# DynamicSurf: Dynamic Neural RGB-D Surface Reconstruction with an Optimizable Feature Grid

Mirgahney Mohamed
Department of Computer Science
University College London

Lourdes Agapito
Department of Computer Science
University College London

Figure 1. **Reconstruction result on two dynamic sequences:** DynamicSurf takes as input a monocular RGB-D sequence of a deforming object and recovers high-fidelity surface reconstructions.

## Abstract

*We propose DynamicSurf, a model-free neural implicit surface reconstruction method for high-fidelity 3D modelling of non-rigid surfaces from monocular RGB-D video. To cope with the lack of multi-view cues in monocular sequences of deforming surfaces, one of the most challenging settings for 3D reconstruction, DynamicSurf exploits depth, surface normals, and RGB losses to improve reconstruction fidelity and optimisation time. DynamicSurf learns a neural deformation field that maps a canonical representation of the surface geometry to the current frame. We depart from current neural non-rigid surface reconstruction models by designing the canonical representation as a learned feature grid which leads to faster and more accurate surface reconstruction than competing approaches that use a single MLP. We demonstrate DynamicSurf on public datasets and show that it can optimize sequences of varying frames with $6\times$ speedup over pure MLP-based approaches while achieving comparable results to the state-of-the-art methods.* [1]

---

[1] Project is available at https://mirgahney.github.io//DynamicSurf.io/.

## 1. Introduction

Reconstructing non-rigid surfaces from sequences acquired from a single static viewpoint arguably poses the most challenging scenario for 3D geometry acquisition. The small amount of effective multi-view supervision signal present in monocular sequences of deforming scenes, results in ambiguities that are usually tackled by making use of scene or deformation priors. While tremendous progress has taken place in the area of model-based approaches to 3D reconstruction of deformable categories such as human bodies, faces or hands (*i.e.* SMPL [42], 3DMM [3, 11, 35, 74], MANO [51]) which fit a pre-trained model to image or depth observations, model-free approaches have received less attention.

The seminal method DynamicFusion [46] was the first to demonstrate model-free real-time reconstruction of 3D surfaces from a live RGB-D sequence. The strategy was to solve for a dense 3D deformation field that maps the current scene geometry to a canonical representation. This idea of decomposing a non-rigid scene into a canonical space, and a per frame deformation field was extremely successful and spurred on follow-up work [26, 54] with different loss terms or hand-crafted deformation priors.

Recent years have seen the development of learning-based approaches to 3D reconstruction. Neural radiance fields [44] encode the density and appearance of scene coordinate points in the weights of a fully connected neural network that is trained in a self-supervised way, purely from images and their associated camera poses. While the focus of neural radiance fields has been their application to novel view synthesis with spectacular success, neural scene representations have also been used for implicit surface reconstruction [19, 62, 68, 69] by expressing volume density as a function of the underlying 3D surface, leading to improved geometry estimation. Some recent approaches have also exploited the combination of RGB losses with depth cues from RGB-D sequences [2, 61] or general-purpose monocular predictors [71] to successfully resolve reconstruction ambiguities in less-observed or texture-less areas.

While the original NeRF formulation encodes the scene using a multi-layer perceptron, taking advantage of the smoothness and coherence priors inherently encoded in its architecture, sparse or dense optimizable feature grids have recently become a powerful alternative to overcome the long training times required by MLP-based approaches[29, 45, 70].

Neural scene representations have also been extended to model dynamic scenes. The most successful formulations have inherited the canonical space and deformation field decomposition [48–50] and model scene geometry and appearance by predicting density and colour for each spatial location. NDR [10], like us, relies on RGB-D inputs and applies depth losses to minimize the discrepancies between the rendered and input RGB and depth images. However, all approaches proposed so far for neural non-rigid reconstruction exploit coordinate-based MLPs instead of optimizable feature grids.

In this paper, we present DynamicSurf a dynamic neural surface reconstruction method from monocular RGB-D input with a static viewpoint. To better model the exact location of the deforming surface, we adopt an SDF-based representation and use a differentiable surface rendering pipeline to render both RGB and depth images to supervise surface reconstruction with photometric, depth and surface smoothness losses. To avoid the slow training times of pure-MLP non-rigid neural RGB-D surface reconstruction [10], DynamicSurf utilizes message passing [31, 60] on dense feature grid [61] which we call *Geometric-feature grid*, and represents the canonical geometry and texture of the deformable objects, using this *Geometric-feature grid*. We combine the canonical grid representation with a topology-aware network [49] to address topological changes. Our efficient representation and architecture allow for $(3\text{-}6)\times$ speed gain to the state-of-the-art methods [10, 39] while maintaining the same level of details and sometimes better. To the best of our knowledge, DynamicSurf is the first method to bring learnable feature grids to dynamic SDF reconstruction from monocular RGB-D sequences.

## 2. Related Work

**Dynamic RGB reconstruction.** Template-based methods [3, 11, 35, 42, 51, 74] are statistical category-specific models learned from high-quality 3D scans. They provide low-dimensional representation that disentangles shape and appearance. Utilizing 3D morphable models [3, 11, 35] some methods [4, 12, 21, 25, 58] learn to reconstruct heads and faces from RGB information. While others [18, 24] in addition to using morphable models they use the recent neural rendering methods [44] to model faces from RGB inputs. Some work [5, 22, 23, 28, 63, 73] recover digital avatars from monocular 2D input with the help of human parametric models [1, 42]. Despite their success, it still remains challenging to extend model-based methods to general objects with limited 3D scans, such as animals, and humans with diverse clothes and articulated objects.

