# Forward Flow for Novel View Synthesis of Dynamic Scenes

Xiang Guo[1], Jiadai Sun[1,3], Yuchao Dai[†1], Guanying Chen[2], Xiaoqing Ye[3], Xiao Tan[3], Errui Ding[3], Yumeng Zhang[3], Jingdong Wang[3]

[1]Northwestern Polytechnical University, [2]FNii and SSE, CUHK-Shenzhen, [3]Baidu Inc.

## Abstract

*This paper proposes a neural radiance field (NeRF) approach for novel view synthesis of dynamic scenes using forward warping. Existing methods often adopt a static NeRF to represent the canonical space, and render dynamic images at other time steps by mapping the sampled 3D points back to the canonical space with the learned* backward flow *field. However, this backward flow field is non-smooth and discontinuous, which is difficult to be fitted by commonly used smooth motion models. To address this problem, we propose to estimate the* forward flow *field and directly warp the canonical radiance field to other time steps. Such forward flow field is smooth and continuous within the object region, which benefits the motion model learning. To achieve this goal, we represent the canonical radiance field with voxel grids to enable efficient forward warping, and propose a differentiable warping process, including an average splatting operation and an inpaint network, to resolve the many-to-one and one-to-many mapping issues. Thorough experiments show that our method outperforms existing methods in both novel view rendering and motion modeling, demonstrating the effectiveness of our forward flow motion modeling. Project page:* [https://npucvr.github.io/ForwardFlowDNeRF](https://npucvr.github.io/ForwardFlowDNeRF).

## 1. Introduction

Novel view synthesis (NVS) is a challenging and long-standing problem in computer vision and graphics, which has many applications in virtual reality, augmented reality, data augmentation, image editing, *etc*. Recently, differentiable neural rendering [26, 30, 59] has been introduced into this area. In particular, the neural radiance field (NeRF) [26]

---

† Corresponding author (daiyuchao@nwpu.edu.cn). The first three authors also with the Shaanxi Key Laboratory of Information Acquisition and Processing.
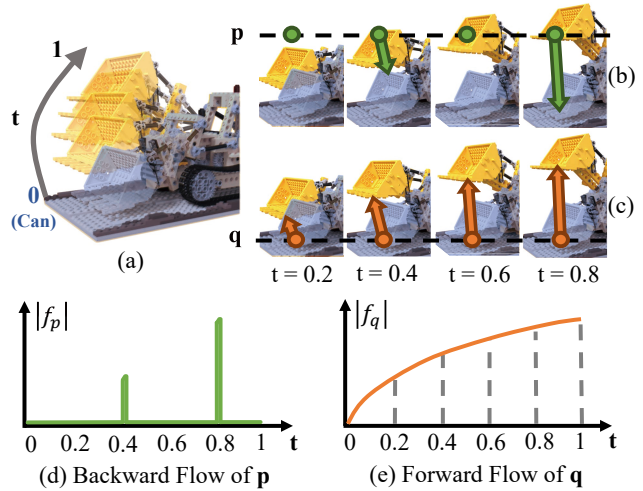


Figure 1. **Comparison of backward flow and forward flow.** This figure shows an example of backward and forward flow changes. **(a)** An example of dynamic scene. **(b)** With the bucket lifting up, different types of points cover the green point **p**, which needs very different backward flows to map this point back to canonical space. **(d)** shows the norm changes of the backward flow, which is not smooth. **(c)** On the other hand, the forward flow of position **q**, which maps the constant object point from canonical space to other times, is smooth and continuous. **(e)** shows the norm changes of the forward flow.

promotes this area significantly and attracts lots of interest within a short time. NeRF [26] produces realistic images by representing the 3D world with a multi-layer perceptron (MLP), which maps the input 3D coordinates and 2D view direction to target density and color.

While the original NeRF [26] can only model static scenes, a series of works extend the NeRF-based framework from static to dynamic scenes [9, 13, 19, 32, 36, 49, 51, 55]. One of the promising directions is using a canonical space representation [15, 36, 49]. This representation sets one of the time steps as canonical time and models the static scene with a canonical radiance field. To render images at other time steps, a deformation field is used to estimate the

*backward flow* for moving the 3D points from the current time step back to the canonical time step. Although the canonical-based representation with the backward flow is easy to implement, the backward flow field is non-smooth and discontinuous. As shown in Fig. 1(b), for a fixed 3D position along the timeline, there will be different types of points covering the position **p**, which needs discontinuous flows to map them back to canonical space (Fig. 1(d)). So the backward flow cannot be fitted well with commonly used smooth motion models (MLP, for example). Also, distortions are introduced to the static canonical space geometry due to the failure of the motion model, as shown in Fig. 8.

To solve the problem of backward flow, we propose using *forward flow* as the deformation model. Instead of warping the sampled points on image rays at other time steps back to the canonical time and rendering at the canonical space, we propose to warp the whole canonical radiance field from the canonical time to other time steps using forward deformation flow and render at corresponding time steps. In this way, the forward flow estimated by the deformation model will be smooth and continuous for the same 3D position along the timeline (Fig. 1(c) and (e)). Note that SNARF [7] has also used forward warping based on an inverse skinning model, but it is designed for dynamic human modeling and cannot be used in general scenes. In this paper, we aim to achieve forward warping for general scenes, which means we must warp the whole space.

However, introducing forward warping into the canonical space based NeRF methods is not straightforward as there are three main problems to be solved. First, the traditional canonical radiance field in existing methods cannot be warped explicitly, since the radiance field is represented as a continuous function parameterized by an MLP. To solve this problem, we propose to use the voxel grid to represent the canonical radiance field as it is finite and discrete. Recent voxel-based methods [27, 44, 62] have proven the effectiveness of this representation. The other two problems are the *many-to-one* and *one-to-many* mapping issues brought by the inherent property of the forward warping operation. To address them, we propose a *differentiable forward warping* method consisting of an *average splatting* operation and an *inpaint network* to solve the *many-to-one* and *one-to-many* issue, respectively. Extensive experiments have been conducted to verify the effectiveness of our method.

Our key contributions can be summarized as follows:

- To the best of our knowledge, we are the first to investigate forward warping in dynamic view synthesis for general scenes. We propose a novel canonical based NeRF with forward flow motion modeling for dynamic view synthesis. Thanks to the forward flow field, our method can better represent the object motions, and explicitly recover the trajectory of a surface point.

- We introduce voxel grid based canonical radiance field to enable reasonable computation of forward warping, and propose a differentiable forward warping method, including an average splatting operation and an inpaint network, to solve the many-to-one and one-to-many issues of forward warping.

- Experiments on multiple datasets show that our method outperforms existing methods on the D-NeRF [36] dataset, Hypernerf [33] dataset, NHR [54] dataset and our proposed dataset.

## 2. Related Works

Novel View Synthesis (NVS) is a long-standing task in both computer vision and graphics [4, 6, 14, 18], and surveys of recent methods can be found in [42, 46, 47]. Some methods require to reconstruct an explicit 3D model to represent the scene, such as point clouds [1], voxels or meshes [16, 40, 41, 48]. Then novel images from arbitrary views can be rendered from this geometry. Another methods try to estimate depth map through multi-view geometry, and then aggregate features from multiple frames through the co-visibility, such as [8, 11, 17, 34, 40, 41, 57]. Recently, neural implicit representations have shown great promise for novel view synthesis and 3D modeling. In this section, we focus more on the neural implicit rendering methods and mainly summarize these schemes according to whether they can handle dynamic scenes.

**NVS for Static Scenes**  NeRF [26], as a seminal work, uses MLPs to model a 5D radiance field, which can render impressive view synthesis for static scenes captured. Numerous subsequent works have extended NeRF to kinds of scenarios, such as larger and unbounded scenes [25, 39, 45, 56, 64], relighting [3, 43, 58, 66], incorporating anti-aliasing for multi-scale rendering [2], and generalization ability [5, 50, 52, 63]. In addition, some methods are devoted to more efficient neural rendering and optimization in NeRF-like framework, such as [20, 22, 24, 28, 35, 62] focus more on efficient sampling along each ray for color accumulation, [37, 38] subdivide the scene into multiple cells for efficient processing, and [27, 44, 61] exploit voxel-grid representation to speed up the optimization of radiance field. However, these methods are mainly applicable to static scenes, leaving out the scenes with dynamic objects, which are actually more extensive and practical.

**NVS for Dynamic Scenes**  There are several works that extend NeRF from static scenes to dynamic scenes with non-rigid deformable objects. One feasible way is to build a 4D spatial-temporal representation. For example, Yoon *et al*. [60] combine single-view depth and depth from multi-view stereo to render virtual views with 3D warping. Gao *et al*. [13] use a time-invariant model (static) and a time-varying model (dynamic) to represent the scenes, and reg-
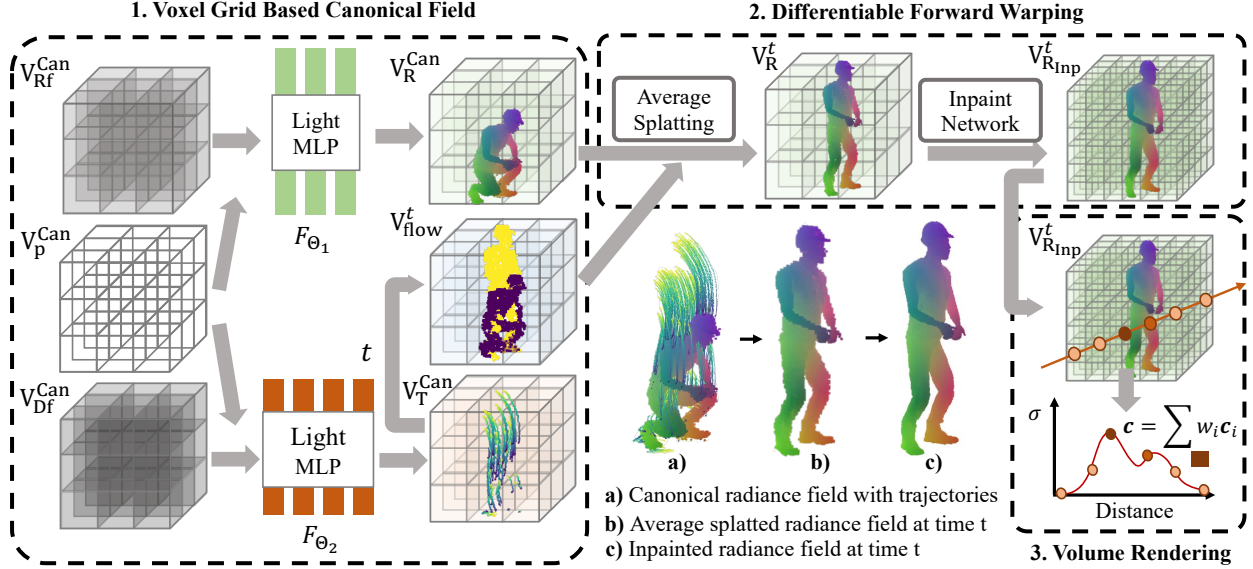
**1. Voxel Grid Based Canonical Field**

