

CAMM: Building Category-Agnostic and Animatable 3D Models from Monocular Videos

Tianshu Kuai Akash Karthikeyan Yash Kant Ashkan Mirzaei Igor Gilitschenski
University of Toronto

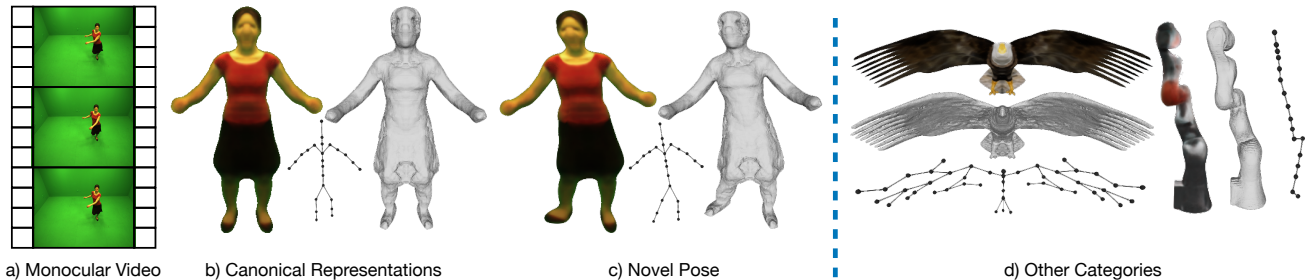


Figure 1. Given monocular videos of an articulated object, our method builds its canonical representation, including 3D shape, appearance, and a corresponding animatable 3D kinematic chain for direct pose manipulations. Our approach does not rely on any information on the object’s shape and underlying structure. In c) we show example of re-posing the human in b) to novel pose. We show the learned canonical representations of other object categories in d).

Abstract

Animating an object in 3D often requires an articulated structure, e.g. a kinematic chain or skeleton of the manipulated object with proper skinning weights, to obtain smooth movements and surface deformations. However, existing models that allow direct pose manipulations are either limited to specific object categories or built with specialized equipment. To reduce the work needed for creating animatable 3D models, we propose a novel reconstruction method that learns an animatable kinematic chain for any articulated object. Our method operates on monocular videos without prior knowledge of the object’s shape or underlying structure. Our approach is on par with state-of-the-art 3D surface reconstruction methods on various articulated object categories while enabling direct pose manipulations by re-posing the learned kinematic chain. Our project page: <https://camm3d.github.io/>.

1. Introduction

Building 3D representations of real-world objects from images has long been studied in 3D Computer Vision. Realistic 3D models allow us to synthesize new images of objects from arbitrary viewpoints and to animate these objects for applications such as virtual and augmented real-

ity. Prior works on building 3D representations for static scenes [2, 6, 10, 14, 28, 35–37, 46, 64, 81, 83] assume the objects in the scene to be rigid to obtain accurate correspondences, which cannot be generalized to deformable objects or scenes with drastic motions. Recent works [21, 23, 27, 41, 45, 49, 57, 61, 73–75] tackle the problem of reconstructing dynamic scenes and 3D deformable objects, and have achieved impressive performance on novel view synthesis. However, these methods can only render objects with the body poses that exist in the given training data because they “mimic” the movements and behaviours of the objects. Therefore, their representations of 3D objects are not applicable for direct pose manipulations or animations.

To fulfill such needs, tremendous efforts have been made to build parametric shape models for humans [31, 42, 68, 69] and quadruped animals [85, 86]. These shape models are often built via capturing large amounts of data using specialized scanners and sensors. They serve as base templates for methods that focus on reconstructing and animating 3D humans and animals [7, 24, 29, 40, 43, 44, 60, 63, 65]. Although these template-based models can offer high-fidelity reconstruction and animation results, they can only be applied to limited object categories. Recent efforts [39, 76] aim to build articulated 3D object models without templates or category-specific priors. However, they either require synchronized multi-view videos which are hard to acquire, or a manually created skeleton as initialization.

In this work, we propose a novel approach to build an animatable 3D model for any articulated object using only monocular videos. We do not rely on prior knowledge of the object shape and structure or manual annotations as initialization. Thus, our pipeline significantly reduces the amount of work needed to create animatable 3D models and eliminates the need for synchronized multi-view observations or specialized sensors. Our model can be directly used for pose manipulations and animations in 3D, which we refer to in this paper as *animatable model*.