Non-rigid structure from motion (NR-SFM) algorithms [9, 17, 32, 33, 53] reconstruct class-agnostic objects from 2D trajectories but depend heavily on reliable point trajectories throughout observed sequences [52, 57]. Sidhu *et al.* [53] introduce the first dense neural non-rigid structure from motion (NR-SfM) approach, which is trained end-to-end in an unsupervised manner using an auto-decoder as a deformation model while imposing subspace constraints on the latent space. LASR [64] and ViSER [65] recover articulated 3D shapes and dense 3D trajectories from monocular videos using differentiable rendering [41]. BANMo [67]
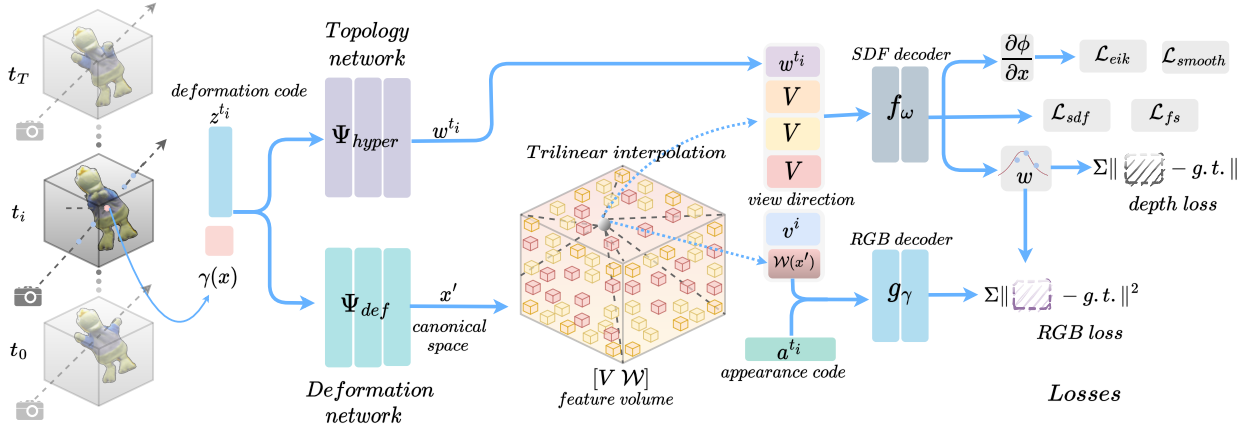
Figure 2. **Overall architecture of DynamicSurf.** Given a monocular sequence of RGB-D frames and segmentation masks, we learn a deformation network to map points to canonical space and a topology network to model topological changes. Points in canonical space are then queried via trilinear interpolation with our feature grid. These features and the hyper-space features (output of the topology network) are concatenated and decoded with shallow MLPs to predict signed distance values (SDF) of the surface and RGB colors.

merges volumetric neural rendering (NeRFs) with invertible deformation through neural blend skinning, and learns dense correspondences via canonical embeddings.

**Dynamic RGB-D reconstruction.** The seminal method of dynamic RGB-D object reconstruction DynamicFusion [46] proposes to model dynamic scenes by estimating a dense 6D motion field that warps the state of the scene into a canonical frame. VolumeDeform [26] reduces the inherent drift during motion tracking by using sparse color features SIFT [43]. KillingFusion [54] is also a templet-free geometry-driven model, estimating a deformation field that aligns signed distance fields (SDFs) representation of the shape. While enforcing a local rigidity constraint and approximating a killing vector field. Guo *et al.* [20] propose to use shading information to leverage appearance information. Then use albedo fusing and geometry to incorporate information for multiple frames. Similar to KillingFusion [54] SobolevFusion [55] is based on a level set variation method where two SDF fields are aligned with a warping field defined using gradient flow in Sobolev space.

With the success of deep learning, DeepDeform [8], based on Siamese networks [34], proposes a network that predicts the probability heat map of point correspondences, which are later used to enhance tracking for complex motion capturing. Bozic *et al.* [7] use convolutional neural networks (CNNs) to predict dense correspondences, which are later used to constrain the optimization using as-rigid-as-possible (ARAP) [56] priors. OcclusionFusion [39] proposes an LSTM-involved graph neural network (GNN) to infer the motion of occluded regions by exploiting visibility and temporal information.

**Dynamic neural radiance fields.** Following the success of neural radiance fields NeRF [44] for novel view synthesis, recent works have extended it to handle dynamic scenes by learning a warping field to deform observations into a shared canonical space [48–50, 59], or by modelling a spatiotemporal 4D radiance field [16, 36, 37]. Nerfies [48] augment NeRF with a **SE**(3) deformation field that warps observation to a canonical space. They propose a coarse-to-fine optimization strategy and inspired by physical simulation they also propose an elastic regularization that further improves optimization robustness. Exploiting a level-set methodology, HyperNeRF [49] extends Nerfies [48] by introducing an ambient dimension to express topological changes. Instead of deforming each sampled point, Tretschk *et al.* [59] model the deformation as ray blending, and a rigid network to model the static background. Both ray blending and rigid networks are trained without explicit supervision. TiNeuVox [15] uses a feature grid to encode density and colour, and augments their representation with temporal information encoding additional time features with an MLP network. Unlike our approach, all the methods above take RGB sequences as input and focus on novel view synthesis, while in DynamicSurf we take RGB-D input and focus on accurate surface reconstruction.

To the best of our knowledge, NDR [10] is the only existing method to tackle the problem of surface reconstruction from a monocular RGB-D camera. However, they propose a fully MLP-based representation, which provides smooth reconstructions at the expense of slow training times. Instead, we propose the use of a feature volume, optimized in a coarse-to-fine fashion, which results in an increase of $6\times$ speedup.

## 3. Methodology

Given a sequence of RGB-D images $\{(I^t, D^t), t = 1, ..., T\}$, captured by a static monocular RGB-D camera and segmentation masks $\mathcal{M}^t$ obtained with an off-the-shelf video segmentation approach [13, 38], we formulate dynamic surface reconstruction as an optimisation problem. Our model learns to map the input sequence into a canonical hyper-space $\mathbb{R}^{3+m}$ composed of 3D canonical and topology networks similar to [10, 49], where $m$ represents the dimensionality of the topology space. Scene geometry and colour are represented using feature grids, which are decoded into SDF and RGB values with two shallow MLP decoders shared across grids [61]. Unlike other approaches, we do not utilize any structural priors such as 2D annotations [6], estimated normal maps [27] or optical flow [66]. Fig. 2 shows an overview of our architecture.