**2. Differentiable Forward Warping**

Average Splatting

Inpaint Network

a) Canonical radiance field with trajectories
b) Average splatted radiance field at time t
c) Inpainted radiance field at time t

$$\boldsymbol{c} = \sum w_i \boldsymbol{c}_i$$

Distance

**3. Volume Rendering**

Figure 2. **Overview of our proposed method. a)** We represent a static scene at canonical time with a voxel grid based radiance field for density&color and a voxel grid based trajectory field for deformations; **b)** We propose to first forward warp canonical radiance field using the forward flow by average splatting; **c)** We then inpaint the warped radiance field using a inpaint network; Specifically, **1. Voxel Grid Based Canonical Field** contains two models. The canonical radiance field $\mathbf{V}_R^{Can}$ is estimated by a Light MLP which takes canonical radiance feature $\mathbf{V}_{Rf}^{Can}$ and corresponding 3D coordinates $\mathbf{V}_p^{Can}$ as input. The canonical trajectory field $\mathbf{V}_T^{Can}$ is estimated by another Light MLP which takes deformation feature and coordinates as input. The deformation flow $\mathbf{V}_{flow}^t$ from canonical to time $t$ can then be obtained; **2. Differential Forward Warping** first warp $\mathbf{V}_R^{Can}$ to get radiance field $\mathbf{V}_R^t$ at time $t$. Then, the $\mathbf{V}_R^t$ is inpainted by a inpaint network, which is $\mathbf{V}_{R_{Inp}}^t$; **3. Volume Rendering** render colors of rays at time $t$ based on $\mathbf{V}_{R_{Inp}}^t$

ularize the dynamic model by scene flow estimation. NeR-Flow [9] learns a 4D spatial-temporal representation of a dynamic scene from a set of RGB images. Xian *et al.* [55] build a 4D space-time irradiance field to map a spatial-temporal location to the emitted color and volume density. Similarly, NSFF [19] models the dynamic scene as a time-variant continuous function of appearance, geometry, and 3D scene motion. DCT-NeRF [51] uses the Discrete Cosine Transform (DCT) to capture dynamic motion, *i.e.* learning smooth and stable trajectories over time for each point in space.

On the other hand, D-NeRF [36], Nerfies [32], HyperN-eRF [33] and NR-NeRF [49] use a static canonical radiance field to capture geometry and appearance, and then learn a deformation/displacement field at each time step w.r.t. the canonical space. Specifically, to render an image at an arbitrary time step, a deformation field is used to estimate backward scene flow, moving 3D points from the current time step back to the canonical step. However, for the same 3D location along the timeline, the backward flow field is not guaranteed to be smooth and continuous. As a result, the canonical geometry usually has distortions and resembles the mean shape of a moving object. We focus on solving the problem of backward flow in this paper.

Along with the two main directions, there is a trend to speed up the training of dynamic NeRF, which is based on voxel grid representation. TiNeuVox [10] models the de-

formation using a tiny MLP and uses multi-distance interpolation to get the feature for the radiance network which estimates the density and color. V4D [12] uses the 3D feature voxel to model the 4D radiance field with additional time dimension concatenated and proposes look-up tables for pixel-level refinement. Although V4D mainly focuses on improving image quality, the training speed is not significant compared to TiNeuVox. DeVRF [21] also builds on voxel-grid representation, which proposes to use multi-view data to overcome the nontrivial problem of the monocular setup. Multi-view data simplifies the learning of motion and geometry compared with others using monocular images.

## 3. Method

**Motivation** Backward-warping based methods [15, 36] propose a network $f_{t\to Can} = F(p, t)$ to estimate the deformation flow $f_{t\to Can}$ which moves the point at the position $p$ from other time steps $t$ back to canonical time Can. However, for the same position $p$, at different time steps $t$s, there could be different object points or even the empty point at this position, as shown in Fig. 1. This means the deformation flow $f_{t\to Can}$ is non-smooth and discontinuous along the timeline for a fixed position $p$. This could introduce difficulties for motion learning and produce distortion in the canonical radiance field, because the backward flow network has limited capacities in learning a correct non-smooth and dis-
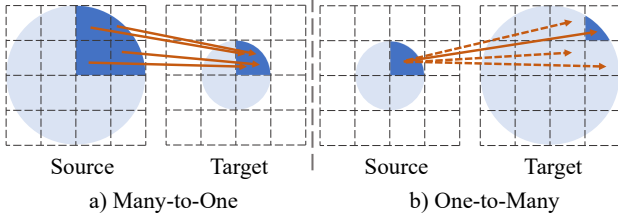
a) Many-to-One        b) One-to-Many

Figure 3. **Two issues of the forward warping.** Many-to-one: multiple positions in the source correspond to the same position in the target. One-to-many: one position in the source corresponds to multiple positions in the target, causing holes if warping naively.

continuous deformation field. On the other hand, our proposed forward warping strategy estimates the deformation $f_{\text{Can}\rightarrow t}$ from the canonical time to another time step of a 3D point in canonical space. The set of the deformation flows of one position along the timeline is actually the trajectory of this point. This guarantees the flows to be smooth and continuous along the timeline as we assume the motions in reality have these properties.

**Method Overview** We model a static scene at the canonical time using the voxel grid representation with a canonical radiance field and a canonical trajectory field. To synthesize dynamic images, we propose to forward warp the canonical radiance field to corresponding time steps and render the images using volume rendering based on the warped radiance field. Fig. 2 shows the overview of our method. In the following, we will introduce three main components of our method: voxel grid based canonical field, differentiable forward warping, and volume rendering. Finally, the model optimization is represented, including the proposed loss functions and training strategy.

## 3.1. Voxel Grid Based Canonical Field

**Canonical Radiance Field** To warp a static scene in canonical time to other time steps using forward warping, the radiance field should be represented by finite 3D points. The original heavy MLP in NeRF [26] is not practical, as we cannot query infinite 3D points in the canonical space. Inspired by recent works [27,44,61] on voxel grid presentations, we propose to use a learnable voxel radiance feature grid $\mathbf{V}_{\text{Rf}}^{\text{Can}}$ and a lightweight MLP network $F_{\theta_1}$ to model the radiance of the static scene as shown in Fig. 2. The canonical radiance field is defined as

$$\mathbf{V}_{\text{R}}^{\text{Can}} = F_{\theta_1}(\mathbf{V}_{\text{p}}^{\text{Can}}, \mathbf{V}_{\text{Rf}}^{\text{Can}}), \tag{1}$$

where $\mathbf{V}_{\text{R}}^{\text{Can}} = \{\mathbf{V}_{\sigma}^{\text{Can}}, \mathbf{V}_{\text{cf}}^{\text{Can}}\}$ consists of density voxel grid $\mathbf{V}_{\sigma}^{\text{Can}}$ and color feature voxel grid $\mathbf{V}_{\text{cf}}^{\text{Can}}$. $\mathbf{V}_{\text{p}}^{\text{Can}}$ are the 3D coordinates of the voxel grid in canonical space, embedded according to NeRF [26].

**Canonical Trajectory Field** We propose to use Discrete Cosine Transform (DCT) [51] to represent the trajectory of

a 3D point, which ensures the smoothness of the motion. Similar to canonical radiance field, we also use a deformation feature grid $\mathbf{V}_{\text{Df}}^{\text{Can}}$ and a lightweight MLP network $F_{\theta_2}$ to model the canonical trajectory field, which is defined as

$$\mathbf{V}_{T}^{\text{Can}} = F_{\theta_2}(\mathbf{V}_{\text{p}}^{\text{Can}}, \mathbf{V}_{\text{Df}}^{\text{Can}}), \tag{2}$$

where $\mathbf{V}_{T}^{\text{Can}}$ contains the DCT coefficients of the voxels. Given time step $t$, we could get the deformation flow of these voxels from canonical to time step $t$ by

$$\mathbf{V}_{\text{flow}}^{t} = f_{\text{DCT}^{-1}}(\mathbf{V}_{T}^{\text{Can}}, t) - f_{\text{DCT}^{-1}}(\mathbf{V}_{T}^{\text{Can}}, \text{Can}), \tag{3}$$

where $f_{\text{DCT}^{-1}}$ is the inverse DCT transform of the DCT coefficients as described in TrajectoryNeRF [51].

## 3.2. Differentiable Forward Warping

Forward warping has two major issues: many-to-one and one-to-many. Fig. 3 shows a simple example, where we aim to warp the source to the target using forward warping: a) Many-to-one: if a circle shrinks, there will be multiple positions corresponding to one position from the source to the target. b) One-to-many: if a circle expands, multiple positions in the target will correspond to the same position in the source, which leaves holes of the target if we warp the source to target naively. For the many-to-one issue, we propose to use the average splatting to fuse multiple values into one. For the one-to-many issue, we use an inpaint network to inpaint the missing positions of the warped grids. More details are in the following sections.

**Average Splatting** We propose to fuse possible multiple values from the source grid that are mapped into the same voxel position of the target grid, motivated by Softmax-splatting [31]. Specifically, we propose a simple yet effective method that calculates the "average" of these values with a trilinear kernel. Formally, suppose that we need to warp the source grid $\mathbf{V}^{\text{S}}$ to target grid $\mathbf{V}^{\text{T}}$ by the flow $f_{\text{S}\rightarrow\text{T}}$, and $\mathbf{p}$, $\mathbf{q}$ are indexes of a voxel grid. We define $\mathbf{V}^{\text{T}} = F_{\text{warp}}(\mathbf{V}^{\text{S}}, f_{\text{S}\rightarrow\text{T}})$ as follows,

$$\mathbf{V}^{\text{T}}[\mathbf{p}] = \frac{\sum_{\forall \mathbf{q} \in \mathbf{V}^{\text{S}}} b[\mathbf{u}] \cdot \mathbf{V}^{\text{S}}[\mathbf{q}]}{\sum_{\forall \mathbf{q} \in \mathbf{V}^{\text{S}}} b[\mathbf{u}]}, \tag{4}$$

$$b[\mathbf{u}] = \prod \max(0, 1 - |\mathbf{u}_i|), i \in \{x, y, z\}, \tag{5}$$

$$\mathbf{u} = (\mathbf{q} + f_{\text{S}\rightarrow\text{T}}[\mathbf{q}]) - \mathbf{p}, \tag{6}$$

where $x, y, z$ are three axes of the voxel grid, and $\mathbf{u}_i \in \mathbb{R}$ is one element of the vector $\mathbf{u} \in \mathbb{R}^3$.

So, we warp the canonical radiance field to time step $t$ by

$$\mathbf{V}_{\text{R}}^{t} = F_{\text{warp}}(\mathbf{V}_{\text{R}}^{\text{Can}}, \mathbf{V}_{\text{flow}}^{t}). \tag{7}$$

**Inpaint Network** Another issue of forward warping is one-to-many, which leaves holes in the voxel grid after
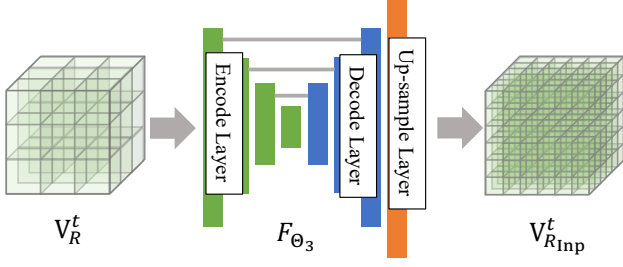
Figure 4. **Inpaint network.** It consists of a 3D U-Net structure and an up-sampling layer, which could inpaint and upsample the input voxel grid.

warping. To resolve this issue, we propose building an inpaint network to fill these holes. As shown in Fig. 4, we modify a 3D U-Net structure network [53] followed by an up-sampling layer, which is $\mathbf{V}_{R_{Inp}}^t = F_{\theta_3}(\mathbf{V}_R^t)$, where $\mathbf{V}_{R_{Inp}}^t$ is the inpainted and up-sampled voxel grid.

