blue fonts indicate the scores obtained by the baselines (PAC-NeRF-3v/3v<sup>†</sup>). Red fonts indicate the scores obtained by the proposed methods (+LPO and +LPO<sup>4</sup>). Given the initial estimation by the baseline (b)(e), +LPO first corrects the geometric structures (including appearance and shape) (c)(f). By repeatedly conducting physical identification and geometric correction via Algorithm 1, +LPO<sup>4</sup> reidentifies physical properties and recorrects geometric structures (d)(g). In the Torus scene, the baselines (b)(e) have difficulty correctly capturing color and shape. They are improved by applying +LPO (c)(f), and the fine details are also improved by utilizing +LPO<sup>4</sup> (d)(g). Also, in the Bird scene, the baselines (b)(e) fail to capture color and shape correctly. The shape (e.g., the directions of the tail) is first corrected by +LPO (c)(f), and then the color is corrected by +LPO<sup>4</sup> (d)(g).

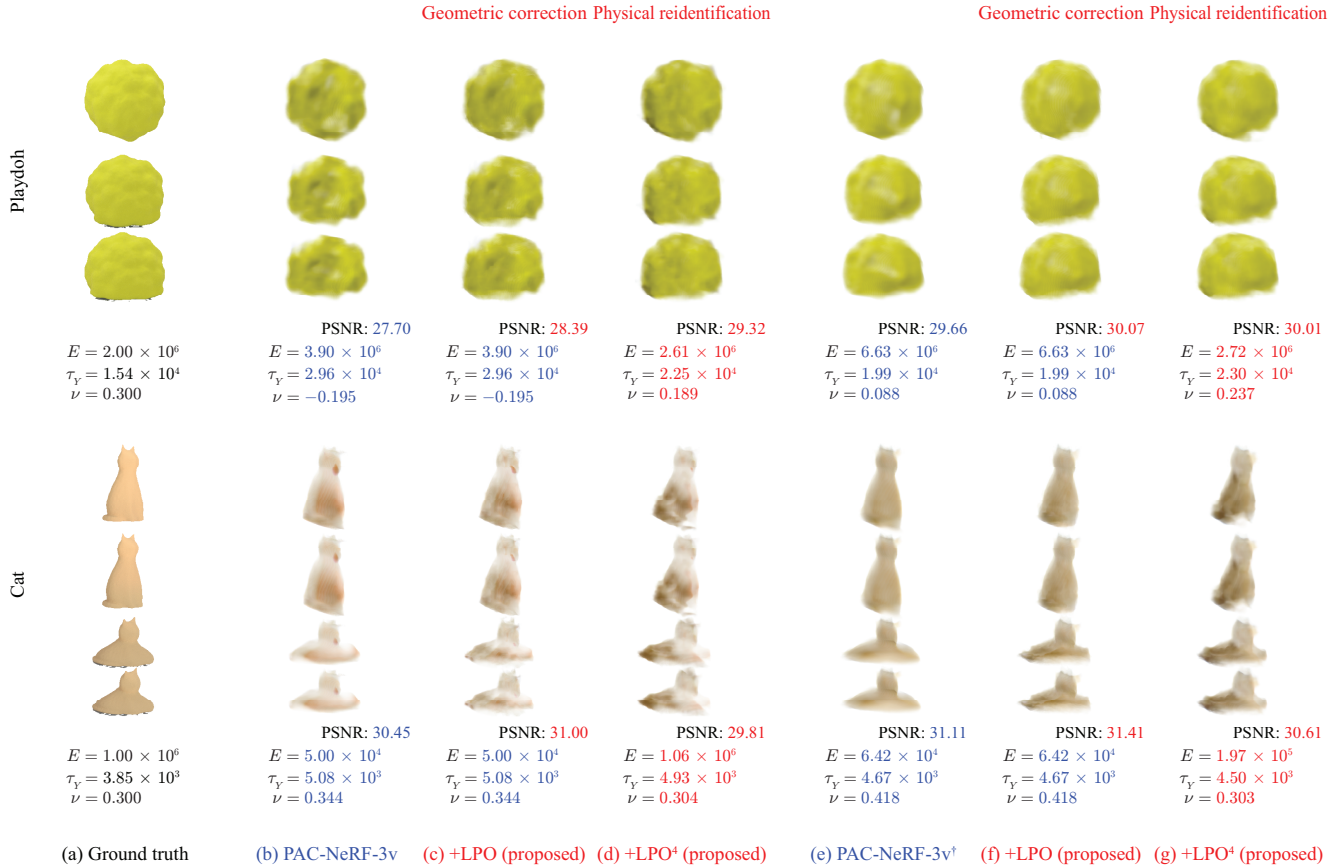


Figure 6. Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on plasticine (Playdoh and Cat). Blue fonts indicate the scores obtained by the baselines (PAC-NeRF-3v/3v<sup>†</sup>). Red fonts indicate the scores obtained by the proposed methods (+LPO and +LPO<sup>4</sup>). Given the initial estimation by the baseline (b)(e), +LPO first corrects the geometric structures (including appearance and shape) (c)(f). By repeatedly conducting physical identification and geometric correction via Algorithm 1, +LPO<sup>4</sup> reidentifies physical properties and recorrects geometric structures (d)(g). In the Playdoh scene, the baselines (b)(e) make the playdoh a little crushed compared to the ground truth (a). These geometric failures are gradually alleviated by applying +LPO (c)(f) and +LPO<sup>4</sup> (d)(g). In the Cat scene, the baselines (b)(e) fail to capture the tail of the cat in the lower left corner. +LPO (c)(f) and +LPO<sup>4</sup> (d)(g) struggle to recover it.

