# SceNeRFlow: Time-Consistent Reconstruction of General Dynamic Scenes

Edith Tretschk[1]  Vladislav Golyanik[1]  Michael Zollhöfer[2]
Aljaž Božič[2]  Christoph Lassner[2]  Christian Theobalt[1]

[1]Max Planck Institute for Informatics, Saarland Informatics Campus  [2]Meta Reality Labs Research

## Abstract

*Existing methods for the 4D reconstruction of general, non-rigidly deforming objects focus on novel-view synthesis and neglect correspondences. However, time consistency enables advanced downstream tasks like 3D editing, motion analysis, or virtual-asset creation. We propose SceNeRFlow to reconstruct a general, non-rigid scene in a time-consistent manner. Our dynamic-NeRF method takes multi-view RGB videos and background images from static cameras with known camera parameters as input. It then reconstructs the deformations of an estimated canonical model of the geometry and appearance in an online fashion. Since this canonical model is time-invariant, we obtain correspondences even for long-term, long-range motions. We employ neural scene representations to parametrize the components of our method. Like prior dynamic-NeRF methods, we use a backwards deformation model. We find non-trivial adaptations of this model necessary to handle larger motions: We decompose the deformations into a strongly regularized coarse component and a weakly regularized fine component, where the coarse component also extends the deformation field into the space surrounding the object, which enables tracking over time. We show experimentally that, unlike prior work that only handles small motion, our method enables the reconstruction of studio-scale motions.*

## 1. Introduction

One criterion for the proper reconstruction of a deforming scene is time consistency, *i.e.* correspondences across time. Reconstructing general dynamic objects from RGB input in a time-consistent manner is a little explored [9, 54], but challenging and highly relevant research direction. Establishing correspondences is equivalent to factorizing the reconstruction into time-varying deformations and a time-invariant geometry model. High-level tasks that go beyond novel-view synthesis benefit immensely from such a deeper analysis of the scene, namely a reconstruction with long-range, long-term dense 3D correspondences. For example,
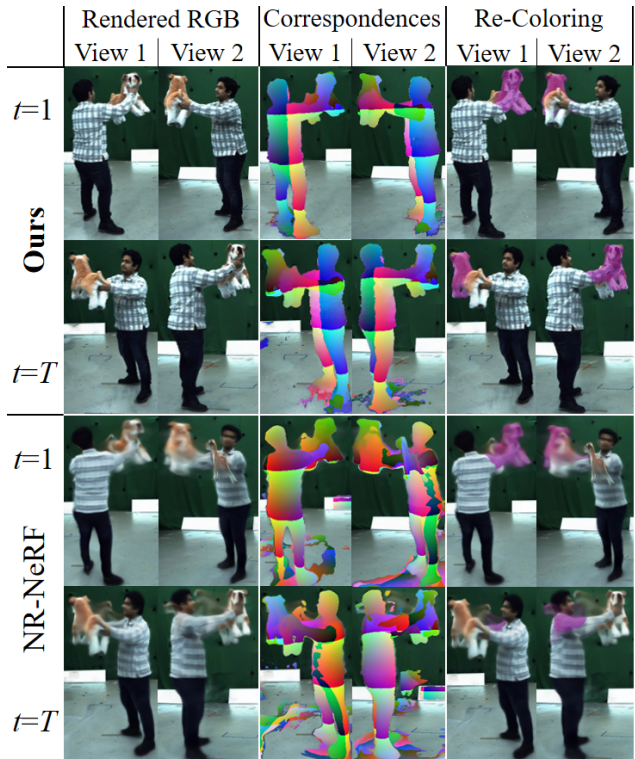


Figure 1. **SceNeRFlow.** Our NeRF-based method reconstructs a general non-rigid scene from multi-view videos *with time consistency*. Here, a person with a plush dog rotates 180° from time $t=1$ until $t=T$. Novel view 1 at $t=1$ is consistent with novel view 2 at $t=T$ (placed opposite of novel view 1) for our method, but not for NR-NeRF [78]. This enables, *e.g.*, time- and view-consistent re-coloring. We color the correspondences according to 3D positions in static canonical space, which differs between both methods.

having access to such a virtual dynamic 3D object is a crucial step towards sophisticated 3D editing or estimating a motion model for virtual-asset creation [90]. Almost all existing work on reconstructing general dynamic objects either only obtains time consistency for small motion or relaxes it to very small time windows. In this work, we explore the setting of time consistency for large motions.

Prior work [16, 18, 42, 59, 60, 63, 78, 85] for non-

rigid 3D reconstruction based on NeRF [50] focuses on novel-view synthesis and does not aim for long-range correspondences. Almost all dynamic-NeRF papers (except for [44, 63, 78], which only handle small motions) weaken the long-term consistency of the reconstruction by design by having a time-varying geometry and/or appearance model. Thus, only rather simple tasks like replaying the input scene under novel views are straightforward. Furthermore, as we will show in our experiments, this time dependency improves the novel-view rendering quality but loosens long-term correspondences. In contrast, we explore the other extreme, where *only* deformations are time-variant.

Classical, non-NeRF-based work on 4D reconstruction only handles small motion [5, 33, 36, 66, 67] or is not time-consistent [3, 15, 55, 88]. Similarly, scene flow [91] focuses on the deformations and not a full reconstruction, and exhibits correspondence drift [29], similar to optical flow.

In this paper, we propose *SceNeRFlow* (Scene Flow + NeRF) to tackle time-consistent reconstruction of a general, non-rigidly deforming scene; see Fig. 1. Our method is NeRF-based, trained per scene, and by design estimates each timestamp's deformation of a time-invariant canonical model, thereby counteracting correspondence drift. As is common for dynamic NeRFs [59, 63, 78], we employ a backward deformation model. Standard reconstruction techniques like online optimization, coarse-and-fine deformation decomposition, and rigidity regularization turn out to be insufficient for large motion. While category-specific priors like an articulated human skeleton would help by normalizing out large motion, this work tackles the general setting, where such prior knowledge is not available. Instead, we crucially find that carefully "extending the deformation field" is the minimally necessary change to make backwards deformation modeling work for large non-rigid motion.

Since ours is the first NeRF-based work with time consistency for large motion, we want to focus on the core problem of long-term correspondences. We therefore use multi-view input to avoid the need for correctly modeling the deformations of occluded geometry that arises in the monocular setting. Importantly, our scenes exhibit significantly larger motion than any prior time-consistent general reconstruction method could handle. Our experiments show how SceNeRFlow enables the time-consistent reconstruction even of studio-scale motion without any category-specific priors, *e.g.* without a human skeleton.

In summary, our **contributions** are (1) an end-to-end differentiable, time-consistent 4D reconstruction method for general dynamic scenes from multi-view RGB input from static cameras, which is built around (2) a general approach to make backward deformation models work for larger motion via "extending the deformation field" (Sec. 3.2).

## 2. Related Work

**3D Correspondence Estimation.** For temporal data, Vedula *et al.*'s seminal work [81] introduced scene flow as the 3D equivalent of optical flow, which led to many approaches tackling the problem [6, 28, 64, 75, 82, 87], as a recent survey [91] summarizes. Lately, learning-based methods [24, 41, 45, 51] focus on 3D point clouds as input. More generally, non-rigid 3D point-cloud registration [14] estimates correspondences of provided 3D point clouds. We go beyond just estimating 3D correspondences to create a full reconstruction (including appearance) in 3D space from only 2D input data, in an end-to-end differentiable pipeline.

**4D Reconstruction of General Non-Rigid Scenes.** Multiple lines of research target 4D reconstruction without relying on NeRF. As a recent survey [79] discusses, monocular RGB input has traditionally been tackled by Non-Rigid Structure-from-Motion [20, 22, 36, 58, 66] and Shape-from-Template methods [5, 12, 33, 56, 67], although learning-based approaches [34, 35, 40, 84, 86, 88] have grown in popularity over the last years. Due to the restrictive input setting, they either only handle small motions or lack time consistency. Slightly larger motions are possible with a single RGB-D camera [98], although these methods [7, 8, 10, 25, 26, 30, 43, 55, 69, 70, 97] all update [13] both previously unseen *and previously seen* parts of the canonical model over time during their online optimization, thereby losing time consistency (except for the template-based [97]). Multi-view input has seen comparatively less work over the years [21, 38, 80, 92]. Fusion4D [15, 57] operates in real time but also modifies the canonical model over time. Mustafa *et al.* [54] track a mesh through a temporal sequence according to optical flow and refine its topology for each timestamp, but do not reconstruct the appearance. Unlike such multi-view mesh-tracking approaches [9, 95], our method is end-to-end differentiable and can easily integrate recent advances in neural scene representations [77]. Bansal *et al.* [3] obtain impressive but time-inconsistent novel-view results.

**Dynamic NeRFs.** 2D neural rendering [76] and 3D neural scene representations [77] based on NeRF [50] have in recent years seen great success. Apart from early work [46], current methods [16, 18, 42, 49, 59, 60, 63, 78, 85] for dynamic scene reconstruction with neural representations build on NeRF. Some approaches focus on surface [31], fast [17, 27], generalizable [83], or depth-supported [1] reconstruction. Our method is most related to this field of work but differs in its goal: we do not primarily aim for novel-view synthesis but rather for a *time-consistent* reconstruction. A few prior methods [44, 63, 78] do obtain time consistency but either only in synthetic [63] or small-motion [44, 78] settings. Furthermore, most works use monocular input [19] and only a few [39, 44, 71, 83] explore multi-view input, where the latter forego long-term

time consistency or only handle small motion [44].

## 3. Method

SceNeRFlow takes as input multi-view RGB images of size $h \times w$ over $T$ consecutive timestamps from $C$ static cameras with known extrinsics and intrinsics, *i.e.* $\{\mathbf{I}^{c,t} \in [0,1]^{h \times w \times 3}\}_{c=1,t=1}^{C,T}$, and associated background images $\{\mathbf{B}^c\}_c$. In a time-consistent manner, it then reconstructs the general dynamic scene as a time-invariant, NeRF-style canonical model of the geometry and appearance, with time-dependent deformations. The optimization proceeds in an online manner: we build a canonical model from the first timestamp (Sec. 3.1) and then track it frame-by-frame through the temporal input sequence (Sec. 3.2). Fig. 2 provides an overview of the pipeline.

### 3.1. Constructing the Canonical Model

**Canonical Model.** The canonical model $m$ encodes the geometry and appearance in a time-independent manner. Specifically, we use a *HashMLP* [53] to represent opacity $\sigma$ and RGB color $\mathbf{c} \in [0,1]^3$ for any 3D point $\mathbf{x}$: $(\sigma, \mathbf{c}) = m(\mathbf{x})$. A HashMLP combines a *hash grid* $v$ with a subsequent shallow MLP $M$: $m(\mathbf{x}) = M(v(\mathbf{x}))$, which is significantly faster than NeRF's [50] pure MLP. A hash grid consists of about a dozen voxel grids of increasing resolution, each containing learnable features. To evaluate, each grid is queried via trilinear interpolation and the resulting features from all grids are concatenated. Crucially, each voxel grid is implemented via hashed indexing into an array of feature vectors rather than as a dense voxel grid.

**Rendering.** Since we take 2D images as input, we need to render $m$ into 2D. To this end, we follow the quadrature discretization of volumetric rendering from NeRF [50] for a ray $\mathbf{r}(s) = \mathbf{o} + s\mathbf{d}$ with origin $\mathbf{o} \in \mathbb{R}^3$ and direction $\mathbf{d} \in \mathbb{R}^3$:

$$\mathbf{C}(\mathbf{r}) = \sum_{i=1}^{S} w_i \mathbf{c}_i, \text{ with } w_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)(1 - \exp(-\sigma_i \delta_i)),$$

(1)

where $i$ indexes $S$ discrete samples $\{\mathbf{r}(s_i)\}_i$ along the ray; $(\sigma_i, \mathbf{c}_i) = m(\mathbf{r}(s_i))$ are the opacity and color of the $i$-th sample, respectively; and $\delta_i = s_{i+1} - s_i$. Like NeRF [50], we use stratified sampling of $S$ evenly sized intervals between the near plane and far plane to obtain $\{s_i\}_i$.