This structure helps encoding layers to learn information from local neighboring voxels and decoding layers to recover the original resolution with filled content. We only feed the warped voxel grid to the inpaint network, which has no access to temporal information. Thus, the inpaint network can not learn to warp, limiting its ability only to inpainting. Since warping operations and the inpaint network are relatively expensive to compute, we propose to attach an up-sample layer after the U-Net structure. In this way, we could warp the grids with a low resolution to save time and memory, while rendering with a higher resolution for better image quality. The up-sample layer consists of a interpolation layer, convolution layers, and activation layers.

### 3.3. Volume Rendering

After we get the radiance voxel grid at time $t$, the pixel colors of the image rays can be rendered using volume rendering techniques [26]. Given a ray $\mathbf{r}(w) = \mathbf{o} + w\mathbf{d}$ emitted from the camera center $\mathbf{o}$ with view direction $\mathbf{d}$ through a given pixel on the image plane, we render the corresponding pixel color $\mathbf{C}_{Inp}(\mathbf{r}) = F_{render}(\mathbf{V}_{R_{Inp}}^t, \mathbf{r})$. To do this, we get all 3D points $\mathbf{p}$ that the ray intersects with voxel grids. Next, trilinear interpolation is applied to obtain the density and color feature of each 3D point.

$$(\sigma, \mathbf{c_f}) = F_{inter}(\mathbf{V}_{R_{Inp}}^t, \mathbf{p}). \tag{8}$$

Then, we get the color of this 3D point $\mathbf{p}$ by

$$\mathbf{c} = \begin{cases} \mathbf{c_f} & \text{if } \mathbf{c_f} \in \mathbb{R}^3, \\ F_{\theta_4}(\mathbf{c_f}, \mathbf{d}) & \text{otherwise,} \end{cases} \tag{9}$$

where $F_{\theta_4}$ is a small MLP network that is used to produce colors that depend on the direction of the ray.

Finally, the pixel color could be rendered by

$$\mathbf{C}(\mathbf{r}) = \sum_{k=1}^{K} T(w_k)\,\alpha\left(\sigma(w_k)\delta_k\right)\mathbf{c}(w_k), \tag{10}$$

$$T(w_k) = \exp\left(-\sum_{j=1}^{k-1} \sigma(w_j)\delta_j\right), \tag{11}$$

where $K$ is the number of sampled points along the ray, $\delta_k$ is the distance between adjacent samples along the ray, and $\alpha\left(\sigma(w_k)\delta_k\right) = 1 - \exp(-\sigma(w_k)\delta_k)$.

Note that the warped radiance grid $\mathbf{V}_R^t$ should also render reasonable images. So we directly up-sample $\mathbf{V}_R^t$ to get $\mathbf{V}_{R_{Up}}^t$ without inpainting, and render the corresponding pixel color $\mathbf{C}_{Up}(\mathbf{r}) = F_{render}(\mathbf{V}_{R_{Up}}^t, \mathbf{r})$.

### 3.4. Model Optimization

To optimize the model, we design a series of loss functions. **Photometric Loss** First, the most important loss is the photometric loss, which is the mean square error (MSE) of the rendered color and ground truth color $\mathbf{C}_{gt}(\mathbf{r})$.

$$\begin{aligned} \mathcal{L}^{photo} = &\frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\mathbf{C}_{Inp}(\mathbf{r}) - \mathbf{C}_{gt}(\mathbf{r})\|_2^2 \\ &+ \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \|\mathbf{C}_{Up}(\mathbf{r}) - \mathbf{C}_{gt}(\mathbf{r})\|_2^2, \end{aligned} \tag{12}$$

where $\mathcal{R}$ is the set of rays sampled in one batch.

Following DVGO [44], we use $\mathcal{L}^{ptc}$ to directly supervise the color of sampled points. The intuition is that sampled points with larger weights contribute more to the rendered color. Also, we use the background entropy loss $\mathcal{L}^{bg}$ to encourage the densities concentrating on either the foreground or the background.

**Inpaint Network Loss** As the inpaint network is used to complete the warped voxel grids, the output of the inpaint network should be close to the input. Therefore, we propose to use the L1 norm of the difference between the input and output of the inpaint network.

$$\mathcal{L}^{vdiff} = \left\| \mathbf{V}_{R_{Inp}}^t - \mathbf{V}_{R_{Up}}^t \right\|_1. \tag{13}$$

**Regularization Terms** Since modeling the appearance and motion of a dynamic scene from monocular images is a non-trivial problem, we propose a series of regularization terms.

First, we encourage most 3D points in canonical space to be static by introducing $\mathcal{L}^{flow}$ to be the L1 norm of $\mathbf{V}_{flow}^t$. Second, denoting $\mathcal{L}^{tv}$ as the total variation function. We compute $\mathcal{L}^{tv}(\mathbf{V}_\sigma^{Can})$ to ensure the spatial smoothness of the density in canonical space, and introduce $\mathcal{L}^{tv}(\mathbf{V}_{flow}^t)$ to encourage the spatial smoothness of the motions. Finally, inspired by RegNeRF [29], we render depths $\mathbf{D}$ of image patches sampled from random views and minimize their total variations $\mathcal{L}^{tv}(\mathbf{D})$.

Table 1. **Quantitative comparison on D-NeRF Dataset.** Comparison of our method with others on LPIPS and PSNR/SSIM. The **Red** text indicates the best and **blue** text is the second best result.

| Methods | Type | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|---|
| T-NeRF [36] | N | 29.51 | 0.95 | 0.08 |
| TiNeuVox-S $(100^3)$ [10] | NPC | 30.75 | **0.96** | 0.07 |
| TiNeuVox-B $(160^3)$ [10] | NPC | **32.67** | **0.97** | **0.04** |
| D-NeRF [36] | PC | 30.50 | 0.95 | 0.07 |
| NDVG $(160^3)$ [15] | PC | 30.54 | **0.96** | **0.05** |
| Ours $(80^3)$ | PC | **32.68** | **0.97** | **0.04** |

**Overall Loss**   The overall loss function used for optimization defined as follow, where $w_{1\sim7}$ are weights to balance each component in the final coarse loss.

$$\mathcal{L} = \mathcal{L}^{\text{photo}} + w_1\mathcal{L}^{\text{ptc}} + w_2\mathcal{L}^{\text{bg}} + w_3\mathcal{L}^{\text{flow}} + w_4\mathcal{L}^{\text{vdiff}}$$
$$+ w_5\mathcal{L}^{\text{tv}}(\mathbf{V}_\sigma^{\text{Can}}) + w_6\mathcal{L}^{\text{tv}}(\mathbf{V}_{\text{flow}}^t) + w_7\mathcal{L}^{\text{tv}}(\mathbf{D}), \quad (14)$$

**Training Strategy**   First, we propose *progressive training*, which starts training with images close to the canonical time and progressively adds images with farther time steps until all images are added. Second, we set up a *coarse-to-fine training strategy*. In the coarse stage, we set the color feature to be the color itself, without considering the ray directions. We compute a smaller bounding box with the proxy geometry learned from the coarse stage, which could filter a large portion of empty space for the fine stage training. The model trained in the fine stage is our final model, and we set the color features to be high dimensional features and model the ray direction dependency with $F_{\theta_4}$.

## 4. Experiments

### 4.1. Baselines and Evaluation Datasets

**Baselines**   The methods compared in this paper are classified into three types. First, we compare methods which are *non-canonical based* (**N**), including NeRF[26], and T-NeRF[36]. Second, *physical canonical based methods* (**PC**) contain D-NeRF [36], NDVG[15] and Ours. These methods set their canonical space as one frame of the whole timeline by explicitly giving zeros to estimated flows at the canonical time. The canonical space of these methods should have reasonable physical 3D reconstructions of the scene at the canonical time. Third, *non-physical canonical based methods* (**NPC**) consists of NV[23], Nerfies[32], HyperNeRF[33] and TiNeuVox[10]. Their canonical space geometries do not necessarily have physical meanings.

**D-NeRF Dataset**   We evaluate our method on 8 dynamic scenes of D-NeRF [36] dataset. We report several common metrics for the evaluation: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity (SSIM) and Learned Perceptual Image Patch Similarity (LPIPS) [65].



GT    Ours    D-NeRF   TiNueVox-S  TiNueVox-B

Figure 5. **Qualitative comparison on D-NeRF Dataset.** We show some novel view synthesized images on the selected test set of the dataset. Comparing ours with ground truth, D-NeRF [36] and TiNeuVox [10]. Our model yields cleaner images with more details.

**HyerNeRF Dataset**   Besides synthetic datasets, we also conduct experiments on real scenes, proposed by HyperNeRF [33]. This dataset uses a multi-view camera rig consisting of 2 phones to capture real unbounded scenes with challenging rigid and non-rigid deformations.

**NHR Dataset**   NHR [54] dataset capture dynamic human data using a multi-camera dome system with up to 80 cameras arranged on a cylinder. We conduct experiments on four scenes where the performers are in different clothing and perform different actions. We test with 100 frames of each scene, selecting 90% views for training and 10% views for testing.

**Lego Complete Dataset**   The D-NeRF [36] dataset only contains 20 test images for each scene, and the evaluations of D-NeRF do not contain images in the canonical space. However, for *physical canonical based methods* (**PC**), canonical space geometry is important, as it affects images quality rendered at the canonical time and reflects the quality of flow estimation. Hence, we build a new dataset, named *Lego Complete Dataset*, which animates the object *LEGO* in D-NeRF dataset with 3 different motion patterns. For each scene, the test set is split into 3 categories to evaluate three abilities: *space interpolation* (rand views for each train time step), *time interpolation*(times between trained time steps), and *canonical interpolation* (rand views for canonical time step) abilities. For more details, please refer to our *supp*.

Besides the image quality, we also consider geometry precision by evaluating the depth of each test image to analyze the reconstructed geometry. For depth metrics, we use the relative error $\delta = \max\left(\frac{d}{d_{\text{gt}}}, \frac{d_{\text{gt}}}{d}\right)$, where $d$ is the estimated depth and $d_{\text{gt}}$ is the ground truth depth. We report the percentage of the pixels with $\delta < \text{Threshold}$.