Specifically, our approach uses Neural Radiance Fields [37] and Signed Distance Functions to represent the appearance and shape of the object. The body pose of the object is represented by a *kinematic chain*. Our kinematic chain is initialized based on the initial estimate of the object’s shape, and further optimized jointly with the object’s shape, appearance, and deformation parameters. We adopt skinning weights mechanisms from [5, 47, 75] and propose a novel way to use our optimized kinematic chain to drive pose changes and deformations. Due to the under-constrained nature of this problem, we leverage the foreground masks and optical flow predicted by off-the-shelf models as robust visual cues to learn the shape and underlying structure of the object. Inspired by [76], we utilize the 2D image features from DINO-ViT [3] that are trained in a self-supervised manner, to establish long-range correspondences and consistencies at the object parts level [1, 58] between different frames. This is achieved by matching canonical features at the object surface to the pre-trained DINO-ViT image features. Examples of rendered objects in canonical representations and novel poses, along with the corresponding 3D meshes and kinematic chains are shown in Figure 1. Our contributions are summarized as follows:

- Our work is the first to build an animatable 3D model for objects of arbitrary categories from monocular videos without any template or prior knowledge of the object’s shape and underlying structure.
- We propose a novel optimization technique to iteratively refine the kinematic chain and its associated deformation parameters. Users can directly manipulate the optimized kinematic chain to animate the object.
- We achieve similar surface reconstruction results to state-of-the-art 3D surface reconstruction methods on various articulated and deformable object categories.

2. Related Work

3D deformable shape reconstruction from images. Inspired by NeRF [37], recent works [21, 27, 41, 45, 57, 61] are able to obtain robust 3D reconstructions on dynamic scenes or deformable objects from a collection of multi-view images. Some works [15, 20, 26, 79] use additional supervi-

sions from human annotated data or shape templates to recover 3D shapes. These methods suffer from large performance degradation when the input has large deformations and self-occlusions. Recent works [23, 49, 82] explore structure from motion approaches to reconstruct non-rigid 3D scenes using video sequences. Similar to us, LASR [73], ViSER [74], and BANMo [75] utilize monocular videos for reconstructing deformable and articulated 3D shapes and achieve great surface reconstruction results. However, these approaches do not suffice the needs of direct pose manipulations of the reconstructed 3D objects as they reconstruct by memorizing the poses and movements in training data.

Category-specific animatable models. A group of works [4, 30, 34, 80, 84] tackle the problem of animating humans in 2D using annotated data with human poses. To enable animations in 3D space, significant efforts from the community have been put into creating 3D human and animal shape models [31, 42, 68, 69]. Many works [7, 12, 13, 24, 29, 40, 43, 44, 60, 63, 65–67, 70] utilize these template shapes or pose priors from shape models to recover 3D shapes and perform animations. Recent methods [25, 51, 52] do not use shape templates, but rely on pre-defined skeletons or poses predicted by models [19, 20] trained on human annotated data, and these methods require synchronized multi-view inputs. These pre-defined skeletons or poses are used as priors for learning the shape and articulation for animations. However, the annotations are usually expensive to get and category-specific, which cannot be applied to other object categories in the real world and do not work for out-of-distribution articulation modes or poses.

Category-agnostic animatable models. Several methods have recently emerged for building 3D animatable model for any articulated object. Watch it Move [39] builds a 3D articulated structure by identifying physically meaningful joints in an unsupervised manner from calibrated multi-view videos and corresponding foreground masks. Recent work [55] learns deformation behaviours from a given 3D mesh of the object to allow users to define desired deformations at regions of interest. However, synchronized multi-videos and 3D meshes are not easily accessible. Similar to our goal, LASSIE [76] recovers a 3D shape and skeleton from image ensembles without pre-defined shape priors, but it requires a manually annotated 3D skeleton as input. In contrast, we aim to simplify the process of building animatable 3D models for any articulated object using only collections of monocular videos that can be easily collected.