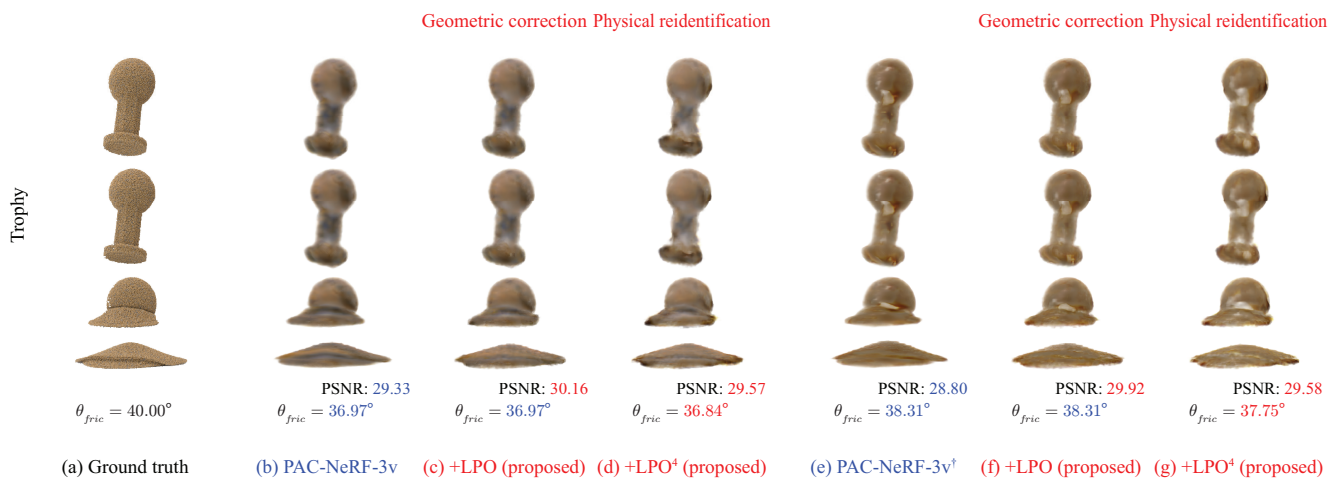


Figure 7. Qualitative comparisons among PAC-NeRF-3v/3v<sup>†</sup>, +LPO, and +LPO<sup>4</sup> on granular media (Trophy). Blue fonts indicate the scores obtained by the baselines (PAC-NeRF-3v/3v<sup>†</sup>). Red fonts indicate the scores obtained by the proposed methods (+LPO and +LPO<sup>4</sup>). Given the initial estimation by the baseline (b)(e), +LPO first corrects the geometric structures (including appearance and shape) (c)(f). By repeatedly conducting physical identification and geometric correction via Algorithm 1, +LPO<sup>4</sup> reidentifies physical properties and recorrects geometric structures (d)(g). In the Trophy scene, the baselines (b)(e) achieve good performance in terms of physical identification (comparable with PAC-NeRF with full-view supervision). However, they tend to make the trophy darker than the ground truth (a). This misestimation is corrected by +LPO and +LPO<sup>4</sup>.