Table 2. **Quantitative comparison on real scene HyperNeRF dataset.** Comparison of our method with others on PSNR and MS-SSIM.

| Methods | Type | 3D Printer | | Broom | | Chicken | | Peel Banana | | Mean | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | MS-SSIM↑ | PSNR↑ | MS-SSIM↑ | PSNR↑ | MS-SSIM↑ | PSNR↑ | MS-SSIM↑ | PSNR↑ | MS-SSIM↑ |
| NeRF [26] | N | 20.7 | 0.780 | 19.9 | 0.653 | 19.9 | 0.777 | 20.0 | 0.769 | 20.1 | 0.745 |
| NV [23] | NPC | 16.2 | 0.665 | 17.7 | 0.623 | 17.6 | 0.615 | 15.9 | 0.380 | 16.9 | 0.571 |
| Nerfies [32] | NPC | 20.6 | 0.830 | 19.2 | 0.567 | 26.7 | 0.943 | 22.4 | 0.872 | 22.2 | 0.803 |
| HyperNeRF [33] | NPC | 20.0 | 0.821 | 19.3 | 0.591 | 26.9 | **0.948** | 23.3 | **0.896** | 22.4 | 0.814 |
| TiNeuVox-S ($100^3$) [10] | NPC | **22.7** | 0.836 | **21.9** | **0.707** | 27.0 | 0.929 | 22.1 | 0.780 | 23.4 | 0.813 |
| TiNeuVox-B ($160^3$) [10] | NPC | **22.8** | **0.841** | 21.5 | 0.686 | **28.3** | **0.947** | **24.4** | **0.873** | **24.3** | **0.837** |
| NDVG ($160^3$) [15] | PC | 22.4 | 0.839 | **21.5** | 0.703 | 27.1 | 0.939 | 22.8 | 0.828 | 23.3 | 0.823 |
| Ours ($70^3$) | PC | **22.8** | **0.845** | **21.9** | **0.715** | **28.0** | 0.944 | **24.3** | 0.865 | **24.2** | **0.842** |



GT     Ours     NDVG     TiNueVox-B     TiNueVox-S       GT     Ours     NDVG     TiNueVox-B     TiNueVox-S
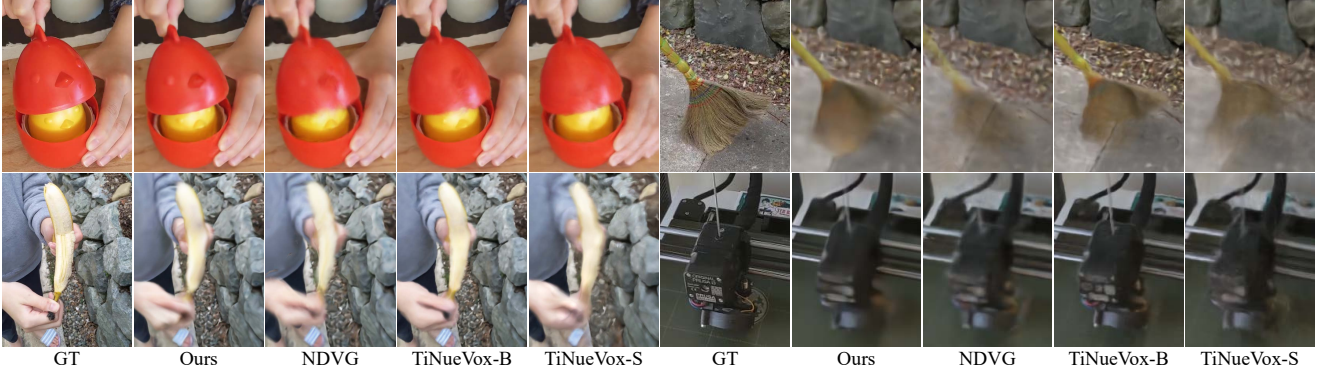
Figure 6. **Qualitative comparison on HyperNeRF Dataset.** Our results are closer to ground truth than other methods.

## 4.2. Experimental Results

**Evaluation on D-NeRF Dataset** We compared our method with the canonical-based backward flow methods D-NeRF [36], NDVG [15] and TiNeuVox [10] on D-NeRF Dataset. As shown in Table 1, our method achieves significant improvements in terms of all three metrics for *physical canonical based methods*. On average, our method improves PSNR by 2.18 compared with D-NeRF and 2.14 compared with NDVG. For *non-physical canonical based methods*, our method also achieves the best result. Ours show obviously better qualitative rendering quality compared with TiNeuVox-S ($100^3$) [10] with even smaller voxel resolution. More detailed statistics are provided in *supp*.

We also provide some visual comparisons in Fig. 5. Ours could recover accurate and detailed images, *e.g.*, the helmet and arm in the top scene, and could also produce cleaner boundaries, *e.g.* the hand and feet in the bottom scene.

**Evaluation on HyperNeRF Dataset** We further compare our method with some highly related works on real scene dataset proposed by [33]. As shown in Table 2, our method achieves consistently better performance among *physical canonical based methods*. Our method has a voxel grid resolution limitation, due to we need to warp the whole voxel grid for rendering. Compared with other voxel grid based methods, our voxel grid resolution is limited to $70^3$, along with a $160^3$ static voxel grid for background modeling. However, the performance of our method is still com-



GT    Ours    T-B [57]   T-S [57]   NDVG [12]

Figure 7. **Visual comparisons on NHR dataset.**

parable with SOTA of *non-physical canonical based methods*, proving the effectiveness of our forward flow design.

In Fig. 13, we show visual comparisons with NDVG [15], and TiNueVox [10]. The forward flow helps to recover the correct structure of dynamic objects, like the eye and edges of the chicken (top left) and the broom (top right). The relatively low resolution may harm the ability to recover very fine details, like the patterns on the 3D printer compared with TiNueVox-B [10] (bottom right).

**Evaluation on NHR Dataset** We also test our method on NHR dataset. Quantitative results are shown in Table 8, which shows ours achieving best results consistently. We show qualitative comparison in Fig. 11, and our method could render clean and detailed images.

Table 3. **Quantitative comparison on NHR dataset.** The <span style="color:red">red</span> text indicates the best and <span style="color:blue">blue</span> text is the second best result.

| Methods | type | Sport 1 | | | Sport 2 | | | Sport 3 | | | Bacsketball | | | Mean | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| TiNeuVox-S [10] | NPC | 26.06 | **0.93** | **0.10** | 25.98 | 0.93 | 0.11 | 25.90 | **0.93** | **0.11** | 23.75 | 0.91 | 0.14 | 25.42 | 0.92 | 0.12 |
| TiNeuVox-B [10] | NPC | **26.44** | **0.93** | **0.10** | **26.68** | **0.94** | **0.10** | **26.09** | **0.93** | **0.11** | **25.06** | **0.92** | **0.12** | **26.07** | **0.93** | **0.11** |
| NDVG [15] | PC | 23.66 | 0.89 | 0.15 | 24.43 | 0.91 | 0.13 | 22.54 | 0.88 | 0.16 | 22.55 | 0.89 | 0.17 | 23.29 | 0.89 | 0.15 |
| Ours | PC | <span style="color:red">27.71</span> | <span style="color:red">0.95</span> | <span style="color:red">0.08</span> | <span style="color:red">27.89</span> | <span style="color:red">0.95</span> | <span style="color:red">0.08</span> | <span style="color:red">27.57</span> | <span style="color:red">0.94</span> | <span style="color:red">0.08</span> | <span style="color:red">24.85</span> | <span style="color:red">0.93</span> | <span style="color:red">0.11</span> | <span style="color:red">27.00</span> | <span style="color:red">0.94</span> | <span style="color:red">0.09</span> |

Table 4. **Quantitative comparison on Lego Complete Dataset.** We report the average PSNR, LPIPS and the average relative depth error.

| Methods | Type | Space Interpolation | | | | Time Interpolation | | | | Canonical Interpolation | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | LPIPS↓ | $\delta < 1.25$ | $\delta < 1.25^2$ | PSNR↑ | LPIPS↓ | $\delta < 1.25$ | $\delta < 1.25^2$ | PSNR↑ | LPIPS↓ | $\delta < 1.25$ | $\delta < 1.25^2$ |
| D-NeRF [36] | PC | 23.98 | 0.15 | 98.720 | 99.328 | 24.18 | 0.15 | 98.854 | 99.378 | 17.36 | 0.23 | 97.341 | 98.868 |
| NDVG ($160^3$) [15] | PC | 28.12 | **0.06** | 99.655 | 99.853 | 28.19 | **0.06** | 99.749 | 99.953 | 19.90 | 0.11 | 94.828 | 97.560 |
| Ours_notv ($80^3$) [1] | PC | 27.58 | 0.07 | 99.758 | 99.958 | 27.63 | 0.07 | 99.788 | 99.966 | 23.56 | 0.13 | 99.238 | 99.820 |
| Ours_noup ($80^3$) [2] | PC | 26.11 | 0.09 | 99.596 | 99.981 | 26.19 | 0.09 | 99.636 | 99.983 | 24.64 | 0.10 | 99.230 | 99.925 |
| Ours_noinp ($80^3$) [3] | PC | 27.21 | 0.08 | 99.783 | 99.979 | 27.27 | 0.08 | 99.797 | 99.980 | 23.91 | 0.09 | 99.586 | 99.963 |
| Ours ($80^3$) | PC | **28.18** | **0.06** | **99.818** | **99.985** | **28.22** | **0.06** | **99.839** | **99.987** | **25.63** | **0.07** | **99.654** | **99.968** |

[1] not use the all total variation losses   [2] not up-sample the voxel grid   [3] not use the inpaint network

## 4.3. Method Analysis

**Flow Estimation and Canonical Space Geometry**  Since evaluating the estimated flow with ground truth is not practical, we choose to compare the reconstructed canonical geometry of *physical canonical based methods*. The idea is that if the canonical space is one frame of the whole time-line (like D-NeRF [36] and NDVG [15]), better flow estimation will result in better canonical scene geometry. We test these methods on Lego Complete Dataset, and all set the canonical space as the first frame of the image sequence.

In Table 4, D-NeRF can not reconstruct the canonical radiance field well, since both the image and depth quality in canonical interpolation are significantly worse than those in space interpolation and time interpolation. NDVG [15] could recover much better images and geometry for other time steps, but the differences between canonical interpolation with the other two are even more significant. This means NDVG [15] has difficulty estimating correct backward flows. To render good quality images at other time steps, the correct geometries (depth) at these time steps are warped back into distorted canonical geometry due to incorrect backward flows.

This is the proof of our claim that the non-smooth and discontinuous nature of the backward flow deformation field makes it difficult to fit with smooth functions, especially for NDVG [15], which has limited MLP capacity for fast training. However, our method performs significantly better for both image quality and depth, and there are fewer variations between canonical, space, and temporal interpolations. This demonstrates the benefit of forward warping, which makes it easier for the deformation model to learn deformation flows from canonical time to other times.

We compare the canonical image of our method with D-NeRF [36] in Fig. 8. Our method could recover the correct geometry and image details of the canonical frame. For example, the arms of the bucket produced by D-NeRF are at the "mean" positions of the whole trajectory, and ours recover the correct status.

**Trajectory Visualization**  We show the learned DCT trajectories generated by our canonical trajectory field in Fig. 11 and Fig. 9. The canonical trajectory field is capable of recovering reasonable object point trajectories. Backward warping designs cannot track the motions of the same object point, so this is not a viable solution. With this feature, geometry constraints, motion models, and prior knowledge could be introduced in future work.

**Regularization terms**  We study the effect of all regularization terms we proposed, including $\mathcal{L}^{\text{flow}}$, $\mathcal{L}^{\text{vdiff}}$ and $\mathcal{L}^{\text{tv}}$ in Table 5. We could observe that all these three regularization terms has positive effect on the performance, but the improvement is minor. This proves the improvement of our method compared with others come from the forward warping desgin we proposed in the paper. Also, backward flow based method NDVG [15] use similar regularization terms with ours, and our method has clear advantage compared with NDVG [15].

**Canonical setup**  In our method, we set canonical time to be the first frame for D-NeRF dataset to compare with *physical canonical based method* and the middle frame for HyperNeRF dataset for better performance. Setting canonical time to be the middle frame helps improving the performance as the state of the middle frame geometry is closer to other time steps compared with the first frame. To prove this, we set canoical time to be middel frame in Table 5 (can-time t at mid), and the PSNR improves slightly.