### 3.1. Deformation Field

For a 3D point sampled from $t$-th frame, we learn a continuous deformation field to map the points to a canonical space with our deformation network $\Psi_{def}$. We formulate $\Psi_{def}$ as a **SE**(3) field conditioned on a learnable deformation code $z^t$ [48]. More formally, given an input point $x_i^t \in \mathbb{R}^3$ – it is important to note $x_i^t$ can represent both surface and free-space points. We encode the fields with a rotation around an anchor point $a$ using log-quaternion $q$ and a displacement $d$ [48]. We use an MLP architecture to model the deformation field $\Psi_{def} : (x, z^t) \rightarrow (r, a, d)$. The deformation network predicts rotation vector $r$, anchor point $a$, and a displacement vector $d$. We use log-quaternion $(0, r)$ and represent valid rotations using the exponent which is guaranteed to be a unit quaternion:

$$q = exp\left(\frac{cos\|r\|}{\frac{r}{\|r\|}sin\|r\|}\right) \tag{1}$$

Finally, the deformation is given by the following equation:

$$x_i'^t = q(x_i^t - a)q^{-1} + a + d \tag{2}$$

Dynamic scenes may undergo topological changes which can make it hard for the deformation network alone to learn. Therefore, we augment our deformation field with a topology-aware network [49]. Our topology network takes as input 3D point $x_i^t$ and the learnable deformation code $z^t$ and learns to map it into $w_i^t$ in the hyper-space $\mathbb{R}^m$. The corresponding coordinate of $x_i^t$ in the canonical hyper-space is:

$$p = [x_i'^t, w_i^t] = [\Psi_{def}(x_i^t, z^t), \Psi_{hyper}(x_i^t, z^t)] \in \mathbb{R}^{3+m} \tag{3}$$

where $m$ is the dimension of the topological hyper-space.

### 3.2. Canonical Grid Representation

While most feature-grid-based approaches choose multi-resolution feature volumes [14, 61], we use a single fea-

ture grid to encode the canonical scene geometry and shared shallow MLP decoders, but we employ a coarse-to-fine strategy on the grid resolution. The canonical scene geometry a one-level feature grid $\mathcal{V}_\theta = \{V^l\}$, where $l \in \{0, 1, 2\}$ represents the current level used within the coarse-to-fine strategy.

Given a deformed point $p$ in canonical space $\mathbb{R}^{3+m}$, we obtain its geometric feature by trilinearly interpolating its 3D location $x_i'^t$ features in our feature grid and concatenating its hyper-space coordinate $w_i^t$ to make what we denote *hyper-features*. The hyper-features are then decoded into an SDF value $\varphi(\mathbf{x})$ via the geometry MLP $f_\omega(\cdot)$:

$$\varphi(\mathbf{x}_i^t) = f_\omega([V^l(\mathbf{x}_i'^t), w_i^t]) \tag{4}$$

To decode colour information, following [10] we first map the viewing direction into the canonical space using the Jacobian matrix $J_x(x_i^t) = \frac{\partial x_i'^t}{\partial x_i^t}$ of the 3D canonical point $x_i'^t$ w.r.t observed point $x_i^t$. Then we encode the color using a separate feature grid $\mathcal{W}_\beta$ and decoder $g_\gamma(\cdot)$ conditioned on the canonical viewing direction, surface normals $n_i^t$ and a global appearance latent code $a^t$ as in [48]:

$$\mathbf{c}_i^t = g_\gamma(\mathcal{W}_\beta(\mathbf{x}_i'^t), \mathbf{d}_c, n_i^t, a^t) \tag{5}$$

where $\mathbf{d}_c = J_x(x_i^t)\mathbf{d}$ is the viewing direction in canonical space, $n_i^t = \nabla\varphi(x_i^t)$ is the surface normal at point $x_i^t$.

**Coarse-to-Fine strategy on the grid resolution.** Unlike MLP-based methods which enjoy natural smoothness, grid-based methods can suffer from noise and surface artefacts and are prone to fall into suboptimal local minima [40]. To circumvent these problems and to better enhance feature grid smoothness, we use a coarse-to-fine optimization approach [40], progressively increasing the feature grid resolution from $35^3$ to $140^3$. The lowest resolution models the overall surface, and during training, we increase the grid resolution via interpolation to capture finer surface details

$$V^l = interp(., V^{l-1}) \tag{6}$$

where . represents the 3D coordinates of points at the previous resolution level and $V^{l-1}$ and $V^l$ denote the features at the previous and current resolution respectively.

### 3.3. Depth and Color Rendering

Inspired by the recent work on learnable feature grid [61] and volume rendering [44], we adopt rendering equations to account for dynamic scenes. For frame $t$, we sample points $\{\mathbf{x}_i^t | \mathbf{x}_i^t = \mathbf{o} + d_i\mathbf{r}, i = 1, \dots, N\}$ along ray $r$ parameterised by camera centre $\mathbf{o}$ and ray direction $d_i$. We use unbiased and occlusion-aware weights $w_i^t = T_i^t \alpha_i^t$ [62], and $T_i^t = \prod_{j=1}^{i-1}(1 - \alpha_j^t)$ represents the *accumulated transmittance* at
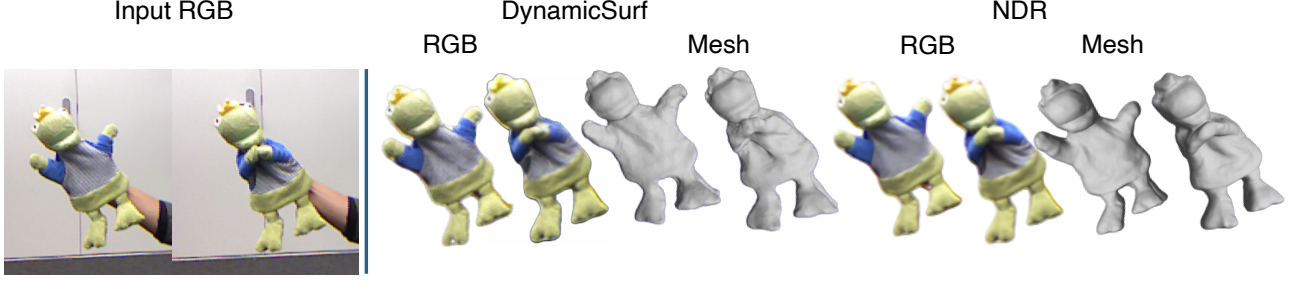
Figure 3. Qualitative comparison with NDR [10]. The result shows our method is on par with NDR [10] while enjoying $6\times$ speed gain.