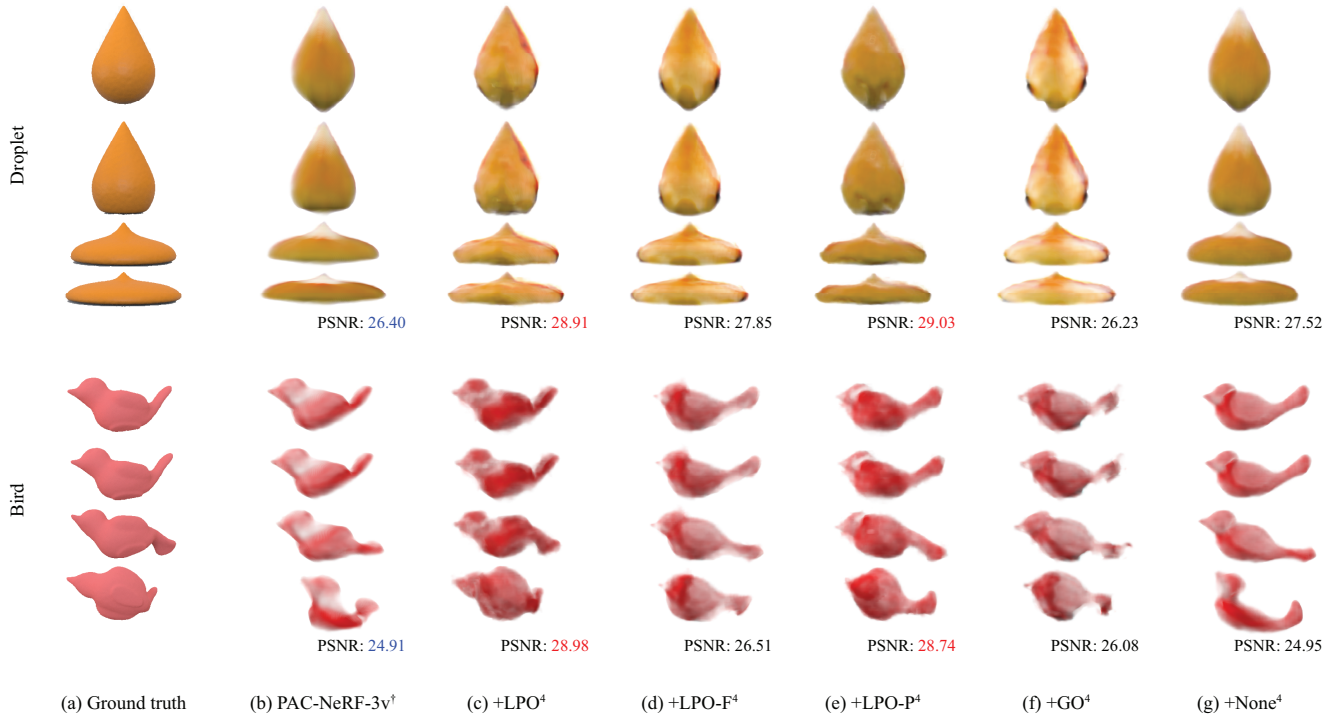


Figure 8. Qualitative comparisons among  $+LPO^4$ ,  $+LPO-F^4$ ,  $+LPO-P^4$ ,  $+GO^4$ , and  $+None^4$  in the scenes where  $+LPO-P^4$  outperformed  $+LPO-F^4$  in terms of PSNR. In the  $+LPO-F^4$ , position (that is, shape) optimization was ablated, and only feature (that is, appearance) optimization was conducted. In the  $+LPO-P^4$ , feature (that is, appearance) optimization was ablated, and only position (that is, shape) optimization was conducted. In the above scenes, large geometric gaps exist between the ground truth (a) and PAC-NeRF-3v<sup>†</sup> (b). For example, in the second row of the Droplet scene, the bottom of the droplet is flat in the ground truth (a), while that is bulging in PAC-NeRF-3v<sup>†</sup> (b). As shown in (d), only appearance correction by  $+LPO-F^4$  is insufficient to correct this geometry failure estimation, and the bottom of the droplet is still bulging. Instead,  $+LPO-F^4$  attempts to solve this problem by changing the appearance, making the colors overcorrected. In contrast, shape correction by  $+LPO-P^4$  effectively addresses this failure, and the bottom of the droplet is flat in (e). The same correction was also conducted in  $+LPO^4$  (c), a combination of  $+LPO-F^4$  and  $+LPO-P^4$ . Similarly, in the Bird scene, the directions of the tail are adequately corrected in  $+LPO^4$  (c) and  $+LPO-P^4$  (e). In contrast, those are not sufficiently conducted in  $+LPO-F^4$  (d). Also, in this case,  $+LPO-F^4$  (d) attempts to solve this problem by overcorrecting the colors.  $+GO^4$  (f), which also only corrects appearance, has the same difficulty as  $+LPO-F^4$  (d). In the Droplet scene, the bottom of the droplet is bulging; in the Bird scene, the bird’s tail is corrupted.  $+None^4$  (g), which performs Algorithm 1 without geometry correction, does not have a sufficient correction ability. In the Droplet scene, shape and appearance are almost identical to those in PAC-NeRF-3v<sup>†</sup> (b). In the Bird scene, the pose of the bird is not corrected.