**Photmetric loss for $\mathbf{V_{R_{Up}}}$**  We use photometric terms on $\mathbf{V_{R_{Up}}}$ to make sure the warped grid before inpainting could already render reasonable images. This make sure the UNet
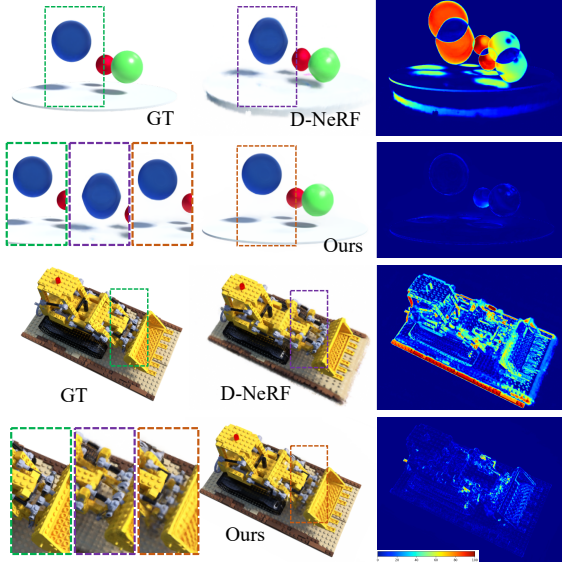
**Figure 8. Canonical qualitative comparison.** We show canonical radiance field comparison with D-NeRF [36]. Given the error map between the ground truth and rendered images, we can see that the canonical frame yielded by ours is closer to the ground truth. The re...
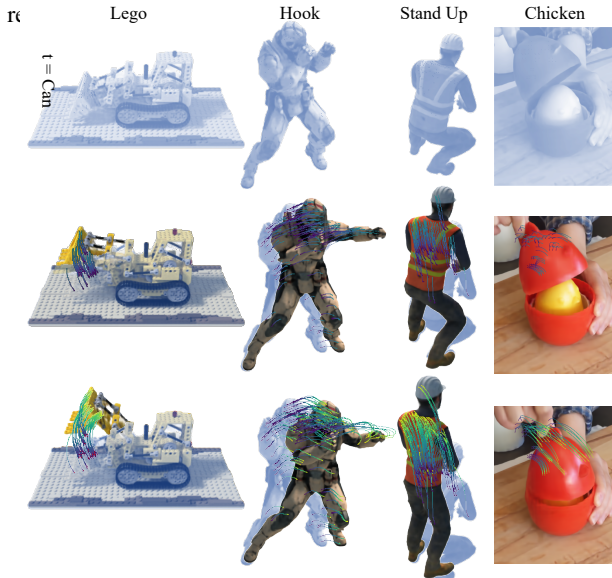


**Figure 9. Trajectory learned by the canonical trajectory field.** Light blue is the canonical frame, the curve represents the historical motion trajectory. More videos can be found in *supp*.

Table 5. **Ablations.** Mean of Hell Warrior, Mutant, Hook and Bouncing Balls in D-NeRF dataset.

| Method | Ours | w/o $\mathcal{L}^{\text{tv}}$ | w/o $\mathcal{L}^{\text{flow}}$ | w/o $\mathcal{L}^{\text{vdiff}}$ | can-time t at mid | w/o $\mathbf{V}_{R_{Up}}$ photo |
|---|---|---|---|---|---|---|
| PSNR | 33.75 | 33.30 | 33.74 | 33.57 | 34.09 | 31.56 |
| SSIM | 0.979 | 0.974 | 0.978 | 0.977 | 0.980 | 0.966 |
| LPIPS | 0.040 | 0.048 | 0.039 | 0.041 | 0.037 | 0.057 |

Table 6. **Ablation of losses on two dataset.**

| Methods | D-NeRF | | HyperNeRF | |
|---|---|---|---|---|
| | PSNR↑ | LPIPS↓ | PSNR↑ | MS-SSIM↑ |
| Ours_notv [1] | 33.09 | 0.046 | 23.84 | 0.831 |
| Ours_noup [2] | 32.03 | 0.043 | 23.72 | 0.822 |
| Ours_noinp [3] | 31.78 | 0.044 | 23.25 | 0.814 |
| Ours | 32.68 | 0.037 | 24.23 | 0.842 |

plies smoothness over the deformation and radiance space, but may harm the high-frequency details (PSNR) for some of the scenes in D-NeRF dataset. For HyperNeRF dataset, the total variation losses improve the results, as the camera settings are more challenging than D-NeRF dataset. However, this improvement is minor compared with the other two factors. Since the resolution is critical for voxel grid based method, up-sampling operation of inpaint network is vital to improve the performance. Finally, the inpaint network plays the most essential role among these three factors in our method. More ablations are provided in our *supp*.

## 5. Conclusion

This paper presents a canonical-based representation with forward warping for novel view synthesis of dynamic scenes. Our method models a static scene at the canonical field and forward warp this whole field to other time steps for dynamic scene rendering. To address the many-to-one and one-to-many mapping difficulties, we present a differentiable forward warping with an average splatting process and an inpaint network. Our proposed forward warping pipeline achieves SOTA performance on the public datasets and our newly built dataset, especially for canonical frames.

**Limitations and Future Directions** Our current implementation is relatively memory-consuming, especially for real scenes. Also, the training speed is relatively slow (one day for each scene). Since we have a smooth trajectory field thanks to the forward warping, additional constraints and motion models could be introduced to learn better trajectories in future works.

is actually doing 'inpainting'. Also, we do observe inpainting and upsampling could refine grids. In Fig. 11 and Fig. 9, the trajectories are reasonable which means we do learn trajectories without inpainting overfitting. We also test w/o photometric terms of $\mathbf{V}_{R_{Up}}$ in Table 5 (w/o $\mathbf{V}_{R_{Up}}$ photo) and the performance drops as there is no direct supervise signal for trajectory learning.

**More Ablations** We study the effects of the total variation losses, the up-sample strategy, and the inpaint network we proposed in Table 6. The total variation losses help to improve visual quality (LPIPS) in D-NeRF dataset as it ap-

Table 7. **Preliminaries in the paper.**

| Variables | Description |
|---|---|
| $\mathbf{V}_{\text{Rf}}^{\text{Can}}$ | voxel grid contains canonical radiance feature |
| $\mathbf{V}_{\text{p}}^{\text{Can}}$ | voxel grid contains coordinate of each voxel |
| $F_{\theta_1}$ | light weight MLP network estimating canonical radiance field |
| $\mathbf{V}_{\text{R}}^{\text{Can}}$ | canonical radiance field estimated by $F_{\theta_1}$ |
| $\mathbf{V}_{\sigma}^{\text{Can}}$ | canonical density voxel grid in $\mathbf{V}_{\text{R}}^{\text{Can}}$ |
| $\mathbf{V}_{\text{cf}}^{\text{Can}}$ | canonical color feature voxel grid in $\mathbf{V}_{\text{R}}^{\text{Can}}$ |
| $\mathbf{V}_{\text{Df}}^{\text{Can}}$ | voxel grid contains canonical deformation feature |
| $F_{\theta_2}$ | light weight MLP network estimating trajectory of each voxel |
| $\mathbf{V}_{T}^{\text{Can}}$ | trajectory voxel grid, containing DCT params of each voxel |
| $\mathbf{V}_{\text{flow}}^{t}$ | deformation flow of each voxel from canonical to time $t$ |
| $F_{\text{warp}}$ | average splatting operation |
| $\mathbf{V}_{\text{R}}^{t}$ | radiance field at time $t$ warped from $\mathbf{V}_{\text{R}}^{\text{Can}}$ |
| $F_{\theta_3}$ | Inpaint Network |
| $\mathbf{V}_{\text{R}_{\text{Inp}}}^{t}$ | inpainted radiance field at time $t$ by $F_{\theta_3}$ |
| $\mathbf{V}_{\text{R}_{\text{Up}}}^{t}$ | upsampled radiance field at time $t$ |
| $F_{\text{render}}$ | volume rendering function |
| $\mathbf{C}_{\text{Inp}}(\mathbf{r})$ | color of ray $r$ rendered from field $\mathbf{V}_{\text{R}_{\text{Inp}}}^{t}$ |
| $\mathbf{C}_{\text{Up}}(\mathbf{r})$ | color of ray $r$ rendered from field $\mathbf{V}_{\text{R}_{\text{Up}}}^{t}$ |

## A. More Implementation Details

### A.1. Preliminaries

We provide preliminaries in Table 7

### A.2. Network Architecture

There are four networks in proposed network: canonical radiance network $F_{\theta_1}$, canonical deformation network $F_{\theta_2}$, view dependent color network $F_{\theta_4}$ and inpaint network $F_{\theta_3}$.

Canonical radiance network $F_{\theta_1}$ is a 3 layer MLP, with width set to 128. The input contains the radiance feature sampled from radiance feature grid $\mathbf{V}_{\text{Rf}}^{\text{Can}}$ with dimension 12, and embedded position with dimension 33. The output contains density whose dimension is 1, and color feature whose dimension is 3 or 12, depending the training stage.

Canonical deformation network $F_{\theta_2}$ is a 4 layer MLP, with width set to 64. Similar with canonical radiance network, the input of $F_{\theta_2}$ contains the deformation feature and embedded position with dimension 33. We set the deformation feature dimension to be 0 and the number of dct bases to be 15 for D-NeRF dataset. For Hypernerf dataset, we set the deformation feature dimension to be 12 and the number of dct bases to be 25.

View dependent color network $F_{\theta_4}$ is a simpler MLP with only 2 layers and width is 64. For input, the dimension of embedded view direction is 27 and color feature is 12. The output is the rgb color with dimension 3.

Inpaint network $F_{\theta_3}$ consist of a UNet structure and a up-sample layer. The UNet structure has 4 encode layers and 3 decode layers. Each encode layer consists of one max pooling layer and two convolution layers. For one convolution layer, there is one instance norm, followed by convolution

and ReLU activation. Each decode layer consists of one up-sample interpolation layer and two convolution layers, which are the same with encode layer. The up-sample layer has the same structure with the decode layer.

### A.3. Average Splatting vs. Softmax Splatting

We use average splatting in this paper, rather then more complex splatting methods in [31], like softmax splatting. We tried to predict weights for softmax splatting, and there is no obvious improvement compared with average splatting. It may introduce too much complexity and we find the UNet could refine the voxel grid to some extend. Average splatting is enough in our setting.