When a background color $\mathbf{c}_{back}$ is provided (*e.g.* from $\mathbf{B}^c$), we composite it at the end of the ray:

$$\mathbf{C}(\mathbf{r}, \mathbf{c}_{back}) = \mathbf{C}(\mathbf{r}) + \left(1 - \sum_i w_i\right)\mathbf{c}_{back}. \qquad (2)$$

**Losses.** To obtain the canonical model, we iteratively optimize it on images $\{\mathbf{I}^{c,1}\}_c$ from $t=1$ for $20k$ iterations. Each iteration uses a batch of $R$ rays with an $\ell_1$ reconstruction loss w.r.t. the ground-truth color $\mathbf{I}^{c,1}(\mathbf{r}_r)$:

$$\mathcal{L}_{\text{rec}} = \frac{1}{R} \sum_{r=1}^{R} \|\mathbf{C}(\mathbf{r}_r, \mathbf{c}_{back,r}) - \mathbf{I}^{c,1}(\mathbf{r}_r)\|_1. \qquad (3)$$

However, using only $\mathcal{L}_{\text{rec}}$ leads to a canonical model with floating geometry artifacts. We remove them by steering the model to commit to foreground *or* background by encouraging $\sum_i w_i$ to be 1 or 0 via the beta distribution [46]:

$$\mathcal{L}_{\text{back}} = \frac{1}{R} \sum_r \log\left(\sum_i w_{r,i}\right) + \log\left(1 - \sum_i w_{r,i}\right), \quad (4)$$

where $w_{r,i}$ is $w_i$ of ray $r$. We also discourage transparent geometry by encouraging each $w_i$ to be 0 or 1 via a mixture of two Laplacian distributions [65]:

$$\mathcal{L}_{\text{hard}} = -\frac{1}{RS} \sum_r \sum_i \log\left(e^{-w_{r,i}} + e^{-(1-w_{r,i})}\right). \quad (5)$$

We note that these losses do not use foreground masks; the canonical model needs to learn on its own whether it can be transparent and use the background image $\mathbf{B}$ or not.

The total loss for the canonical model is thus:

$$\mathcal{L}_{\text{canon}} = \mathcal{L}_{\text{rec}} + \lambda_{\text{back}} \mathcal{L}_{\text{back}} + \lambda_{\text{hard}} \mathcal{L}_{\text{hard}}, \qquad (6)$$

where $\lambda_{\text{back}}, \lambda_{\text{hard}} \in \mathbb{R}$ are loss weights. Since we want the canonical model to be time-consistent, we keep its parameters fixed after constructing it from the first timestamp.

### 3.2. Optimizing per Timestamp

Given the canonical model from $t=1$, we next reconstruct its deformations for the remaining timestamps $t>1$.

**Space Warping.** We model the deformations at time $t$ like prior dynamic-NeRF work: we use backwards space warping $d_c(\mathbf{x}; \theta_c^t) = \mathbf{x} + \Delta_c(\mathbf{x}; \theta_c^t) = \mathbf{x}'$, where $\mathbf{x}$ is a 3D point in deformed world space, $\Delta_c$ outputs a coarse offset (fine offsets $\Delta_f$ will be introduced later), $\mathbf{x}'$ is in undeformed canonical space, and $\theta_c^t$ are the time-dependent parameters of $d_c/\Delta_c$. We use a HashMLP to parametrize $\Delta_c$, *i.e.* there is one HashMLP per timestamp for coarse deformations. To query the canonical model from world space at $t>1$, we first undo the deformation: $(\sigma(\mathbf{x},t), \mathbf{c}(\mathbf{x},t)) = m(d_c(\mathbf{x}; \theta_c^t))$. When rendering, this leads to view-consistent ray bending: $\{\mathbf{r}(s_i)\}_i$ in world space becomes $\{d_c(\mathbf{r}(s_i); \theta_c^t)\}_i$ in canonical space. We can then apply Eq. 1 or Eq. 2 to the bent ray to render, which in turn enables us to use $\mathcal{L}_{\text{rec}}$ to optimize for the deformation parameters at time $t$ using $\{\mathbf{I}^{c,t}\}_c$.

**Frame-Wise Tracking.** However, naïvely reconstructing the entire scene by optimizing all timestamps at once does not converge to a recognizable canonical model when the scene contains large motion. Furthermore, keeping all $\{\mathbf{I}^{c,t}\}_{c,t}$ in memory at once is expensive in practice.
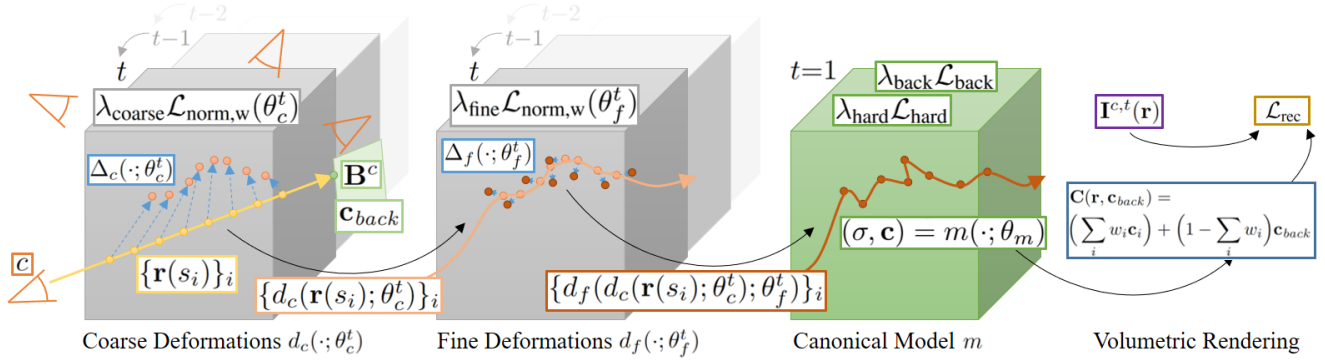
Figure 2. **SceNeRFlow Overview.** *Input:* We take as input a multi-view RGB video $\{\mathbf{I}^{c,t}\}_{c,t}$ of a general dynamic scene from $C$ cameras with background images $\{\mathbf{B}^c\}_c$ and known camera parameters. *Output:* From this, we build a NeRF-based [50], time-consistent 4D reconstruction. *Model:* To render a ray $\mathbf{r}$ from camera $c$, we first sample discrete points $\{\mathbf{r}(s_i)\}_i$ along the straight ray (with background color $\mathbf{c}_{back}$ at the end), then coarsely bend the ray with the coarse deformations $d_c$, and then finely bend the resulting ray with the fine deformations $d_f$. We then query the *canonical model $m$* at the resulting positions. $m$ represents geometry (opacity $\sigma$) and appearance (color $\mathbf{c}$) volumetrically for any 3D point, in a time-invariant manner. Finally, we use NeRF-style volumetric rendering to obtain the ray's color $\mathbf{C}$. We parametrize all three components with HashMLPs [53]. *Training:* We construct $m$ from time $t=1$ using a reconstruction loss ($\mathcal{L}_{rec}$ w.r.t. the ground-truth color $\mathbf{I}^{c,t}(\mathbf{r})$) and geometric regularizers ($\mathcal{L}_{back}$ and $\mathcal{L}_{hard}$), and keep $m$'s parameters $\theta_m$ fixed for future timestamps, which leads to time consistency. For $t>1$, we split the deformations into a strongly regularized coarse component $\Delta_c$ and a weakly regularized fine component $\Delta_f$. Our online optimization performs frame-wise tracking to handle large motion and we thus initialize the deformation parameters $\theta_c^t$ and $\theta_f^t$ at time $t$ with $\theta_c^{t-1}$ and $\theta_f^{t-1}$. To handle a peculiarity of backward deformation models in this tracking setting, we regularize $\Delta_c$ to "extend the deformation field" ($\mathcal{L}_{norm,w}$).
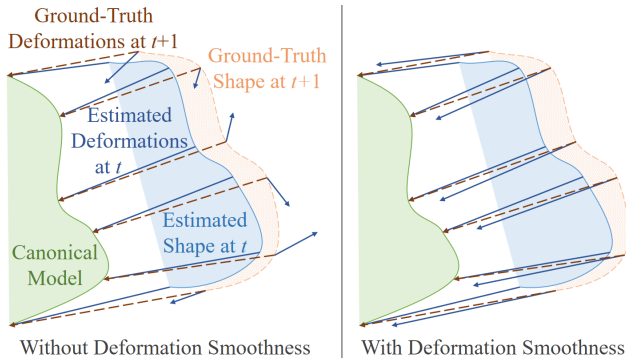


Figure 3. **Extending the Deformation Field for Tracking.** Without smoothness, the estimated backwards deformations at time $t$ give a bad initialization for $t+1$. Smoothness initializes the deformations closer to the ground truth.

We thus propose online, timestamp-by-timestamp tracking. Since $t=1$ has, by construction, zero offsets, we set the last layer of $\Delta_c(\cdot;\theta_c^1)$ to zeros [78] and do not optimize the deformations at $t=1$. After reconstructing time $t$, we fix $\theta_c^t$ and proceed with $t+1$, where we initialize $\theta_c^{t+1}$ with $\theta_c^t$.

**Extending the Deformation Field.** Unfortunately, naïvely employing a HashMLP for $d_c$ fails to reconstruct the scene for any deformed state that significantly differs from the canonical model. This is because, at time $t+1$, the dynamic object resides in a slightly different place in world space, part of which was empty space at time $t$ and is hence

not well initialized by $d_c(\cdot;\theta_c^t)$. This becomes especially prevalent when the object has undergone large motion. We propose to mitigate this failure of the backward model by extending the deformations into the area surrounding the dynamic object at time $t$, as Fig. 3 illustrates.

To this end, we employ two means: (1) We find that reducing the resolution of the coarsest grid of the deformation hash grid to $32^3$ stabilizes tracking. We hypothesize that when the object enters a previously empty voxel in a fine grid, the uninitialized latent codes of this voxel are too quickly too influential in the trilinear interpolation of the latent codes. These uninitialized latent codes thus do not obtain the right values before being relevant but rather negatively impact the deformations, leading to artifacts and a lack of large-scale smoothness. (2) However, this structural change gives an unpredictable, badly-controlled smoothness, similar to an MLP. We encourage well-behavedness via a smoothness loss, which uses the inherent smoothness of the coarse grid to propagate its influence.

**Smoothness Loss.** To stabilize the tracking, we propose to impose a smoothness loss on the deformations. Because our deformation model is continuous and fully differentiable, we do not need to discretize the loss. Instead, we influence the local behavior directly via the (spatial) Jacobian $\mathbf{J} \in \mathbb{R}^{3\times3}$ of the deformations: $\mathbf{J} = \mathbf{J_x} = \frac{\partial d_c(\mathbf{x};\theta_t)}{\partial \mathbf{x}}$. As is common for general reconstruction methods [79], we assume the object to deform locally in an as-rigid-as-possible manner [72]. Specifically, we take inspiration from Ner-

4

fies's elastic loss on $\mathbf{J}$ [59]. However, their loss is computationally expensive because it needs to compute all rows of $\mathbf{J}_{\mathbf{r}_r(s_i)}$ for each sample $\mathbf{r}_r(s_i)$ (which takes three backward passes during the forward pass of the loss computation) and then performs an SVD of $\mathbf{J}_{\mathbf{r}_r(s_i)}$. Denoting the identity matrix by $\mathbf{I}$, we can avoid the SVD and allow for trivial computation of gradients by relaxing the constraint from the rotation group $SO(3) = \{\mathbf{A} \in \mathbb{R}^{3\times3} | \mathbf{A}^\top \mathbf{A}{=}\mathbf{I}, \det \mathbf{A}{=}1\}$ to the orthogonal group $O(3) = \{\mathbf{A} \in \mathbb{R}^{3\times3} | \mathbf{A}^\top \mathbf{A}{=}\mathbf{I}\}$, which additionally allows reflections since $\det \mathbf{A}{=} \pm 1$ for $\mathbf{A} \in O(3)$. We can then encourage rigidity via:

$$\mathcal{L}_{\text{rigid}} = \frac{1}{9RS} \sum_r \sum_i \sum_{j,k} \left| (\mathbf{J}_{\mathbf{r}_r(s_i)}^\top \mathbf{J}_{\mathbf{r}_r(s_i)} - \mathbf{I})_{j,k} \right|, \quad (7)$$

where $\mathbf{A}_{j,k} \in \mathbb{R}$ is the entry of matrix $\mathbf{A}$ at index $(j, k)$.

