

# REMATCHING DYNAMIC RECONSTRUCTION FLOW

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Reconstructing dynamic scenes from image inputs is a fundamental computer vision task with many downstream applications. Despite recent advancements, existing approaches still struggle to achieve high-quality reconstructions from unseen viewpoints and timestamps. This work introduces the ReMatching framework, designed to improve generalization quality by incorporating deformation priors into dynamic reconstruction models. Our approach advocates for velocity-field-based priors, for which we suggest a matching procedure that can seamlessly supplement existing dynamic reconstruction pipelines. The framework is highly adaptable and can be applied to various dynamic representations. Moreover, it supports integrating multiple types of model priors and enables combining simpler ones to create more complex classes. Our evaluations on popular benchmarks involving both synthetic and real-world dynamic scenes demonstrate a clear improvement in reconstruction accuracy of current state-of-the-art models.

## 1 INTRODUCTION

This work addresses the challenging task of novel-view dynamic reconstruction. That is, given a set of images of a dynamic scene evolving over time, the task objective is to render images from any novel view or intermediate point in time. Despite significant progress in dynamic reconstruction (Lombardi et al., 2021; Fridovich-Keil et al., 2023; Yunus et al., 2024), effectively learning dynamic scenes still remains an open challenge. The main hurdle arises from the typically sparse nature of multi-view inputs, both temporally and spatially. While tackling sparsity often involves incorporating some form of prior knowledge into the dynamic reconstruction model - either from a physical prior such as rigidity (Sorkine & Alexa, 2007), or learnable priors derived from large foundation models (Ling et al., 2024; Wang et al., 2024) - the optimal scheme for integrating these priors without compromising the fidelity of model reconstructions remains unclear.

To address this issue, this paper presents the ReMatching framework, a novel approach for designing and integrating deformation priors into dynamic reconstruction models. The ReMatching framework has three core goals: i) suggest an optimization objective that aims at achieving a reconstruction solution that closely satisfies the prior regularization **without compromising fidelity**; ii) ensure applicability to various model functions, including time-dependent rendered pixels or particles representing scene geometry; and, iii) provide a flexible design of deformation prior classes, allowing more complex classes to be built from simpler ones.

To support the usage of rich deformation prior classes, we advocate for priors expressed through velocity fields. A velocity field is a mathematical object that describes the instantaneous change in time the deformation induces. As such, a velocity field can potentially provide a simpler characterization of the underlying *flow* deformation. For example, the complex class of volume-preserving flow deformations is characterized by the condition of being generated by divergence-free velocity fields (Eisenberger et al., 2019). However, representing a deformation through its generating velocity field typically necessitates numerical simulation for integration, **a procedure that can be computationally expensive and time-consuming**. Nevertheless, recent progress in flow-based generative models (Ben-Hamu et al., 2022; Lipman et al., 2022; Albergo et al., 2023) supports simulation-free flow training, inspiring this work to explore simulation-free training for flow-based dynamic reconstruction models. Therefore, our framework is specifically designed to integrate with dynamic reconstruction models that represent dynamic scenes directly through time-dependent reconstruction functions (Pumarola et al., 2021; Yang et al., 2023).

Exploiting the simplicity offered by velocity-field-based deformation prior classes, we observe that the *projection* of a time-dependent reconstruction function onto a velocity-field prior class can be framed as a flow-matching problem, solvable analytically. The opportunity to access the projected flow is reminiscent of the Alternating Projections Method (APM) (Deutsch, 1992), a greedy algorithm *guaranteed* in finding the closest points between two sets. Therefore, we suggest an optimization objective aimed at re-projecting back onto the set of reconstruction flows. This corresponds to a flow-matching loss that we term the *ReMatching* loss. Our hypothesis is that by mimicking the APM, this optimization would converge to solutions that not only meet the reconstruction objective, but also reach the *closest* possible alignment to the required prior class. By doing so, we achieve the desired goal of improving generalization without compromising solutions’ fidelity levels.

We instantiate our framework with a dynamic model based on the popular Gaussian Splats (Kerbl et al., 2023) rendering model. We explore several constructions for deformation prior classes including piece-wise rigid and volume-preserving deformations. Additionally, we demonstrate our framework’s usability for two different types of time-dependent functions: rendered image intensity values, and particle positions representing scene geometry. Lastly, we evaluate our framework on standard dynamic reconstruction benchmarks, involving both synthetic and real-world scenes, and showcase clear improvement in generalization quality.

**Our contributions.** In summary, the main contributions of this paper are:

1. We propose the ReMatching framework, which controls the optimization of dynamic reconstruction models to converge to solutions that closely align with a predefined prior class of deformations, without strictly enforcing membership in the prior class, thereby improving the ability to achieve high-fidelity reconstructions.
2. The framework unifies different types of model functions, including geometry representations and image rendering functions, under a single cohesive approach, ensuring wide applicability and making future advancements within this framework relevant to many models.
3. The framework allows for the combination of multiple prior classes, enabling users to design the method for their specific reconstruction problem, enhancing adaptability across varied scenarios.

## 2 RELATED WORK

**Flow-based 3D dynamics.** There is an extensive body of works utilizing flow-based deformations for 3D related problems. For shape interpolation, (Eisenberger et al., 2019) considers volume-preserving flows. For dynamic geometry reconstruction, (Niemeyer et al., 2019) suggests learning neural parametrizations of velocity fields. This representation is further improved by augmenting it with a canonicalized object space parameterization (Rempe et al., 2020; Ren et al., 2021) or by simultaneously optimizing for 3D reconstruction and motion flow estimation (Vu et al., 2022). Similarly to (Niemeyer et al., 2019), (Du et al., 2021) suggests flow-based representation of dynamic rendering model based on a neural radiance field (Mildenhall et al., 2020). More recently, (Chu et al., 2022; Yu et al., 2023) explores combining a time-aware neural radiance field with a velocity field for modelling fluid dynamics. In contrast to our framework they focus exclusively on recovering the deformation of specific fluids i.e. smoke and not on reconstructing generic non-rigid objects.

**Dynamic novel-view rendering models.** Neural Radiance Fields (NeRF) (Mildenhall et al., 2020) is a popular image rendering model combining an implicit neural network with volumetric rendering. Several follow-up works (Pumarola et al., 2021; Park et al., 2021a; Tretschk et al., 2021) explore using NeRF for non-rigid reconstruction, by optimizing for time-dependent deformations. More recently, several works (Fridovich-Keil et al., 2023; Cao & Johnson, 2023; Wu et al., 2023; Song et al., 2023; Guo et al., 2023) try to address the training and inference inefficiencies of continuous volumetric representations by incorporating planes and grids into a spatio-temporal NeRF. An alternative to NeRF, suggesting an explicit scene representation, is the Gaussian Splatting (Kerbl et al., 2023) rendering model. Several works incorporate dynamics with Gaussian Splatting. (Yang et al., 2023) introduce a time-conditioned local deformation network. Similarly, (Wu et al., 2023) also relies on a canonical representation of a scene but further improves efficiency by considering a deformation

108 model based on on  $k$ -planes (Fridovich-Keil et al., 2023). (Lu et al., 2024) propose the integration of  
 109 a global deformation model.

### 111 3 METHOD

112  
 113 Given a collection of images,  $F_t = \{I_i^t\}_{i=1}^M$ , captured at  $T$  time steps, from  $M \geq 1$  viewing directions,  
 114 we seek to develop an image-based model for novel-view synthesis that can effectively render new  
 115 images from unseen viewpoints in any direction  $\mathbf{d} \in \mathcal{S}^2$  and any time  $t \in [t_1, t_T]$ . Since we aim to  
 116 support several time-dependent elements in a dynamic reconstruction model, we employ a general  
 117 notation for a dynamic image model. That is,  
 118

$$119 \quad t \mapsto \Psi_t = \{\psi(t) \mid \psi : \mathbb{R}_+ \rightarrow V\}, \quad (1)$$

120  
 121 with  $\Psi_t$  representing the evaluation at time  $t$  of all of the model components. Each element function  
 122  $\psi : \mathbb{R}_+ \rightarrow V$ , where  $V$  is a vector space, can specify any time-dependent quantity specified by the  
 123 model.  $V$  denotes a different vector space depending on the definition of  $\psi$ . For instance, if  $\psi$  models  
 124 time-dependent image intensity values,  $V = C^1(\mathbb{R}^d) = \{f \mid f : \mathbb{R}^d \rightarrow \mathbb{R}, \nabla f \text{ exists and continuous}\}$   
 125 with  $d = 2$ . Whereas, if  $\psi$  models the time-dependent position of  $n$  particles representing the  
 126 underlying scene geometry,  $V = \mathbb{R}^{n \times d}$  with  $d = 3$ . Lastly, in what follows, we interchangeably switch  
 127 between the notations  $\psi(t)$  and  $\psi_t$ .

128 We defer the specific details of the time-dependent reconstruction function  $\Psi_t$  to Section 5 and begin  
 129 by describing our proposed framework for incorporating priors via velocity fields.

#### 131 3.1 VELOCITY FIELDS

132  
 133 We consider a velocity field to be a time-dependent function of the form:

$$134 \quad v : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^d, \quad (2)$$

135  
 136 where usually  $d = 3$  or  $d = 2$ . A velocity field defines a time-dependent deformation in space  
 137  $\phi_t : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , also known as a *flow*, via an Ordinary Differential Equation (ODE):

$$138 \quad \begin{cases} \frac{\partial}{\partial t} \phi_t(\mathbf{x}) = v(\phi_t(\mathbf{x}), t) \\ \phi_0(\mathbf{x}) = \mathbf{x}. \end{cases} \quad (3)$$

139  
 140  
 141 Flow-based deformations are an ubiquitous modeling tool (Rezende & Mohamed, 2015; Chen et al.,  
 142 2018) that has been extensively used in various dynamic reconstruction tasks (Niemeyer et al., 2019;  
 143 Du et al., 2021). In a dynamic reconstruction model, a flow deformation can be incorporated by  
 144 defining a time-dependent function  $\psi_t : \mathbb{R}^d \rightarrow \mathbb{R}$  as a push-forward of some reference function  
 145  $\psi_0$ , i.e.,  $\psi_t = \phi_{t*} \psi_0$ . One key advantage of flow-based deformations is that they enable simple  
 146 parametrizations for the velocity field, in turn facilitating the integration of priors into the model.  
 147 For example, restricting  $\phi_t$  to be volume-preserving can be achieved by imposing the condition  
 148  $\text{div}(v) = 0$  (Eisenberger et al., 2019).

149  
 150 However, recovering  $\psi_t$  values in the case  $\psi_t = \phi_{t*} \psi_0$  is not explicit. Typically, this is achieved by  
 151 solving the continuity equation <sup>1</sup>

$$152 \quad \frac{\partial}{\partial t} \psi_t(\mathbf{x}) + \text{div}(\psi_t(\mathbf{x})v_t(\mathbf{x})) = 0, \forall \mathbf{x} \in \mathbb{R}^d, \quad (4)$$

153  
 154 which necessitates a numerical simulation. This introduces challenges for training flow-based models,  
 155 as errors in the numerical simulation can destabilize the optimization process. Therefore, to overcome  
 156 this hurdle, our framework assumes a reconstruction model consisting of functions  $\psi_t$  that are  
 157 simulation-free, i.e., each evaluation of  $\psi_t$  requires only a single step. Coupling  $\psi_t$  to a prior class  
 158 stemming from velocity-field-based formulation is the core issue our framework aims to address,  
 159 described in the following section.

160  
 161 <sup>1</sup>Assuming  $\psi_t$  obeys a conservation law, where  $v$  continuously deforms  $\psi_t$ .

### 3.2 FLOW REMATCHING

We assume that for a time-dependent reconstruction function,  $\psi_t \in \Psi_t$ , there exists an underlying flow  $\phi_t$  such that  $\psi_t$  can be described as a push-forward by  $\phi_t$ . We refer to such  $\phi_t$  as *reconstruction flow* and denote its velocity field by  $v_t$ . Under our assumption that  $\psi_t$  is simulation-free, neither  $\phi_t$  nor  $v_t$  is directly accessible. Nevertheless, let us assume that we can work with an element  $v_t \in \mathcal{V}$ , where  $\mathcal{V}$  represents the set of all the possible reconstruction generating velocity fields. Let  $\mathcal{P} \subset \{u_t | u : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}^d\}$  be a prior class of velocity fields to which  $v$  should belong. In Section 4, we discuss different choices for a class  $\mathcal{P}$ .

In some of the choices for  $\mathcal{P}$ , requiring  $v \in \mathcal{P}$  could be over-restrictive, conflicting with the fact that  $v$  also adheres to `generate` the reconstruction flow. Hence, an appealing objective would be to optimize  $v$  so that it is the closest element to  $\mathcal{P}$  out of the set  $\mathcal{V}$ . We suggest an optimization procedure mimicking the alternating projections method (APM) (Deutsch, 1992). The APM is an iterative procedure where alternating orthogonal projections are performed between two closed Hilbert sub-spaces  $V$  and  $P$ . Specifically,  $v_{k+1} = \text{proj}_V(\text{proj}_P(v_k))$  guarantees the convergence of  $v_k$  to  $\text{dist}(V, P)$ . Following this concept, our next step is to find a suitable notion for defining the projection operator for reconstruction generating velocity fields.

Since  $v$  is unknown, in our case, we utilize the continuity equation (4), which provides both a sufficient and a *necessary* condition for the generating velocity field of  $\phi_t$  in terms of  $\psi_t$  and its partial derivatives. In particular, we propose a projection procedure corresponding to the following matching optimization problem on the reconstruction flow:

$$u(\cdot, t) = \arg \min_{u_t \in \mathcal{P}} \rho(u_t, \psi_t), \quad (5)$$

where,

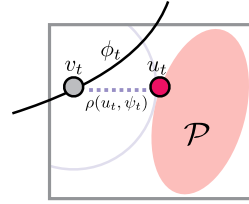
$$\rho(u_t, \psi_t) = \int \left| \frac{\partial}{\partial t} \psi_t(\mathbf{x}) + \text{div}(\psi_t(\mathbf{x})u_t(\mathbf{x})) \right|^2 d\mathbf{x}. \quad (6)$$

This procedure is illustrated in the right inset, where  $u_t$  (red dot) is the closest point to  $v_t$  on  $\mathcal{P}$ .

Following the alternating projections concept, the matched  $u_t$  should be projected back onto  $\mathcal{V}$  to propose a better candidate for  $v$ . This corresponds to a flow matching problem in  $u_t$ . We refer to this procedure as ReMatching and introduce the flow ReMatching loss,  $L_{\text{RM}}$ , a matching loss striving for the reconstruction flow to match  $u_t$ . That is,

$$L_{\text{RM}}(\theta) = \mathbb{E}_{t \sim U[0,1]} \rho(u_t, \psi_t) \quad (7)$$

where  $\theta$  denotes *solely* the parameters of  $\psi_t$ .




---

#### Algorithm 1 ReMatching loss

---

**Require:** Solver for 5, times  $\{t_l\}$   
 $L_{\text{RM}} = 0$   
**for**  $t \in \{t_l\}$  **do**  
 $u_t(\cdot) \leftarrow \text{solve}(\rho, \psi_t(\cdot))$   
 $L_{\text{RM}} \leftarrow L_{\text{RM}} + \rho(u_t(\cdot), \psi)$   
**end for**  
**Return:**  $L_{\text{RM}}$

---

The ReMatching loss is designed to supplement a reconstruction loss  $L_{\text{REC}}$  on  $\psi$  parameters  $\theta$ . Thus, our framework’s final loss for dynamic reconstruction `training` is

$$L(\theta) = L_{\text{REC}}(\theta) + \lambda L_{\text{RM}}(\theta) \quad (8)$$

where  $\lambda > 0$  is a hyper-parameter. In practice, for the ReMatching procedure to be seamlessly incorporated into a reconstruction training process, it is essential that problem 5 can be solved efficiently. Additionally, the integral in equation 7 is approximated by a sum using random samples  $\{t_l\} \sim U[0, 1]$ .