point $\mathbf{x}_i^t$, and $\alpha_i^t$ is the *opacity value* defined by:

$$\alpha_i^t = \max\left(\frac{\mathcal{T}_\lambda(\varphi(\mathbf{x}_i^t)) - \mathcal{T}_\lambda(\varphi(\mathbf{x}_{i+1}^t))}{\mathcal{T}_\lambda(\varphi(\mathbf{x}_i^t))}, 0\right) \quad (7)$$

where $\mathcal{T}_\lambda(x) = (1 + e^{-\lambda x})^{-1}$ is the Sigmoid function modulated by a learnable parameter $\lambda$ which controls surface transition. The expected values of predicted colors $\mathbf{c}_r^t$ and sampled depth $d_r^t$ become:

$$\hat{\mathbf{c}}_r^t = \sum_{i=1}^N w_i^t \hat{\mathbf{c}}_{ri}^t, \quad \hat{d}_r^t = \sum_{i=1}^N w_i^t \hat{d}_{ri}^t \quad (8)$$

The RGB and depth per-ray rendering losses are:

$$\ell_{rgb}^r = \|\mathbf{c}_r^t - \hat{\mathbf{c}}_r^t\|, \quad \ell_d^r = |d_r^t - \hat{d}_r^t| \quad (9)$$

where $\mathbf{c}_r^t$ and $d_r^t$ represent the ray $r$ RGB value in the image and the corresponding depth map value respectively.

### 3.4. SDF Supervision and Regularization

**SDF supervision.** Similar to [2, 47, 61], we approximate the ground truth SDF value based on distance to the observed depth along ray direction $d_r^t$. We define the bound $b_r(\mathbf{x}_i^t) = d_r^t - d(x_i^t)$ and divide the set of sampled points into two disjoint sets: near-surface points $S_{tr}^r = \{x_i^t | b_r(x_i^t) <= \epsilon\}$ and free-space points (far from the surface) $S_{fs}^r = \{x_i^t | b_r(x_i^t) > \epsilon\}$. The truncation threshold $\epsilon$ is a hyper-parameter. For the set of near-surface points along the ray $S_{tr}^r$ we apply the following SDF loss:

$$\mathcal{L}_{sdf}^r = \frac{1}{|S_{tr}^r|} \sum_{s \in S_{tr}} |\varphi(\mathbf{x}_s^t) - b_r(\mathbf{x}_s^t)| \quad (10)$$

For the set of points far from the surface $S_{fs}$ we apply the free space loss as in [47, 61] to encourage free space prediction and provide more direct supervision than the rendering terms in 9.

$$\mathcal{L}_{fs}^r = \frac{1}{|S_{fs}^r|} \sum_{s \in S_{fs}} \max\left(0, e^{-\alpha\varphi(\mathbf{x}_s^t)} - 1, \varphi(\mathbf{x}_s^t) - b_r(\mathbf{x}_s^t)\right) \quad (11)$$

An exponential penalty is applied for negative SDF values, a linear penalty for positive SDF values larger than the bound and no penalty is applied if it is smaller.

**SDF regularization.** To encourage valid SDF values, especially in areas without direct supervision we employ the Eikonal regularisation term $\ell_{eik}$, which encourages a uniform increase in the absolute value of SDF as we move far from the surface [19, 47, 62]. More formally for any query point $x_i'^t$ in the canonical space $\mathbb{R}^3$, we encourage the gradient of the SDF value w.r.t. the 3D observed point to have unit length:

$$\mathcal{L}_{eik}^r = \frac{1}{|S_{fs}^r|} \sum_{s \in S_{fs}} \left(1 - \|\nabla\varphi(\mathbf{x}_s'^t)\|\right)^2 \quad (12)$$

**Surface smoothness regularization.** To further enhance surface smoothness we enforce nearby points to have similar normals. Unlike [61], where they sample uniformly inside the grid, we sample surface points only $x_s^t \in S_{surf}$ which reduces the computation cost drastically.

$$\mathcal{L}_{smooth} = \frac{1}{R} \sum_{s \in S_{surf}} \|\nabla\varphi(x_s^t) - \nabla\varphi(x_s^t + \delta)\|^2 \quad (13)$$

where $x_s^t$ is back-projected using depth maps, $\delta$ is a small perturbation sampled from a uniform Gaussian distribution with standard deviation $\delta_{std}$, and $R$ is the total number of sampled rays.

### 3.5. Optimization

From a randomly sampled frame $t$ we sample a batch of rays $R$ from all pixels across the image and we minimize the following loss function $\mathcal{L}$:

$$\mathcal{L} = \lambda_{rgb}\mathcal{L}_{rgb} + \lambda_d\mathcal{L}_d + \lambda_{sdf}\mathcal{L}_{sdf} + \lambda_{fs}\mathcal{L}_{fs} + \\ \lambda_{eik}\mathcal{L}_{eik} + \lambda_{smooth}\mathcal{L}_{smooth} + \lambda_{mask}\mathcal{L}_{mask} \quad (14)$$

The RGB rendering loss $\mathcal{L}_{rgb}$ measures the difference between the ground truth ray colour and predicted colour over
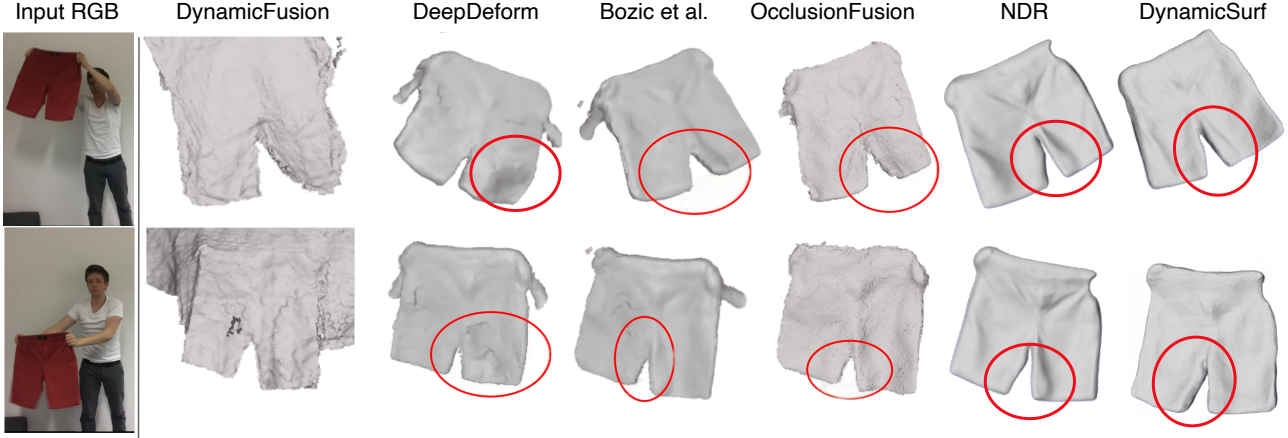
5