We can compute this prior even faster by reducing the need for three backward passes to just one, lowering overall training time by a factor of $2-3$. To this end, we first note that norm preservation is an equivalent definition of $O(3)$:

$$\mathbf{A} \in O(3) \iff \forall \mathbf{e} \in \mathbb{R}^3 : \|\mathbf{A}\mathbf{e}\|_2 = \|\mathbf{e}\|_2. \quad (8)$$

We next exploit the fact that automatic differentiation [23, 61, 73] computes Jacobian-vector products $\mathbf{J}^\top \mathbf{e}$ in a single backward pass, where $\mathbf{e}$ is an arbitrary vector. Since $\mathbf{J}^\top = \mathbf{J}$ for $\mathbf{J} \in O(3)$, we can compute the norm of $\mathbf{J}\mathbf{e}$ as $\|\mathbf{J}\mathbf{e}\|_2 = \|\mathbf{J}^\top \mathbf{e}\|_2$. The following loss then encourages norm preservation and only requires one backward pass:

$$\mathcal{L}_{\text{norm}} = \frac{1}{RS} \sum_r \sum_i \mathbb{E}_{\mathbf{e}} \left[ \left| \|\mathbf{J}_{\mathbf{r}_r(s_i)}^\top \mathbf{e}\|_2 - 1 \right| \right], \quad (9)$$

where $\mathbf{e}$ is distributed uniformly on the unit sphere and hence $\|\mathbf{e}\|_2 = 1$. We note that it is sufficient to only consider vectors on the unit sphere due to linearity. In practice, we approximate this expectation with one random sample.

**Weighting the Smoothness Loss.** To extend the deformation field as discussed earlier, we could naïvely encourage smoothness uniformly everywhere. However, this restricts empty space too much, leading it to push back against object deformations. We thus focus on the object by weighting $\mathcal{L}_{\text{norm}}$ by $\hat{\sigma}_{r,i} = \exp(-\sigma_{r,i}\delta)$, where $\sigma_{r,i}$ is $\sigma_i$ of the $r$-th ray and we do not backpropagate into the opacity [78].

We now need to regularize the space *around* the object. However, this space is hard to locate efficiently and we approximate it by max-pooling over $\{\hat{\sigma}_{r,i}\}_i$ along ray $r$, with a window size of $1\%$ of the ray length.

This still makes the space around the object too stiff. We thus weaken the regularization by dividing the resulting weight by $u$ if $\hat{\sigma}_{r,i}$ and the weight after max-pooling differ by at least a factor of $u$. We empirically set $u{=}10$.

Finally, regularizing very small offsets, which are widespread during early timestamps, tends to collapse the

network. We hence only regularize deformations that are not very small. We denote this weighted loss as $\mathcal{L}_{\text{norm,w}}(\theta_c^t)$. The supplement contains a full mathematical description.

**Fine Deformations.** In addition to enabling tracking by extending the deformation field, we also use the smoothness loss to stabilize the surface by strongly regularizing its deformations $\Delta_c$. However, this leads to a loss of detail because $\Delta_c$ can now only represent *coarse* deformations. To counteract this, we add *fine* deformations $\Delta_f(\cdot; \theta_f^t)$ on top, after normalizing out (*i.e.*, applying) the coarse deformations: $\mathbf{x}'' = d_f(\mathbf{x}'; \theta_f^t) = \mathbf{x}' + \Delta_f(\mathbf{x}'; \theta_f^t)$, where $\mathbf{x}' = d_c(\mathbf{x}; \theta_c^t)$ for any 3D point $\mathbf{x}$ in world space and we query $m$ at $\mathbf{x}''$. Crucially, we allow for details by using a much weaker weight $\lambda_{\text{fine}}$ for $\mathcal{L}_{\text{norm,w}}(\theta_f^t)$, where the Jacobian is w.r.t. $\mathbf{x}'$. We parametrize $\Delta_f$ with a HashMLP.

**Frame-Wise Tracking Revisited.** We apply tracking to the coarse and fine deformations as follows: At $t{=}1$, we initialize the last layers of $\Delta_c$ and $\Delta_f$ to zeros and do not optimize for the deformations. At any $t{>}1$, we initialize $\theta_c^t$ with the final $\theta_c^{t-1}$, and optimize for $\Delta_c(\cdot; \theta_c^t)$ for $5k$ iterations while setting $\Delta_f = \mathbf{0}$. We then fix $\theta_c^t$, initialize $\theta_f^t$ with the final $\theta_f^{t-1}$, and optimize for $\Delta_f(\cdot; \theta_f^t)$ for $5k$ iterations.

With $\lambda_{\text{coarse}}, \lambda_{\text{fine}} \in \mathbb{R}$, the total loss for any $t{>}1$ is:

$$\mathcal{L}_{\text{time}} = \mathcal{L}_{\text{rec}} + \lambda_{\text{coarse}} \mathcal{L}_{\text{norm,w}}(\theta_c^t) + \lambda_{\text{fine}} \mathcal{L}_{\text{norm,w}}(\theta_f^t). \quad (10)$$

### 3.3. Implementation Details

We next describe how we can speed up our method and correct vignetting effects, and provide optimization details. Fig. 4 visualizes some of these methods.

**Foreground-Focused Batches.** To speed up training, we focus the batches on the foreground. We first use background subtraction w.r.t. $\mathbf{B}^c$ to get a rough foreground mask $\mathbf{F}^{c,t}$ for each $\mathbf{I}^{c,t}$. We then pick $80\%$ of the rays in the batch from the foreground and the rest from the background.

**Pruning.** To speed up rendering, we prune any sample $\mathbf{r}(s_i)$ in world space that does not contain any foreground. To determine this for time $t$, we query a binary voxel grid $g_t$ that we overlay over the scene. $g_t$ is 1 for voxels that are potentially in the foreground and 0 otherwise. We fill $g_t$ via space carving [37] with $\{\mathbf{F}^{c,t}\}_c$. To be conservative and to let the deformation field extend into the surrounding area, we dilate both the foreground masks and the foreground in the resulting voxel grid. To clear out geometry artifacts in empty space, we do not prune when constructing the canonical model, only when optimizing for $t{>}1$. Pruning removes $\sim 80\%$ of samples, making our method four times faster.

**Vignetting Correction.** Cameras collect less light around the image border, which leads to darkening in the input images. We model this vignetting with a radial model [2, 74]:

$$\mathbf{C}_{vig}(\mathbf{r}) = \mathbf{C}(\mathbf{r})(1 + k_1 p(\mathbf{r}) + k_2 p(\mathbf{r})^2 + k_3 p(\mathbf{r})^3), \quad (11)$$

where $p(\mathbf{r}) \in \mathbb{R}$ is the squared distance to the camera center
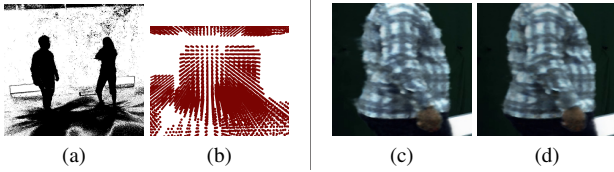
Figure 4. **Implementation.** (a) Foreground $\mathbf{F}^{c,t}$. (b) Pruning grid for (a). (c) W/o and (d) w/ vignetting correction.

in pixel space (we divide the distance by $w$ for resolution invariance), and $k_1, k_2, k_3 \in \mathbb{R}$ are correction parameters. We initialize these parameters to zeros, optimize for them when constructing the canonical model, and keep them fixed for $t>1$. We share the same set of parameters across all cameras and timestamps. We use $\mathbf{C}_{vig}$ in place of $\mathbf{C}$ in Eq. (2). The quality improves since the canonical model no longer needs to use artifacts to account for vignetting effects.

**Optimizer.** We use AdamW [47] with weight decay of $0.01$ to stabilize the training. Auto-decoded parameters like hash grids get sparsely non-zero gradients in any given training iteration, which degrades the momentum accumulation. We thus use a modified version of AdamW: rather than treating all parameters $\{\theta_i \in \mathbb{R}\}_i$ in a tensor the same, we separately keep track of AdamW's parameters $\theta_i^{opt}$ (*e.g.* the number of iterations) for each individual parameter $\theta_i$. In any given iteration, only $\theta_i$ whose derivative is non-zero have their AdamW parameters $\theta_i^{opt}$ updated and AdamW applied. The supplement describes the learning rates in detail.

**Hyperparameters.** All scenes use the same settings: a batch size of $R=1024$, $S=3072$ samples per ray, and loss weights $\lambda_{\text{back}}=0.001$, $\lambda_{\text{hard}}=1$, $\lambda_{\text{coarse}}=1000$, and $\lambda_{\text{fine}}=30$. We use LeakyReLUs [48] for the deformations and ReLUs for the canonical model. The supplement has more details.

**Code.** We use PyTorch [62] and tiny-cuda-nn [52] via its Python wrappers for its fast implementation of hash grids.

## 4. Experiments

We evaluate the reconstruction quality and time consistency of our method, perform ablations, show simple scene editing, and discuss limitations and future work. The supplemental video contains many additional results.

**Prior Work.** We compare with all prior time-consistent general NeRF methods: NR-NeRF [78], D-NeRF [63], and DeVRF [44]. To evaluate time consistency, we compare against PREF [71], a NeRF-based method to estimate correspondences.

**Variants.** We evaluate reconstruction quality by using novel-view synthesis *as a proxy*. Unlike most works, we target *time-consistent* reconstruction and not novel-view synthesis, which leads to a trade-off with novel-view quality.

Most dynamic-NeRF papers condition the canonical model on time. Therefore, we evaluate two variants of our method that do the same and, hence, synthesize better novel
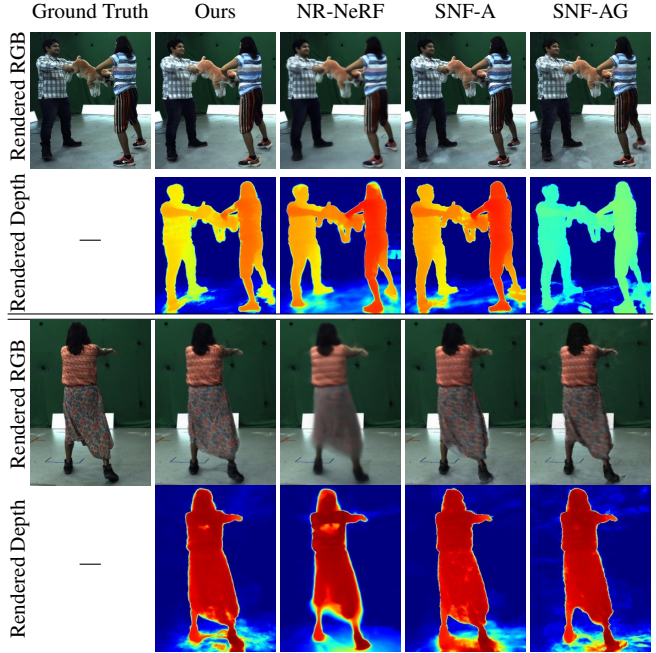


Figure 5. **Novel-View Synthesis.** (Top) Seq. 1 at $t=T$. (Bottom) Seq. 5 at $t=38$. The depth color differences between methods are due to normalization.

views at the cost of correspondences. The first variant, *SNF-A*, only makes the appearance time-varying, while the second variant, *SNF-AG*, makes the appearance and canonical geometry time-varying. The supplement has further details.