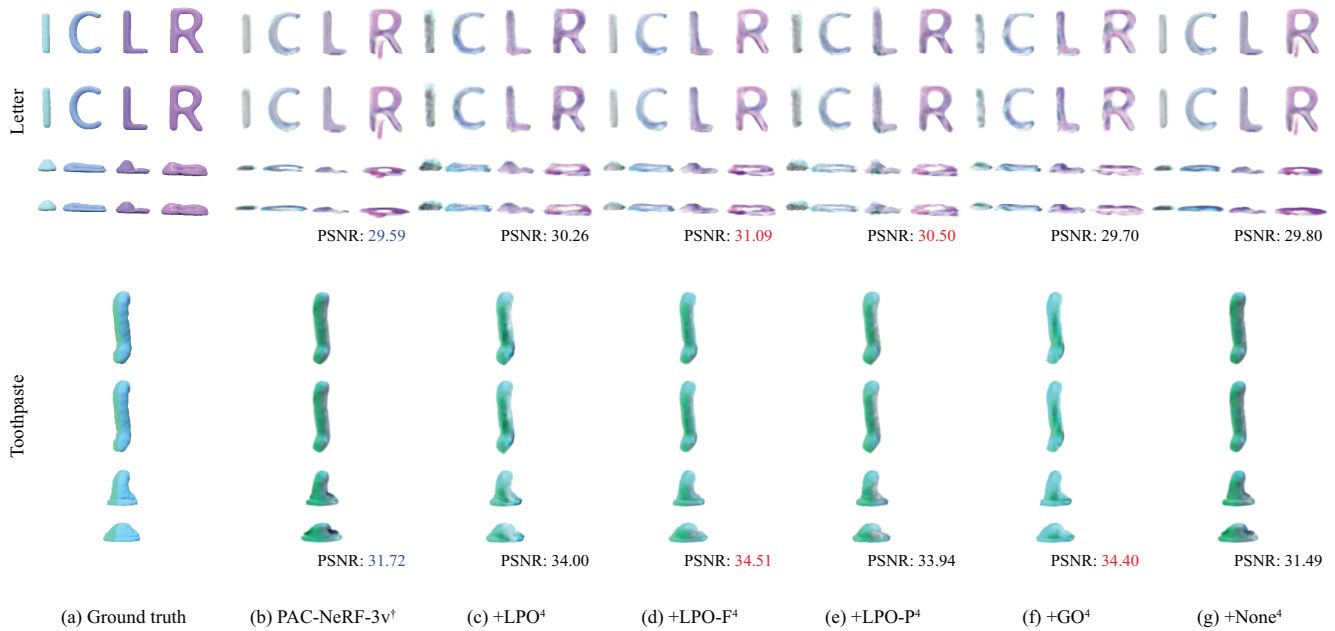


Figure 9. Qualitative comparisons among +LPO<sup>4</sup>, +LPO-F<sup>4</sup>, +LPO-P<sup>4</sup>, +GO<sup>4</sup>, and +None<sup>4</sup> in the scenes where +LPO-F<sup>4</sup> outperformed +LPO-P<sup>4</sup> in terms of PSNR. In the +LPO-F<sup>4</sup>, position (that is, shape) optimization was ablated, and only feature (that is, appearance) optimization was conducted. In the +LPO-P<sup>4</sup>, feature (that is, appearance) optimization was ablated, and only position (that is, shape) optimization was conducted. As discussed in Figure 8, +LPO-F<sup>4</sup> is unsuitable for significantly correcting the geometric shape because it focuses on correcting appearance. However, as shown in the above scenes, when PAC-NeRF-3v<sup>†</sup> (b) captures the geometric structure relatively well, +LPO-F<sup>4</sup> (d) also works well because appearance correction is more important than shape correction. For example, in the Letter scene, +LPO-F<sup>4</sup> (d) succeeds in eliminating artifacts existing in the vicinity of the left line of the letter “R.” In the Toothpaste scene, +LPO-F<sup>4</sup> (d) succeeds in making the color of the material brighter and closer to the ground truth (a). Notably, in both scenes, +None<sup>4</sup> (g) fails to do so and produces almost the same results as those in PAC-NeRF-3v<sup>†</sup> (b). These results indicate that appearance correction by +LPO-F<sup>4</sup> is essential (simple iterative updates in Algorithm 1 are insufficient) to address these problems. Another interesting finding is that +LPO-P<sup>4</sup> (e), which focuses on correcting