Algorithm 1 summarises the details of computing 7. Note that calculating  $\nabla_{\theta} L_{\text{RM}}$  does *not* necessarily require the cumbersome calculation of  $\nabla_{\theta} \arg \min_{u_t \in \mathcal{P}} \rho(u(\cdot, t), \psi_t)$ , since according to Danskin’s theorem (Madry, 2017),  $\nabla_{\theta} \rho(u_t, \psi_t) = \nabla_{\theta} \min_{u_t \in \mathcal{P}} \rho(u_t, \psi_t)$ . Additional implementation details regarding the losses can be found in the Appendix.

## 4 FRAMEWORK INSTANCES

This section presents several instances of the ReMatching framework discussed in this work. One notable setting is when  $V = \mathbb{R}^{n \times d}$ , i.e.,  $\psi_t = (\gamma_t^1, \dots, \gamma_t^n)^T$ , where each  $\gamma^i : \mathbb{R}_+ \rightarrow \mathbb{R}^d$ . In this case,

equation 6 becomes:

$$\rho(u_t, \psi_t) = \sum_{i=1}^n \left\| u_t(\gamma_t^i) - \frac{d}{dt} \gamma_t^i \right\|^2. \quad (9)$$

For the settings where  $V = C^1(\mathbb{R}^d)$ , equation 6 involves the computation of a spatial integral, which can be approximated by sampling a set of points  $\{\mathbf{x}_i\}_{i=1}^n$ . Moreover, taking into account that all prior classes [incorporated in this work](#) are divergence-free, equation 6 becomes:

$$\rho(u_t, \psi_t) = \sum_{i=1}^n \left| \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i) + \langle \nabla \psi_t(\mathbf{x}_i), u_t(\mathbf{x}_i) \rangle \right|^2, \quad (10)$$

since  $\text{div}(\psi_t(\mathbf{x})u_t(\mathbf{x})) = \langle \nabla_x \psi_t(\mathbf{x}), u_t(\mathbf{x}) \rangle + \psi_t(\mathbf{x}) \text{div} u_t(\mathbf{x})$ .

We now formulate several useful prior classes of velocity fields  $\mathcal{P}$ . A key feature of all the following constructions is their reliance on linear parameterizations, capitalizing on the fact that linear subspaces are sufficiently expressive to represent the velocity-based prior classes considered. This approach enables the use of efficient solvers for problem 5, reducing the computational task to solving a system of  $d$  linear equations, with a run-time complexity of at most  $O(n)$ .

#### 4.1 PRIOR DESIGN

**Directional restricted deformation.** In certain scenarios, it is safe to assume that the reconstruction flow can only deform along specific directions. For example, in an indoor scene, where furniture is placed on the floor, deformations would typically occur only in directions parallel to the floor plane. Let  $\mathbf{v} \in \text{span}\{\mathbf{v}_1, \dots, \mathbf{v}_l\}$ ,  $1 \leq l \leq d$  and  $\{\mathbf{v}_1, \dots, \mathbf{v}_l\}$  is a predefined orthonormal basis in which the flow remains static. Then, the prior class becomes:

$$\mathcal{P}_I = \{u_t \mid \langle u(\mathbf{x}, t), \mathbf{v}_m \rangle = 0, \forall m \in [l]\}. \quad (11)$$

When considering the matching minimization problem 5 in the settings of equation 9, we get:

$$\min_{u \in \mathcal{P}_I} \sum_{i=1}^n \left\| u_t(\gamma_t^i) - \frac{d}{dt} \gamma_t^i \right\|^2 = \sum_{i=1}^n \left\| \mathbf{V}^T \frac{d}{dt} \gamma_t^i \right\|^2, \quad (12)$$

where  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_l]$ . For the settings involving equation 10, the matching minimization problem is solved by:

$$\min_{u_t \in \mathcal{P}_I} \sum_{i=1}^n \left| \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i) + \langle \nabla \psi_t(\mathbf{x}_i), u_t(\mathbf{x}_i) \rangle \right|^2 = \sum_{i=1}^n \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i)^2 \left( 1 - \frac{\langle \nabla \psi_t(\mathbf{x}_i), \mathbf{V}_* \nabla \psi_t(\mathbf{x}_i) \rangle}{\|\nabla \psi_t(\mathbf{x}_i)\|^2} \right)^2, \quad (13)$$

where  $\mathbf{V}_* = (I - \mathbf{V}\mathbf{V}^T)$ .

**Rigid deformation.** One widely used prior in the dynamic reconstruction literature is rigidity, i.e., objects in a scene can only be deformed by a rigid transformation consisting of a translation and an orthogonal transformation. In a simple case, where it is assumed that the underlying dynamics consists of *one* rigid motion, the reconstruction flow would be of the form

$$\gamma(t) = \mathbf{R}(t)\mathbf{x}_0 + \mathbf{b}(t) \quad (14)$$

with  $\mathbf{R}(t) \in O(3)$  and  $\mathbf{b}(t) \in \mathbb{R}^3$ . Differentiating  $\gamma$  and solving for  $\mathbf{x}_0$  yields that

$$\frac{d}{dt} \gamma(t) = \dot{\mathbf{R}}(t)\mathbf{R}^T(t)(\gamma(t) - \mathbf{b}(t)) + \dot{\mathbf{b}}(t). \quad (15)$$

Since  $\dot{\mathbf{R}}(t)\mathbf{R}^T(t)$  is a skew-symmetric matrix, we suggest the following natural parameterization for the prior class

$$\mathcal{P}_{II} = \{u_t \mid u(\mathbf{x}, t) = \mathbf{A}_t \mathbf{x} + \mathbf{b}_t, \mathbf{A}_t \in \mathbb{R}^{d \times d}, \mathbf{A}_t = -\mathbf{A}_t^T, \mathbf{b}_t \in \mathbb{R}^d\}. \quad (16)$$

Substituting  $\mathcal{P}_{II}$  in problem 5 using equation 9 yields the following minimization problem:

$$\min_{(\mathbf{A}_t, \mathbf{b}_t)} \sum_{i=1}^n \left\| \mathbf{A}_t \gamma_t^i + \mathbf{b}_t - \frac{d}{dt} \gamma_t^i \right\|^2 \text{ s.t. } \mathbf{A}_t = -\mathbf{A}_t^T. \quad (17)$$

For the settings involving equation 10, the minimization problem 5 becomes:

$$\min_{(\mathbf{A}_t, \mathbf{b}_t)} \sum_{i=1}^n \left| \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i) + \langle \nabla \psi_t(\mathbf{x}_i), \mathbf{A}_t \mathbf{x}_i + \mathbf{b}_t \rangle \right|^2 \quad \text{s.t. } \mathbf{A}_t = -\mathbf{A}_t^T. \quad (18)$$

Importantly, both 17 and 18 are constrained least-squares problems. Thus, as detailed in Lemma 1, they enjoy an analytic solution that can be computed efficiently.

**Volume-preserving deformation.** So far we have only covered prior classes that may be too simplistic for capturing complex real-world dynamics. To address this, a reasonable assumption would be to include deformations that preserve the volume of any subset of the space. Notably, the rigid deformations prior class discussed earlier strictly falls within this class as well. Interestingly, volume-preserving flows are characterized by being generated via a divergence-free velocity field, i.e.,  $\text{div } u = 0$ . To this end, we propose the following prior class:

$$\mathcal{P}_{III} = \left\{ u_t | u_t(\mathbf{x}) = \sum_{j=1}^k \beta_j b_j(\mathbf{x}), \beta = [\beta_1, \dots, \beta_k]^T \in \mathbb{R}^k \right\}, \quad (19)$$

where for each basis  $b_j : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , we assume that  $\text{div}(b_j) = 0$ . Clearly,  $\text{div}(u_t) = 0$  for any choice of  $\beta \in \mathbb{R}^k$ . Taking into account that  $\text{div } \text{curl } u = 0$ , we follow (Eisenberger et al., 2019), and incorporate the following basis functions:

$$b_j(\mathbf{x}) \in \left\{ \text{curl}(\phi_j(\mathbf{x})\mathbf{e}_1^T), \dots, \text{curl}(\phi_j(\mathbf{x})\mathbf{e}_d^T) \right\}, \quad (20)$$

where  $\phi_j : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $\phi_j(\mathbf{x}) = \prod_{l=1}^d \sin(j_l \pi \mathbf{e}_l^T \mathbf{x})$  with  $j_l \in \mathbb{N}$  denoting the frequency for the  $l^{\text{th}}$  coordinate of the  $j^{\text{th}}$  basis function. Combining this prior with equation 9, yields the following minimization problem:

$$\min_{\beta} \sum_{i=1}^n \left\| \sum_{j=1}^k \beta_j b_j(\gamma_t^i) - \frac{d}{dt} \gamma_t^i \right\|^2. \quad (21)$$

Similarly, for the case of equation 10, we get:

$$\min_{\beta} \sum_{i=1}^n \left| \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i) + \left\langle \nabla \psi_t(\mathbf{x}_i), \sum_{j=1}^k \beta_j b_j(\mathbf{x}_i) \right\rangle \right|^2. \quad (22)$$

In particular, both minimization problems of 21 and 22 correspond to a standard least-squares problem and have an analytic solution.

A key decision involved in using the prior class  $\mathcal{P}_{III}$  is to select the number of basis functions  $k$ . However, setting  $k$  equal to a large value would make  $\mathcal{P}_{III}$  overly permissive, effectively neutralizing the ReMatching loss. To address this, we propose an additional procedure for constructing more complex prior classes, based on an adaptive choice of complexity level.

**Adaptive-combination of prior classes.** To address the challenge of setting the complexity level of the prior class, we introduce an adaptive (learnable) construction scheme for a prior class. Let  $w_j(\mathbf{x}, t) : \mathbb{R}^d \times \mathbb{R}_+ \rightarrow \mathbb{R}$ ,  $1 \leq j \leq k$ , be learnable functions, which are part of the reconstruction model, i.e.,  $w_j(\cdot, t) \in \Psi_t$  and  $w_j$  are normalized, i.e.,  $\sum_{j=1}^k w_j(\mathbf{x}, t) = 1$ . The details of  $w_j$  architecture are left to Section 5. We can construct a complex prior class by assigning simpler prior classes to different parts of the space, according to the weights  $w_j$ . For example, let us consider a *piece-wise* rigid deformation prior class defined as:

$$\mathcal{P}_{IV} = \left\{ u_t | u(\mathbf{x}, t) = \sum_{j=1}^k w_j(\mathbf{x}, t) u_j(\mathbf{x}, t), u_j \in \mathcal{P}_{II} \text{ for } 1 \leq j \leq k, \sum_{j=1}^k w_j(\mathbf{x}, t) = 1 \right\}. \quad (23)$$

In a similar manner, we can also combine  $\mathcal{P}_I$  with rigid deformations and derive a prior class defined as:

$$\mathcal{P}_V = \left\{ u_t | u(\mathbf{x}, t) = \sum_{j=1}^k w_j(\mathbf{x}, t) u_j(\mathbf{x}, t), u_1 \in \mathcal{P}_I, u_j \in \mathcal{P}_{II} \text{ for } 2 \leq j \leq k, \sum_{j=1}^k w_j(\mathbf{x}, t) = 1 \right\}. \quad (24)$$

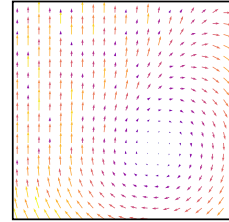


Figure 1: A vector field in  $\mathcal{P}_V$ .

Figure 1 illustrates an element in  $\mathcal{P}_V$ , with weights  $w_j$  dividing the plane to a restricted up direction deformation above the diagonal, and a rigid deformation below the diagonal. Note that directly substituting an adaptive-combination prior class in 5 would no longer yield a linear problem. Therefore, we propose to use a linear problem that upper bounds the matching optimization problem of 5. For example, in the case of equation 9 with  $\mathcal{P}_{IV}$ , we can solve:

$$\min_{\{(\mathbf{A}_{jt}, \mathbf{b}_{jt})\}} \sum_{i=1}^n \sum_{j=1}^k w_j(\gamma_t^i, t) \left\| \mathbf{A}_{jt} \gamma_t^i + \mathbf{b}_{jt} - \frac{d}{dt} \gamma_t^i \right\|^2 \text{ s.t. } \mathbf{A}_{jt} = -\mathbf{A}_{jt}^T. \quad (25)$$

Using Jensen’s inequality, it can be seen that 25 upper bounds the matching optimization from 5. Again the minimization problem of 25 can be solved efficiently, as it corresponds to a weighted least squares problem that is solvable independently for each  $j \in [k]$ , similarly to problem 17.

Lastly, as incorporating  $\mathcal{P}_{II}$  in 5 involves a non standard least-squares problem which includes a constraint, we formulate the analytic solutions for  $\mathcal{P}_{IV}$  in the next lemma, covering problems 17 and 18 as a special case.

**Lemma 1.** For the prior class  $\mathcal{P}_{IV}$ , the solutions  $(\mathbf{A}_{jt}, \mathbf{b}_{jt})$  to the minimization problem 25 are given by,

$$\begin{bmatrix} \text{vech}(\mathbf{A}_{jt}) \\ \mathbf{b}_{jt} \end{bmatrix} = \mathbf{P}_{jt}^{-1} \begin{bmatrix} \text{vec}(\dot{\mathbf{\Gamma}}_t^T \mathbf{W}_{jt} \mathbf{\Gamma}_{jt} - \mathbf{\Gamma}_{jt}^T \mathbf{W}_{jt} \dot{\mathbf{\Gamma}}_t) \\ \frac{1}{\mathbf{1}^T \mathbf{W}_{jt} \mathbf{1}} \mathbf{1}^T \mathbf{W}_{jt} \dot{\mathbf{\Gamma}}_t \end{bmatrix}$$

where  $\mathbf{\Gamma}_{jt} = [\gamma_t^1, \dots, \gamma_t^n]^T \in \mathbb{R}^{n \times d}$ ,  $\dot{\mathbf{\Gamma}}_t = [\frac{d}{dt} \gamma_t^1, \dots, \frac{d}{dt} \gamma_t^n]^T \in \mathbb{R}^{n \times d}$ ,  $\mathbf{W}_{jt} = \sum_{i=1}^n w_j(\gamma_t^i, t) \mathbf{e}_i \mathbf{e}_i^T$  with  $\{\mathbf{e}_i\}$  as the standard basis in  $\mathbb{R}^n$ ,  $\text{vech}(\mathbf{A}_{jt}) \in \mathbb{R}^{\frac{d(d-1)}{2}}$  denotes the half-vectorization of the anti-symmetric matrix  $\mathbf{A}_{jt}$ , and the matrix  $\mathbf{P}_{jt}^{-1}$  depends solely on  $\mathbf{\Gamma}_{jt}$ ,  $\dot{\mathbf{\Gamma}}_t$ , and  $\mathbf{W}_{jt}$ .