We emphasize that we do not claim that these variants are competitive with state-of-the-art 4D reconstruction methods that neglect correspondences and focus solely on novel-view synthesis. Our design is not tailored to this setting.

**Data.** For ease of recording, we evaluate our *general* method on real-world scenes of one or two people. Crucially, the recordings also contain a plush dog and loose clothing. Furthermore, our method reconstructs scenes with multiple people *without any knowledge that there are multiple entities in the scene nor that they are human*. These aspects demonstrate its generality. We use a total of $C=117$ static cameras in a studio to record six scenes of $4-5$sec and one of $12$sec at $25$fps. For the test sets, we pick the same two cameras for all scenes. We train at an input resolution of $h=1504$ and $w=2056$. In addition, we consider the short multi-view scene from NR-NeRF [78], which contains 16 camera pairs.

**Runtime.** On an Nvidia A100 GPU, the canonical model trains at $8$iter/sec, the coarse deformations at $16$iter/sec, and the fine deformations at $14$iter/sec. For efficiency, we render all images at half the input resolution, at $10-15$sec/image.

### 4.1. Qualitative Results

We first qualitatively evaluate reconstructions at any one point in time, and then examine their consistency over time.

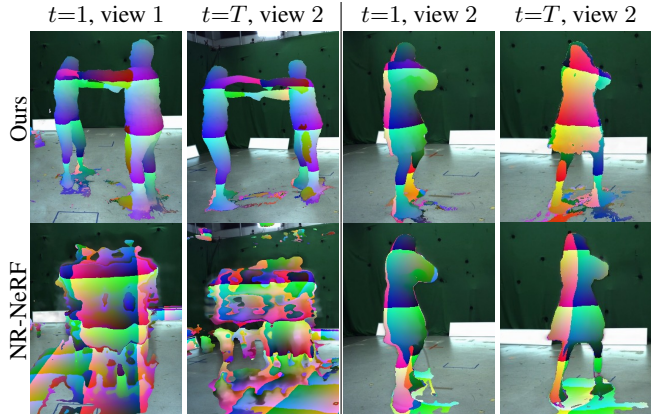| | $t=1$, view 1 | $t=T$, view 2 | $t=1$, view 2 | $t=T$, view 2 |

Figure 6. **Correspondences.** (Left) Seq. 2, opposite views. NR-NeRF fails. (Right) Seq. 4, same view. Only our correspondences follow the $90°$ left turn; see orange at $t=1$ & $T$.

### 4.1.1 Volumetric 3D Scene Reconstruction

Fig. 5 shows novel views of the reconstruction and its geometry via depth maps. While NR-NeRF gives blurry results, our method yields high-quality reconstructions. Although blurrier, SNF-A matches the pattern on the dress better than ours. SNF-AG reconstructs the non-rigid geometry more accurately as it loosens the correspondences.

### 4.1.2 Time Consistency

**Correspondence Visualization.** Time consistency enables 3D correspondences over time: Fig. 6 shows that our estimated correspondences are temporally stable. However, on scenes with large motion, NR-NeRF fails to converge to a recognizable model since this requires simultaneous correspondence estimation in this highly challenging setting.

**Comparisons with Prior Work.** For further evaluation, we follow PREF [71], which estimates world-space 3D human joint positions over time: $\{\tilde{\mathbf{p}}_j^t \in \mathbb{R}^3\}_{t,j}$, where $j$ indexes the $J=23$ joints. Off-the-shelf commercial systems [11] exploit strong human-specific priors to reliably estimate 3D joints, which makes for excellent pseudo-ground-truth long-term 3D correspondences $\{\hat{\mathbf{p}}_j^t\}_{t,j}$. (The joints are used only for evaluation; our method does not use human priors.)

To track the joints with our method, we first initialize the estimated positions $\{\tilde{\mathbf{p}}_j^1\}_j$ with the ground-truth joints $\{\hat{\mathbf{p}}_j^1\}_j$ at $t=1$, which are the positions in the canonical model. We then need to invert the backwards deformation field for $t>1$. To this end, for each joint $j$ at time $t$, we evaluate $d_f(d_c(\cdot; \theta_c^t); \theta_f^t)$ on a voxel grid (with a resolution of $128^3$ and a side length of $40$cm) centered on the estimated position at $t-1$, $\tilde{\mathbf{p}}_j^{t-1}$, and pick the voxel center that lands closest to the ground-truth joint $\hat{\mathbf{p}}_j^1$ in the canonical model as the estimated world-space joint position $\tilde{\mathbf{p}}_j^t$. The supple-
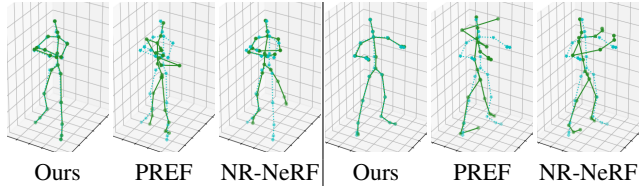


Figure 7. **Time Consistency.** (Left) Seq. 4. (Right) Seq. 6. The solid skeleton is the tracking estimate at $t=T$. The dotted skeleton is the pseudo-ground truth at $t=T$.



| Ours at $t_1$ & $t_2$ | SNF-A at $t_1$ | SNF-A at $t_2$ | Ours at $t_1$ & $t_2$ | SNF-AG at $t_1$ | SNF-AG at $t_2$ |

Figure 8. **Canonical Model.** Ours uses a static canonical model, while those of SNF-A(G) vary over time.

| Sequence | 1 | 2 | 3 | 4 | 5 | 6 | 7 | Mean |
|---|---|---|---|---|---|---|---|---|
| Ours | **0.9** | **0.9** | **1.8** | **1.8** | **2.8** | **0.9** | 1.4 | **1.5** |
| PREF | 11.0 | 13.9 | 47.8 | 8.9 | 13.4 | 12.1 | 3.9 | 15.9 |
| NR-NeRF | 2.2 | 46.4 | 110.9 | 6.1 | 8.5 | 14.8 | **1.2** | 27.2 |
| D-NeRF | 40.7 | 46.2 | 110.1 | 12.2 | 81.9 | 35.3 | 70.0 | 56.6 |

Table 1. **Time Consistency.** We report per-scene and mean MPJPE in cm. Lower is better.

ment contains tracking details for NR-NeRF and PREF.

Fig. 7 shows results at $t=T$. Our method yields stable correspondences, while NR-NeRF does not. PREF's correspondences drift strongly over time due to error accumulation from chaining frame-to-frame correspondences.

**Variants.** We contrast the time consistency of our method with its variants. To this end, Fig. 8 visualizes the canonical model at two different timestamps. Unlike SceNeRFlow, the canonical model of the variants changes over time, *i.e.*, they lack inherent correspondences. However, as results in the video show, these additional degrees of freedom improve the variants' novel-view synthesis, showing a trade-off between time consistency and novel-view synthesis.

### 4.2. Quantitative Results

**Time Consistency.** Like PREF [71], we measure time consistency via 3D joint tracking. We report the mean per-joint position error (MPJPE) [68] over all timestamps $t>1$ and the $J$ joints: $\text{MPJPE} = \frac{1}{(T-1)J} \sum_{t=2}^{T} \sum_{j=1}^{J} \|\hat{\mathbf{p}}_j^t - \tilde{\mathbf{p}}_j^t\|_2$. Lower is better. For scenes with two people, we report the average across both. We evaluate all scenes except for the 5-frame Seq. 8 used in NR-NeRF [78]. Tab. 1 contains the results. SceNeRFlow has the lowest error with only marginal drift. PREF's frame-wise approach leads to large drift. D-NeRF cannot handle large motion.

| | | Ours | NR-NeRF | SNF-A | SNF-AG | Background |
|---|---|---|---|---|---|---|
| Unmasked | PSNR ↑ | 30.32 | 28.77 | **30.88** | 30.34 | 21.69 |
| | SSIM ↑ | 0.939 | 0.922 | **0.940** | 0.929 | 0.920 |
| | LPIPS ↓ | **0.054** | 0.099 | 0.057 | 0.089 | 0.101 |
| Masked | PSNR ↑ | 32.38 | 30.80 | **33.31** | 33.22 | — |
| | SSIM ↑ | 0.976 | 0.965 | **0.978** | 0.977 | — |
| | LPIPS ↓ | 0.018 | 0.041 | **0.016** | 0.018 | — |

Table 2. **Novel-View Synthesis.** Mean PSNR, SSIM, and LPIPS across scenes. SNF-A(G) are variants of ours.
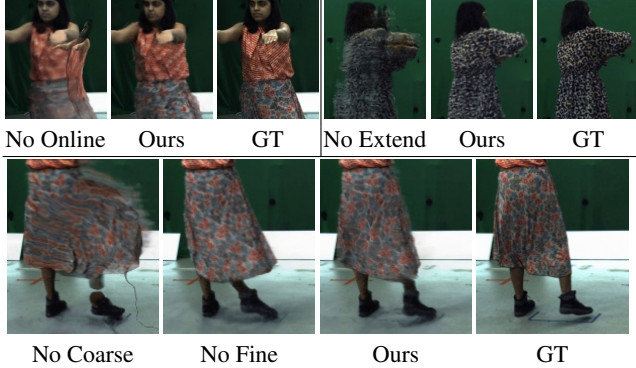


Figure 9. **Ablations.** (Top left) Online optimization, (top right) extending the deformations, (bottom) coarse and fine deformations.

**Reconstruction.** We quantify the reconstruction quality by using novel-view quality as a proxy. Like NeRF [50], we measure the latter with PSNR and SSIM [96] (for both, higher is better), and the neural perceptual metric LPIPS [32] (lower is better). We focus on the dynamic scene part by (1) also providing these metrics w.r.t. the static background image, and (2) reporting *masked scores* where we mask out the renderings with foreground masks estimated from the ground truth. The results are in Tab. 2. Our method outperforms NR-NeRF. The variants have artifacts in empty space and only outperform our method on the masked scores. The supplement has per-scene results and further comparisons with D-NeRF and DeVRF.

## 4.3. Ablations

We ablate the core components of our method. We investigate the relevance of (1) optimizing in an online manner (by optimizing all timestamps at once), (2) "extending the deformation field" (by using a coarse resolution of $512^3$ and $\mathcal{L}_{\text{norm,w}}$ without max-pooling), (3) the coarse deformation model (by setting $\Delta_c = \mathbf{0}$ at training time), and (4) the fine deformation model (by setting $\Delta_f = \mathbf{0}$ at test time). Fig. 9 shows the results, the supplement confirms these quantitatively: (1) Optimizing online simplifies implicit correspondence estimation, avoiding ghosting artifacts; (2) extending the deformation field helps the online optimization by initializing the deformations for the next timestamp well; (3) the coarse deformations stabilize the reconstruction; and (4) the fine deformations add details.
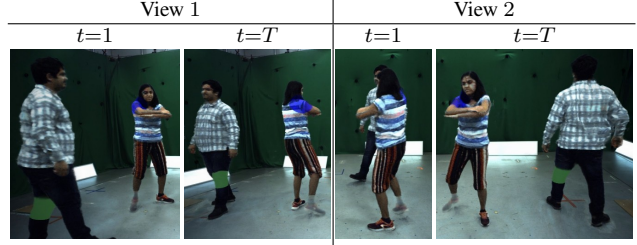


Figure 10. **Scene Editing.** (1) Coloring the left knee of person 1 green, (2) blending the right shoulder of person 2 with blue, and (3) making the right foot of person 2 transparent.

## 4.4. Simple Editing

The time consistency enables simple editing of the geometry and appearance. As a proof of concept, Fig. 10 shows that we can re-color scene parts (*e.g.,* of Seq. 3) or make them transparent in a straightforward manner—by modifying the static canonical model. The deformations then consistently propagate these changes to all timestamps.