Figure 4. We compare aganist NDR [10], OcclusionFusion [39], Bozic *et al.* [7], DeepDeform [8] and DynamicFusion [46] on the *shorts* sequence from DeepDeform [8] dataset. As the code is not available the results of DeepDeform [8] and Bozic [7] *et al.* are taken from the video of Bozic *et al.* [7] The qualitative comparison shows our method achieves comparable results to NDR [10] and better results than all the other baselines.

all sampled rays, while the depth rendering loss $\mathcal{L}_d$ measures the difference between ground-truth depth and predicted depth values over rays with valid depth $R_d$.

$$\mathcal{L}_{rgb} = \frac{1}{|R_{rgb}|} \sum_{r \in R_{rgb}} \mathcal{M}_r^t \ell_{rgb}^r, \quad \mathcal{L}_d = \frac{1}{|R_d|} \sum_{r \in R_d} \mathcal{M}_r^t \ell_d^r \tag{15}$$

$\mathcal{M}_r$ is the object mask, which we also utilize as a mask loss to better focus on the object of interest.

$$\mathcal{L}_{mask} = H(\mathcal{M}_r^t, \hat{\mathcal{M}}_r^t) \tag{16}$$

where $\hat{\mathcal{M}}_r^t = \sum T_i^t \sigma(x_i^t)$ is the density accumulated along ray $r$ and $H$ is the binary cross entropy function. The SDF loss $\mathcal{L}_{sdf}$ is applied on sampled points inside the truncation region $S_{tr}$.

$$\mathcal{L}_{sdf} = \frac{1}{|R_d|} \sum_{r=1}^{R} \mathcal{L}_{sdf}^r \tag{17}$$

The free-space and Eikonal losses $\mathcal{L}_{fs}$, $\mathcal{L}_{eik}$ are applied to the rest of the points $S_{fs}$:

$$\mathcal{L}_{fs} = \frac{1}{|R_d|} \sum_{r=1}^{R} \mathcal{L}_{fs}^r \tag{18}$$

$$\mathcal{L}_{eik} = \frac{1}{|R_d|} \sum_{r=1}^{R} \mathcal{L}_{eik}^r \tag{19}$$

**Geometric Initialisation.** We found it very important to initialize our hyper-feature grid and geometry decoder to predict a sphere [19, 61] centred at volume origin and with a radius proportional to the scale of the object of interest.

## 4. Experimental Evaluation

**Implementation details.** We train DynamicSurf using Adam optimizer [30] with a learning rate $5e^{-4}$. And step based weighting for the RGB, depth and other regularization loss, where we give more weight to the RGB and other regularization losses at the beginning of the training and we reduce that gradually toward the end. First, we randomly sample an image, then we sample 1024 rays per batch with 128 points along each ray. Similar to [10, 62] after we sample uniformly 64 points, we iteratively use importance sampling 4 times for 16 points each for a total of 128 points. At the beginning of the training, we start with $35^3$ feature grid resolution and 32-dim features. We double the grid size twice during the training using trilinear interpolation as discussed in Sec. 3.2. Out of the 32 feature dimensions we use 6 dimensions for colour decoder and the remaining for the SDF decoder with an ambient dimension of 2 for KIllingFusion [54] sequences and 8 for DeepDeform [8] sequences. We use a two-layered MLP network for both colour and SDF with 64 neurons each. And for the topology-aware network with follow the same architecture as in HyperNeRF [49]. We also utilize a coarse-to-fine positional encoding strategy, as in Nerfies [48]. Following [10] we leaverage off-the-shelf segmentation methods for humans [38] and objects [13] and we apply Robust ICP [72] for per-frame pose initialization. All experiments were trained on a single NVIDIA Tesla V100 GPU for a total of 80K for DeepDeform [8] sequences and 60K for KIllingFusion [54] sequences.

**Baseline methods.** We compare with DynamicFusion [46], two recent learning-based methods [7, 8] both
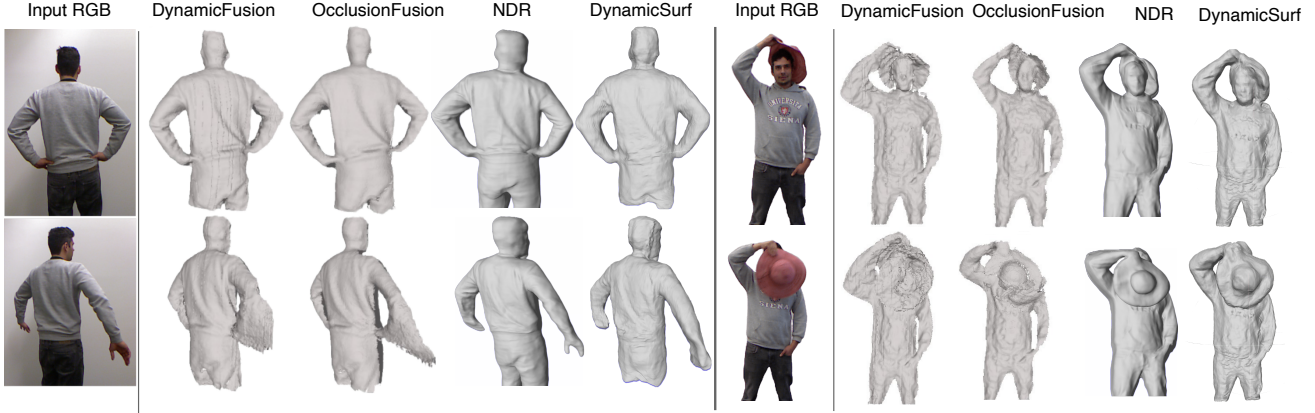
Figure 5. Qualitative comparison with NDR [10], OcclusionFusion [39], and DynamicFusion [46] on two sequences from the KillingFusion [54] dataset. The Qualitative comparison shows our method achieves comparable results to the baselines while being $6.7\times$ faster.
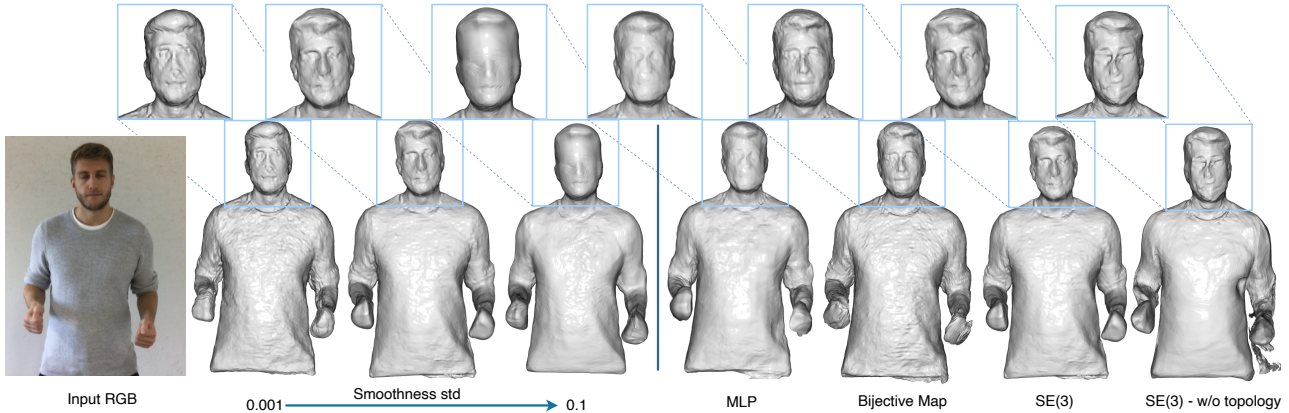