3. Method

Given a collection of monocular videos, our goal is to build an animatable 3D model for the articulated object in the scene. Our method optimizes the canonical representation (Section 3.1) of the object, where its shape and appearance are modeled implicitly using Multi-layer Per-

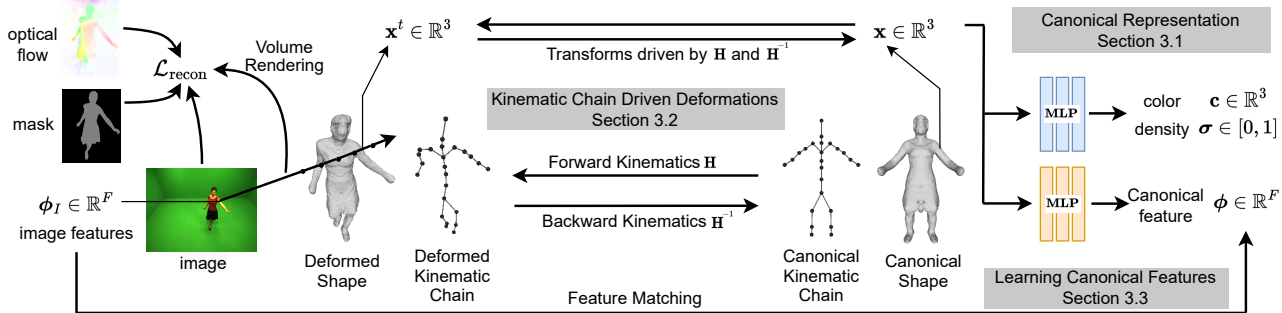


Figure 2. **Method overview.** Our method optimizes the canonical representation, including the object’s shape, appearance, and kinematic chain (Section 3.1), to enable direct pose manipulations. It uses the kinematic chain to transform 3D points between the canonical and deformed space (Section 3.2), and render 2D observation predictions of colors, foreground masks and optical flows to match the actual 2D observations. In addition, canonical feature embeddings (Section 3.3) are learned by matching the corresponding 2D image features to establish the object parts level correspondences across different frames within a collection of monocular videos.

ceptrons (MLPs). Its canonical body pose is represented by a kinematic chain. Our approach consists of two optimization stages: the first stage (Section 3.4.1) optimizes the canonical shape and deformation parameters in each frame without a kinematic chain to get an initial estimate of the 3D shape. We then apply RigNet [71], an off-the-shelf category-agnostic skeleton estimator, on the initial estimate of the 3D mesh to obtain an initial kinematic chain. In the second stage (Section 3.4.2), the kinematic chain adapts to object shape to enable kinematic chain driven deformations (Section 3.2). We supervise the learning by matching the rendered images, foreground masks, and optical flows to the actual 2D observations in both stages. In addition, we learn canonical feature embeddings (Section 3.3) to obtain long-range correspondences between video frames, supervised by pre-trained 2D image features from [3].

3.1. Canonical Representation

Canonical shape. Similar to [37, 75], we represent the shape and appearance of an object implicitly in the canonical space. A 3D point $\mathbf{x} \in \mathbb{R}^3$ has two properties: its color $\mathbf{c} \in \mathbb{R}^3$ and density $\sigma \in [0, 1]$. The color \mathbf{c} is given by a Multilayer Perceptron (MLP) network, and the density σ is given by the cumulative of a Laplacian distribution with zero mean and learnable scale on the learned Signed Distance Function (SDF) [62, 77] in the canonical space. The value of SDF at any 3D point is given by a separate MLP network, and the canonical mesh can be extracted by finding the zero-level set of the SDF [77] with the marching cubes algorithm. Please see our supplementary material for additional details about the MLPs.

Canonical kinematic chain. We use a kinematic chain to represent the body pose of any articulated object. The canonical kinematic chain is defined as a set of connected 3D joints $\mathbf{P} = \{\mathbf{p}_i \mid i = 1, \dots, n_p\}$, where $\mathbf{p}_i \in \mathbb{R}^3$ denotes joint positions in the canonical space, with a set of

$n_j - 1$ links $\mathbf{L} = \{\ell_{jk} = (\mathbf{p}_j, \mathbf{p}_k)\}^{(n_j-1)}$. The joints are connected in a hierarchical manner with a pre-defined root joint to form a tree structure. As a result, there is no cycle, and a unique path exists between every two joints.