## 4.5. Limitations and Future Work

While our method design makes a few assumptions, we exploit them as little as possible, *e.g.* we do not use accurate segmentation masks, depth estimates, or 3D estimates. This makes it easier to extend our work to other input settings.

Adjacent work offers promising remedies for our assumptions. Most dynamic-NeRF methods are monocular (but not time-consistent for large motion) and some static-NeRF works [77, 89] use only few cameras. Adapting these ideas might reduce the number of cameras required by our method. Since our core problem lies in the dynamic foreground, we exclude the background in a soft manner. The background could be included via a static NeRF as some existing dynamic-NeRF approaches do. Naïvely modeling time-varying appearance loosens correspondences. Sophisticated regularization from Shape-from-Shading [4, 93, 94] might help. To preserve time consistency, we do not update the canonical model with newly visible geometry from $t>1$. However, multi-view input mitigates the effect of occlusions at $t=1$. Finally, reconstructing topology changes *in a time-consistent manner* remains an open challenge in the field and we do not currently aim to address it.

## 5. Conclusion

SceNeRFlow demonstrates the great potential of modern neural scene representations for time-consistent reconstruction of general dynamic objects. In particular, we showed how backwards deformation models can be adapted to tracking large motion. As a proof of concept, we demonstrated how SceNeRFlow enables simple edits that are consistent over time. We also showed how conditioning the

canonical model on time trades off novel-view and correspondence quality. As discussed in Sec. 4.5, we make some simplifying assumptions but exploit them as little as possible, paving the way for future extensions to weaker inputs.

# References

[1] Benjamin Attal, Eliot Laidlaw, Aaron Gokaslan, Changil Kim, Christian Richardt, James Tompkin, and Matthew O'Toole. Törf: Time-of-flight radiance fields for dynamic scene view synthesis. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2

[2] Artur Bal and Henryk Palus. Image vignetting correction using a deformable radial polynomial model. *Sensors*, 2023. 5

[3] Aayush Bansal, Minh Vo, Yaser Sheikh, Deva Ramanan, and Srinivasa Narasimhan. 4d visualization of dynamic events from unconstrained multi-view videos. In *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[4] Jonathan T Barron and Jitendra Malik. Shape, illumination, and reflectance from shading. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2014. 8

[5] Adrien Bartoli, Yan Gérard, François Chadebecq, Toby Collins, and Daniel Pizarro. Shape-from-template. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2015. 2

[6] Tali Basha, Yael Moses, and Nahum Kiryati. Multi-view scene flow estimation: A view centered variational approach. *International Journal of Computer Vision (IJCV)*, 2013. 2

[7] Aljaz Bozic, Pablo Palafox, Michael Zollhöfer, Angela Dai, Justus Thies, and Matthias Nießner. Neural non-rigid tracking. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[8] Aljaz Bozic, Michael Zollhofer, Christian Theobalt, and Matthias Nießner. Deepdeform: Learning non-rigid rgb-d reconstruction with semi-supervised data. In *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[9] Cedric Cagniart, Edmond Boyer, and Slobodan Ilic. Probabilistic deformable surface tracking from multiple videos. In *European Conference on Computer Vision (ECCV)*, 2010. 1, 2

[10] Hongrui Cai, Wanquan Feng, Xuetao Feng, Yan Wang, and Juyong Zhang. Neural surface reconstruction of dynamic scenes with monocular rgb-d camera. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

[11] The Captury. Captury Studio. https://captury.com/, 2023. 7

[12] David Casillas-Perez, Daniel Pizarro, David Fuentes-Jimenez, Manuel Mazo, and Adrien Bartoli. The isowarp: the template-based visual geometry of isometric surfaces. *International Journal of Computer Vision (IJCV)*, 2021. 2

[13] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proc. Conference on Computer Graphics and Interactive Techniques*, 1996. 2

[14] Bailin Deng, Yuxin Yao, Roberto M Dyke, and Juyong Zhang. A survey of non-rigid 3d registration. In *Computer Graphics Forum (Eurographics State of the Art Reports)*, 2022. 2

[15] Mingsong Dou, Sameh Khamis, Yury Degtyarev, Philip Davidson, Sean Ryan Fanello, Adarsh Kowdle, Sergio Orts Escolano, Christoph Rhemann, David Kim, Jonathan Taylor, et al. Fusion4d: Real-time performance capture of challenging scenes. *ACM Transactions on Graphics*, 2016. 2

[16] Yilun Du, Yinan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *International Conference on Computer Vision (ICCV)*, 2021. 1, 2

[17] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, 2022. 2

[18] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *International Conference on Computer Vision (ICCV)*, 2021. 1, 2

[19] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular dynamic view synthesis: A reality check. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

[20] Ravi Garg, Anastasios Roussos, and Lourdes Agapito. Dense variational reconstruction of non-rigid surfaces from monocular video. In *Computer Vision and Pattern Recognition (CVPR)*, 2013. 2

[21] Bastian Goldluecke and Marcus Magnor. Space-time isosurface evolution for temporally coherent 3d reconstruction. In *Computer Vision and Pattern Recognition (CVPR)*, 2004. 2

[22] Stella Graßhof and Sami Sebastian Brandt. Tensor-based non-rigid structure from motion. In *Winter Conference on Applications of Computer Vision (WACV)*, 2022. 2

[23] Andreas Griewank and Andrea Walther. *Evaluating derivatives: principles and techniques of algorithmic differentiation*. Society for Industrial and Applied Mathematics (SIAM), 2008. 5

[24] Xiuye Gu, Yijie Wang, Chongruo Wu, Yong Jae Lee, and Panqu Wang. Hplflownet: Hierarchical permutohedral lattice flownet for scene flow estimation on large-scale point clouds. In *Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[25] Kaiwen Guo, Feng Xu, Yangang Wang, Yebin Liu, and Qionghai Dai. Robust non-rigid motion tracking and surface reconstruction using l0 regularization. In *International Conference on Computer Vision (ICCV)*, 2015. 2

[26] Kaiwen Guo, Feng Xu, Tao Yu, Xiaoyang Liu, Qionghai Dai, and Yebin Liu. Real-time geometry, albedo, and motion reconstruction using a single rgb-d camera. *ACM Transactions on Graphics*, 2017. 2

[27] Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. Neural deformable voxel grid for fast optimization of dynamic view synthesis. In *Asian Conference on Computer Vision (ACCV)*, 2022. 2

[28] Frédéric Huguet and Frédéric Devernay. A variational method for scene flow estimation from stereo sequences. In *International Conference on Computer Vision (ICCV)*, 2007. 2

[29] Chun Ho Hung, Li Xu, and Jiaya Jia. Consistent binocular depth and scene flow with chained temporal profiles. *International Journal of Computer Vision (IJCV)*, 2013. 2

[30] Matthias Innmann, Michael Zollhöfer, Matthias Nießner, Christian Theobalt, and Marc Stamminger. VolumeDeform: Real-time Volumetric Non-rigid Reconstruction. 2016. 2

[31] Erik C.M. Johnson, Marc Habermann, Soshi Shimada, Vladislav Golyanik, and Christian Theobalt. Unbiased 4d: Monocular 4d reconstruction with a neural deformation model. *arXiv:2206.08368*, 2022. 2

[32] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision (ECCV)*, 2016. 8

[33] Navami Kairanda, Edith Tretschk, Mohamed Elgharib, Christian Theobalt, and Vladislav Golyanik. $\phi$-sft: Shape-from-template with a physics-based deformation model. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[34] Angjoo Kanazawa, Shubham Tulsiani, Alexei A. Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *European Conference on Computer Vision (ECCV)*, 2018. 2

[35] Filippos Kokkinos and Iasonas Kokkinos. Learning monocular 3d reconstruction of articulated categories from motion. In *Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[36] Suryansh Kumar, Anoop Cherian, Yuchao Dai, and Hongdong Li. Scalable dense non-rigid structure-from-motion: A grassmannian perspective. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[37] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *International Journal of Computer Vision (IJCV)*, 2000. 5

[38] E Scott Larsen, Philippos Mordohai, Marc Pollefeys, and Henry Fuchs. Temporally consistent reconstruction from multiple video streams using enhanced belief propagation. In *International Conference on Computer Vision (ICCV)*, 2007. 2

[39] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. Neural 3d video synthesis from multi-view video. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[40] Xueting Li, Sifei Liu, Shalini De Mello, Kihwan Kim, Xiaolong Wang, Ming-Hsuan Yang, and Jan Kautz. Online adaptation for consistent mesh reconstruction in the wild. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[41] Xueqian Li, Jhony Kaesemodel Pontes, and Simon Lucey. Neural scene flow prior. *Advances in Neural Information Processing Systems (NeurIPS)*, 2021. 2

[42] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2

[43] Wenbin Lin, Chengwei Zheng, Jun-Hai Yong, and Feng Xu. Occlusionfusion: Occlusion-aware motion estimation for real-time dynamic 3d reconstruction. 2022. 2

[44] Jia-Wei Liu, Yan-Pei Cao, Weijia Mao, Wenqiao Zhang, David Junhao Zhang, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Devrf: Fast deformable voxel radiance fields for dynamic scenes. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2, 3, 6

[45] Xingyu Liu, Charles R Qi, and Leonidas J Guibas. Flownet3d: Learning scene flow in 3d point clouds. *Computer Vision and Pattern Recognition (CVPR)*, 2019. 2

[46] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Transactions on Graphics*, 2019. 2, 3

[47] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*, 2019. 6

[48] Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *International Conference on Machine Learning (ICML)*, 2013. 6

[49] Willi Menapace, Stéphane Lathuilière, Aliaksandr Siarohin, Christian Theobalt, Sergey Tulyakov, Vladislav Golyanik, and Elisa Ricci. Playable environments: Video manipulation in space and time. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[50] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision (ECCV)*, 2020. 2, 3, 4, 8

[51] Himangi Mittal, Brian Okorn, and David Held. Just go with the flow: Self-supervised scene flow estimation. In *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[52] Thomas Müller. tiny-cuda-nn, 2021. 6

[53] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics*, 2022. 3, 4, 17

[54] Armin Mustafa, Hansung Kim, Jean-Yves Guillemaut, and Adrian Hilton. Temporally coherent 4d reconstruction of complex dynamic scenes. In *Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2

[55] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Computer Vision and Pattern Recognition (CVPR)*, 2015. 2

[56] Dat Tien Ngo, Sanghyuk Park, Anne Jorstad, Alberto Crivellaro, Chang D. Yoo, and Pascal Fua. Dense image registration and deformable surface reconstruction in presence of occlusions and minimal texture. In *International Conference on Computer Vision (ICCV)*, 2015. 2

10

[57] Sergio Orts-Escolano, Christoph Rhemann, Sean Fanello, Wayne Chang, Adarsh Kowdle, Yury Degtyarev, David Kim, Philip L Davidson, Sameh Khamis, Mingsong Dou, et al. Holoportation: Virtual 3d teleportation in real-time. In *Annual Symposium on User Interface Software and Technology*, 2016. 2

[58] Shaifali Parashar, Mathieu Salzmann, and Pascal Fua. Local non-rigid structure-from-motion from diffeomorphic mappings. In *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[59] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 5, 13

[60] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: a higher-dimensional representation for topologically varying neural radiance fields. *ACM Transactions on Graphics*, 2021. 1, 2

[61] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. *Advances in Neural Information Processing Systems (NeurIPS) Workshops*, 2017. 5

[62] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (NeurIPS)*. 2019. 6

[63] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 6

[64] Julian Quiroga, Thomas Brox, Frédéric Devernay, and James Crowley. Dense semi-rigid scene flow estimation from rgbd images. In *European Conference on Computer Vision (ECCV)*, 2014. 2

[65] Daniel Rebain, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. Lolnerf: Learn from one look, 2022. 3

[66] Chris Russell, João Fayad, and Lourdes Agapito. Dense non-rigid structure from motion. In *International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission*, 2012. 2