The solutions  $(\mathbf{A}_{jt}, \mathbf{b}_{jt})$  to the minimization problem 5 with 10 are given by,

$$\begin{bmatrix} \text{vech}(\mathbf{A}_{jt}) \\ \mathbf{b}_{jt} \end{bmatrix} = \mathbf{P}_{jt}^{-1} \begin{bmatrix} \sum_{i=1}^n w_j(\mathbf{x}_i, t) \text{vec}(s_i(\mathbf{x}_i [\mathbf{g}_t^i]^T - \mathbf{g}_t^i \mathbf{x}_i^T)) \\ \mathbf{G}_t^T \mathbf{W}_{jt} \mathbf{s} \end{bmatrix}$$

where  $\mathbf{g}_t^i = [\nabla \psi_t(\mathbf{x}_i)]^T$ ,  $\mathbf{G}_t = [\mathbf{g}_t^1, \dots, \mathbf{g}_t^n]^T \in \mathbb{R}^{n \times d}$ ,  $s_i = \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i)$ ,  $\mathbf{s} = [s_t^1, \dots, s_t^n]^T \in \mathbb{R}^n$ .

For the proof of lemma 1, including the details of  $\mathbf{P}_{jt}^{-1}$  computation, we refer the reader to the Appendix.

## 5 IMPLEMENTATION DETAILS

In this section, we provide additional details about the dynamic image model  $\Psi_t$  employed in this work, based on Gaussian Splatting (Kerbl et al., 2023). We provide an overview of this image model, followed by details about the dynamic model used in the experiments.

**Gaussian Splatting image model.** The Gaussian Splatting image model is parameterized by a collection of  $n$  3D Gaussians augmented with color and opacity parameters. That is,  $\theta = \{\boldsymbol{\mu}^i, \boldsymbol{\Sigma}^i, \mathbf{c}^i, \alpha^i\}_{i=1}^n$  with  $\boldsymbol{\mu}^i \in \mathbb{R}^3$  denoting the  $i^{\text{th}}$  Gaussian mean,  $\boldsymbol{\Sigma}^i \in \mathbb{R}^{3 \times 3}$  its covariance matrix,  $\mathbf{c}^i \in \mathbb{R}^3$  its color, and  $\alpha^i \in \mathbb{R}$  its opacity. To render an image, the 3D Gaussians are projected to the image plane to form a collection of 2D Gaussians parameterized by  $\{\boldsymbol{\mu}_{2D}^i, \boldsymbol{\Sigma}_{2D}^i\}$ . Given  $K, E$  denoting the intrinsic and extrinsic camera transformations, the image plane Gaussians parameters are calculated using the point rendering formula:

$$\boldsymbol{\mu}_{2D}^i = K \frac{E \boldsymbol{\mu}^i}{(E \boldsymbol{\mu}^i)_z}, \quad (26)$$

and,

$$\boldsymbol{\Sigma}_{2D}^i = J E \boldsymbol{\Sigma}^i E^T J^T, \quad (27)$$

where  $J$  denotes the Jacobian of the affine transformation of 26. Lastly, an image pixel  $I(\mathbf{p})$  is obtained by alpha-blending the ordered by depth visible Gaussians:

$$I(\mathbf{p}) = \sum_{i=1}^n c^i \alpha^i \sigma^i(\mathbf{p}) \prod_{j=1}^{i-1} (1 - \alpha^j \sigma^j(\mathbf{p})), \quad (28)$$

where  $\sigma^i(\mathbf{p}) = \exp\left(-\frac{1}{2}(\mathbf{p} - \boldsymbol{\mu}_{2D}^i)^T (\boldsymbol{\Sigma}_{2D}^i)^{-1} (\mathbf{p} - \boldsymbol{\mu}_{2D}^i)\right)$ .

**Dynamic image model.** We utilize the Gaussian Splatting image model to construct our dynamic model as:

$$\Psi_t = \left\{ \boldsymbol{\mu}^i + \boldsymbol{\mu}^i(t), \boldsymbol{\Sigma} + \boldsymbol{\Sigma}^i(t), \mathbf{c}^i, \alpha^i, w_{ij}(t) \right\}_{i=1}^n, \quad (29)$$

where  $\boldsymbol{\mu}^i(t) = f_{\boldsymbol{\mu}}(\boldsymbol{\mu}^i, t)$ ,  $\boldsymbol{\Sigma}^i(t) = f_{\boldsymbol{\Sigma}}(\boldsymbol{\mu}^i, t)$ ,  $w_{ij}(t) = e_j^T \text{softmax}(f_w(\boldsymbol{\mu}^i + \boldsymbol{\mu}^i(t), \boldsymbol{\mu}^i, t))$ . We follow (Yang et al., 2023) and each of the functions:  $f_{\boldsymbol{\mu}} : \mathbb{R}^3 \times \mathbb{R} \rightarrow \mathbb{R}^3$ ,  $f_{\boldsymbol{\Sigma}} : \mathbb{R}^3 \times \mathbb{R} \rightarrow \mathbb{R}^6$ ,  $f_w : \mathbb{R}^3 \times \mathbb{R}^3 \times \mathbb{R} \rightarrow \mathbb{R}^k$  is a Multilayer perceptron (MLP). For more details regarding the MLP architectures, we refer the reader to the Appendix. Note that the model element  $w_{ij}(t)$  is only relevant to instances where the adaptive-combination prior class is assumed. Lastly, in our experiments we apply the ReMatching loss for  $\boldsymbol{\mu}^i + \boldsymbol{\mu}^i(t)$ , and for time-dependent rendered images  $I_t$ .

**Training details.** We follow the training protocol of (Yang et al., 2023). We initialize the model using  $n = 100K$  3D Gaussians. Training is done for 40K iterations, where for the first 3K iterations, only  $\{\boldsymbol{\mu}^i, \boldsymbol{\Sigma}^i, \mathbf{c}^i, \alpha^i\}_{i=1}^n$  are optimized. In instances where the adaptive-combination prior class is applied, we supplement the ReMatching optimization objective with an entropy loss on the weights  $w_{ij}$  as follows:

$$L_{\text{entropy}} = \frac{1}{k} \sum_{j=1}^k \frac{1}{n} \sum_{i=1}^n w_{ij} \log \left( \frac{1}{n} \sum_{i=1}^n w_{ij} \right). \quad (30)$$

Lastly, for all the experiments considered in this work, we set the ReMatching loss weight  $\lambda = 0.001$ . Additional details are provided in the Appendix.

## 6 EXPERIMENTS

We evaluate the ReMatching framework on benchmarks involving synthetic and real-world video captures of deforming scenes. For quantitative analysis in both cases, we report the PSNR, SSIM (Wang et al., 2004) and LPIPS (Zhang et al., 2018) metrics.

**D-NeRF synthetic.** D-NeRF dataset (Pumarola et al., 2021) comprises of 8 scenes, each consisting from 100 to 200 frames, hence providing a dense multi-view coverage of the scene. We follow D-NeRF’s evaluation protocol and use the same train/validation/test split at  $800 \times 800$  image resolution with a black background. In terms of baseline methods, we consider recent state-of-the-art dynamic models, including **Deformable 3D Gaussians** (D3G) (Yang et al., 2023), 3D Geometry-aware Deformable Gaussians (**GA3D**) (Lu et al., 2024), Neural Parametric Gaussians (**NPG**) (Das et al., 2024), and K-Planes (Fridovich-Keil et al., 2023). Note that some of these baselines incorporate prior regularization losses such as local rigidity and smoothness to their optimization procedure. Table 1 summarizes the average image quality results for unseen frames in each scene. We include two variants of our framework: i) Using the divergence-free prior  $\mathcal{P}_{III}$ ; and ii) Using the adaptive-combination prior class  $\mathcal{P}_{IV}$  or the class  $\mathcal{P}_V$  specifically for scenes that include a floor component.

Method	Bouncing Balls			Hell Warrior			Hook			JumpingJacks		
	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑
<b>K-Planes</b> (Fridovich-Keil et al., 2023)	0.0242	37.78	0.9929	0.1074	32.57	0.9316	0.0655	29.46	0.9481	0.0417	31.73	0.9715
D3G (Yang et al., 2023)	0.0089	41.52	0.9978	0.0261	41.28	0.9928	0.0165	37.03	0.9906	0.0137	37.59	0.9930
GA3D (Liu et al., 2024)	0.0093	40.76	0.9950	0.0210	41.30	0.9871	0.0124	37.78	0.9887	0.0121	37.00	0.9887
NPG (Das et al., 2024)				0.0537	38.68	0.9780	0.0460	33.39	0.9735	0.0345	33.97	0.9828
Ours - $\mathcal{P}_{III}$	0.0087	41.84	0.9979	0.0244	41.59	0.9932	0.0161	37.19	0.9909	0.0134	37.72	0.9931
Ours - $\mathcal{P}_{IV}$ or $\mathcal{P}_V$	0.0089	41.61	0.9978	0.0245	41.69	0.9977	0.0158	37.39	0.9911	0.0131	38.01	0.9934
Method	Lego			Mutant			Stand Up			T-Rex		
	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑
<b>K-Planes</b> (Fridovich-Keil et al., 2023)	0.0472	25.15	0.9431	0.0215	35.30	0.9825	0.0211	36.55	0.9831	0.0284	30.41	0.9778
D3G (Yang et al., 2023)	0.0453	24.93	0.9537	0.0066	42.09	0.9966	0.0083	43.85	0.9970	0.0105	37.89	0.9956
GA3D (Liu et al., 2024)	0.0446	24.87	0.9420	0.0050	42.39	0.9951	0.0062	43.96	0.9948	0.0100	37.70	0.9929
NPG (Das et al., 2024)	0.0716	24.63	0.9312	0.0311	36.02	0.9840	0.0257	38.20	0.9889	0.0310	32.10	0.9959
Ours - $\mathcal{P}_{III}$	0.0503	24.89	0.9522	0.0067	42.13	0.9966	0.0085	43.99	0.9969	0.0105	38.07	0.9958
Ours - $\mathcal{P}_{IV}$ or $\mathcal{P}_V$	0.0456	24.95	0.9537	0.0065	42.40	0.9968	0.0081	44.31	0.9971	0.0103	38.38	0.9961

Table 1: Image quality evaluation on unseen frames for the D-NeRF dataset (Pumarola et al., 2021).



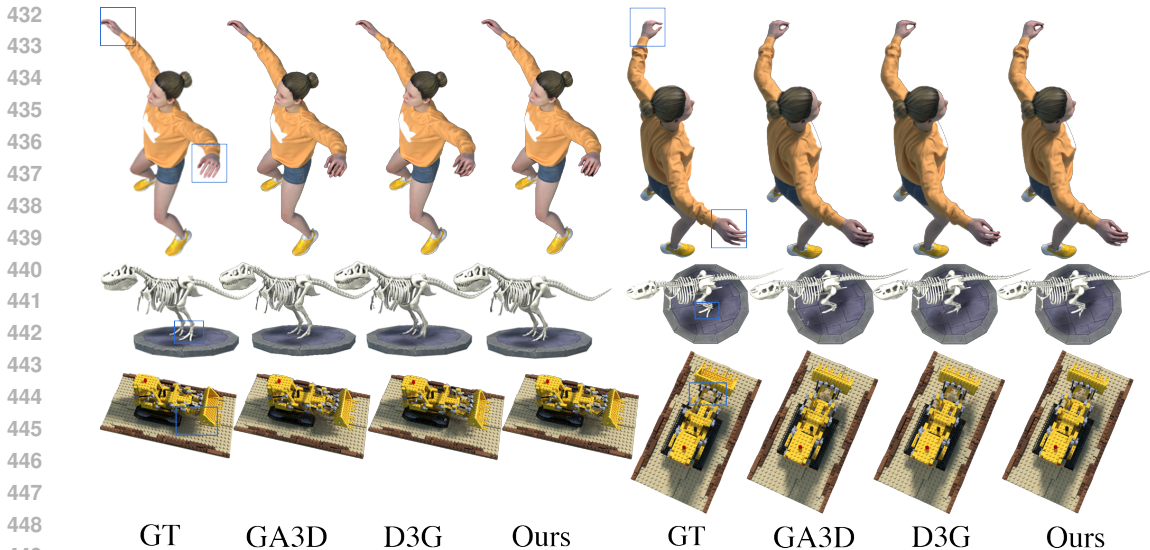


Figure 2: Qualitative comparison of baselines and our model on the D-NeRF dataset (Pumarola et al., 2021). We note that our framework consistently produces high fidelity reconstructions, accurately capturing fine-grained details, as highlighted in the blue boxes.

Figure 2 provides a qualitative comparison of rendered test frames, highlighting the improvements of our approach, which: i) produces plausible reconstructions that avoid unrealistic distortions, e.g., the human fingers in the jumping jacks scene; ii) reduces rendering artifacts of extraneous parts, especially in moving parts such as the leg in the T-Rex scene.

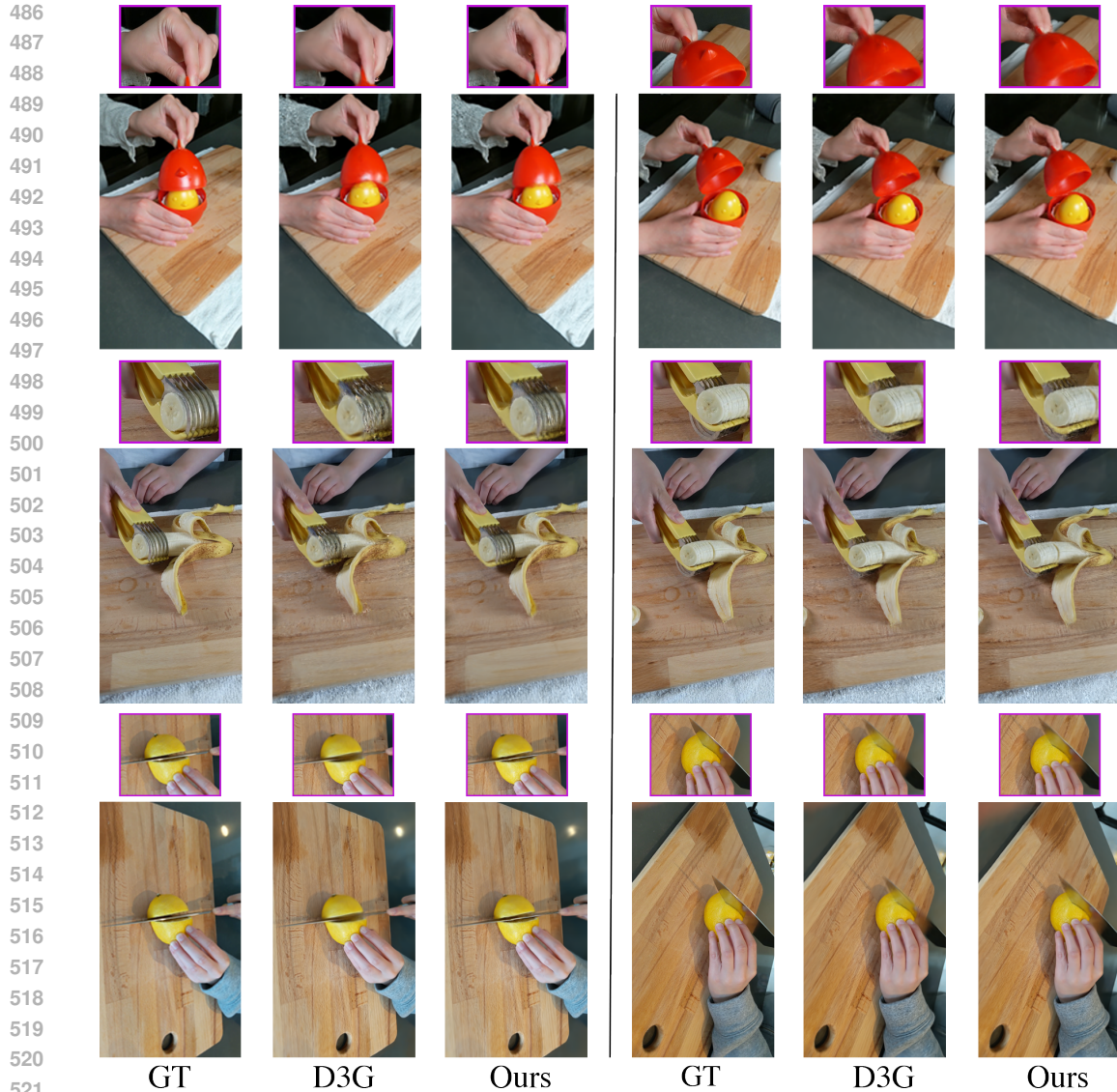
**HyperNeRF real-world.** The HyperNeRF dataset (Park et al., 2021b) consists of real-world videos capturing a diverse set of human activities involving interactions with common objects. We follow the evaluation protocol provided with the dataset, and use the same train/test split. In table 2 we report image quality results for unseen frames on 5 scenes from the dataset: Slice Banana, Chicken, Lemon, Torch, and Split Cookie. Figure 3 shows qualitative comparison to the baseline D3G (Yang et al., 2023). Our approach demonstrates similar types of improvements as noticed in the synthetic case providing more realistic reconstructions, especially in areas involving deforming parts.