3.2. Kinematic Chain Driven Deformations

3.2.1 Neural Blend Skinning Deformations

Inspired by [73–75], we define a set of learnable *deformation anchors* $\mathbf{A} = \{\mathbf{a}_i \mid i = 1, \dots, n_a\}$, where $\mathbf{a}_i \in \mathbb{R}^3$ denotes the anchor positions in the canonical space. We compute surface deformations with respect to anchors via linear blend skinning with learned weights. At time t , let $\mathbf{C}^t \in SE(3)$ be the object’s root pose with respect to the canonical root pose, $\mathbf{T}_i^t \in SE(3)$ be the transformation of anchor \mathbf{a}_i from canonical space to the deformed space, and $\mathbf{w}_{a_i}^t$ be the skinning weight of \mathbf{x} with respect to anchor \mathbf{a}_i . Given any 3D point $\mathbf{x} \in \mathbb{R}^3$ in the canonical space, its corresponding point in the deformed space \mathbf{x}^t at time t is given by the weighted average of anchor’s transformation \mathbf{T}_i^t by:

$$\mathbf{x}^t = \mathbf{C}^t \sum_{i=1}^{n_a} (\mathbf{w}_{a_i} \mathbf{T}_i^t) \mathbf{x} \quad (1)$$

where we define $\mathbf{W} = [\mathbf{w}_{a_1}, \dots, \mathbf{w}_{a_{n_a}}] \in \mathbb{R}^{n_a}$ to be the forward skinning weights for canonical space point \mathbf{x} at time t , and the weights are determined by the Euclidean distances between point \mathbf{x} and all anchors \mathbf{A} in the canonical space as:

$$\mathbf{W} = \sigma \left(-\|\mathbf{x} - \mathbf{A}\|_2^2 \right) \quad (2)$$

where σ is softmax to normalize the skinning weights and its temperature τ is a learnable parameter.

3.2.2 Kinematic Chain Driven Deformations

As the kinematic chain represents the body pose of the object, the deformation anchors’ transformations must be driven by the kinematic chain to obtain properly deformed surfaces corresponding to the actual body pose. To achieve this, we associate each anchor \mathbf{a}_i to its closest kinematic chain link ℓ_{jk} based on Euclidean distances in the canonical space. Example associations are shown as dashed lines in Figure 3. We let the intersection between the dashed line of anchor \mathbf{a}_i and the associated link ℓ_{jk} be \mathbf{m}_{ijk} , where the parent joint and child joint of link ℓ_{jk} are \mathbf{p}_j and \mathbf{p}_k . Then this association between anchor \mathbf{a}_i and link ℓ_{jk} can be represented by the following terms:

$$\alpha_i = \|\mathbf{m}_{ijk} - \mathbf{p}_j\|_2, \beta_i = \|\mathbf{a}_i - \mathbf{m}_{ijk}\|_2,$$

$$\mathbf{G}_i = \text{Rot}(\mathbf{m}_{ijk} - \mathbf{p}_j, \mathbf{a}_i - \mathbf{m}_{ijk}), \quad (3)$$

where α_i is the Euclidean distance between the parent joint \mathbf{p}_j and the intersection \mathbf{m}_{ijk} , and β_i is the Euclidean distance between the intersection \mathbf{m}_{ijk} and associated anchor \mathbf{a}_i . The rotation matrix $\mathbf{G}_i \in SO(3)$ aligns vectors $\mathbf{m}_{ijk} - \mathbf{p}_j$ and $\mathbf{a}_i - \mathbf{m}_{ijk}$ about the intersection \mathbf{m}_{ijk} . We refer to these three terms as the *association parameters* for anchor \mathbf{a}_i , and they are kept constant in *any* kinematic chain configuration to ensure stable surface deformations.

With defined association parameters, we can express an anchor’s position using its associated link and association parameters. As anchor \mathbf{a}_i is associated to link ℓ_{jk} in the canonical space, its universal position $\tilde{\mathbf{a}}_i$ in *any* kinematic chain configuration can be computed as:

$$\tilde{\mathbf{a}}_i = \tilde{\mathbf{p}}_j + \underbrace{\alpha_i \frac{\tilde{\mathbf{p}}_k - \tilde{\mathbf{p}}_j}{\|\mathbf{p}_k - \mathbf{p}_j\|_2}}_{\text{parent joint to intersection}} + \underbrace{\beta_i \mathbf{G}_i \frac{\tilde{\mathbf{p}}_k - \tilde{\mathbf{p}}_j}{\|\mathbf{p}_k - \mathbf{p}_j\|_2}}_{\text{intersection to anchor}}, \quad (4)$$

where $\tilde{\mathbf{p}}_j$ and $\tilde{\mathbf{p}}_k$ are given by:

$$\tilde{\mathbf{p}}_{j,k} = \begin{cases} \mathbf{p}_{j,k} & \text{in canonical space} \\ \mathbf{H}_{j,k}^t \mathbf{p}_{j,k} & \text{in deformed space at time } t \end{cases} \quad (5)$$

Note that $\tilde{\mathbf{a}}_i = \mathbf{a}_i$ in the canonical space. $\mathbf{H}_j^t \in SE(3)$ and $\mathbf{H}_k^t \in SE(3)$ are the rigid transformations of joint \mathbf{p}_j and joint \mathbf{p}_k from the canonical space to the deformed space without breaking the kinematic chain, which we refer as *forward kinematics* on the kinematic chain.