### A.4. Losses and Hyper-parameter Settings

Following DVGO [44], we use $\mathcal{L}^{\text{ptc}}$ to directly supervise the color of sampled points. The intuition is that sampled points with bigger weights contribute more to the rendered color.

$$\mathcal{L}^{\text{ptc}} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \mathcal{L}(\mathbf{C}_{\text{Inp}}(\mathbf{r})) + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \mathcal{L}(\mathbf{C}_{\text{Up}}(\mathbf{r})), \quad (15)$$

$$\mathcal{L}(\mathbf{C}(\mathbf{r})) = \frac{1}{K} \sum_{k=1}^{K} A_{\text{accum}}(k) \|\mathbf{c}(w_k) - \mathbf{C}_{\text{gt}}(\mathbf{r})\|_2^2, \quad (16)$$

$$A_{\text{accum}}(k) = T(w_k)\,\alpha\,(\sigma(w_k)\delta_k). \quad (17)$$

Also, we use the background entropy loss $\mathcal{L}^{\text{bg}}$ to encourage the densities concentrating on either the foreground or the background.

$$\mathcal{L}^{\text{bg}} = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \mathcal{L}(\mathbf{A}_{\text{Inp}}(\mathbf{r})) + \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} \mathcal{L}(\mathbf{A}_{\text{Up}}(\mathbf{r})), \quad (18)$$

$$\begin{aligned} \mathcal{L}(\mathbf{A}(\mathbf{r})) = &- \frac{1}{K-1} \sum_{k=1}^{K-1} A_{\text{accum}}(k) log(A_{\text{accum}}(k)) \\ &+ (1 - A_{\text{accum}}(k))log(1 - A_{\text{accum}}(k)), \end{aligned} \quad (19)$$

As show in Eq. (15) of the paper, the overall loss can be written as

$$\begin{aligned} \mathcal{L} = &\mathcal{L}^{\text{photo}} + w_1 \mathcal{L}^{\text{ptc}} + w_2 \mathcal{L}^{\text{bg}} + w_3 \mathcal{L}^{\text{flow}} + w_4 \mathcal{L}^{\text{vdiff}} \\ &+ w_5 \mathcal{L}^{\text{tv}}(\mathbf{V}_{\sigma}^{\text{Can}}) + w_6 \mathcal{L}^{\text{tv}}(\mathbf{V}_{\text{flow}}^{t}) + w_7 \mathcal{L}^{\text{tv}}(\mathbf{D}), \end{aligned} \quad (20)$$

where $w_1$, $w_2$, $w_3$, $w_4$, $w_5$, $w_6$ and $w_7$ are weights to balance each component in the final coarse loss. In experiments, we set $w_1 = 1e - 1$, $w_2 = 1e - 2$, $w_3 = 1e - 5$, $w_4 = 0.$, $w_5 = 1e - 6$, $w_6 = 1e - 3$ and $w_7 = 1e - 1$ in coarse stage for all datasets, and set $w_1 = 1e - 2$,

$w_2 = 1e-3$, $w_3 = 1e-5$, $w_4 = 1e-5$, $w_5 = 1e-6$, $w_6 = 1e-3$ and $w_7 = 1e-1$ in fine stage for all datasets.

For progressive training in fine stage, we first train with 10 images with closest time steps with canonical step, and progressively add image with the closest time step every 60 iterations.

### A.5. Training Strategy and Settings

We propose a coarse-to-fine training strategy. For the coarse stage, we set the expected voxel number to $67^3$, the scale factor of the inpaint network is 1.5, and the iteration number is 20k. As the purpose of the coarse stage is learning a proxy geometry to calculate the bounding box for the fine stage, we do not use inpaint network during the coarse stage. For the fine stage, the expected voxel number is set to $80^3$ for D-NeRF dataset and $70^3$ for HyperNeRF dataset. The scale factor of the inpaint network is 2.0 and the iteration number is 100k. We use Adam optimizaer and set learning rate 1e-1 for voxel grids and 1e-3 for networks. The training process of a scene takes around 1 day on a GeForce RTX 3090 GPU. Our final model size on average is 260M for D-NeRF dataset and 440M for HyperNeRF dataset. The rendering peed is 7s per image for D-NeRF dataset and 15s per image for HyperNeRF dataset.

### A.6. Lego Complete Dataset

We build a new dataset, named *Lego Complete Dataset*, that animates the object *LEGO* with three different motion patterns. For each scene, the test set is split into three categories to evaluate three abilities: *space interpolation*, *time interpolation*, and *canonical interpolation* abilities. For space interpolation ability, we test four random views for each training time step. Also, we interpolate three time steps between two near training time steps to test the model's generalization ability over unseen times. Finally, to evaluate the learned canonical radiance field, we test 50 random views in canonical space. This results in $200 + 197 + 50 = 447$ test images in the test set.

## B. More Results

### B.1. Results on NHR dataset

We also test our method on NHR dataset [54]. We test all four scenes with 100 frames selecting 90% views for train and 10% views for test. Quantitative results are shown in Table 8, which shows ours achieving best results consistently. We show qualitative comparison in Figure 10, and our method could render clean and detailed images. We also visualize the learnt trajectories in Figure 11. Use Acrobat to view animations.



Figure 10. **NHR dataset qualitative comparison.** We show some synthesized images on NHR dataset.

Figure 11. **Trajectory visualization Use Acrobat to view animations.**

### B.2. Quantitative Results

**Detailed results for D-NeRF Dataset** We show more detailed results in Table 9 for D-NeRF dataset. As show in Table 9, the inpaint network plays a relative important role. Also, the up-sample layer could improves the performance,

Table 8. **Quantitative comparison on NHR dataset.** The red text indicates the best and blue text is the second best result.

| Methods | type | Sport 1 | | | Sport 2 | | | Sport 3 | | | Bacsketball | | | Mean | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| TiNeuVox-S [10] | NPC | 26.06 | 0.93 | 0.10 | 25.98 | 0.93 | 0.11 | 25.90 | 0.93 | 0.11 | 23.75 | 0.91 | 0.14 | 25.42 | 0.92 | 0.12 |
| TiNeuVox-B [10] | NPC | 26.44 | 0.93 | 0.10 | 26.68 | 0.94 | 0.10 | 26.09 | 0.93 | 0.11 | 25.06 | 0.92 | 0.12 | 26.07 | 0.93 | 0.11 |
| NDVG [15] | PC | 23.66 | 0.89 | 0.15 | 24.43 | 0.91 | 0.13 | 22.54 | 0.88 | 0.16 | 22.55 | 0.89 | 0.17 | 23.29 | 0.89 | 0.15 |
| Ours | PC | 27.71 | 0.95 | 0.08 | 27.89 | 0.95 | 0.08 | 27.57 | 0.94 | 0.08 | 24.85 | 0.93 | 0.11 | 27.00 | 0.94 | 0.09 |

Table 9. **Quantitative comparison.** Comparison of our method with others on LPIPS (lower is better) and PSNR/SSIM (higher is better) on eight dynamic scenes of the D-NeRF dataset.

| Methods | type | Hell Warrior | | | Mutant | | | Hook | | | Bouncing Balls | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| TNeRF [36] | N | 23.19 | 0.93 | 0.08 | 30.56 | 0.96 | 0.04 | 27.21 | 0.94 | 0.06 | 32.01 | 0.97 | 0.04 |
| TiNeuVox-S ($100^3$) [10] | NPC | 27.00 | 0.95 | 0.09 | 31.09 | 0.96 | 0.05 | 29.30 | 0.95 | 0.07 | 39.05 | 0.99 | 0.06 |
| TiNeuVox-B ($160^3$) [10] | NPC | 28.17 | 0.97 | 0.07 | 33.61 | 0.98 | 0.03 | 31.45 | 0.97 | 0.05 | 40.73 | 0.99 | 0.04 |
| DNeRF [36] | PC | 25.03 | 0.95 | 0.07 | 31.29 | 0.98 | 0.03 | 29.26 | 0.97 | 0.12 | 38.93 | 0.99 | 0.10 |
| NDVG [15] | PC | 25.53 | 0.95 | 0.07 | 35.53 | 0.99 | 0.01 | 29.80 | 0.97 | 0.04 | 34.58 | 0.97 | 0.11 |
| Ours | PC | 27.71 | 0.97 | 0.05 | 34.97 | 0.98 | 0.03 | 32.29 | 0.98 | 0.04 | 40.02 | 0.99 | 0.04 |
| Ours_notv [1] | PC | 27.96 | 0.97 | 0.05 | 35.26 | 0.98 | 0.03 | 29.57 | 0.96 | 0.06 | 40.40 | 0.99 | 0.05 |
| Ours_noup [2] | PC | 27.60 | 0.96 | 0.06 | 34.15 | 0.98 | 0.04 | 31.51 | 0.97 | 0.04 | 38.89 | 0.99 | 0.05 |
| Ours_noinp [3] | PC | 27.15 | 0.96 | 0.06 | 33.98 | 0.98 | 0.03 | 31.77 | 0.97 | 0.04 | 38.22 | 0.99 | 0.04 |

| Methods | type | Lego | | | T-Rex | | | Stand Up | | | Jumping Jacks | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| TNeRF [36] | N | 23.82 | 0.90 | 0.15 | 30.19 | 0.96 | 0.13 | 31.24 | 0.97 | 0.02 | 32.01 | 0.97 | 0.03 |
| TiNeuVox-S ($100^3$) [10] | NPC | 24.35 | 0.88 | 0.13 | 29.95 | 0.96 | 0.06 | 32.89 | 0.98 | 0.03 | 32.33 | 0.97 | 0.04 |
| TiNeuVox-B ($160^3$) [10] | NPC | 25.02 | 0.92 | 0.07 | 32.70 | 0.98 | 0.03 | 35.43 | 0.99 | 0.02 | 34.23 | 0.98 | 0.03 |
| DNeRF [36] | PC | 21.64 | 0.84 | 0.17 | 31.76 | 0.98 | 0.04 | 32.80 | 0.98 | 0.02 | 32.80 | 0.98 | 0.04 |
| NDVG [15] | PC | 25.23 | 0.93 | 0.05 | 30.15 | 0.97 | 0.05 | 34.05 | 0.98 | 0.02 | 29.45 | 0.96 | 0.08 |
| Ours | PC | 25.27 | 0.94 | 0.05 | 30.71 | 0.96 | 0.04 | 36.91 | 0.99 | 0.02 | 33.55 | 0.98 | 0.03 |
| Ours_notv [1] | PC | 24.33 | 0.89 | 0.11 | 35.02 | 0.99 | 0.02 | 37.01 | 0.99 | 0.02 | 35.14 | 0.98 | 0.03 |
| Ours_noup [2] | PC | 25.20 | 0.93 | 0.06 | 30.24 | 0.97 | 0.05 | 35.47 | 0.98 | 0.02 | 33.14 | 0.98 | 0.04 |
| Ours_noinp [3] | PC | 25.42 | 0.93 | 0.06 | 30.00 | 0.97 | 0.04 | 36.46 | 0.99 | 0.02 | 31.24 | 0.98 | 0.04 |

[1] not use the three total variation losses    [2] not up-sample the voxel grid    [3] not use the inpaint network

which proves this layer learns to recover the details of the 3D voxel grid. This point gives some insight for image super-resolution direction, which is working in 3D dimension may benefit the 2D image tasks. Finally, the total variation losses helps the training to be more stable and get cleaner images.

**Resolutions** We report performance with different resolutions on HyperNeRF dataset in Figure 12. According to Figure 12, the resolution of voxel grid plays an important role when it is relative small and the improvement of the performance decrease when resolution increasing.