[67] Mathieu Salzmann, Vincent Lepetit, and Pascal Fua. Deformable surface tracking ambiguities. In *Computer Vision and Pattern Recognition (CVPR)*, 2007. 2

[68] Leonid Sigal, Alexandru O Balan, and Michael J Black. Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision (IJCV)*, 2010. 7

[69] Miroslava Slavcheva, Maximilian Baust, Daniel Cremers, and Slobodan Ilic. Killingfusion: Non-rigid 3d reconstruction without correspondences. In *Computer Vision and Pattern Recognition (CVPR)*, 2017. 2

[70] Miroslava Slavcheva, Maximilian Baust, and Slobodan Ilic. Sobolevfusion: 3d reconstruction of scenes undergoing free non-rigid motion. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[71] Liangchen Song, Xuan Gong, Benjamin Planche, Meng Zheng, David Doermann, Junsong Yuan, Terrence Chen, and Ziyan Wu. Pref: Predictability regularized neural motion fields. In *European Conference on Computer Vision (ECCV)*, 2022. 2, 6, 7, 17

[72] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry Processing (SGP)*, 2007. 4

[73] Bert Speelpenning. *Compiling fast partial derivatives of functions given by algorithms*. University of Illinois at Urbana-Champaign, 1980. 5

[74] Jessi Stumpfel, Andrew Jones, Andreas Wenger, Chris Tchou, Tim Hawkins, and Paul Debevec. Direct hdr capture of the sun and sky. In *ACM SIGGRAPH Courses*. 2006. 5

[75] Zachary Teed and Jia Deng. Raft-3d: Scene flow using rigid-motion embeddings. In *Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[76] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, Rohit Pandey, Sean Fanello, Gordon Wetzstein, Jun-Yan Zhu, Christian Theobalt, Maneesh Agrawala, Eli Shechtman, Dan B Goldman, and Michael Zollhöfer. State of the art on neural rendering. *Computer Graphics Forum (Eurographics State of the Art Reports)*, 2020. 2

[77] A. Tewari, J. Thies, B. Mildenhall, P. Srinivasan, E. Tretschk, W. Yifan, C. Lassner, V. Sitzmann, R. Martin-Brualla, S. Lombardi, T. Simon, C. Theobalt, M. Nießner, J. T. Barron, G. Wetzstein, M. Zollhöfer, and V. Golyanik. Advances in Neural Rendering. *Computer Graphics Forum (Eurographics State of the Art Reports)*, 2022. 2, 8

[78] 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 *International Conference on Computer Vision (ICCV)*, 2021. 1, 2, 4, 5, 6, 7, 13, 17

[79] Edith Tretschk, Navami Kairanda, Mallikarjun B R, Rishabh Dabral, Adam Kortylewski, Bernhard Egger, Marc Habermann, Pascal Fua, Christian Theobalt, and Vladislav Golyanik. State of the art in dense monocular non-rigid 3d reconstruction, 2023. 2, 4

[80] Tony Tung, Shohei Nobuhara, and Takashi Matsuyama. Complete multi-view reconstruction of dynamic scenes from probabilistic fusion of narrow and wide baseline stereo. In *International Conference on Computer Vision (ICCV)*, 2009. 2

[81] Sundar Vedula, Simon Baker, Peter Rander, Robert Collins,

and Takeo Kanade. Three-dimensional scene flow. In *International Conference on Computer Vision (ICCV)*, 1999. 2

[82] Christoph Vogel, Konrad Schindler, and Stefan Roth. 3d scene flow estimation with a piecewise rigid scene model. *International Journal of Computer Vision (IJCV)*, 2015. 2

[83] Liao Wang, Jiakai Zhang, Xinhua Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenoctrees for dynamic radiance field rendering in real-time. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[84] Shangzhe Wu, Tomas Jakab, Christian Rupprecht, and Andrea Vedaldi. DOVE: Learning deformable 3d objects by watching videos. *arXiv preprint arXiv:2107.10844*, 2021. 2

[85] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2

[86] Gengshan Yang, Minh Vo, Neverova Natalia, Deva Ramanan, Vedaldi Andrea, and Joo Hanbyul. Banmo: Building animatable 3d neural models from many casual videos. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[87] Jae Shin Yoon, Ziwei Li, and Hyun Soo Park. 3d semantic trajectory reconstruction from 3d pixel continuum. In *Computer Vision and Pattern Recognition (CVPR)*, 2018. 2

[88] Jae Shin Yoon, Kihwan Kim, Orazio Gallo, Hyun Soo Park, and Jan Kautz. Novel view synthesis of dynamic scenes with globally coherent depths from a monocular camera. In *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[89] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *Computer Vision and Pattern Recognition (CVPR)*, 2021. 8

[90] Hsiao yu Chen, Edith Tretschk, Tuur Stuyck, Petr Kadlecek, Ladislav Kavan, Etienne Vouga, and Christoph Lassner. Virtual elastic objects. In *Computer Vision and Pattern Recognition (CVPR)*, 2022. 1

[91] Mingliang Zhai, Xuezhi Xiang, Ning Lv, and Xiangdong Kong. Optical flow and scene flow estimation: A survey. *Pattern Recognition*, 2021. 2

[92] Li Zhang, Brian Curless, and Steven M Seitz. Spacetime stereo: Shape recovery for dynamic scenes. In *Computer Vision and Pattern Recognition (CVPR)*, 2003. 2

[93] Ruo Zhang and Ping-Sing Tsai. Shape-from-shading: a survey. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 1999. 8

[94] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nerfactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics*, 2021. 8

[95] Fuqiang Zhao, Yuheng Jiang, Kaixin Yao, Jiakai Zhang, Liao Wang, Haizhao Dai, Yuhui Zhong, Yingliang Zhang, Minye Wu, Lan Xu, et al. Human performance modeling and rendering via neural animated mesh. *ACM Transactions on Graphics (Proceedings of SIGGRAPH Asia)*, 2022. 2

[96] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 2004. 8

[97] Michael Zollhöfer, Matthias Nießner, Shahram Izadi, Christoph Rehmann, Christopher Zach, Matthew Fisher, Chenglei Wu, Andrew Fitzgibbon, Charles Loop, Christian Theobalt, and Marc Stamminger. Real-time non-rigid reconstruction using an rgb-d camera. *ACM Transactions on Graphics*, 2014. 2

[98] Michael Zollhöfer, Patrick Stotko, Andreas Görlitz, Christian Theobalt, Matthias Nießner, Reinhard Klein, and Andreas Kolb. State of the art on 3d reconstruction with rgb-d cameras. 2018. 2

# Supplementary Material

We provide details on how we visualize correspondences in Sec. A, a description of each scene in the dataset in Sec. B, per-scene quantitative novel-view-synthesis results in Sec. C, novel-view-synthesis comparisons with D-NeRF and DeVRF in Sec. D, more qualitative joint-tracking and novel-view results in Sec. E, quantitative ablation results in Sec. F, further architecture and training details in Sec. G, details on how we adapt the baselines to joint tracking in Sec. H, and more details on the foreground masks used for evaluation in Sec. I.

## A. Correspondence Visualization Details

We follow NR-NeRF's [78] visualization and replace the appearance in the canonical model with a voxel grid of RGB cubes. Like prior work [59, 78], we pick that sample $i'$ of the ray as the surface point $\mathbf{r}(s_{i'})$ that is closest to an accumulated transmittance $\sum_{i=1}^{i'} w_i$ of 0.5. For better visualization, we filter rays with a total accumulated transmittance below 0.4 (*i.e.* those hitting the background) and visualize them as transparent. Fig. 11 shows an example.

## B. Data Details

Tab. 3 contains a description and the length in frames of each scene. We record at 25fps.

## C. Per-Scene Quantitative Results

We collect the quantitative results per scene for novel-view synthesis in Tab. 4.

## D. Novel-View-Synthesis Comparison with Prior Work

Tab. 5 contains further comparisons, without Seq. 8 due to memory constraints. SceNeRFlow outperforms D-NeRF and DeVRF, which both fail on large motion.

## E. More Qualitative Joint-Tracking and Novel-View Results

Fig. 12 contains the qualitative joints estimated at $t=T$ for all sequences not shown in the main paper. Fig. 13 and Fig. 14 contain more novel-view results of all methods for four scenes at $t=\frac{T}{2}$.

## F. Quantitative Ablation Results

Tab. 6 confirms the qualitative ablation results quantitatively.



| With Filtering | Without Filtering |

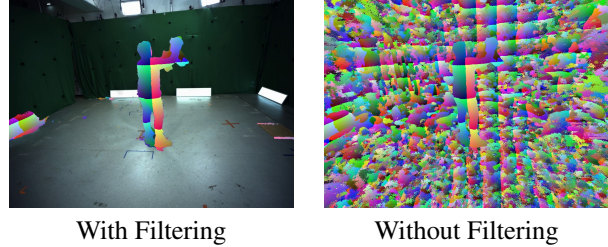Figure 11. **Filtering for Correspondence Visualization.** For better visualization, we filter out rays with an accumulated transmittance below 0.4.

| Name | Description | Frames |
|---|---|---|
| Seq. 1 | two people playing with a plush dog | 125 |
| Seq. 2 | two people holdings hands | 125 |
| Seq. 3 | one person rotating, one person walking | 300 |
| Seq. 4 | a person dancing while wearing a dark dress | 100 |
| Seq. 5 | a person walking like a zombie while wearing a light dress | 125 |
| Seq. 6 | a person playing with a plush dog | 125 |
| Seq. 7 | a person standing with a brown dress | 125 |
| Seq. 8 | a person doing squads (from NR-NeRF [78]) | 5 |

Table 3. **Dataset Description.**

## G. Further Architecture and Training Details

For all methods, we normalize the scene into the unit cube by tightly fitting an axis-aligned bounding box to all near and far plane samples of all images.

### G.1. Ours

**Architecture.** Both the coarse and fine deformation fields use the same architecture. The hash grid consists of 16 levels, with two feature dimensions per level. The coarsest level has a resolution of $32^3$, and each subsequent level has a 1.3819 times higher resolution. The hashmap has a size of $2^{20}$. The shallow MLP has one hidden layer with 64 hidden dimensions.

The canonical model uses a hash grid with 13 levels, two feature dimensions per level, a coarsest resolution of $128^3$, and a scaling factor of 1.3819. The shallow MLP that outputs the opacity has one hidden layer with 64 hidden dimensions. A second MLP outputs the appearance and takes as input a 15-dimensional vector additionally regressed by the first shallow MLP. The second MLP has two hidden layers with 64 hidden dimensions.