**Adaptive-combination prior class.** Employing the adaptive-combination prior classes  $\mathcal{P}_{IV}$  and  $\mathcal{P}_V$  with learnable parts assignments  $\{w_{ij}\}$  raises the question of whether the learning process successfully produced assignments  $\{w_{ij}\}$  that align with the scene segmentation based on its deforming parts. Figure 4 shows our results for test frames from the Bouncing-Balls and Lego synthetic scenes (left), and the Chicken real-world scene (right). For comparison, we include the results of the Segment Anything Model (SAM) (Kirillov et al., 2023), which tends to over-segment the scene, mostly influenced by color variations and unable to capture the underlying geometry effectively.

**ReMatching time-dependent image.** In this experiment we validate the applicability of the ReMatching loss for controlling model solutions via rendered images. To that end, we apply our framework with the  $\mathcal{P}_{III}$  prior class to the Jumping Jacks scene from D-NeRF on a single specific front view through time. The qualitative comparison to D3G (Yang et al., 2023), as shown in the Appendix, supports the benefits of prior integration in this case as well, demonstrating more plausible reconstructions in areas involving moving parts.

Scene		LPIPS ↓	PSNR ↑	SSIM ↑
Slice Banana	D3G	0.3692	24.87	0.7935
	GA3D	0.4160	25.34	0.6722
	Ours - $\mathcal{P}_{III}$	0.3829	25.08	0.7992
	Ours - $\mathcal{P}_{IV}$	0.3673	25.28	0.8025
Chicken	D3G	0.3030	26.66	0.8813
	GA3D	0.4721	25.13	0.7555
	Ours - $\mathcal{P}_{III}$	0.2987	26.74	0.8836
	Ours - $\mathcal{P}_{IV}$	0.3044	26.80	0.8835
Lemon	D3G	0.2858	28.65	0.8873
	GA3D	0.3252	28.37	0.7596
	Ours - $\mathcal{P}_{III}$	0.2760	27.91	0.8842
	Ours - $\mathcal{P}_{IV}$	0.2675	28.30	0.8883
Torch	D3G	0.2340	25.41	0.9207
	GA3D	0.3278	23.79	0.8174
	Ours - $\mathcal{P}_{III}$	0.2221	26.00	0.9251
	Ours - $\mathcal{P}_{IV}$	0.2260	25.62	0.9229
Split Cookie	D3G	0.0971	32.61	0.9657
	GA3D	0.1144	32.28	0.9290
	Ours - $\mathcal{P}_{III}$	0.1097	31.31	0.9600
	Ours - $\mathcal{P}_{IV}$	0.0937	32.67	0.9667

Table 2: Unseen frames evaluation for the HyperNeRF dataset (Park et al., 2021b).



522 Figure 3: Qualitative comparison of our method to D3G (Yang et al., 2023) on the HyperNeRF  
523 dataset (Park et al., 2021b). Our framework yields more accurate reconstructions, in particular around  
524 moving parts.

## 527 7 CONCLUSIONS

528  
529  
530 We presented the ReMatching framework for integrating priors into dynamic reconstruction models. Our experimental  
531 results align with our hypothesis that the proposed ReMatching loss can induce solutions that match the required prior while  
532 achieving high fidelity reconstruction. We believe that the generality with which the framework was formulated would  
533 enable broader applicability to various dynamic reconstruction models. An interesting research venue is the construction of  
534 velocity-field-based prior classes emerging from video generative models, possibly utilizing our ReMatching formulation for time-dependent image intensity  
535 values. Another potential direction is the design of richer prior classes to handle more complex  
536 physical phenomena, such as ones including liquids and gases.  
537  
538  
539

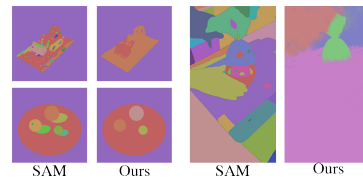


Figure 4: Part assignments for the adaptive-combination prior class.

## REFERENCES

- 540  
541  
542 Michael S. Albergo, Nicholas M. Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying  
543 framework for flows and diffusions, 2023. URL <https://arxiv.org/abs/2303.08797>.
- 544 Heli Ben-Hamu, Samuel Cohen, Joey Bose, Brandon Amos, Aditya Grover, Maximilian Nickel,  
545 Ricky TQ Chen, and Yaron Lipman. Matching normalizing flows and probability paths on  
546 manifolds. *arXiv preprint arXiv:2207.04711*, 2022.
- 547 Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of*  
548 *the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 130–141, 2023.
- 549  
550 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary  
551 differential equations. *Advances in neural information processing systems*, 31, 2018.
- 552 Mengyu Chu, Lingjie Liu, Quan Zheng, Erik Franz, Hans-Peter Seidel, Christian Theobalt, and  
553 Rhaleb Zayer. Physics informed neural fields for smoke reconstruction with sparse data. *ACM*  
554 *Trans. Graph.*, 2022.
- 555 Devikalyan Das, Christopher Wewer, Raza Yunus, Eddy Ilg, and Jan Eric Lenssen. Neural parametric  
556 gaussians for monocular non-rigid object reconstruction. In *Proceedings of the IEEE/CVF*  
557 *Conference on Computer Vision and Pattern Recognition*, pp. 10715–10725, 2024.
- 558 Frank Deutsch. The method of alternating orthogonal projections. In *Approximation theory, spline*  
559 *functions and applications*, pp. 105–121. Springer, 1992.
- 560  
561 Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance  
562 flow for 4d view synthesis and video processing. In *2021 IEEE/CVF International Conference on*  
563 *Computer Vision (ICCV)*, pp. 14304–14314. IEEE Computer Society, 2021.
- 564 Marvin Eisenberger, Zorah Löhner, and Daniel Cremers. Divergence-free shape correspondence by  
565 deformation. In *Computer Graphics Forum*, volume 38, pp. 1–12. Wiley Online Library, 2019.
- 566  
567 Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo  
568 Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *IEEE/CVF*  
569 *Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada,*  
570 *June 17-24, 2023*, pp. 12479–12488. IEEE, 2023.
- 571  
572 Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo  
573 Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *CVPR*, 2023.
- 574 Haoyu Guo, Sida Peng, Yunzhi Yan, Linzhan Mou, Yujun Shen, Hujun Bao, and Xiaowei Zhou.  
575 Compact neural volumetric video representations with dynamic codebooks. In *Advances in Neural*  
576 *Information Processing Systems 36: Annual Conference on Neural Information Processing Systems*  
577 *2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.
- 578 Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting  
579 for real-time radiance field rendering. *ACM Transactions on Graphics*, 42(4), July 2023. URL  
580 <https://repo-sam.inria.fr/fungraph/3d-gaussian-splatting/>.
- 581  
582 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete  
583 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick.  
584 Segment anything. *arXiv:2304.02643*, 2023.
- 585 Huan Ling, Seung Wook Kim, Antonio Torralba, Sanja Fidler, and Karsten Kreis. Align your  
586 gaussians: Text-to-4d with dynamic 3d gaussians and composed diffusion models. In *IEEE/CVF*  
587 *Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada,*  
588 *June 17-24, 2023*. IEEE, 2024.
- 589 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching  
590 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 591  
592 Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhöfer, Yaser Sheikh, and Jason M.  
593 Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Trans. Graph.*, 40  
(4):59:1–59:13, 2021.

- 594 Zhicheng Lu, Xiang Guo, Le Hui, Tianrui Chen, Min Yang, Xiao Tang, Feng Zhu, and Yuchao Dai.  
595 3d geometry-aware deformable gaussian splatting for dynamic view synthesis. *arXiv preprint*  
596 *arXiv:2404.06270*, 2024.
- 597 Aleksander Madry. Towards deep learning models resistant to adversarial attacks. *arXiv preprint*  
598 *arXiv:1706.06083*, 2017.
- 600 Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and  
601 Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020.
- 602 Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Occupancy flow: 4d  
603 reconstruction by learning particle dynamics. In *Proceedings of the IEEE/CVF International*  
604 *Conference on Computer Vision*, pp. 5379–5389, 2019.
- 606 Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman, Steven M.  
607 Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *2021 IEEE/CVF*  
608 *International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17,*  
609 *2021*, 2021a.
- 610 Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B. Goldman,  
611 Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation  
612 for topologically varying neural radiance fields. *CoRR*, abs/2106.13228, 2021b. URL <https://arxiv.org/abs/2106.13228>.
- 614 Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural  
615 radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer*  
616 *Vision and Pattern Recognition*, pp. 10318–10327, 2021.
- 618 Davis Rempe, Tolga Birdal, Yongheng Zhao, Zan Gojcic, Srinath Sridhar, and Leonidas J. Guibas.  
619 Caspr: Learning canonical spatiotemporal point cloud representations. In *Advances in Neural*  
620 *Information Processing Systems 33: Annual Conference on Neural Information Processing Systems*  
621 *2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 622 Zhongzheng Ren, Xiaoming Zhao, and Alex Schwing. Class-agnostic reconstruction of dynamic  
623 objects from videos. *Advances in Neural Information Processing Systems*, 34:509–522, 2021.
- 624 Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International*  
625 *conference on machine learning*, pp. 1530–1538. PMLR, 2015.
- 627 Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas  
628 Geiger. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance  
629 fields. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2732–2742, 2023.
- 630 Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry*  
631 *processing*, volume 4, pp. 109–116. Citeseer, 2007.
- 632 Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and  
633 Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a  
634 dynamic scene from monocular video. In *2021 IEEE/CVF International Conference on Computer*  
635 *Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, 2021.
- 637 Tuan-Anh Vu, Duc Thanh Nguyen, Binh-Son Hua, Quang-Hieu Pham, and Sai-Kit Yeung. Rfnet-4d:  
638 Joint object reconstruction and flow estimation from 4d point clouds. In Shai Avidan, Gabriel J.  
639 Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (eds.), *Computer Vision -*  
640 *ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part*  
641 *XXIII*, 2022.
- 642 Chaoyang Wang, Peiye Zhuang, Aliaksandr Siarohin, Junli Cao, Guocheng Qian, Hsin-Ying Lee,  
643 and Sergey Tulyakov. Diffusion priors for dynamic view synthesis from monocular videos. *CoRR*,  
644 abs/2401.05583, 2024.
- 646 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from  
647 error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612,  
2004.

648 Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian,  
649 and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint*  
650 *arXiv:2310.08528*, 2023.

651 Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. De-  
652 formable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint*  
653 *arXiv:2309.13101*, 2023.

654 Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of  
655 dynamic scenes with globally coherent depths from a monocular camera. *CoRR*, abs/2004.01294,  
656 2020. URL <https://arxiv.org/abs/2004.01294>.

657 Hong-Xing Yu, Yang Zheng, Yuan Gao, Yitong Deng, Bo Zhu, and Jiajun Wu. Inferring hybrid  
658 neural fluid fields from videos. In *Advances in Neural Information Processing Systems 36: Annual*  
659 *Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA,*  
660 *USA, December 10 - 16, 2023*, 2023.

661 Raza Yunus, Jan Eric Lenssen, Michael Niemeyer, Yiyi Liao, Christian Rupprecht, Christian Theobalt,  
662 Gerard Pons-Moll, Jia-Bin Huang, Vladislav Golyanik, and Eddy Ilg. Recent trends in 3d  
663 reconstruction of general non-rigid scenes. In *Computer Graphics Forum*, pp. e15062. Wiley  
664 Online Library, 2024.

665 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable  
666 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on*  
667 *computer vision and pattern recognition*, pp. 586–595, 2018.