Then the rigid transformation \mathbf{T}_i^t of anchor \mathbf{a}_i at time t can be directly inferred based on its new position in the deformed space given by Eq. (4):

$$\mathbf{T}_i^t = [\mathbb{I} \mid \tilde{\mathbf{a}}_i - \mathbf{a}_i], \quad (6)$$

where $\mathbb{I} \in \mathbb{R}^{3 \times 3}$ is the identity matrix.

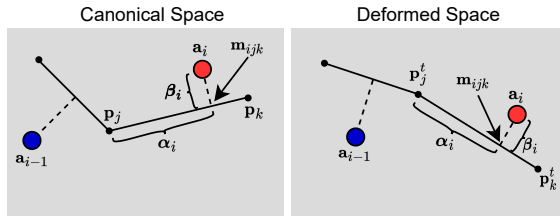


Figure 3. Simplified 2D example on how anchors (red and blue dots) move along with associated kinematic chain links. Associations can be visualized at dashed lines.

Given optimized kinematic chain and anchors, we can perform explicit re-posing of the object by directly applying user defined forward kinematics \mathbf{H} to the kinematic chain.

3.3. Learning Semantically Consistent Features

Large deformations and self-occlusions often cause 3D reconstruction methods to fail or perform poorly. This is usually due to the lack of strong correspondence cues or the lack of long-range correspondences between different frames in long videos [48, 54]. Recent works [1, 58] show that the pre-trained image features from DINO-ViT [3, 8] can provide robust correspondences at the object parts level between different 2D images even under large appearance variations and viewing angle changes.

To better optimize the canonical shape and deformations, we follow [74, 75] to learn surface feature embeddings $\phi(\mathbf{x}) \in \mathbb{R}^F$ for points on the canonical shape’s surface. The feature is given by a separate MLP network that takes in any 3D point in the canonical space and it is trained in a self-supervised manner such that the surface feature embeddings match the pre-trained 2D image features at its corresponding pixel given by the camera projection model. Compared to BANMo [75] where pre-trained DensePose CSE features [38] specifically designed for humans and animals are used for learning surface embedding, DINO-ViT [3] features have shown to be meaningful and rich in downstream tasks for a wide variety of objects [1, 3]. This is because DINO-ViT is pre-trained in a self-supervised manner on large-scale datasets with diverse object categories.

3.4. Two-Stage Optimizations

3.4.1 Initial Optimization

The initial optimization stage aims to build a reasonable object shape with unconstrained anchors to capture deformations in each frame of the training videos. We initialize the MLP network that predicts the canonical shape of the object to approximate the object as a 3D ellipsoid, and initialize all the anchors \mathbf{A} at the global origin in the canonical space. During the initial optimization stage, the anchors are updated without any constraints in terms of their positions. Inspired by [75], the *unconstrained transformations* of the

anchors at time t are given by a separate MLP, \mathcal{F}_A as:

$$\hat{\mathbf{T}}^t = \mathcal{F}_A(\psi_a^t), \quad (7)$$

where $\psi_a^t \in \mathbb{R}^{128}$ is a learnable latent code for time t .

As we define the transformation of any 3D point \mathbf{x} from the canonical space to the deformed space as weighted rigid transformation of anchors in Eq. (1), we can directly infer the backward transformation for any 3D point \mathbf{x}^t in the deformed space by taking the inverse of the rigid transformations as:

$$\mathbf{x} = \left(\sum_{i=1}^{n_a} (\mathbf{w}_{a_i}^t (\hat{\mathbf{T}}_i^t)^{-1}) \right) (\mathbf{C}^t)^{-1} \mathbf{x}^t$$

$$\mathbf{W}^t = \sigma \left(-\|\mathbf{x}^t - \mathbf{A}^t\|_2^2 \right), \quad (8)$$

where the backward skinning weights $\mathbf{W}^t = [\mathbf{w}_{a_1}^t, \dots, \mathbf{w}_{a_{n_a}}^t] \in \mathbb{R}^{n_a}$ are computed in the same way as in forward skinning in Eq. (2) but based on deformed anchors and the 3D point of interest in the deformed space instead. This allows us to directly transform any 3D point back and forth between the canonical space and each frame’s deformed space.