**Weighting Scheme for Smoothness Loss.** We weigh sample $i$ on ray $r$ depending on its closeness to the object. Mathematically, we start with $\hat{\sigma}_{r,i} = \exp(-\sigma_{r,i}\delta)$, where $\sigma_{r,i}$ is the opacity of the $i$-th sample on ray $r$. We next apply max-pooling with window size $k = \lfloor f \cdot S \rfloor$, where we empirically set $f=0.005$:

$$\hat{\sigma}'_{r,i} = \max_{i' \in [i-k,...,i+k]} \hat{\sigma}_{r,i'}. \qquad (12)$$

We then weaken the regularization on empty space by

Table 4 (left):

| | | | | Ours | NR-NeRF | SNF-A | SNF-AG | Background |
|---|---|---|---|---|---|---|---|---|
| Seq. 1 | Unmasked | PSNR | ↑ | 27.49 | 26.84 | **27.84** | 26.73 | 18.64 |
| | | SSIM | ↑ | **0.928** | 0.915 | **0.928** | 0.906 | 0.894 |
| | | LPIPS | ↓ | **0.074** | 0.117 | 0.076 | 0.138 | 0.139 |
| | Masked | PSNR | ↑ | 29.73 | 29.38 | **30.02** | 29.89 | — |
| | | SSIM | ↑ | 0.970 | 0.966 | **0.971** | 0.970 | — |
| | | LPIPS | ↓ | 0.021 | 0.036 | **0.017** | 0.018 | — |
| Seq. 2 | Unmasked | PSNR | ↑ | 27.93 | 21.64 | **28.24** | 26.99 | 18.57 |
| | | SSIM | ↑ | **0.929** | 0.875 | **0.929** | 0.907 | 0.898 |
| | | LPIPS | ↓ | **0.069** | 0.183 | 0.074 | 0.136 | 0.130 |
| | Masked | PSNR | ↑ | 29.65 | 22.75 | **30.23** | 29.96 | — |
| | | SSIM | ↑ | 0.970 | 0.931 | **0.972** | 0.970 | — |
| | | LPIPS | ↓ | 0.019 | 0.081 | **0.017** | 0.019 | — |
| Seq. 3 | Unmasked | PSNR | ↑ | **27.72** | 21.48 | 27.50 | 26.61 | 18.94 |
| | | SSIM | ↑ | **0.923** | 0.874 | 0.920 | 0.904 | 0.895 |
| | | LPIPS | ↓ | **0.075** | 0.181 | 0.089 | 0.136 | 0.140 |
| | Masked | PSNR | ↑ | **28.95** | 22.25 | 28.90 | 28.47 | — |
| | | SSIM | ↑ | **0.961** | 0.925 | **0.961** | 0.958 | — |
| | | LPIPS | ↓ | **0.029** | 0.089 | 0.029 | 0.033 | — |
| Seq. 4 | Unmasked | PSNR | ↑ | 31.59 | 31.79 | 31.92 | **32.08** | 24.68 |
| | | SSIM | ↑ | 0.950 | 0.948 | **0.951** | 0.949 | 0.943 |
| | | LPIPS | ↓ | **0.034** | 0.055 | 0.036 | 0.044 | 0.070 |
| | Masked | PSNR | ↑ | 33.68 | 33.67 | 34.15 | **34.41** | — |
| | | SSIM | ↑ | 0.980 | 0.978 | **0.981** | **0.981** | — |
| | | LPIPS | ↓ | **0.011** | 0.031 | 0.012 | 0.014 | — |
| Seq. 5 | Unmasked | PSNR | ↑ | 32.25 | 31.73 | 32.58 | **32.77** | 23.33 |
| | | SSIM | ↑ | 0.946 | 0.944 | **0.947** | 0.945 | 0.938 |
| | | LPIPS | ↓ | **0.042** | 0.057 | 0.044 | 0.049 | 0.076 |
| | Masked | PSNR | ↑ | 34.08 | 33.29 | 34.49 | **34.95** | — |
| | | SSIM | ↑ | 0.976 | 0.973 | **0.978** | 0.978 | — |
| | | LPIPS | ↓ | **0.017** | 0.033 | 0.017 | 0.017 | — |
| Seq. 6 | Unmasked | PSNR | ↑ | 28.52 | 26.69 | **29.05** | 27.70 | 19.76 |
| | | SSIM | ↑ | **0.938** | 0.917 | 0.937 | 0.918 | 0.916 |
| | | LPIPS | ↓ | **0.060** | 0.123 | 0.069 | 0.120 | 0.105 |
| | Masked | PSNR | ↑ | 30.24 | 28.85 | **31.19** | 30.99 | — |
| | | SSIM | ↑ | 0.978 | 0.969 | **0.980** | 0.979 | — |
| | | LPIPS | ↓ | 0.015 | 0.035 | **0.013** | 0.013 | — |
| Seq. 7 | Unmasked | PSNR | ↑ | 34.38 | 34.79 | 35.32 | **35.43** | 26.08 |
| | | SSIM | ↑ | 0.959 | 0.958 | **0.960** | 0.956 | 0.950 |
| | | LPIPS | ↓ | **0.026** | 0.032 | **0.026** | 0.036 | 0.057 |
| | Masked | PSNR | ↑ | 37.88 | 38.01 | 39.33 | **39.83** | — |
| | | SSIM | ↑ | 0.991 | 0.990 | **0.992** | 0.992 | — |
| | | LPIPS | ↓ | 0.005 | 0.007 | **0.003** | 0.004 | — |
| Seq. 8 | Unmasked | PSNR | ↑ | 32.67 | **35.18** | 34.58 | 34.38 | 23.53 |
| | | SSIM | ↑ | 0.937 | 0.942 | **0.944** | 0.941 | 0.926 |
| | | LPIPS | ↓ | 0.046 | **0.041** | 0.042 | 0.047 | 0.088 |
| | Masked | PSNR | ↑ | 34.82 | **38.17** | 38.13 | 37.26 | — |
| | | SSIM | ↑ | 0.980 | 0.984 | **0.986** | 0.984 | — |
| | | LPIPS | ↓ | 0.022 | 0.017 | **0.013** | 0.018 | — |

Table 4. **Novel-View Synthesis.** We report per-scene PSNR, SSIM, and LPIPS.

| | | | SceNeRFlow | D-NeRF | DeVRF |
|---|---|---|---|---|---|
| Unmasked | PSNR | ↑ | **29.98** | 24.89 | 25.24 |
| | SSIM | ↑ | **0.939** | 0.913 | 0.905 |
| | LPIPS | ↓ | **0.054** | 0.107 | 0.125 |
| Masked | PSNR | ↑ | **32.03** | 26.13 | 27.66 |
| | SSIM | ↑ | **0.975** | 0.951 | 0.957 |
| | LPIPS | ↓ | **0.017** | 0.060 | 0.041 |

Table 5. **Novel-View Synthesis.** Mean PSNR, SSIM, and LPIPS across scenes, except for Seq. 8 (due to memory limitations from the high resolution).

| | | | SceNeRFlow | No Online | No Extend | No Coarse | No Fine |
|---|---|---|---|---|---|---|---|
| Unmasked | PSNR | ↑ | 30.32 | 30.36 | 28.53 | **30.40** | 28.95 |
| | SSIM | ↑ | **0.939** | 0.929 | 0.927 | 0.936 | 0.933 |
| | LPIPS | ↓ | **0.054** | 0.074 | 0.067 | 0.056 | 0.056 |
| Masked | PSNR | ↑ | 32.38 | **32.61** | 30.20 | 32.45 | 30.45 |
| | SSIM | ↑ | **0.976** | 0.974 | 0.966 | 0.974 | 0.970 |
| | LPIPS | ↓ | **0.018** | 0.020 | 0.027 | 0.020 | 0.021 |
| | MPJPE | ↓ | **1.5** | 33.4 | 15.6 | 10.4 | **1.5** |

Table 6. **Ablations.** Mean PSNR, SSIM, LPIPS, and MPJPE across scenes.
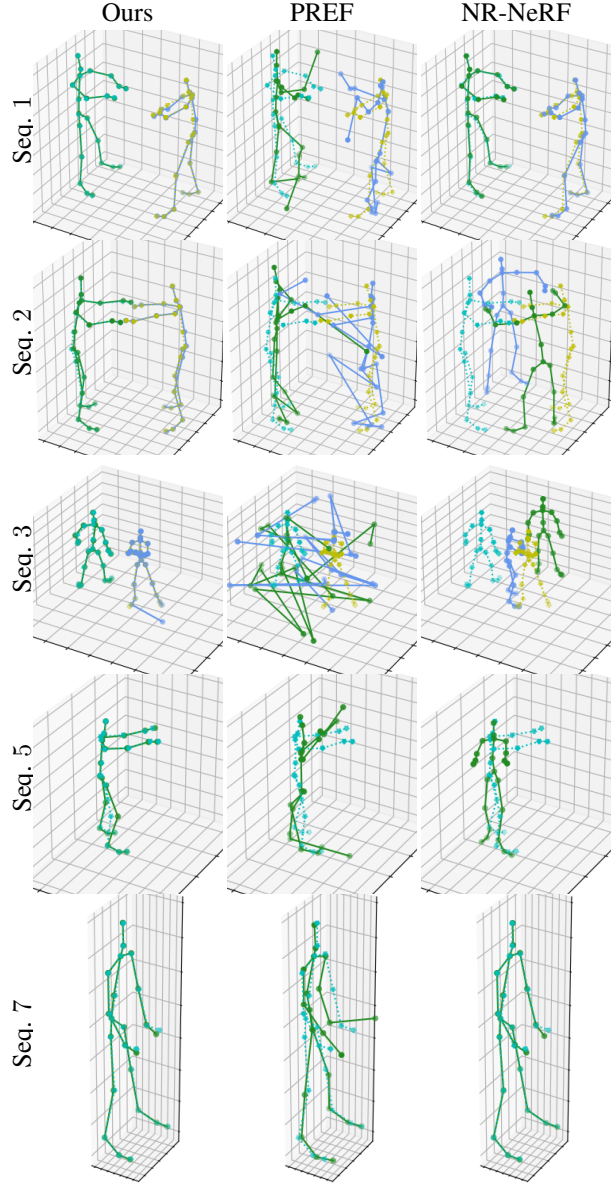


Figure 12. **Time Consistency.** The solid skeleton is the tracking estimate at $t=T$. The dotted skeleton is the pseudo-ground truth at $t=T$.

$u=10$:

$$\hat{\sigma}''_{r,i} = \begin{cases} \frac{1}{u}\hat{\sigma}'_{r,i} & \text{if } \hat{\sigma}'_{r,i} > u \cdot \hat{\sigma}_{r,i} \\ \hat{\sigma}'_{r,i} & \text{else.} \end{cases} \tag{13}$$

Figure 13. **Novel-View Synthesis.** (First row) Seq. 3 at $t=\frac{T}{2}$. (Second row) Seq. 4 at $t=\frac{T}{2}$.

Finally, we weaken the regularization on very small offsets $\Delta \in \mathbb{R}^3$ with a soft threshold of $s_t=0.001$:

$$\hat{\sigma}_{r,i}'''(\Delta) = \text{sig}\left(\frac{4\|\Delta\|_2}{s_t} - 2\right)\hat{\sigma}_{r,i}'', \qquad (14)$$

where $\text{sig}(x)=\frac{1}{1+\exp(-x)}$ is the sigmoid function. We thus have:

$$\mathcal{L}_{\text{norm,w}} = \frac{1}{RS}\sum_r\sum_i \hat{\sigma}_{r,i}'''(\Delta)\mathbb{E}_{\mathbf{e}}\left[\left|\|\mathbf{J}_{\mathbf{r}_r(s_i)}^\top\mathbf{e}\|_2 - 1\right|\right], \qquad (15)$$

where $\Delta=\Delta_c(\mathbf{r}_r(s_i))$ when regularizing the coarse deformations, and $\Delta=\Delta_f(d_c(\mathbf{r}_r(s_i)))$ when regularizing the fine deformations.

**Learning Rates.** We use exponential decay for the learning rates when constructing the canonical model and for each timestamp $t>1$. For $t=1$, we decay by a factor of 0.01

for all parameters. For $t>1$, we decay the coarse deformations by 0.1 when they get optimized, and the fine deformations by 0.1 when they get optimized. All parameters of the canonical model have an initial learning rate of $10^{-2}$. All deformation parameters have an initial learning rate of $10^{-3}$. The vignetting parameters have an initial learning rate of $10^{-2}$ and are the only parameters with no weight decay.

During the first 1000 iterations when building the canonical model, we warm up the learning rates with an additional factor that exponentially increases from 0.01 to 1.