## 671 8 APPENDIX

### 672 8.1 PROOFS

#### 673 8.1.1 PROOF OF LEMMA 1

674 *Proof.* (Lemma 1)

675 Let  $\Gamma_t, \dot{\Gamma}_t, \mathbf{W}_{jt}$  be given. Without loss of generality, we show the proof only for individuals  $j \in [k]$   
676 and  $t$ . So in order to ease the notation, in what follows we omit the subscripts  $t$  and  $j$ . Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$   
677 and  $\mathbf{b} \in \mathbb{R}^d$ , and define  $[w_1, \dots, w_n]^T = \mathbf{W}\mathbf{1}$ . First, we show that  $\sum_{i=1}^n w_i \|\mathbf{A}\gamma^i + \mathbf{b} - \dot{\gamma}^i\|^2$  can be  
678 reformulated as a weighted norm-squared minimization problem in  $\mathbf{A}$  and  $\mathbf{b}$ . That is,

$$679 \sum_{i=1}^n w_i \|\mathbf{A}(\gamma^i - \mathbf{c}) + \mathbf{b} - \dot{\gamma}^i\|^2 = \sum_{i=1}^n \text{tr} \mathbf{A} \sqrt{w_i} \gamma^i \sqrt{w_i} \gamma^{iT} \mathbf{A}^T - \quad (31)$$

$$680 2\text{tr} \sqrt{w_i} (\dot{\gamma}^i - \mathbf{b}) \sqrt{w_i} \gamma^i \mathbf{A}^T + \text{tr} \sqrt{w_i} (\dot{\gamma}^i - \mathbf{b}) \sqrt{w_i} (\dot{\gamma}^i - \mathbf{b})^T \quad (32)$$

$$681 = \text{tr} \mathbf{A} \Gamma^T \mathbf{W} \Gamma \mathbf{A}^T - 2\text{tr} (\dot{\Gamma} - \mathbf{1b}^T)^T \mathbf{W} \Gamma \mathbf{A}^T + \quad (33)$$

$$682 \text{tr} (\dot{\Gamma} - \mathbf{1b}^T)^T \mathbf{W} (\dot{\Gamma} - \mathbf{1b}^T) \quad (34)$$

$$683 = \left\| \sqrt{\mathbf{W}} (\Gamma \mathbf{A}^T - (\dot{\Gamma} - \mathbf{1b}^T)) \right\|^2. \quad (35)$$

684 Next, we consider the following optimization problem:

$$685 \min_{\mathbf{A}, \mathbf{b}} \left\| \sqrt{\mathbf{W}} (\Gamma \mathbf{A}^T - (\dot{\Gamma} - \mathbf{1b}^T)) \right\|^2 \text{ s.t. } \mathbf{A} = -\mathbf{A}^T. \quad (36)$$

686 Use the fact that  $\mathbf{A} = -\mathbf{A}^T$  to define the following Lagrangian,

$$687 \mathcal{L}(\mathbf{A}, \mathbf{b}, \Lambda) = \left\| \sqrt{\mathbf{W}} (\Gamma \mathbf{A} - \mathbf{1b}^T + \dot{\Gamma}) \right\|^2 + \text{tr} \Lambda^T (\mathbf{A} + \mathbf{A}^T). \quad (37)$$

688 Then,

$$689 \frac{\partial \mathcal{L}}{\partial \mathbf{A}} = 2\Gamma^T \mathbf{W} (\Gamma \mathbf{A} - \mathbf{1b}^T + \dot{\Gamma}) + \Lambda + \Lambda^T.$$

Thus,  $\frac{\partial \mathcal{L}}{\partial \mathbf{A}} = 0$  yields that  $\Gamma^T \mathbf{W}(\Gamma \mathbf{A} - \mathbf{1} \mathbf{b}^T + \dot{\Gamma})$  is symmetric. Then, using again the fact that  $\mathbf{A} = -\mathbf{A}^T$ , we get that,

$$\Gamma^T \mathbf{W} \Gamma \mathbf{A} + \mathbf{A} \Gamma^T \mathbf{W} \Gamma + \mathbf{b} \mathbf{1}^T \mathbf{W} \Gamma - \Gamma^T \mathbf{W} \mathbf{1} \mathbf{b}^T = \dot{\Gamma}^T \mathbf{W} \Gamma - \Gamma^T \mathbf{W} \dot{\Gamma}. \quad (38)$$

Now, taking the derivative w.r.t. to  $\mathbf{b}$  gives,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = -2 \mathbf{1}^T \mathbf{W}(\Gamma \mathbf{A} - \mathbf{1} \mathbf{b}^T + \dot{\Gamma})$$

and,  $\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = 0$ , yields that,

$$-\widehat{\mathbf{w}}^T \Gamma \mathbf{A} + \mathbf{b}^T = \widehat{\mathbf{w}}^T \dot{\Gamma}, \quad (39)$$

where  $\widehat{\mathbf{w}} = \frac{\mathbf{W} \mathbf{1}}{\mathbf{1}^T \mathbf{W} \mathbf{1}}$ . Vectorizing the LHS of 38 gives,

$$(I_d \otimes \Gamma^T \mathbf{W} \Gamma + \Gamma^T \mathbf{W} \Gamma \otimes I_d) D_d \text{vech}(\mathbf{A}) + (\Gamma^T \mathbf{W} \mathbf{1} \otimes I_d - I_d \otimes \Gamma^T \mathbf{W} \mathbf{1}) \mathbf{b} \quad (40)$$

where  $D_d$  is the duplication matrix transforming  $\text{vech}(\mathbf{A})$  to  $\text{vec}(\mathbf{A})$ , with  $\text{vech}(\mathbf{A})$  denoting the half-vectorization of the anti-symmetric matrix  $\mathbf{A}$ . Similarly, vectorizing the LHS of 39 yields,

$$-\frac{1}{\mathbf{1}^T \mathbf{W} \mathbf{1}} (I_d \otimes \mathbf{1}^T \mathbf{W} \Gamma) D_d \text{vech}(\mathbf{A}) + \mathbf{b}. \quad (41)$$

Based on 40 and 41, we can define the following block matrix:

$$\mathbf{P} = \left[ \begin{array}{c|c} \mathbf{Q} = \mathbf{Q}' D_d & \mathbf{R} = \Gamma^T \mathbf{W} \mathbf{1} \otimes I_d - I_d \otimes \Gamma^T \mathbf{W} \mathbf{1} \\ \hline \mathbf{S} = \mathbf{S}' D_d & \mathbf{T} = I_d \end{array} \right] \quad (42)$$

where  $\mathbf{Q}' = I_d \otimes \Gamma^T \mathbf{W} \Gamma + \Gamma^T \mathbf{W} \Gamma \otimes I_d$ , and,  $\mathbf{S}' = -\frac{1}{\mathbf{1}^T \mathbf{W} \mathbf{1}} (I_d \otimes \mathbf{1}^T \mathbf{W} \Gamma)$ . Then, let,

$$\mathbf{U} = (\mathbf{Q} - \mathbf{R} \mathbf{T}^{-1} \mathbf{S})^{-1} = L_d (\mathbf{Q}' - \mathbf{R} \mathbf{S}')^{-1} \quad (43)$$

where  $L_d$  is the matrix satisfying  $D_d L_d = I_{d^2}$ . Consequently,

$$\mathbf{P}^{-1} = \left[ \begin{array}{c|c} \mathbf{U} & -\mathbf{U} \mathbf{R} \\ \hline -\mathbf{S}' (\mathbf{Q}' - \mathbf{R} \mathbf{S}')^{-1} & I_d + \mathbf{S}' (\mathbf{Q}' - \mathbf{R} \mathbf{S}')^{-1} \mathbf{R} \end{array} \right] \quad (44)$$

and,

$$\left[ \begin{array}{c} \text{vech}(\mathbf{A}) \\ \mathbf{b} \end{array} \right] = \mathbf{P}^{-1} \left[ \begin{array}{c} \text{vec}(\dot{\Gamma}^T \mathbf{W} \Gamma - \Gamma^T \mathbf{W} \dot{\Gamma}) \\ \widehat{\mathbf{w}}^T \dot{\Gamma} \end{array} \right]. \quad (45)$$

□

Now, for the second part of the lemma. Let  $\mathbf{g}_i = [\nabla \psi_t(\mathbf{x}_i)]^T$ ,  $s_i = \frac{\partial}{\partial t} \psi_t(\mathbf{x}_i)$ . Consider the following energy,

$$L = \sum_{i=1}^n w_i (\mathbf{g}_i^T (A \mathbf{x}_i + \mathbf{b}) + s_i)^2. \quad (46)$$

Note that,

$$\mathbf{g}_i^T A \mathbf{x}_i = \mathbf{y}_i^T \mathbf{a} \quad (47)$$

where  $\mathbf{a} := \text{vec}(A)$ , and  $\mathbf{y}_i := \mathbf{x}_i \otimes \mathbf{g}_i$ . Then,

$$L = \sum_{i=1}^n w_i (\mathbf{a}^T \mathbf{y}_i \mathbf{y}_i^T \mathbf{a} + \mathbf{b}^T \mathbf{g}_i \mathbf{g}_i^T \mathbf{b} + 2 \mathbf{a}^T \mathbf{y}_i \mathbf{g}_i^T \mathbf{b} + 2 \mathbf{g}_i^T s_i \mathbf{b} + 2 s_i \mathbf{a}^T \mathbf{y}_i + s_i^2). \quad (48)$$

Define the Lagrangian,

$$\mathcal{L}(\mathbf{a}, \mathbf{b}, \lambda) = \mathbf{a}^T \sum_i \mathbf{y}_i w_i \mathbf{y}_i^T \mathbf{a} + \mathbf{b}^T \mathbf{G}^T \mathbf{W} \mathbf{G} \mathbf{b} + 2 \mathbf{a}^T \sum_i \mathbf{y}_i w_i \mathbf{g}_i^T \mathbf{b} + \quad (49)$$

$$2 \mathbf{s}^T \mathbf{W} \mathbf{G} \mathbf{b} + 2 \mathbf{a}^T \sum_i w_i \mathbf{y}_i s_i + \mathbf{t}^T \mathbf{W} \mathbf{t} + \lambda^T (\mathbf{a} + \mathbf{P} \mathbf{a}) \quad (50)$$

where  $P$  is the permutation matrix s.t.  $\text{vec}(A^T) = P\mathbf{a}$ .

Then,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}} = 2 \sum_i w_i \mathbf{y}_i \mathbf{y}_i^T \mathbf{a} + 2 \sum_i w_i \mathbf{y}_i \mathbf{g}_i^T \mathbf{b} + 2 \sum_i w_i \mathbf{y}_i s_i + \lambda + P\lambda \quad (51)$$

Equating the above to 0 and unvectorizing it, yields the following matrix equation,

$$\sum_{i=1}^n w_i (\mathbf{g}_i \mathbf{g}_i^T A \mathbf{x}_i \mathbf{x}_i^T + \mathbf{g}_i \mathbf{g}_i^T \mathbf{b} \mathbf{x}_i^T + s_i \mathbf{g}_i \mathbf{x}_i^T) = \frac{1}{2} (\mathbf{\Lambda} + \mathbf{\Lambda}^T), \quad (52)$$

yielding that the LHS is a symmetric matrix. Therefore,

$$\sum_{i=1}^n w_i (\mathbf{g}_i \mathbf{g}_i^T A \mathbf{x}_i \mathbf{x}_i^T + \mathbf{x}_i \mathbf{g}_i^T \mathbf{b} \mathbf{x}_i^T + s_i \mathbf{g}_i \mathbf{x}_i^T) = \sum_{i=1}^n w_i (\mathbf{x}_i \mathbf{x}_i^T A^T \mathbf{g}_i \mathbf{g}_i^T + \mathbf{x}_i \mathbf{b}^T \mathbf{g}_i \mathbf{g}_i^T + s_i \mathbf{x}_i \mathbf{g}_i^T). \quad (53)$$

Rearranging the above and half-vectorizing both sides yields that,

$$\sum_{i=1}^n w_i (\mathbf{x}_i \mathbf{x}_i^T \otimes \mathbf{g}_i \mathbf{g}_i^T + \mathbf{g}_i \mathbf{g}_i^T \otimes \mathbf{x}_i \mathbf{x}_i^T) D_d \text{vech}(\mathbf{A}) + w_i (\mathbf{x}_i \otimes \mathbf{g}_i \mathbf{g}_i^T - \mathbf{g}_i \mathbf{g}_i^T \otimes \mathbf{x}_i) \mathbf{b} = \quad (54)$$

$$\sum_{i=1}^n w_i \text{vec}(s_i (\mathbf{x}_i \mathbf{g}_i^T - \mathbf{g}_i \mathbf{x}_i^T)). \quad (55)$$

Now,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}} = 0 \quad (56)$$

yields that,

$$\mathbf{G}^T \mathbf{W} \mathbf{G} \mathbf{b} + \sum_i w_i \mathbf{g}_i \mathbf{y}_i^T D_d \text{vech}(\mathbf{A}) = -\mathbf{G}^T \mathbf{W} \mathbf{s}. \quad (57)$$

Therefore,

$$\mathbf{P} = \left[ \begin{array}{c|c} \mathbf{Q} = \mathbf{Q}' D_d & \mathbf{R} = \sum_{i=1}^n w_i (\mathbf{x}_i \otimes \mathbf{g}_i \mathbf{g}_i^T - \mathbf{g}_i \mathbf{g}_i^T \otimes \mathbf{x}_i) \\ \mathbf{S} = \mathbf{S}' D_d & \mathbf{T} = \mathbf{G}^T \mathbf{W} \mathbf{G} \end{array} \right], \quad (58)$$

where  $\mathbf{S}' = \sum_{i=1}^n w_i \mathbf{g}_i \mathbf{y}_i^T$ , and  $\mathbf{Q}' = \sum_{i=1}^n w_i (\mathbf{x}_i \mathbf{x}_i^T \otimes \mathbf{g}_i \mathbf{g}_i^T + \mathbf{g}_i \mathbf{g}_i^T \otimes \mathbf{x}_i \mathbf{x}_i^T)$ . Then, let,

$$\mathbf{U} = (\mathbf{Q} - \mathbf{R} \mathbf{T}^{-1} \mathbf{S})^{-1} = L_d (\mathbf{Q}' - \mathbf{R} \mathbf{T}^{-1} \mathbf{S}')^{-1} \quad (59)$$

where  $L_d$  is the matrix that satisfies  $D_d L_d = I_{d^2}$ . Consequently,

$$\mathbf{P}^{-1} = \left[ \begin{array}{c|c} \mathbf{U} & -\mathbf{U} \mathbf{R} \mathbf{T}^{-1} \\ \hline -\mathbf{T}^{-1} \mathbf{S}' (\mathbf{Q}' - \mathbf{R} \mathbf{T}^{-1} \mathbf{S}')^{-1} & \mathbf{T}^{-1} + \mathbf{T}^{-1} \mathbf{S}' (\mathbf{Q}' - \mathbf{R} \mathbf{T}^{-1} \mathbf{S}')^{-1} \mathbf{R} \mathbf{T}^{-1} \end{array} \right]. \quad (60)$$

### 8.1.2 CONTINUITY EQUATION CONSTRAINT DERIVATION FOR $V = \mathbb{R}^n$

In the main text, we stated that in the case when  $V = \mathbb{R}^n$ , i.e.,  $\psi_t = (\gamma_t^1, \dots, \gamma_t^n)^T$ , equation 6 becomes:

$$\rho(u_t, \psi_t) = \sum_{i=1}^n \left\| u_t(\gamma_t^i) - \frac{d}{dt} \gamma_t^i \right\|^2. \quad (61)$$

To see this formally, let  $\delta(\mathbf{x} - \mathbf{a})$  denote the Dirac delta generalized function concentrated around  $\mathbf{a}$ , satisfying

$$\delta(\mathbf{x} - \mathbf{a}) = 0, \forall \mathbf{x} \neq \mathbf{a}, \quad (62)$$

and,

$$\int \phi(\mathbf{x}) \delta(\mathbf{x} - \mathbf{a}) d\mathbf{x} = \phi(\mathbf{a}), \quad (63)$$

for any test function  $\phi$ . Consider  $\psi_t(\mathbf{x}) = \sum_{i=1}^n \psi_t^i(\mathbf{x})$ , where  $\psi_t^i(\mathbf{x}) = \delta(\mathbf{x} - \gamma_t^i)$ . Note that under this definition of  $\psi_t$ ,  $V$  is in fact the space of generalized functions. Then,

$$\frac{\partial}{\partial t} \psi_t^i = \left\langle \nabla \delta(\mathbf{x} - \gamma_t^i), -\frac{d}{dt} \gamma_t^i \right\rangle, \quad (64)$$

and, using the chain rule as applied in the simplification of equation 10, we have that,

$$\operatorname{div} \psi_t^i u(\mathbf{x}) = \langle \nabla \delta(\mathbf{x} - \gamma_t^i), u(\mathbf{x}) \rangle + \delta(\mathbf{x} - \gamma_t^i) \operatorname{div}(u(\mathbf{x})). \quad (65)$$

Substituting these computations in the continuity equation 4, yields that,

$$0 = \int \left| \frac{\partial}{\partial t} \psi_t(\mathbf{x}) + \operatorname{div}(\psi_t(\mathbf{x}) v_t(\mathbf{x})) \right| d\mathbf{x} \geq \quad (66)$$

$$\left| \int \frac{\partial}{\partial t} \psi_t(\mathbf{x}) + \operatorname{div}(\psi_t(\mathbf{x}) v_t(\mathbf{x})) d\mathbf{x} \right| = \quad (67)$$

$$\left| \int \left\langle \nabla \delta(\mathbf{x} - \gamma_t^i), -\frac{d}{dt} \gamma_t^i \right\rangle + \langle \nabla \delta(\mathbf{x} - \gamma_t^i), u(\mathbf{x}) \rangle + \delta(\mathbf{x} - \gamma_t^i) \operatorname{div}(u(\mathbf{x})) d\mathbf{x} \right| = \quad (68)$$

$$\left| \int \left\langle \nabla \delta(\mathbf{x} - \gamma_t^i), u(\mathbf{x}) - \frac{d}{dt} \gamma_t^i \right\rangle d\mathbf{x} + \int \delta(\mathbf{x} - \gamma_t^i) \operatorname{div}(u(\mathbf{x})) d\mathbf{x} \right|. \quad (69)$$