Figure 6. Ablation on smoothness loss and model architecture. Left: shows the effect of standard deviation $\delta_{std}$ on the surface smoothness. Right: shows the importance of deformation and topology networks in recovering surface details and capturing motion.

are fusion-based methods and utilize neural networks to learn correspondences. and two more state-of-the-art neural reconstruction methods: OcclusionFusion [39] which estimates the motion of occluded regions from visible ones using a temporal graph neural network (GNN). and NDR [10], which models motion and surface reconstruction using neural implicit functions.

**Datasets.** We evaluate our method and baselines on RGB-D videos from the DeepDeform [8] and KillingFusion [54] datasets. The DeepDeform [8] dataset is captured with a Structure Sensor mounted on an iPad. The RGB-D stream is captured and aligned at $640 \times 480$ and 30 frames per second Although the data also provides sparse annotated correspondences and scene flow we did not utilize any of them. The KillingFusion [54] dataset is collected using Kinect v1 and aligned at $640 \times 480$ resolution. We use 6 scenes from DeepDeform [8], and 7 scenes from the KillingFusion [54]

dataset, containing human motion, animals, objects, clothes and toys.

**Evaluation metrics.** Following previous work [7, 8, 10, 39] for quantitative evaluation we computed geometric errors by comparing the reconstructed geometry to the depth values inside the object mask.

## 4.1. Reconstruction Quality

**Quantitative evaluation.** We evaluate DynamicSurf following previous work [7, 8, 10, 39] using the geometric reconstruction error as a metric and the evaluation protocol described in [8]. Tab. 1 shows a comparison with NDR [10], which we trained using their code release and settings, on the same datasets as DynamicSurf. For completeness, we report the results we obtained using their code (marked with * in Tab. 1), and the results from their paper. As shown in Tab. 1, DynamicsSurf achieves reconstruction

errors on par with NDR [10] while being $5 - 6.7\times$ faster as shown in Tab. 2.

| Metric | Method | Dataset | | | |
|--------|--------|---------|---|---|---|
| | | DeepDeform | KillingFusion | | |
| | | Human | Alex | Hat | Frog |
| Mean | OcclusionFusion | - | 3.75 | 6.77 | 1.61 |
| | NDR | - | 4.24 | 4.93 | 1.34 |
| | NDR* | 3.89 | 5.69 | 3.24 | 2.26 |
| | Ours | 2.96 | 6.12 | 2.83 | 1.39 |
| Median | OcclusionFusion | - | 3.41 | 6.47 | 1.59 |
| | NDR | - | 4.11 | 4.66 | 1.33 |
| | NDR* | 3.8 | 4.69 | 3.26 | 2.25 |
| | Ours | 2.89 | 5.21 | 2.75 | 1.27 |

Table 1. Quantitative results on 5 sequences. The geometry error ($\downarrow$) represents the average per-point per-frame 3D error between the reconstructed shape and the point cloud obtained from back-projecting the GT depth values inside the mask. NDR denotes the results from their paper, while NDR* denotes the result obtained by training on NDR's [10] public code release and settings.

**Qualitative results.** We present qualitative results of reconstructed sequences in 3, 4 and 5 which show that DynamicSurf can reconstruct smoother and more complete surfaces than fusion-based methods [39, 46] and learning based methods [7, 8] and achieves comparable reconstruction quality to the SOTA method NDR [10] while being significantly faster Tab. 2.

## 4.2. Ablation studies

To validate our architecture choices and losses, we evaluate three key components of our method based on their contribution to surface reconstruction quality. Please refer to supplementary materials for more extended evaluations.

### 4.2.1 Deformation networks

We evaluate the reconstruction result with different deformation network architectures. We replace it with the Bijective map proposed by NDR [10], and with a pure MLP-based network to model the deformation. To keep the comparison fair, we make sure that all networks have the same number of layers and neurons. Fig. 6 shows that using our **SE**(3) deformation network yields better results than the pure MLP-based method, especially in capturing face details. As for the bijective map [10], while Fig. 6 shows comparable reconstruction quality, the bijective map deformation model suffers from surface artefacts, furthermore, it requires high computation cost. Therefore, we opted for the **SE**(3) due to its lower computational complexity and for providing smoother surfaces.

| Method | Dataset | | | | |
|--------|---------|---|---|---|---|
| | KillingFusion | | | DeepDeform | |
| | Alex | Hat | Frog | Human | Dog |
| NDR | 1 (36) | 1 (37.8) | 1 (36) | 1 (41) | 1 (41.5) |
| Ours | 5.8× | 6.7× | 6.2× | 5.1× | 5× |

Table 2. Optimization time results on 5 sequences. DynamicSurf archives $(5\text{-}7)\times$ speed to NDR. The numbers between parentheses show time in hours.

### 4.2.2 Topology network

We train DynamicSurf without the topology network, such that the decoders only take input features from the feature grid. For a fair comparison, we increase the feature volume dimensions to match those of the topology network. As shown in Fig. 6 the network with the topology can better capture the details of the surface, especially with difficult deformations.