shape, also works well. For example, in the Letter scene, +LPO-P<sup>4</sup> (e) eliminates the artifacts existing in the vicinity of the left line of the letter “R,” and in the Toothpaste scene, it makes the material brighter. This is possible because moving correctly colored particles from other places allows for appearance changes.

## E. Implementation details

### E.1. Dataset

We investigated the benchmark performance on the dataset provided by the original study on PAC-NeRF [42]. This dataset comprised nine scenes and various continuum materials, including the following:

- Newtonian fluids with fluid viscosity  $\mu$  and bulk modulus  $\kappa$ :
  - *Droplet* with  $\mu = 200$  and  $\kappa = 10^5$
  - *Letter* with  $\mu = 100$  and  $\kappa = 10^5$
- Non-Newtonian fluids with shear modulus  $\nu$ , bulk modulus  $\kappa$ , yield stress  $\tau_Y$ , and plasticity viscosity  $\eta$ :
  - *Cream* with  $\mu = 10^4$ ,  $\kappa = 10^6$ ,  $\tau_Y = 3 \times 10^3$ , and  $\eta = 10$
  - *Toothpaste* with  $\mu = 5 \times 10^3$ ,  $\kappa = 10^5$ ,  $\tau_Y = 200$ , and  $\eta = 10$
- Elastic materials with Young’s modulus  $E$  and Poisson’s ratio  $\nu$ :
  - *Torus* with  $E = 10^6$  and  $\nu = 0.3$
  - *Bird* with  $E = 3 \times 10^5$  and  $\nu = 0.3$
- Plasticine with Young’s modulus  $E$ , Poisson’s ratio  $\nu$ , and yield stress  $\tau_Y$ :
  - *Playdoh* with  $E = 2 \times 10^6$ ,  $\nu = 0.3$ , and  $\tau_Y = 1.54 \times 10^4$
  - *Cat* with  $E = 10^6$ ,  $\nu = 0.3$ , and  $\tau_Y = 3.85 \times 10^3$
- Granular media with friction angle  $\theta_{fric}$ :
  - *Trophy* with  $\theta_{fric} = 40^\circ$