Now, under the assumption that  $\operatorname{div}(u) = 0$  almost everywhere, using 63 yields that the second term in the last equation vanishes. Therefore,

$$0 = \left| \int \left\langle \nabla \delta(\mathbf{x} - \gamma_t^i), u(\mathbf{x}) - \frac{d}{dt} \gamma_t^i \right\rangle d\mathbf{x} \right| = \quad (70)$$

$$\left| \int \delta(\mathbf{x} - \gamma_t^i) \left\langle \nabla \log \delta(\mathbf{x} - \gamma_t^i), u(\mathbf{x}) - \frac{d}{dt} \gamma_t^i \right\rangle d\mathbf{x} \right| \geq \quad (71)$$

$$\int \delta(\mathbf{x} - \gamma_t^i) \left\| \nabla \log \delta(\mathbf{x} - \gamma_t^i) \right\| \left\| u(\mathbf{x}) - \frac{d}{dt} \gamma_t^i \right\| d\mathbf{x}, \quad (72)$$

where we applied the Cauchy-Schwarz inequality in the final step. Therefore,

$$\int \delta(\mathbf{x} - \gamma_t^i) \left\| \nabla \log \delta(\mathbf{x} - \gamma_t^i) \right\| \left\| u(\mathbf{x}) - \frac{d}{dt} \gamma_t^i \right\| d\mathbf{x} = 0. \quad (73)$$

Applying property 63 yields that equation 73 can be true only if when  $\mathbf{x} = \gamma_t^i$ , we have that,

$$\left\| u(\mathbf{x}) - \frac{d}{dt} \gamma_t^i \right\| = 0. \quad (74)$$

Utilizing this constraint for each  $i$ , we can derive equation 9.

## 8.2 ADDITIONAL IMPLEMENTATION DETAILS

### 8.2.1 ARCHITECTURE

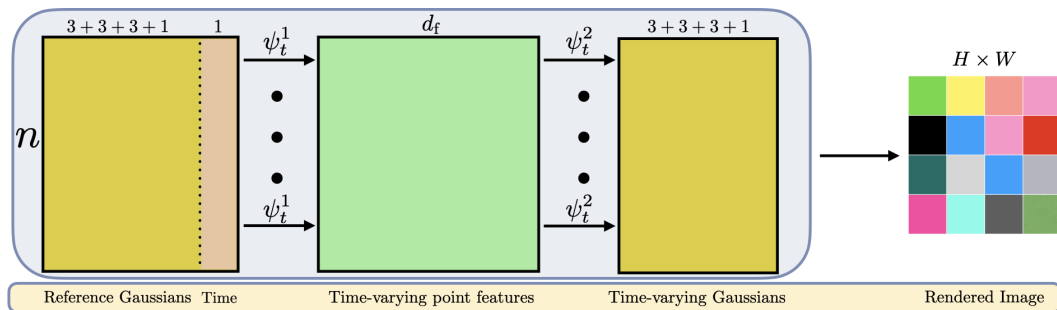


Figure 5: Illustration of the architecture for  $\Psi_t$  used in the experiments, based on (Yang et al., 2023). Reference Gaussians parameters are propagated to time  $t$  through a shared function,  $\psi_t^1$ , implemented as an MLP with positional encoding features to compute time-varying point features of dimension  $d_t$ . These features are then processed by a second shared function,  $\psi_t^2$ , to generate time-varying Gaussians parameters. Finally, given a chosen viewing direction, the Gaussian Splatting rendering model is used to produce a rendered image.

We first describe the construction of the Gaussian Splatting dynamic image model referenced in section 5. An illustration of this model is presented in Figure 5. The time invariant base of



864 the model is optimized throughout training and consists of the following set of parameters  $\theta =$   
 865  $\{\mu^i, \mathbf{S}^i, \mathbf{R}^i, \mathbf{c}^i, \alpha^i\}_{i=1}^n$  with Gaussian mean  $\mu^i \in \mathbb{R}^3$ , scaling  $\mathbf{S}^i \in \mathbb{R}^3$ , rotation quaternion  $\mathbf{R}^i \in \mathbb{R}^4$ ,  
 866 color  $\mathbf{c}^i \in \mathbb{R}^3$  and opacity  $\alpha^i \in \mathbb{R}$ . The covariance matrix  $\Sigma^i$  is calculated during the rendering  
 867 process from the temporally augmented scaling and rotation parameters.  
 868

869 The time dependent deformation model transforms the time invariant Gaussian mean  $\mu^i$  and selected  
 870 time  $t$  into the deformation of the mean, scaling, rotation and the model element  $w$  in the case of the  
 871 adaptive-combination prior.

872 We generate positional embeddings (Mildenhall et al., 2020) of the time and mean inputs, which we  
 873 pass to the deformation model Multilayer perceptrons.

$$874 \text{Emb}_{\text{time}}(t) : \mathbb{R} \rightarrow \mathbb{R}^{d_{\text{time emb}}}$$

$$875 \text{Emb}_{\text{mean}}(\mu^i) : \mathbb{R}^3 \rightarrow \mathbb{R}^{d_{\text{mean emb}}}$$

876 The deformation model is made up of layers of the form:

$$877 \psi(n, d_{\text{in}}, d_{\text{out}}) : \mathbf{X} \mapsto \nu(\mathbf{X}\mathbf{W} + \mathbf{1}\mathbf{b}^T)$$

878 where  $\nu = \text{Softplus}_{\beta}$ , with  $\beta = 100$ .

879 For the deformation of the mean, scaling and rotation, the model takes the same form with minor  
 880 differences in the final layer depending on the deforming parameter.

$$881 \text{Emb}_{\text{time}}(t) \rightarrow \psi(n, d_{\text{time emb}}, 256) \rightarrow \psi(n, 256, d_{\tau}) \rightarrow \tau$$

$$882 [\tau, \text{Emb}_{\text{mean}}(\mu)] \rightarrow \psi(n, d_{\tau} + d_{\text{mean emb}}, 256) \rightarrow \psi(n, 256, 256) \rightarrow$$

$$883 \psi(n, 256, 256) \rightarrow \psi(n, 256, 256) \rightarrow [\tau, \text{Emb}_{\text{mean}}(\mu), \psi(n, 256, 256)] \rightarrow$$

$$884 \psi(n, d_{\tau} + d_{\text{mean emb}} + 256, 256) \rightarrow \psi(n, 256, 256) \rightarrow \psi(n, 256, 256) \rightarrow \omega$$

$$885 \text{Mean} : \omega \rightarrow \psi(n, 256, 3) \rightarrow \mu(t)$$

$$886 \text{Scaling} : \omega \rightarrow \psi(n, 256, 3) \rightarrow \mathbf{S}(t)$$

$$887 \text{Rotation} : \omega \rightarrow \psi(n, 256, 4) \rightarrow \mathbf{R}(t)$$

888 For the prediction of the  $w$  we use a shallower Multilayer perceptron.

$$889 [\tau, \text{Emb}_{\text{mean}}(\mu + \mu(t)), \text{Emb}_{\text{mean}}(\mu)] \rightarrow \psi(n, d_{\tau} + 2 \cdot d_{\text{mean emb}}, 256) \rightarrow$$

$$890 \psi(n, 256, K) \rightarrow \text{Softmax} \rightarrow w(t)$$

## 891 8.2.2 HYPER-PARAMETERS AND TRAINING DETAILS

892 We set  $d_{\text{mean emb}} = 63$ ,  $d_{\text{time emb}} = 13$  and  $d_{\tau} = 30$ . For optimization we use an Adam optimizer with  
 893 different learning rates for the network components, keeping the hyper-parameters of the baseline  
 894 model (Yang et al., 2023).  
 895

896 In the case of the adaptive-combination prior we select  $k$  based on a hyper-parameter search between 1  
 897 and 35. The optimal value for most scenes ranges between 5 and 15, though the number also depends  
 898 on the selected composition of priors. For example, a single volume-preserving class can supervise  
 899 multiple moving objects as opposed to a single rigid deformation class. We use the ReMatching loss  
 900 weight  $\lambda = 0.001$ . When supplementing the ReMatching loss with an additional entropy loss, we use  
 901 0.0001 as the entropy loss weight.  
 902

In calculating the partial derivatives of  $\psi_t$ , we note that the input dimension for predicting  $\psi_t$  is relatively small – 1 for time or  $d$  for spatial coordinates – compared to the output dimension, which can be  $n \times d$  for spatial ReMatching or  $H \times W$  in image-space ReMatching. Given this, forward-mode automatic differentiation proves to be more efficient than backward-mode differentiation for this specific computation, both computationally and in terms of memory usage. Consequently, we utilize forward-mode autodiff to compute the partial derivatives of  $\psi_t$  required for the ReMatching loss. Once the ReMatching loss is incorporated, we employ backward-mode autodiff to compute the gradient of the overall loss with respect to the model parameters.

### 8.2.3 REMATCHING RENDERED IMAGE

The reconstruction model architecture in the case of the image ReMatching is the same as for the other experiments. At initialization we select a fixed viewpoint for the evaluation of the image space loss, which is kept throughout training. At every iteration we sample a random time and evaluate the ReMatching loss from the fixed viewpoint.

For approximating equation 10, we calculate a sample by choosing points that their image value is close to 0 after applying the following transformation:

$$f(x) = -0.1 \cdot \ln(1 - |x|) \cdot \text{sign}(x) \quad (75)$$

on the image.

Next, we compute the image gradient using automatic differentiation and use our single class div-free solver to reconstruct the flow and calculate the loss.



Figure 6: Qualitative comparison of the ReMatching loss applied in the image space. Each group of 3 is showing Ground-Truth (left), Ours (center), and D3G (right).

### 8.3 ADDITIONAL EVALUATION

To further evaluate the efficacy of the ReMatching framework in practical applications, we consider the Dynamic Scenes dataset (Yoon et al., 2020), which captures forward-facing views of real-world scenes exhibiting complex dynamics. To that end, we selected 4 scenes from the Human, Interaction, and Vehicle categories, consisting of monocular videos with approximately 80–180 frames for training and an additional 20 frames reserved for testing. Figure 7 shows qualitative comparison to

the D3G (Yang et al., 2023) baseline, highlighting similar patterns of improvement as observed in earlier experiments. Specifically, our approach better preserves fine details, such as the truck’s front lights (Truck) and the bottom teeth (Dynamic Face). Additionally, it demonstrates a reduction in reconstruction motion artifacts, as evident in the humans in motion (Jumping) and the legs and head of the dinosaur (Balloon). Table 3 presents a quantitative evaluation, comparing two variants from our prior classes,  $\mathcal{P}_{III}$  and  $\mathcal{P}_{IV}$  to D3G. These results correlate with the qualitative improvements discussed above.

Method	Balloon			Truck			Jumping			Dynamic Face		
	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑
D3G (Yang et al., 2023)	0.1584	26.79	0.9349	0.2922	26.01	0.9046	0.2726	23.12	0.8958	0.0806	29.22	0.9756
Ours - $\mathcal{P}_{III}$	0.1592	26.96	0.9348	0.2782	25.49	0.9071	0.2720	22.89	0.8971	0.0794	29.30	0.9761
Ours - $\mathcal{P}_{IV}$	0.1578	26.95	0.9356	0.2533	26.66	0.9197	0.2501	23.63	0.9037	0.0793	29.23	0.9754

Table 3: Unseen frames evaluation for the dynamic scenes dataset (Yang et al., 2023).



Figure 7: Qualitative comparison of our method to D3G (Yang et al., 2023) on the dynamic scenes dataset (Yoon et al., 2020).

### 8.4 RECONSTRUCTION FLOW EVALUATION

In this section, we evaluate the ability of the ReMatching framework to recover the underlying reconstruction flow  $\phi_t$ . Since the reconstruction flow is generally unknown, we use the following simple flow to generate training data:

$$\phi_t(\mathbf{x}) = \mathbf{R}(t)\mathbf{S}\mathbf{x} \tag{76}$$

where  $\mathbf{R}^T(t)\mathbf{R}(t) = \mathbf{I}_d$ , and  $\mathbf{S} = \text{diag}\{s_1, \dots, s_d\}$ . We evaluate our framework in its two settings: i) equation 9, corresponding to  $V = \mathbb{R}^{n \times d}$ ; and, ii) equation 10, corresponding to  $V = C^1(\mathbb{R}^d)$ .

**The  $V = \mathbb{R}^{n \times d}$  case.** We evaluate these settings for  $d = 3$ , following a similar approach to the dynamic image model based on Gaussian Splatting described in Section 5. To construct the training

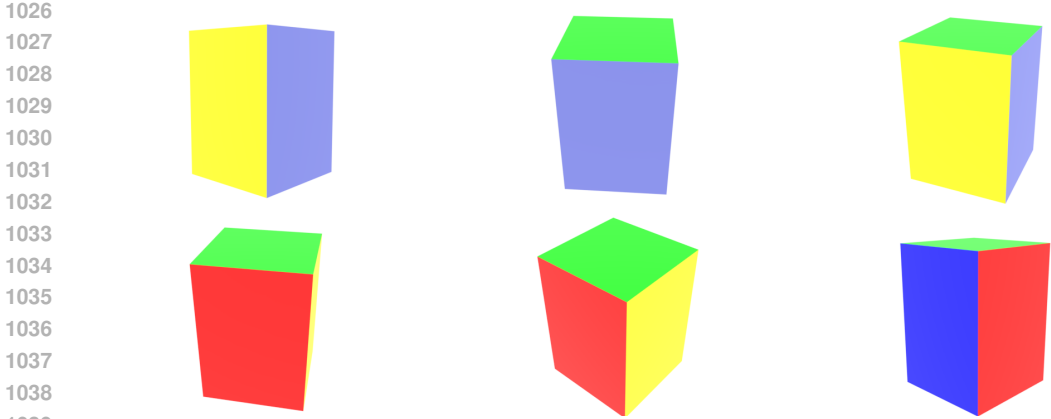


Figure 8: Training frames from the rotating colored box scene.