**Regularization terms** We study the effect of all regularization terms we proposed, including $\mathcal{L}^{\text{flow}}$, $\mathcal{L}^{\text{vdiff}}$ and $\mathcal{L}^{\text{tv}}$ in Table 10. We could observe that all these three regularization terms has positive effect on the performance, but the improvement is minor. This proves the improvement of our method compared with others come from the forward warp-
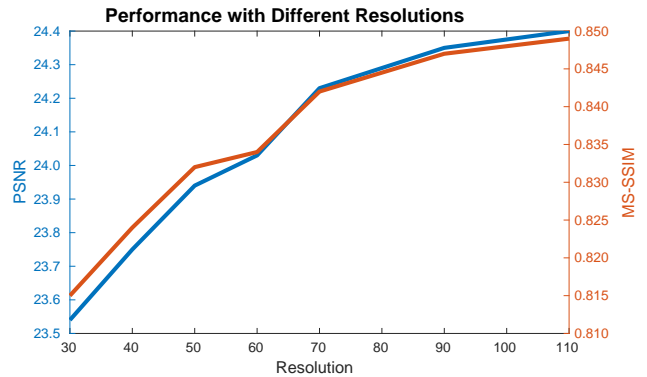


Figure 12. **Performance with different resolutions**

ing desgin we proposed in the paper. Also, backward flow based method NDVG [15] use similar regularization terms with ours, and our method has clear advantage compared with NDVG [15].

Table 10. **More results.** Mean of Hell Warrior, Mutant, Hook and Bouncing Balls in DNeRF dataset.

| Method | ours | w/o $\mathcal{L}^{\text{tv}}$ | w/o $\mathcal{L}^{\text{flow}}$ | w/o $\mathcal{L}^{\text{vdiff}}$ | w/o coarse | can-time t at mid | w/o view dir | w/o $\mathbf{V}_{\mathbf{R}_{\mathbf{U_p}}}$ photo |
|---|---|---|---|---|---|---|---|---|
| PSNR | 33.75 | 33.30 | 33.74 | 33.57 | 31.63 | 34.09 | 22.81 | 31.56 |
| SSIM | 0.979 | 0.974 | 0.978 | 0.977 | 0.967 | 0.980 | 0.925 | 0.966 |
| LPIPS | 0.040 | 0.048 | 0.039 | 0.041 | 0.058 | 0.037 | 0.096 | 0.057 |

**Training strategy** We also set our method without coarse stage training, show in Table 10 (w/o coarse). The performance drops significantly that is reasonable, because the voxel grid covers bigger space without filtering with proxy geometry trained by coarse training.

**Canonical setup** In our method, we set canonical time to be the first frame for D-NeRF dataset to compare with *physical canonical based method* and the middle frame for HyperNeRF dataset for better performance. Setting canonical time to be the middle frame helps improving the performance as the state of the middle frame geometry is closer to other time steps compared with the first frame. To prove this, we set canoical time to be middel frame in Table 10 (can-time t at mid), and the PSNR improves slightly.

**Ray direction modeling** According to Nerfies[32], we need to transfer the current view orientation to the directions in the canonical workspace. But we think using them directly is an acceptable solution for most papers in this area. We leave this as an open question for further study. We test with setting all directions to (0,0,1), the PSNR drops obviously in Table 10 (w/o view dir). This proves the current solution works well to some extend.

**Photometric loss for $\mathbf{V}_{\mathbf{R}_{\mathbf{U_p}}}$** We use photometric terms on $\mathbf{V}_{\mathbf{R}_{\mathbf{U_p}}}$ to make sure the warped grid before inpainting could already render reasonable images. This make sure the UNet is actually doing 'inpainting'. Figure 14 (bottom right) shows an example of how inpainting works. Also, we do observe inpainting and upsampling could refine grids. For videos and Fig. 8 in the main paper, the trajectories are reasonable which means we do learn trajectories without inpainting overfitting. We also test w/o photometric terms of $\mathbf{V}_{\mathbf{R}_{\mathbf{U_p}}}$ in Table 10 (w/o $\mathbf{V}_{\mathbf{R}_{\mathbf{U_p}}}$ photo) and the performance drops as there is no direct supervise signal for trajectory training.

### B.3. Qualitative Results

We show rendered images of our method with different resolutions on HyperNeRF dataset in Figure 13. With bigger resolutions, our method could recover more details, like the pattern of the 3D printer (first row), details of broom (second row) and details of the head in peel-banana (last row). Also, we provide more visual comparison with other methods in Figure 13.
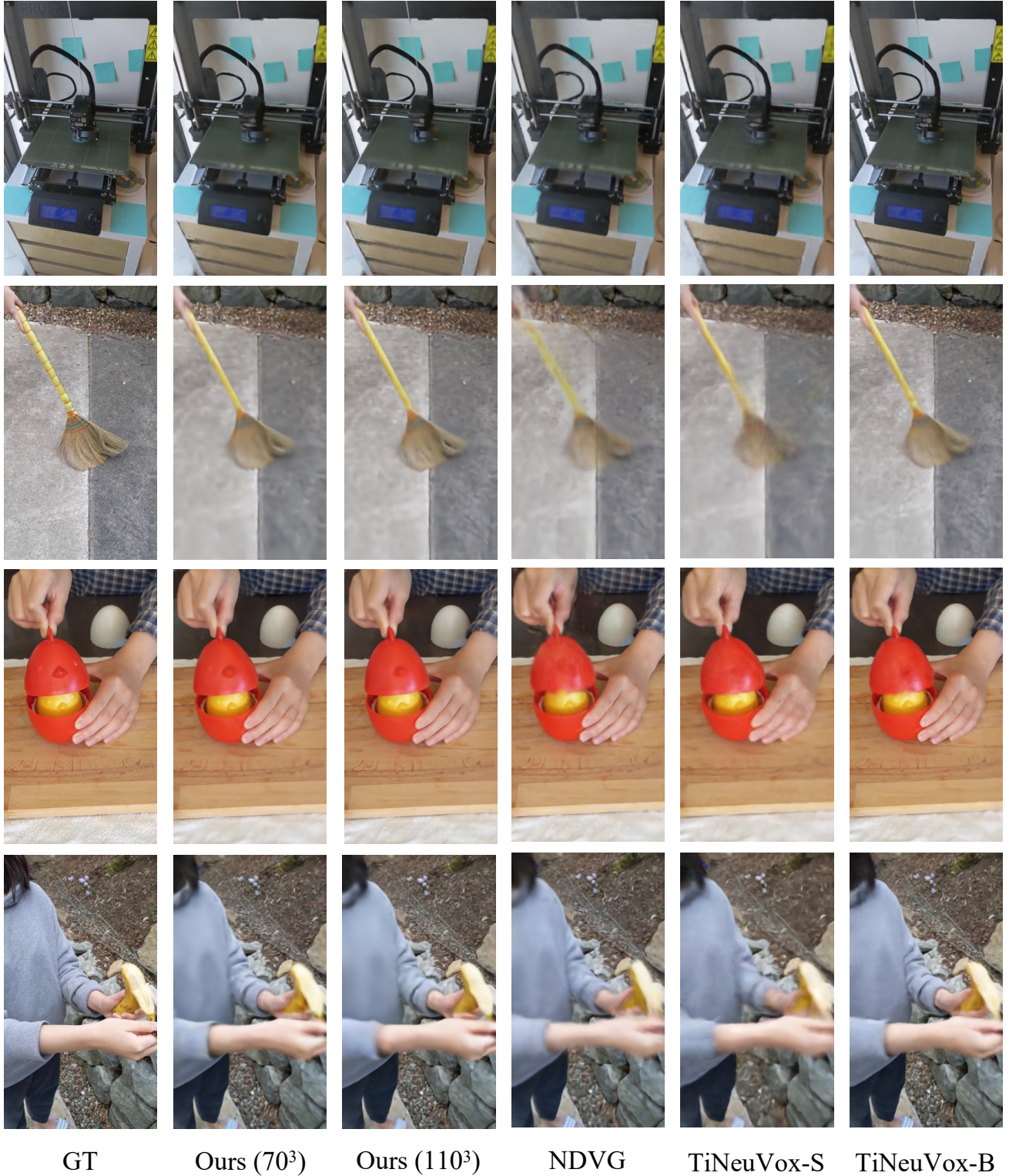
We show the rendered images and depths of our methods and D-NeRF[36] in Figure 14, compared with ground truth. Since we aims to synthesis dynamic scenes from monocular camera, which is a nontrivial problem, the model is highly possible to over-fit the training images. In Figure 14, DNeRF is an example, which produce some clouds in the space which cause artifacts in other views. This is one of the reasons to build lego complete dataset, testing the abilities of the model to interpolate the time and space (including canonical space). Without total variation losses, we could get sharper depth but there may some noise points on the images. Without up-sample layer, the image is blurer (better zoom in for details). Without inpaint network, the rubber band of the lego arms disappears. The rubber band at this time step is stretched and this motion is non-rigid. This non-rigid motion would cause one-to-many issue, compared with rigid motions of mechanical structures of this lego. This proves our inpaint network could handle the one-to-many issue of forward warping.

We show the comparison of canonical image between D-NeRF[36] and proposed method in Figure 15. Figure 15 shows our method could recover correct canonical geometry in the canonical compared with DNeRF[36], which shows the power and potential of the forward warping.

We show more results in our video, including canonical comparison, trajectory visualization and other images render at novel views with different setting. Please refer to the supplement video for more information.

## References

[1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 2

[2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

[3] Mark Boss, Raphael Braun, Varun Jampani, Jonathan T Barron, Ce Liu, and Hendrik Lensch. Nerd: Neural reflectance decomposition from image collections. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

| GT | Ours (70³) | Ours (110³) | NDVG | TiNeuVox-S | TiNeuVox-B |

Figure 13. **HyperNeRF dataset qualitative comparison.** We show some synthesized images on HyperNeRF dataset of our method with different resolutions and other methods.

[4] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, 2001. 2
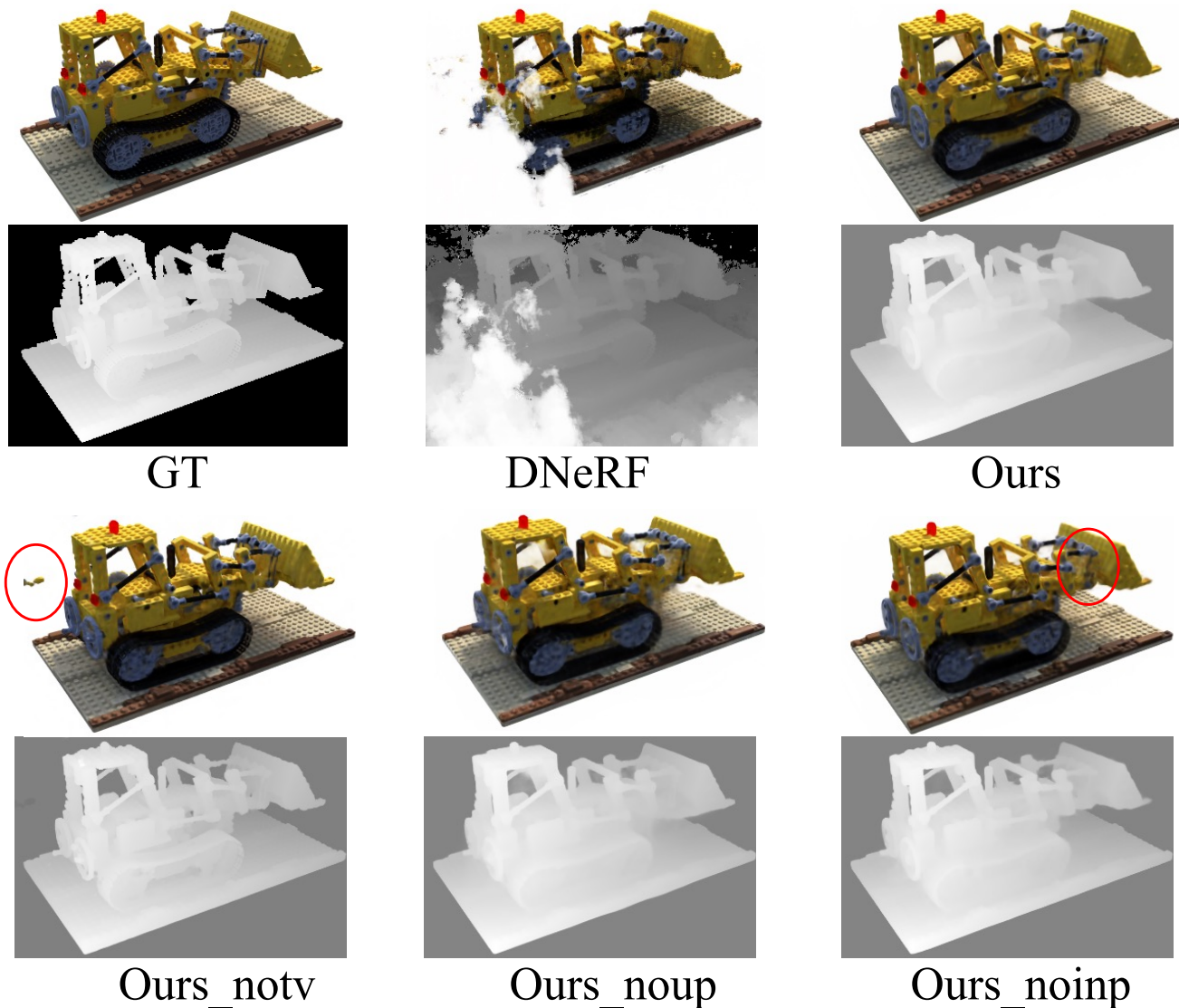
Figure 14. **Lego Complete dataset qualitative comparison.** We show some synthesized images and depth rendered at the same test view selected from one scene of lego complete dataset with different settings.