### 4.2.3 Smoothness loss

Fig. 6 shows the result of varying the standard deviation of the smoothness loss $\delta_{std}$. As expected, as the value increases, the fine details on the surface start to disappear. While having a smaller standard deviation results in noise surface see the chest and face of the human in Fig. 6 left.

## 5. Limitations

DynamicSurf is sequence specific, since it is an optimization approach, and needs to be trained per sequence, which can be impractical. Further, it does not utilize the shared information between identities across sequences and similar motions. Also similar to the baselines, DynamicsSurf requires per-frame pose initialization for sequences with global rotation. Learning the global rotations would bring the method a step closer to online processing. We will investigate how to overcome these limitations in the future. In terms of broader impact, training these models still takes substantial computational and energy resources. Pushing further in the direction of designing light-weight models is an important research direction.

## 6. Conclusions

We have proposed DynamicSurf a template-free method for high-fidelity surface and motion reconstruction of deformable scenes from monocular RGB-D sequences. DynamicSurf maps a canonical representation of the surface geometry to the current frame using a neural deformation field. We represent the canonical space as a learned feature grid, optimized in a coarse-to-fine fashion, and we employ a topology-aware network to handle topology variations.

DynamicSurf achieves comparable reconstruction performance to the state-of-the-art – purely MLP-based – non-rigid surface reconstruction method [10] on different object categories from various public datasets, while being an order of magnitude faster.

## Acknowledgements

## References

[1] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*, pages 408–416. 2005. 2

[2] Dejan Azinović, Ricardo Martin-Brualla, Dan B Goldman, Matthias Nießner, and Justus Thies. Neural rgb-d surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 5

[3] Volker Blanz and Thomas Vetter. A morphable model for the synthesis of 3d faces. *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999. 2

[4] Volker Blanz, Curzio Basso, Tomaso Poggio, and Thomas Vetter. Reanimating faces in images and video. In *Computer graphics forum*, pages 641–650. Wiley Online Library, 2003. 2

[5] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part V 14*, pages 561–578. Springer, 2016. 2

[6] Aljaz Bozic, Michael Zollhöfer, Christian Theobalt, and Matthias Nießner. Deepdeform: Learning non-rigid rgb-d reconstruction with semi-supervised data. *CoRR*, abs/1912.04302, 2019. 4

[7] Aljaz Bozic, Pablo Palafox, Michael Zollhöfer, Angela Dai, Justus Thies, and Matthias Nießner. Neural non-rigid tracking. *Advances in Neural Information Processing Systems*, 33:18727–18737, 2020. 3, 6, 7, 8

[8] Aljaz Bozic, Michael Zollhofer, Christian Theobalt, and Matthias Nießner. Deepdeform: Learning non-rigid rgb-d reconstruction with semi-supervised data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7002–7012, 2020. 3, 6, 7, 8

[9] C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3d shape from image streams. In *Proceedings IEEE Conference on Computer Vision and Pattern Recognition. CVPR 2000 (Cat. No.PR00662)*, pages 690–696 vol.2, 2000. 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 *Thirty-sixth Conference on Neural Information Processing Systems (NeurIPS)*, 2022. 2, 3, 4, 5, 6, 7, 8, 9

[11] Chen Cao, Yanlin Weng, Shun Zhou, Yiying Tong, and Kun Zhou. Facewarehouse: A 3d facial expression database for visual computing. *IEEE Transactions on Visualization and Computer Graphics*, 20(3):413–425, 2013. 2

[12] Chen Cao, Qiming Hou, and Kun Zhou. Displaced dynamic expression regression for real-time facial tracking and animation. *ACM Transactions on graphics (TOG)*, 33(4):1–10, 2014. 2

[13] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Modular interactive video object segmentation: Interaction-to-mask, propagation and difference-aware fusion. In *CVPR*, 2021. 4, 6

[14] 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. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 4

[15] 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. In *SIGGRAPH Asia 2022 Conference Papers*, New York, NY, USA, 2022. Association for Computing Machinery. 3

[16] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5712–5721, 2021. 3

[17] P. F. U. Gotardo and A. M. Martinez. Non-rigid structure from motion with complementary rank-3 spaces. In *Proceedings of the 2011 IEEE Conference on Computer Vision and Pattern Recognition*, page 3065–3072, USA, 2011. IEEE Computer Society. 2

[18] Philip-William Grassal, Malte Prinzler, Titus Leistner, Carsten Rother, Matthias Nießner, and Justus Thies. Neural head avatars from monocular rgb videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18653–18664, 2022. 2

[19] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *International Conference on Machine Learning*, pages 3789–3799. PMLR, 2020. 2, 5, 6

[20] 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 (ToG)*, 36(4):1, 2017. 3

[21] Yudong Guo, Jianfei Cai, Boyi Jiang, Jianmin Zheng, et al. Cnn-based real-time dense face reconstruction with inverse-rendered photo-realistic face images. *IEEE transactions on pattern analysis and machine intelligence*, 41(6):1294–1307, 2018. 2

[22] Marc Habermann, Weipeng Xu, Michael Zollhoefer, Gerard Pons-Moll, and Christian Theobalt. Livecap: Real-time human performance capture from monocular video. *ACM Transactions On Graphics (TOG)*, 38(2):1–17, 2019. 2

[23] Marc Habermann, Weipeng Xu, Michael Zollhofer, Gerard Pons-Moll, and Christian Theobalt. Deepcap: Monocular

human performance capture using weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5052–5063, 2020. 2

[24] Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. Headnerf: A real-time nerf-based parametric head model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20374–20384, 2022. 2

[25] Alexandru Eugen Ichim, Sofien Bouaziz, and Mark Pauly. Dynamic 3d avatar creation from hand-held video input. *ACM Transactions on Graphics (ToG)*, 34(4):1–14, 2015. 2

[26] Matthias Innmann, Michael Zollhöfer, Matthias Nießner, Christian Theobalt, and Marc Stamminger. Volumedeform: Real-time volumetric non-rigid reconstruction. In *European conference on computer vision*, pages 362–379. Springer, 2016. 2, 3

[27] Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. Selfrecon: Self reconstruction your digital avatar from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5605–5615, 2022. 4

[28] Boyi Jiang, Yang Hong, Hujun Bao, and Juyong Zhang. Selfrecon: Self reconstruction your digital avatar from monocular video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5605–5615, 2022. 2