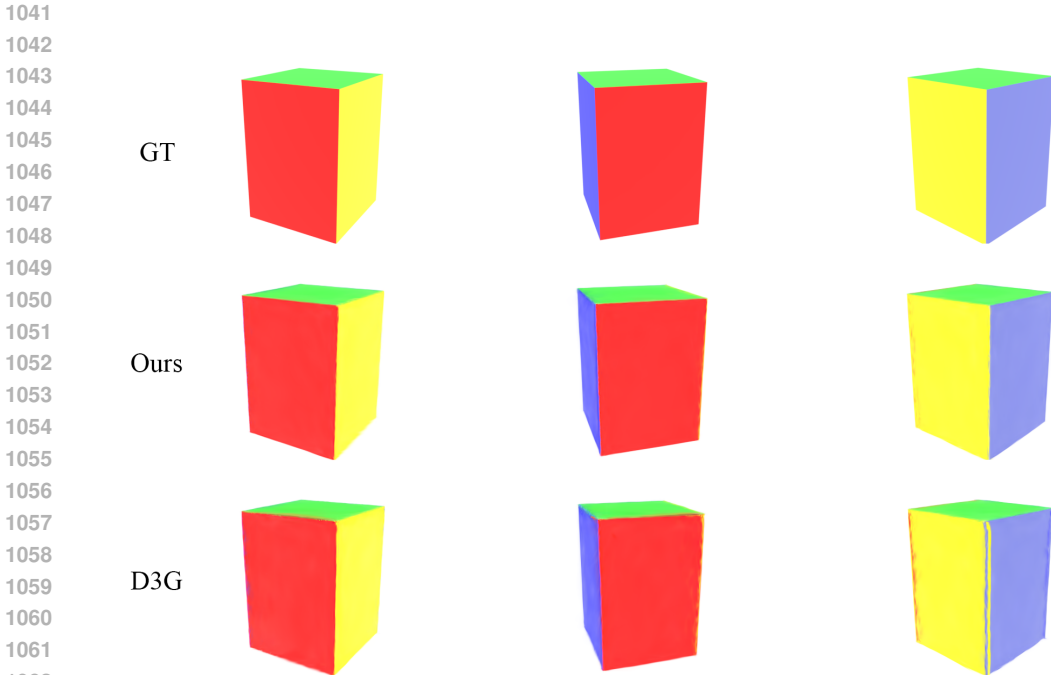
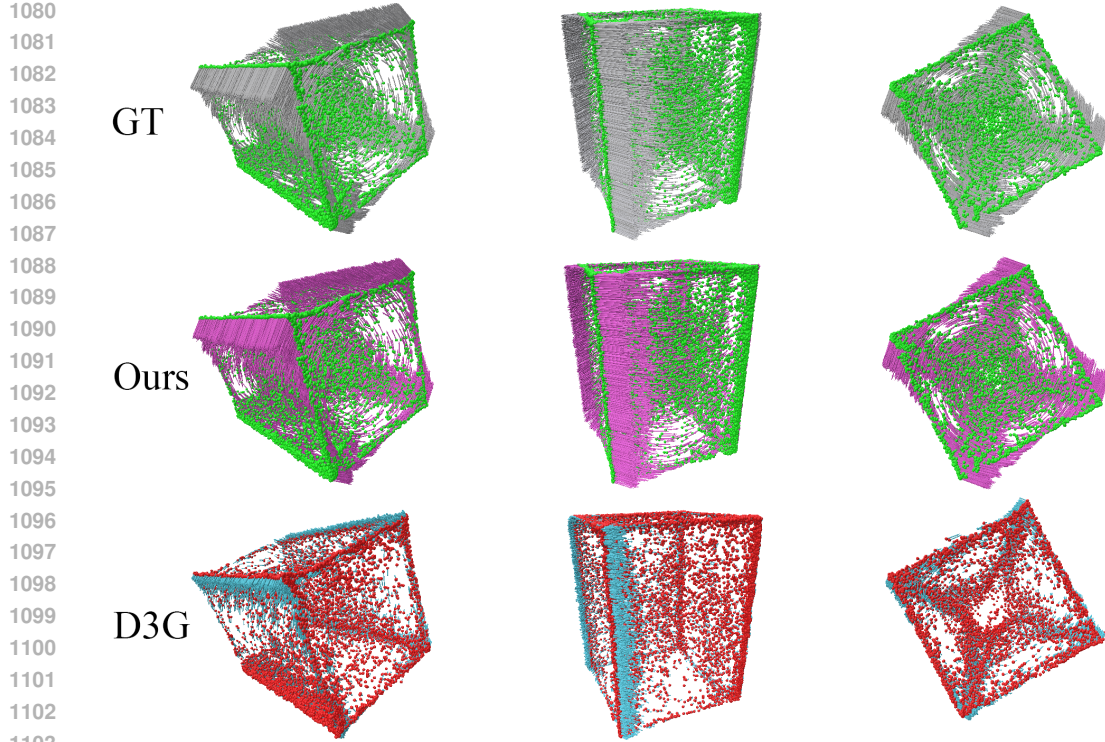
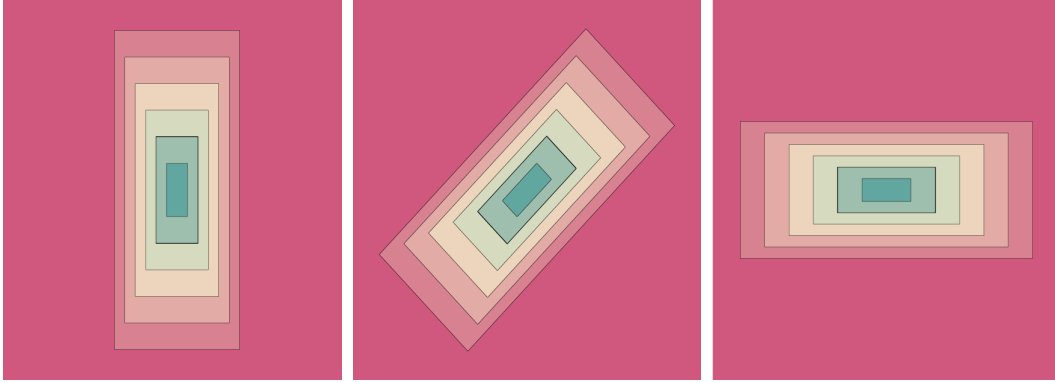


Figure 9: Evaluation of unseen timestamps in the rotating colored box scene.

set, we use a reference scene of a colored box and apply the ground-truth flow  $\phi_t$ , to create a dynamic scene consisting of multi-view captures over 12 time stamps. For the flow  $\phi_t$ , the following

parameters are used:  $S = I_3$ ,  $R = \begin{bmatrix} \cos 2\pi t & -\sin 2\pi t & 0 \\ \sin 2\pi t & \cos 2\pi t & 0 \\ 0 & 0 & 1 \end{bmatrix}$ . Figure 8 shows selected images from

the training set. Two training procedures are considered: (i) a baseline approach using only the reconstruction loss, similar to the D3G model; and, (ii) our approach where both the reconstruction loss and the ReMatching loss are optimized. For the ReMatching loss, we employ the global rigid motion prior class  $\mathcal{P}_{II}$ . Figure 9 compares ground-truth images from time-stamps unseen during training with the model’s predicted renderings. Notably, the ReMatching loss allows the model to generalize in alignment with the ground-truth flow  $\phi_t$ . This is a result of the ReMatching objective’s ability to converge to matched priors  $u_t$  that accurately recover the ground truth velocities  $\frac{\partial}{\partial t} \phi_t$ . To further support this claim, Figure 10 illustrates the velocities of the dynamic Gaussian centers,  $\{\dot{\mu}^i(t)\}$ . These results demonstrate that the ReMatching loss effectively controls  $\{\dot{\mu}^i(t)\}$ , resulting with solutions  $\{\dot{\mu}^i(t)\}$  that match the prior class  $\mathcal{P}_{II}$  and recover the ground-truth velocities  $\frac{\partial}{\partial t} \phi_t$ .

Figure 10: Visualization of  $\{\mu^i(t), \hat{\mu}^i(t)\}$  from multiple timestamps.Figure 11: Visualization of the training set made up of three functions  $\psi_{\text{GT}}(\cdot, t_i)$ .

1120 **The  $V = C^1(\mathbb{R}^d)$  case.** We evaluate these settings for  $d = 2$ . Using the flow described above, we  
1121 define the following ground-truth scalar function,  $\psi_{\text{GT}}: \mathbb{R}^2 \rightarrow \mathbb{R}$ , as:

$$1122 \quad \psi_{\text{GT}}(\mathbf{x}, t) = \|\phi_t^{-1}(\mathbf{x})\|_{\infty} - b. \quad (77)$$

1123 The training data is constructed using three time-stamps, specifically  $t \in \{0.0, 0.25, 0.5\}$ . The  
1124 parameter choices for this procedure are:  $b = 0.2$ , and  $S = \text{diag}\{0.6, 1.4\}$ ,  $\mathbf{R}(t) = \begin{bmatrix} \cos \pi t & \sin \pi t \\ -\sin \pi t & \cos \pi t \end{bmatrix}$ .

1125 Figure 11 visualizes the three distinct functions  $\psi_{\text{GT}}(\cdot, t)$  that constitute the training set. To model  $\psi_t$ ,  
1126 we employ a multi-layer perceptron (MLP) architecture, as described in Section 8.2.1, with the only  
1127 modification of a scalar output dimension in the final layer. For the reconstruction loss, we adopt the  
1128 standard  $L_1$  loss:

$$1129 \quad L_{\text{REC}} = \sum_{i=1}^3 \mathbb{E} \|\psi(\mathbf{x}, t_i) - \psi_{\text{GT}}(\mathbf{x}, t_i)\|. \quad (78)$$

1130 For the ReMatching loss, we employ the global rigid motion prior class  $\mathcal{P}_{II}$ . Two training pro-  
1131 cedures are considered: (1) a baseline approach where only the reconstruction loss is used as the  
1132  
1133

1134 optimization objective, and (2) our suggested approach where both the reconstruction loss and the  
 1135 ReMatching loss are optimized. Figure 12 displays the results of the trained models for times  
 1136  $t \in \{0, 0.0625, 0.125, 0.1875, 0.25, 0.3125, 0.375, 0.4375\}$ . Among these, only  $t \in \{0, 0.25, 0.5\}$ ,  
 1137 shown in the leftmost column, correspond to the training set frames. While both the baseline and our  
 1138 approach perform similarly on the training frames, the results for unseen frames clearly demonstrate  
 1139 the benefits of incorporating the ReMatching loss. Specifically, the ReMatching loss allows the model  
 1140 to recover the ground-truth flow  $\phi_t$ , avoiding the unrealistic distortions observed in the baseline  
 1141 results. To further illustrate this, Figure 13 depicts the matched priors  $u_t$  (white arrows) obtained  
 1142 by solving equation 5, alongside the velocity field of the ground-truth flow  $\frac{\partial}{\partial t}\phi_t$  (green arrows).  
 1143 These comparisons show that the ReMatching loss successfully converges to matched priors  $u_t$  that  
 1144 closely approximate the ground truth. In contrast, the matched priors  $u_t$  obtained with the baseline  
 1145 approach (without the ReMatching loss) deviate significantly from the ground truth. This emphasizes  
 1146 the importance of the reprojection procedure. Not all velocity fields in the prior class are suitable for  
 1147 guiding the optimization process, but the ReMatching loss ensures convergence to an appropriate  
 1148 prior, enabling an accurate recovery of the underlying flow.

## 1149 8.5 RUNTIME AND CONVERGENCE ANALYSIS

1151 We note that our framework is applied solely during the training phase of the algorithm, leaving  
 1152 inference times unaffected. To evaluate computational efficiency, we measured the average time  
 1153 (seconds) for a forward and backward pass over 100 iterations for varying sizes  $n$  of Gaussians sets.  
 1154 Figure 14 presents the results, comparing the computation time of the ReMatching framework to the  
 1155 D3G baseline. The runtime analysis was conducted on a single NVIDIA RTX A6000.

1156 To examine the convergence of the reconstruction model, we compare the loss convergence curves  
 1157 of the D3G (Yang et al., 2023) model and our model. Figure 15 shows that the addition of the  
 1158 ReMatching loss does not affect the convergence behavior of the optimization. We also show the  
 1159 loss curve of the ReMatching loss itself in figure 16. It is important to note that for the ReMatching  
 1160 formulation, the optimal solution does not necessarily achieve 0 loss, similarly to the APM procedure  
 1161 (Deutsch, 1992). Instead, it achieves the lowest loss possible given the reconstruction task and  
 1162 selected prior.

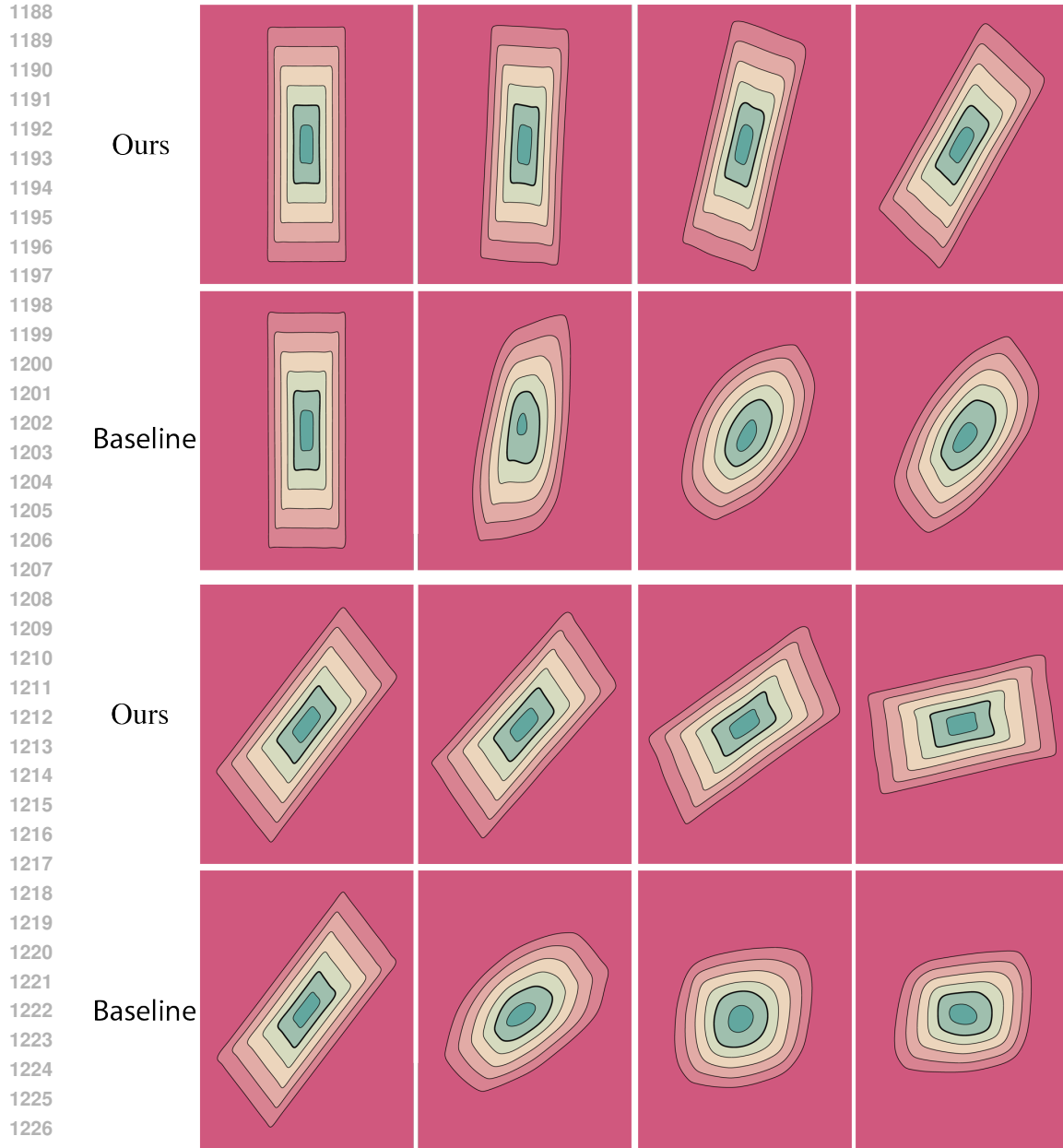
## 1163 8.6 ABLATION OF FRAMEWORK HYPERPARAMETERS

1165 In this section, we present an ablation study on key hyperparameters introduced by the ReMatching  
 1166 framework: i) the weight of the ReMatching loss,  $\lambda$ , as defined in equation 8; ii) maximum number  
 1167 of parts selection,  $k$ , for the adaptive prior class; and iii) the weight of the entropy loss (equation 30),  
 1168 used to optimize the learned part assignments when employing an adaptive prior class.