[5] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. MVSNerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

[6] Shenchang Eric Chen and Lance Williams. View interpolation for image synthesis. In *Proceedings of the 20th annual conference on Computer graphics and interactive techniques*, 1993. 2

[7] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

[8] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H

Kim, and Jan Kautz. Extreme view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019. 2

[9] Yilun Du, Yinan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 3

[10] 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 SIGGRAPH Asia*, 2022. 3, 6, 7, 8, 12

[11] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world's imagery. In *Proceedings of the IEEE Conference*
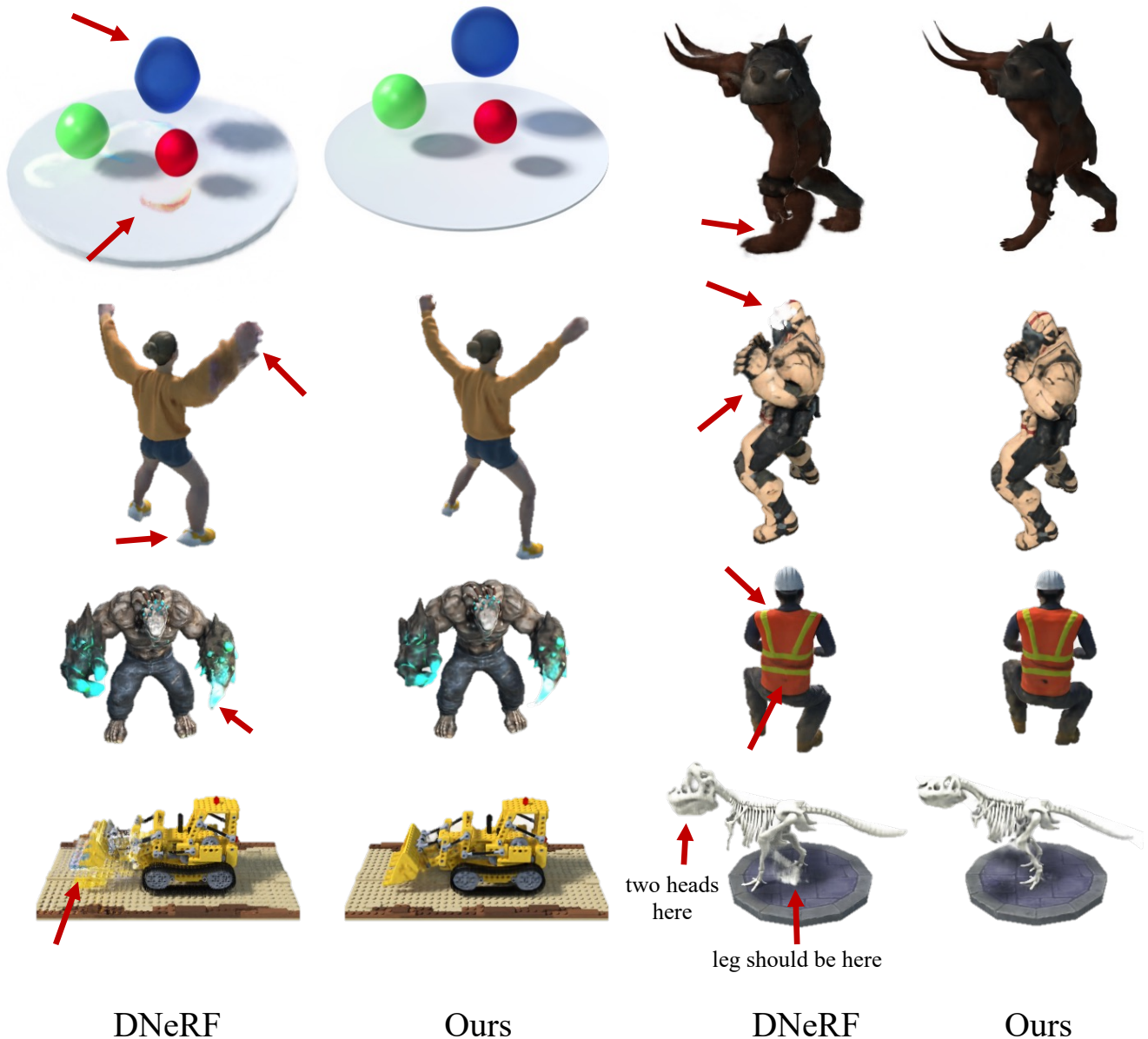
| DNeRF | Ours | DNeRF | Ours |

Figure 15. **Canonical qualitative comparison.** We show canonical comparison of the DNeRF[36] dataset. Note the differences highlighted by the arrows.

*on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2

[12] Wanshui Gan, Hongbin Xu, Yi Huang, Shifeng Chen, and Naoto Yokoya. V4d: Voxel for 4d novel view synthesis. *arXiv preprint arXiv:2205.14332*, 2022. 3

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

[14] Ned Greene. Environment mapping and other applications of world projections. *IEEE Computer Graphics and Applications*, 1986. 2

[15] 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 *Proceedings of the Asian Conference on Computer Vision (ACCV)*, 2022. 1, 3, 6, 7, 8, 12

[16] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Transactions on Graphics (TOG)*, 2018. 2

[17] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Transactions on Graphics (TOG)*, 2016. 2

[18] Marc Levoy and Pat Hanrahan. Light field rendering. In *Pro-*

*ceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. 2

[19] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 3

[20] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[21] 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. *arXiv preprint arXiv:2205.15723*, 2022. 3

[22] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 2

[23] 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 (TOG)*, 38(4):65:1–65:14, July 2019. 6, 7

[24] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)*, 2021. 2

[25] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[26] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 1, 2, 4, 5, 6, 7

[27] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (TOG)*, 2022. 2, 4

[28] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Joerg H. Mueller, Chakravarty R. Alla Chaitanya, Anton S. Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of compact neural radiance fields using depth oracle networks. *Computer Graphics Forum*, 2021. 2

[29] Michael Niemeyer, Jonathan T. Barron, Ben Mildenhall, Mehdi S. M. Sajjadi, Andreas Geiger, and Noha Radwan. Regnerf: Regularizing neural radiance fields for view synthesis from sparse inputs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 5

[30] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 1

[31] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 4, 10

[32] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 3, 6, 7, 13

[33] 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 (TOG)*, 2021. 2, 3, 6, 7

[34] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)*, 2017. 2

[35] Martin Piala and Ronald Clark. Terminerf: Ray termination prediction for efficient neural rendering. In *Proceedings of the International Conference on 3D Vision (3DV)*, 2021. 2

[36] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 2, 3, 6, 7, 8, 9, 12, 13, 16

[37] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[38] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

[39] Konstantinos Rematas, Andrew Liu, Pratul P. Srinivasan, Jonathan T. Barron, Andrea Tagliasacchi, Thomas Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[40] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020. 2

[41] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[42] Harry Shum and Sing Bing Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing*, 2000. 2

[43] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[44] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct Voxel Grid Optimization: Super-fast convergence for radiance fields reconstruction. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 4, 5, 10

[45] Matthew Tancik, Vincent Casser, Xinchen Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[46] Ayush Tewari, Ohad Fried, Justus Thies, Vincent Sitzmann, Stephen Lombardi, Kalyan Sunkavalli, Ricardo Martin-Brualla, Tomas Simon, Jason Saragih, Matthias Nießner, et al. State of the art on neural rendering. In *Computer Graphics Forum*, 2020. 2

[47] Ayush Tewari, O Fried, J Thies, V Sitzmann, S Lombardi, Z Xu, T Simon, M Nießner, E Tretschk, L Liu, et al. Advances in neural rendering. In *ACM SIGGRAPH 2021*, 2021. 2

[48] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 2019. 2

[49] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Nonrigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 1, 3

[50] Alex Trevithick and Bo Yang. Grf: Learning a general radiance field for 3d representation and rendering. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

[51] Chaoyang Wang, Ben Eckart, Simon Lucey, and Orazio Gallo. Neural trajectory fields for dynamic novel view synthesis. *arXiv preprint arXiv:2105.05994*, 2021. 1, 3, 4

[52] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[53] Adrian Wolny, Lorenzo Cerrone, Athul Vijayan, Rachele Tofanelli, Amaya Vilches Barro, Marion Louveaux, Christian Wenzl, Sören Strauss, David Wilson-Sánchez, Rena Lymbouridou, Susanne S Steigleder, Constantin Pape, Alberto Bailoni, Salva Duran-Nebreda, George W Bassel, Jan U Lohmann, Miltos Tsiantis, Fred A Hamprecht, Kay Schneitz, Alexis Maizel, and Anna Kreshuk. Accurate and versatile 3d segmentation of plant tissues at cellular resolution. *eLife*, 9:e57613, 2020. 5

[54] Minye Wu, Yuehao Wang, Qiang Hu, and Jingyi Yu. Multiview neural human rendering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 6, 11

[55] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1, 3

[56] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2

[57] Zexiang Xu, Sai Bi, Kalyan Sunkavalli, Sunil Hadap, Hao Su, and Ravi Ramamoorthi. Deep view synthesis from sparse photometric images. *ACM Transactions on Graphics (TOG)*, 2019. 2

[58] Wenqi Yang, Guanying Chen, Chaofeng Chen, Zhenfang Chen, and Kwan-Yee K Wong. $S^3$-nerf: Neural reflectance field from shading and shadow under a single viewpoint. *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 2

[59] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 1

[60] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2

[61] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 4

[62] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenoctrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021. 2

[63] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2

[64] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. NERF++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 2

[65] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 6

[66] 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 (TOG)*, 2021. 2