### G.2. Variants

We compare against two variants of our method that use time-varying canonical models and thus neglect correspondences.
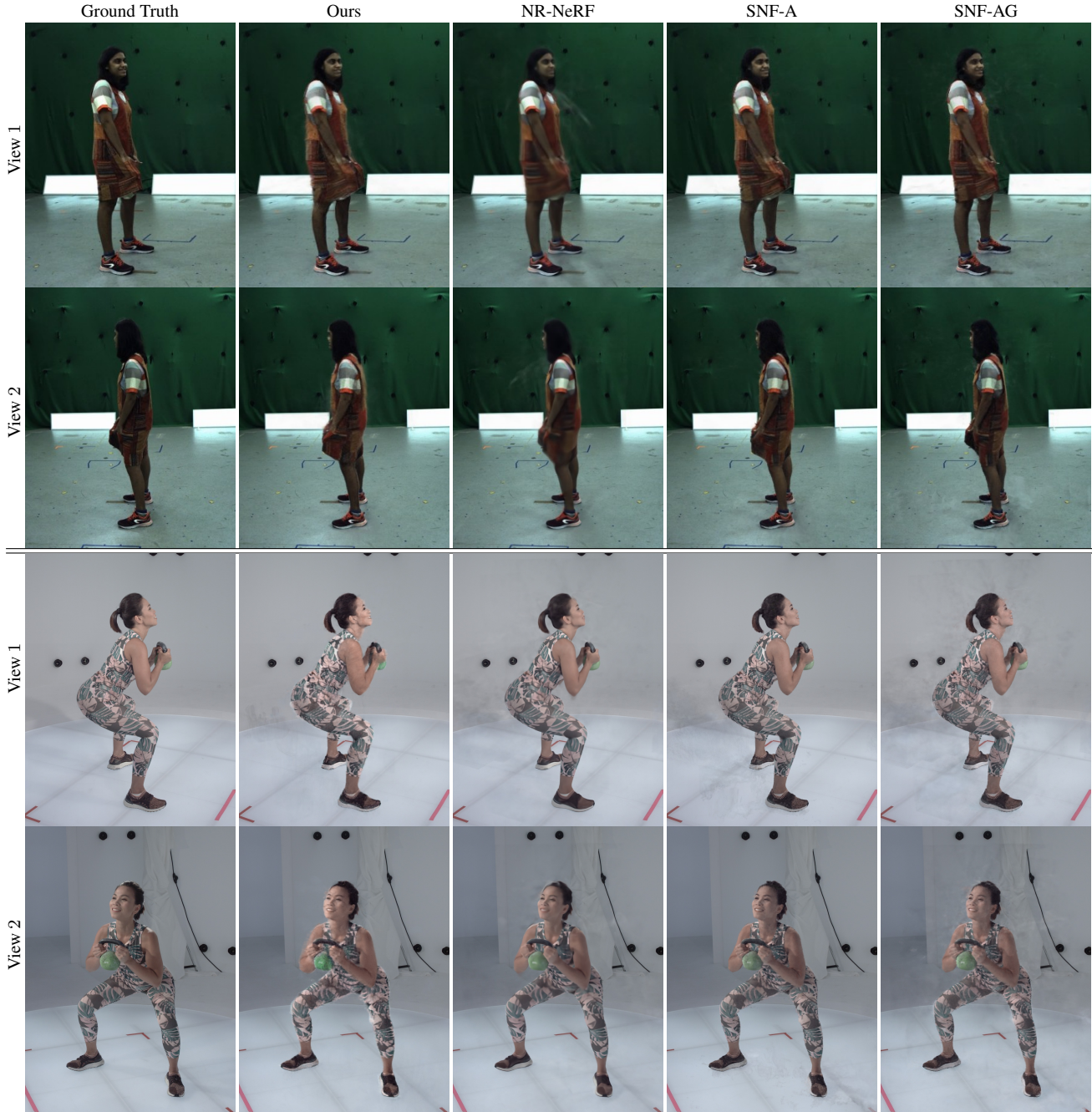
Figure 14. **Novel-View Synthesis.** (First row) Seq. 7 at $t=\frac{T}{2}$. (Second row) Seq. 8 at $t=T$.

The first variant, SNF-A, allows the appearance but not the geometry of the canonical model to change for $t>1$. We implement this via a separate appearance model that has the same HashMLP architecture as the standard canonical model. We can then keep the standard canonical model, which now only predicts the geometry, fixed for $t>1$ while allowing the appearance to vary.

The second variant, SNF-AG, has time-varying appearance and geometry. We thus use the standard canonical model. We apply the geometry regularization losses $\mathcal{L}_{\text{back}}$ and $\mathcal{L}_{\text{hard}}$ to the canonical model at all timestamps $t \geq 1$.

For both variants, we seek to explain as much of the reconstruction as possible via the deformations, such that the canonical model needs to change as little as possible. We thus split the $10k$ iterations per timestamp into three equally long phases: only coarse deformations (as before), then

16

only fine deformations (as before), then only the canonical model. *I.e.*, the canonical model gets optimized only during the third phase, when the deformations are fixed.

We found these variants unstable to train, with frequent divergence, and the following remedies helped: The canonical hash grid(s) use the settings from Instant NGP [53] (a coarsest resolution of 16, with the resolution of each finer level being 1.5 times finer than the previous level, with a total of 16 levels). During the third phase, we use Instant NGP's Huber loss with a threshold of 0.1 as reconstruction loss instead of our $\ell_1$ reconstruction loss. We use an exponentially decaying learning rate for the third phase that goes from $10^{-2}$ to $10^{-3}$ (except for the 300-frame Seq. 3, where we use a ten times lower learning rate for long-term training stability).

### G.3. NR-NeRF

We adapt NR-NeRF [78] to our settings as follows: Since we provide background images to NR-NeRF, we do not use its rigidity network. Furthermore, due to our large-motion setting, we remove its offsets loss that encourages the deformations to remain small. To keep runtime reasonable, we apply pruning. Doing so also enables us to always sample densely and we thus do not apply hierarchical sampling. We also use foreground-focused batches and vignetting correction. Since the recommended [78] training time of $200k$ iterations is only shown for very short scenes, we instead train NR-NeRF for longer on our scenes. Specifically, we extend NR-NeRF's training time to that of our method and thus train NR-NeRF for one third of the number of iterations of our method. We train Seq. 8 for $200k$ iterations, as that is longer.

### G.4. PREF

Like for NR-NeRF, we also supply background images, apply pruning, use foreground-focused batches, and learn the vignetting correction. Since we follow the authors of PREF and split our long scenes into chunks of 25 frames (see Sec. H), we keep the training time the same.

### H. Joint Evaluation Details

We provide details on how we adapt PREF and NR-NeRF to long-term 3D joint tracking.
**PREF.** To train on a longer sequence, PREF [71] splits the sequence into chunks of 25 frames and trains on each chunk independently. We do the same with our scenes. To obtain long-term correspondences across chunks, we overlap the chunks for three frames.
**NR-NeRF.** Unlike SceNeRFlow, NR-NeRF's canonical model does not coincide with the world space at $t{=}1$. We therefore cannot use $\{\hat{\mathbf{p}}_j^1\}_j$ directly as the target canonical positions. Instead, we apply the backward deformation

model at $t{=}1$ to $\{\hat{\mathbf{p}}_j^1\}_j$ to obtain their positions in canonical space. We then have the joint positions in canonical space and the world-space positions at $t{=}1$, which allows us to apply the same tracking procedure as for SceNeRFlow.

### I. Foreground Masks for Evaluation

Since we did not tune the variants beyond making their training stable, they exhibit some significant undesired artifacts in empty space that we did not try to remove. In addition to scores on the full images, we therefore also compute scores that are focused on the actual dynamic object of interest. To this end, we use foreground masks.

During training, we use very coarse and inaccurate foreground masks. However, we use more accurate foreground masks when computing masked scores during evaluation. To also consider reconstructions that are slightly off, these masks are dilated to include the areas surrounding the dynamic foreground in pixel space. Fig. 15 shows example masks.

We use the following procedure to obtain these more accurate foreground masks. We first compute two foreground masks separately: (1) $m_b$ using background subtraction with a threshold $\Delta_t$ followed by a morphological opening (*i.e.*, first erosion, then dilation) of the foreground, and (2) $m_\sigma$. $m_\sigma$ very roughly detects shadows by determining whether the pixel in $\mathbf{I}^{c,t}$ is a scaled version of the background pixel in $\mathbf{B}^c$ (*i.e.*, whether a naïve brightness change of the background has occurred). To this end, we first divide each of the three channels of the pixel in $\mathbf{I}^{c,t}$ by the respective channel of the background pixel. We then take the standard deviation of the resulting three factors and threshold it at $\sigma_t$. We finally apply an opening to obtain the final $m_\sigma$. Since we use a very generous opening for both masks, we combine them into a single foreground mask by considering those pixels foreground that are foreground in both masks. We manually tune $\Delta_t$, $\sigma_t$, and the opening sizes for each sequence.
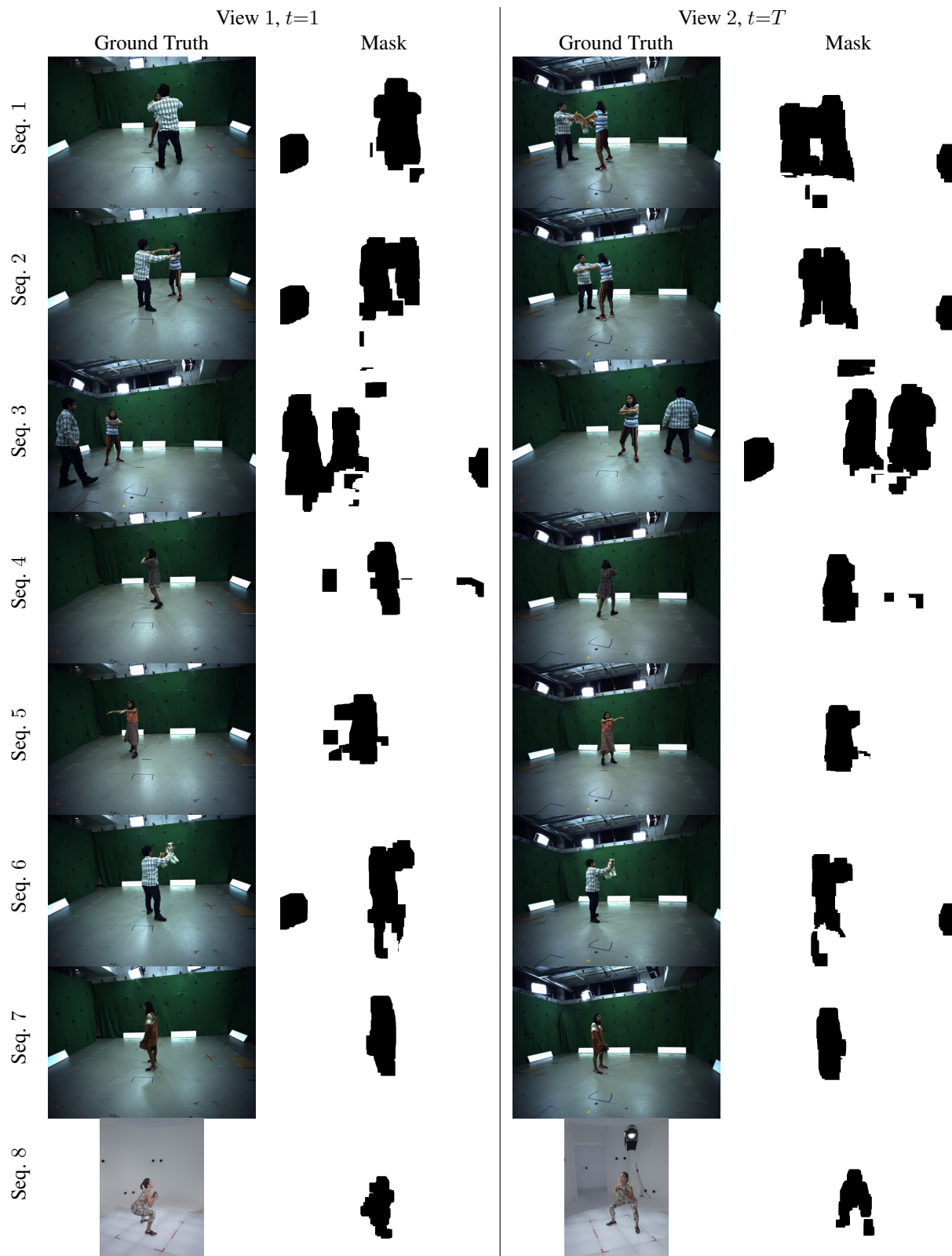
Figure 15. **Foreground Masks for Evaluation.**