1169 **ReMatching loss weight.** We note first that a consistent value of  $\lambda = 0.001$  was used across all  
 1170 scenes experimented with in section 6, already demonstrating the robustness of this parameter. To  
 1171 further test this, we conducted experiments on the Hell Warrior and Lego scenes from the D-NeRF  
 1172 dataset, evaluating how different  $\lambda$  values influence solution quality. Figure 17 shows these findings.  
 1173 We note that for the Hell Warrior scene, we employed the  $\mathcal{P}_{IV}$  prior class, while the Lego scene  
 1174 used the  $\mathcal{P}_V$  class. The results indicate stable improvement within the range of  $\lambda \in [5e-4, 5e-3]$ ,  
 1175 while small values  $\lambda \leq 5e-5$  aligns with the baseline. Larger values,  $\lambda \geq 1e-2$ , may compete with  
 1176 the reconstruction loss, leading to suboptimal solutions.

1178 **Maximum number of parts selection.** To assess the impact of the hyperparameter  $k$ , we selected  
 1179 the Mutant scene, which aligns with the adaptive prior class  $\mathcal{P}_{IV}$ , and the Lego scene, corresponding  
 1180 to the adaptive prior class  $\mathcal{P}_V$  and evaluated how varying  $k$  affects solution quality. Figure 18 presents  
 1181 the results of this analysis. The findings suggest relatively stable performance within the range  $k = 5$   
 1182 to  $k = 15$ , offering flexibility in selecting  $k$  based on leveraging prior knowledge about the expected  
 1183 number of moving parts in the scene.

1185 **Entropy loss weight.** Similar to the  $\lambda$  hyperparameter, the entropy loss weight was kept fixed  
 1186 across all experiments in section 6. To further examine its impact, we evaluated its influence on  
 1187 performance with varying weight values for the Hell Warrior scene. Figure 19 presents the results of  
 this experiment, demonstrating stable performance within the range  $[1e-4, 1e-3]$ .

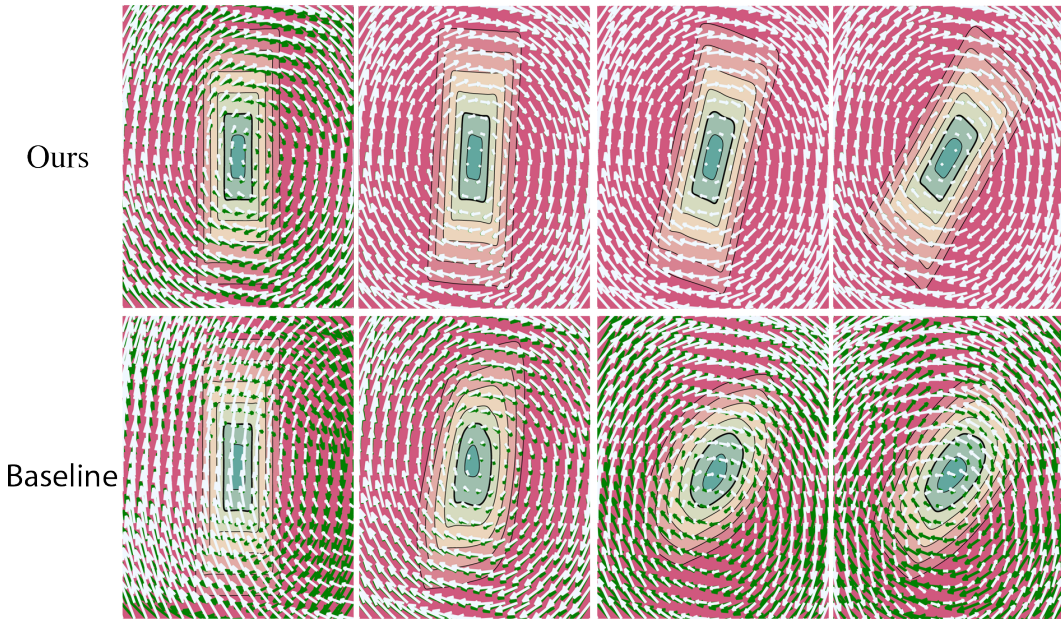


1228 Figure 12: Comparisons of  $\psi_t$  converged solutions between the baseline and our approach, displayed  
1229 in order of increasing time (left to right, top to bottom).

### 1230 8.7 ADDITIONAL QUALITATIVE EVALUATION

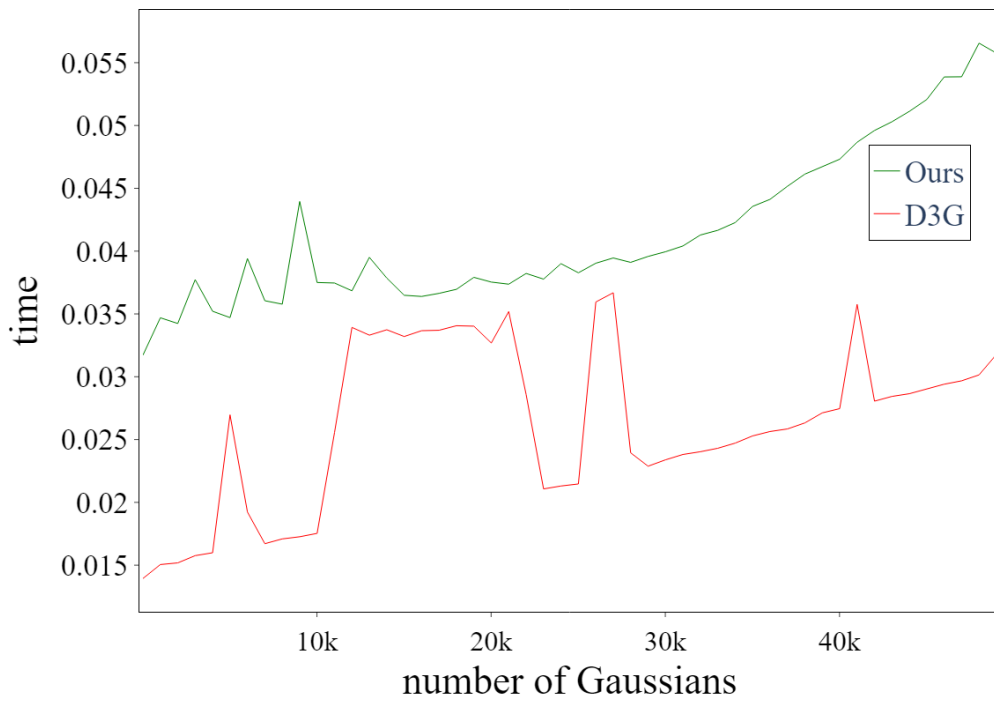
1232 To further support the qualitative results presented in Figures 2 and 3, the supplementary material  
1233 includes additional evidence showcasing novel-view video reconstruction results. These compar-  
1234 isons highlight the performance of our model relative to baseline approaches, providing a more  
1235 comprehensive validation of its efficacy.  
1236  
1237  
1238  
1239  
1240  
1241

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263



1264 Figure 13: Comparisons of the converged  $u_t$  (white arrows) between the baseline and our approach  
1265 that uses the ReMatching loss. The ground-truth velocities,  $\frac{\partial}{\partial t} \phi_t$ , are shown as green arrows. When  
1266 the matched  $u_t$  aligns with the ground truth, the green arrows become indistinguishable.  
1267

1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292



1293 Figure 14: Combined average time (seconds) for a forward and backward pass for varying size of  $n$ .  
1294  
1295



1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

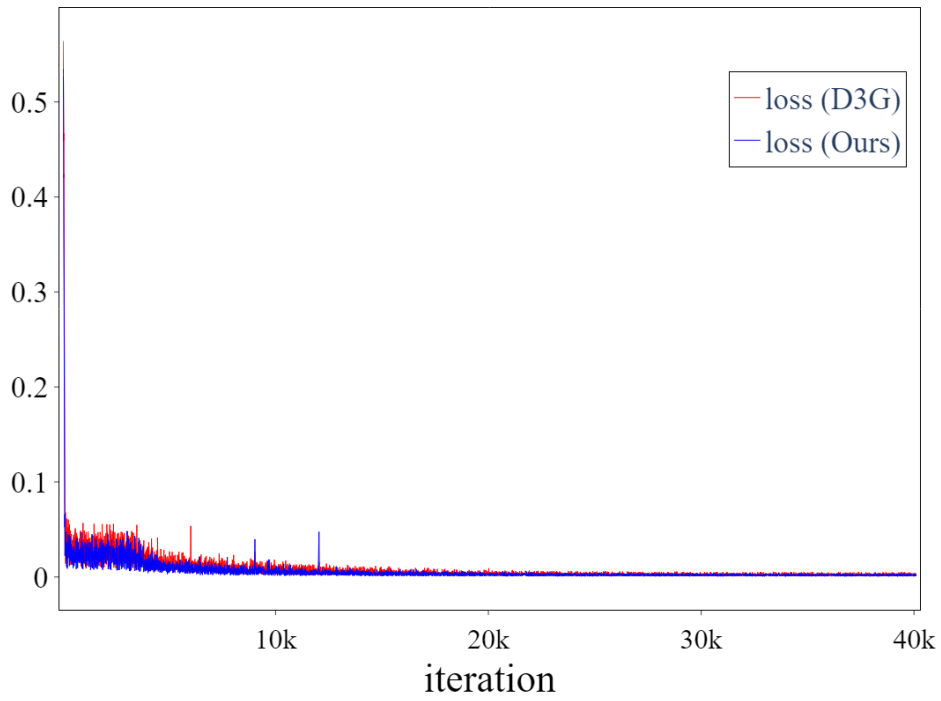


Figure 15: Loss curves report of our model and D3G (Yang et al., 2023) over 40k training iterations.

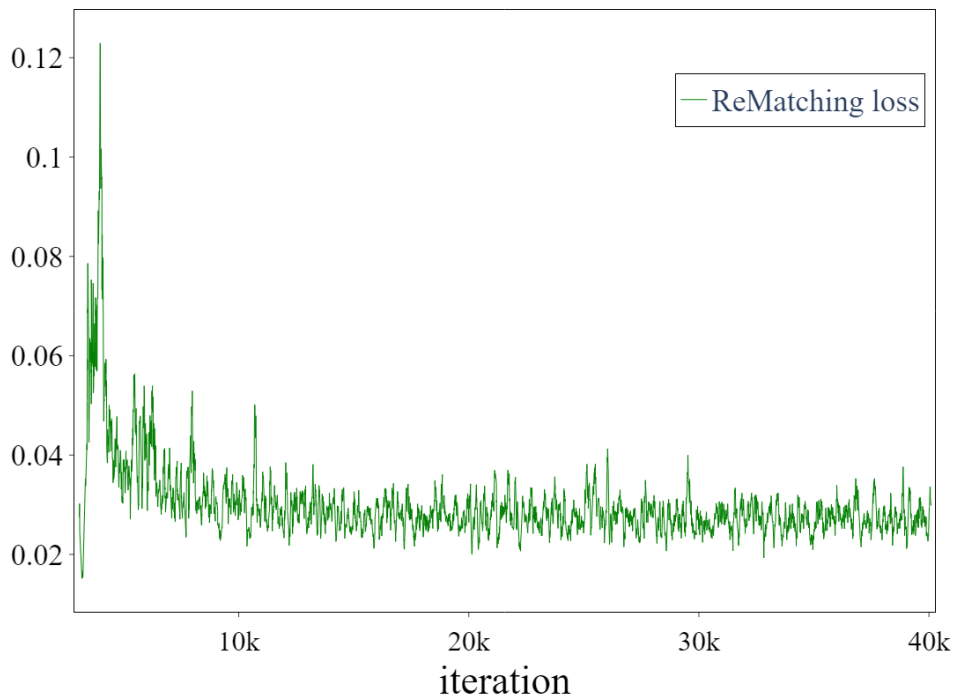


Figure 16: Loss curve report for the ReMatching loss, showing a running average with a window size of 20.

1350  
 1351  
 1352  
 1353  
 1354  
 1355  
 1356  
 1357  
 1358  
 1359  
 1360  
 1361  
 1362  
 1363  
 1364  
 1365  
 1366  
 1367  
 1368  
 1369  
 1370  
 1371  
 1372  
 1373  
 1374  
 1375  
 1376  
 1377  
 1378  
 1379  
 1380  
 1381  
 1382  
 1383  
 1384  
 1385  
 1386  
 1387  
 1388  
 1389  
 1390  
 1391  
 1392  
 1393  
 1394  
 1395  
 1396  
 1397  
 1398  
 1399  
 1400  
 1401  
 1402  
 1403



Figure 17: Effect of the weight parameter  $\lambda$  on the PSNR evaluation metric for the Hell Warrior scene (left) and the Lego scene (right).

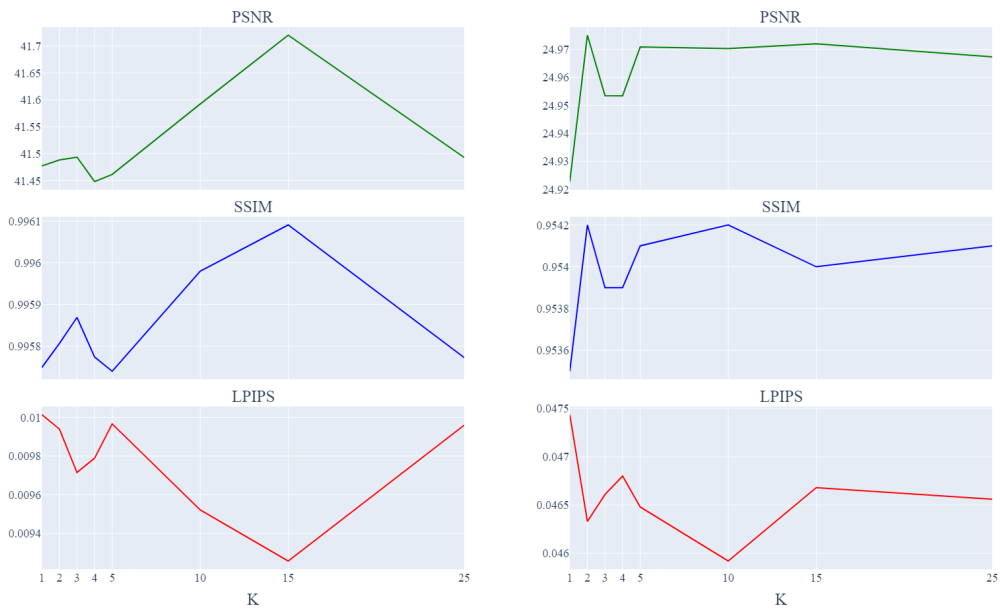


Figure 18: Impact of varying  $k$  values on the PSNR evaluation metric for the Mutant scene (left) and the Lego scene (right).

1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457



Figure 19: Impact of varying entropy loss weights on the PSNR evaluation metric for the Hell Warrior scene.