[29] Animesh Karnewar, Tobias Ritschel, Oliver Wang, and Niloy J Mitra. Relu fields: The little non-linearity that could. *arXiv preprint arXiv:2205.10824*, 2022. 2

[30] Diederik P Kingma and Jimmy Ba. Adam: a method for stochastic optimization. In *Conf. for Learning Representations, San*, 2014. 6

[31] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017. 2

[32] Chen Kong and Simon Lucey. Deep non-rigid structure from motion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 2

[33] Suryansh Kumar. Non-rigid structure from motion: Prior-free factorization method revisited. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2020. 2

[34] Vijay Kumar BG, Gustavo Carneiro, and Ian Reid. Learning local image descriptors with deep siamese and triplet convolutional networks by minimising global loss functions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5385–5394, 2016. 3

[35] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 2

[36] Tianye Li, Mira Slavcheva, Michael Zollhöfer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5521–5531, 2022. 3

[37] 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/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3

[38] Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. Robust high-resolution video matting with temporal guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 238–247, 2022. 4, 6

[39] Wenbin Lin, Chengwei Zheng, Jun-Hai Yong, and Feng Xu. Occlusionfusion: Occlusion-aware motion estimation for real-time dynamic 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1736–1745, 2022. 2, 3, 6, 7, 8

[40] 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. In *Advances in Neural Information Processing Systems*, 2022. 4

[41] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. *The IEEE International Conference on Computer Vision (ICCV)*, 2019. 2

[42] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, 2015. 2

[43] David G. Lowe. Distinctive image features from scale-invariant keypoints. *Int. J. Comput. Vision*, page 91–110, 2004. 3

[44] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2, 3, 4

[45] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 2

[46] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 343–352, 2015. 2, 3, 6, 7, 8

[47] Joseph Ortiz, Alexander Clegg, Jing Dong, Edgar Sucar, David Novotny, Michael Zollhoefer, and Mustafa Mukadam. isdf: Real-time neural signed distance fields for robot perception. *arXiv preprint arXiv:2204.02296*, 2022. 5

[48] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5865–5874, 2021. 2, 3, 4, 6

[49] 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. *arXiv preprint arXiv:2106.13228*, 2021. 2, 3, 4, 6

[50] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 2, 3

[51] Javier Romero, Dimitrios Tzionas, and Michael J. Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)*, 36(6), 2017. 2

[52] Peter Sand and Seth J. Teller. Particle video: Long-range motion estimation using point trajectories. In *CVPR (2)*, pages 2195–2202. IEEE Computer Society, 2006. 2

[53] Vikramjit Sidhu, Edgar Tretschk, Vladislav Golyanik, Antonio Agudo, and Christian Theobalt. Neural dense non-rigid structure from motion with latent space constraints. In *European Conference on Computer Vision*, pages 204–222. Springer, 2020. 2

[54] Miroslava Slavcheva, Maximilian Baust, Daniel Cremers, and Slobodan Ilic. Killingfusion: Non-rigid 3d reconstruction without correspondences. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1386–1395, 2017. 2, 3, 6, 7

[55] Miroslava Slavcheva, Maximilian Baust, and Slobodan Ilic. Sobolevfusion: 3d reconstruction of scenes undergoing free non-rigid motion. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2646–2655, 2018. 3

[56] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Proceedings of the Fifth Eurographics Symposium on Geometry Processing*, page 109–116, Goslar, DEU, 2007. Eurographics Association. 3

[57] Narayanan Sundaram, Thomas Brox, and Kurt Keutzer. Dense point trajectories by gpu-accelerated large displacement optical flow. In *ECCV (1)*, pages 438–451. Springer, 2010. 2

[58] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. 2

[59] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12959–12970, 2021. 3

[60] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018. 2

[61] Jingwen Wang, Tymoteusz Bleja, and Lourdes Agapito. Go-surf: Neural feature grid optimization for fast, high-fidelity rgb-d surface reconstruction. In *Proceedings of the IEEE Conference on 3D Computer Vision (3DV)*, 2022. 2, 4, 5, 6

[62] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *NeurIPS*, 2021. 2, 4, 5, 6

[63] Weipeng Xu, Avishek Chatterjee, Michael Zollhöfer, Helge Rhodin, Dushyant Mehta, Hans-Peter Seidel, and Christian Theobalt. Monoperfcap: Human performance capture from monocular video. *ACM Transactions on Graphics (ToG)*, 37 (2):1–15, 2018. 2

[64] Gengshan Yang, Deqing Sun, Varun Jampani, Daniel Vlasic, Forrester Cole, Huiwen Chang, Deva Ramanan, William T Freeman, and Ce Liu. Lasr: Learning articulated shape reconstruction from a monocular video. In *CVPR*, 2021. 2

[65] Gengshan Yang, Deqing Sun, Varun Jampani, Daniel Vlasic, Forrester Cole, Ce Liu, and Deva Ramanan. Viser: Video-specific surface embeddings for articulated 3d shape reconstruction. In *NeurIPS*, 2021. 2

[66] Gengshan Yang, Minh Vo, Natalia Neverova, Deva Ramanan, Andrea Vedaldi, and Hanbyul Joo. Banmo: Building animatable 3d neural models from many casual videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2863–2873, 2022. 4

[67] Gengshan Yang, Minh Vo, Natalia Neverova, Deva Ramanan, Andrea Vedaldi, and Hanbyul Joo. Banmo: Building animatable 3d neural models from many casual videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2863–2873, 2022. 2

[68] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. 2

[69] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 2

[70] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021. 2

[71] Zehao Yu, Songyou Peng, Michael Niemeyer, Torsten Sattler, and Andreas Geiger. Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction. *arXiv preprint arXiv:2206.00665*, 2022. 2

[72] Juyong Zhang, Yuxin Yao, and Bailin Deng. Fast and robust iterative closest point. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021. 6

[73] Zerong Zheng, Tao Yu, Yebin Liu, and Qionghai Dai. Pamir: Parametric model-conditioned implicit representation for image-based human reconstruction. *IEEE transactions on pattern analysis and machine intelligence*, 44(6): 3170–3184, 2021. 2

[74] Silvia Zuffi, Angjoo Kanazawa, David W Jacobs, and Michael J Black. 3d menagerie: Modeling the 3d shape and pose of animals. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6365–6373, 2017. 2

11