In each scene, the objects fall freely under the influence of gravity and undergo collisions. The ground-truth simulation data were generated using the MLS-MPM framework [31]. A photorealistic simulation engine rendered objects under diverse environmental lighting conditions and ground textures. Each scene was captured from 11 view-points with cameras evenly spaced on the upper hemisphere, including the object. To evaluate our method in sparse-view settings, three views were used for training, and the remaining eight views were used for testing in the main experiments (Section 4) and the experiments described in Appendix C.1. Six views were used for training, and the remaining five were used for testing in the experiments described in Appendix C.2. Data were downloaded from the website<sup>14</sup> provided by the authors of PAC-NeRF [42].

### E.2. Model

The models were implemented based on the official PAC-NeRF source code.<sup>14</sup> For simplicity and fair comparison, we used the default model configurations provided in the source code for the experiments. Specifically, the architecture of a discretized voxel-based NeRF followed direct voxel grid optimization [73], in which a volume density field and color field were represented by voxel grids, and a 2-layer MLP with a hidden dimension of 128 was applied

to the color grid with positional embedding for a view direction  $\mathbf{d}$ . Regarding a differentiable MPM, DiffTaichi [32] was used.

### E.3. Training settings

For simplicity and fair comparison, we conducted Eulerian static voxel grid optimization (Figure 2(1)) and physical property optimization (Figure 2(2)) using the default training settings provided in the source code,<sup>14</sup> in PAC-NeRF, PAC-NeRF-3v, and PAC-NeRF-6v, except that the number of views was changed. In PAC-NeRF-3v<sup>†</sup>, we applied the three modifications described in Appendix A (scheduling of  $\mathcal{L}_{surf}$ , introduction of  $\mathcal{L}_{pixel}^{VI}$ , and adjustment of training length) to the Eulerian static voxel grid optimization and adopted the default training settings for physical property optimization. For the LPO (Figure 2(3)), we trained the features and positions of the particles for 100 iterations using the Adam optimizer [39]. In particular, we optimized the features of the particles at a learning rate of 0.1, which is the default value used for training the feature grids in the Eulerian static voxel grid optimization. The positions of the particles were optimized at a learning rate of  $\frac{dx}{32}$ , where  $dx$  indicates the voxel grid size and differs depending on the scene (set in the configuration files in the source code<sup>14</sup>). In a preliminary experiment, we found that careful setting of this learning rate is vital for stable training because the particles can diverge when the learning rate is exceptionally high. Based on this observation, we used a learning rate adjusted according to  $dx$ . We set the momentum terms of the Adam optimizer,  $\beta_1$  and  $\beta_2$ , to 0.9 and 0.999, respectively.

### E.4. Evaluation metrics

**Evaluation of geometric (re)correction.** We evaluated the performance of the geometric (re)correction using metrics commonly used to assess the performance of novel view synthesis in NeRF studies: the peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM) [79], and learned perceptual image patch similarity (LPIPS) [93]. For PSNR and SSIM, the larger the values, the better the performance. For LPIPS, the smaller the values, the better the performance. In particular, we report the scores averaged over the video sequences in a test set.

**Evaluation of physical identification.** To evaluate the performance of the physical identification, we measured the absolute distance between the ground truth and the estimated physical properties. The values of the ground-truth physical properties are provided in Appendix E.1. For an easy comparison, we calculated these distances after adjusting the scale (i.e., either a logarithmic scale or a linear scale) following the study of PAC-NeRF [42]. The smaller the values, the better the performance.

<sup>14</sup><https://github.com/xuan-li/PAC-NeRF>