Similar to [74, 75], the canonical shape and appearance are learned by minimizing the reconstruction losses on 2D observations: RGB images, foreground masks, and optical flow. The reconstruction losses are formulated in the same way as in [37, 73–75, 78], where we minimize the differences between the rendered and the actual observations:

$$\mathcal{L}_{\text{recon}} = \sum_{\mathbf{x}_I} \|\mathbf{o}_r - \mathbf{o}_{gt}\|^2, \quad (9)$$

where \mathbf{o}_r and \mathbf{o}_{gt} are the pairs of rendered and ground-truth 2D observations (2D images, foreground masks, optical flow) at pixels of interest $\mathbf{x}_I^t \in \mathbb{R}^2$ at time t . To perform volumetric rendering at time t , we first sample a camera ray starting at the pixel of interest on the 2D image at time t , and use Eq. (8) to transform the sampled points along the camera ray back to the canonical space. We then follow the same rendering process in NeRF [37] on the camera ray in the canonical space to get the colour \mathbf{c}^t at the pixel of interest. To render optical flow, we follow [75] to transform canonical space points to consecutive frames’ deformed spaces, and use the camera projection model to find their corresponding positions on the 2D image. We compute the difference in the pixel locations as the optical flow.

The anchors and their transformations are learned from scratch with additional losses to enforce cycle consistency. The consistency losses enforce any 3D point in the deformed space after a backward transformation and a forward transformation to end up at the same position. To supervise canonical feature embeddings, we apply soft argmax

descriptor matching [16, 33, 74, 75] on the 2D pixel of interest to determine the most probable corresponding surface point in the canonical space. The matching is based on cosine similarities between the canonical features and the pre-trained 2D image features. A 2D consistency loss [75] is used to minimize the difference between matched position and the location given by the camera projection model. The details of loss functions are in the supplementary material.

3.4.2 Kinematic chain aware optimization

After the initial optimization stage, we obtain an estimate of the canonical shape and a set of unconstrained anchors. However, the unconstrained anchors introduce unrealistic artifacts in deformations, such as undesired shrinks and stretches on the canonical shape. And the anchors are only optimized to “mimic” object shape in each frame of the videos. A kinematic chain is needed to remove these undesired effects and enable explicit re-posing of the object.

We initialize our canonical kinematic chain with pre-trained RigNet [71] that takes in our canonical mesh from the initial optimization stage and outputs a set of joints and connections. The pre-trained RigNet is able to provide a reasonable estimate of the kinematic chain. However, it requires further optimization such that our kinematic chain can better adapt to the object’s movements in the training videos, and therefore align itself with the object’s actual underlying structure in each frame.

As any point in 3D can be transformed back and forth between the canonical space and the deformed space using the transformations defined by anchors, we utilize the learned unconstrained anchor transformations $\hat{\mathbf{T}}^t$ from the initial optimization stage to transform the canonical kinematic chain joints \mathbf{P} following Eq. (1):

$$\hat{\mathbf{P}}^t = \mathbf{C}^t \sum_{i=1}^{n_a} (\mathbf{w}_{a_i} \hat{\mathbf{T}}_i^t) \mathbf{P}, \quad (10)$$

where $\hat{\mathbf{P}}^t = \{\hat{\mathbf{p}}_i \mid i = 1, \dots, n_p\}$ is the set of *unconstrained kinematic chain joints* in the deformed space at time t . The transformations of anchors are learned without any regularization on how they move or any knowledge on the existence of the kinematic chain. Therefore the transformed joints $\hat{\mathbf{P}}^t$ computed using unconstrained anchors will break the kinematic chain. We can recover the properly deformed kinematic chain as a set of *revised joints* $\tilde{\mathbf{P}}^t = \{\tilde{\mathbf{p}}_i \mid i = 1, \dots, n_p\}$ in the deformed space. It is done by updating each unconstrained joint’s position in a hierarchical order to enforce the length of each kinematic chain link to be constant. Note that the order of joint connections is preserved at all times after initialization to maintain the hierarchical structure of the kinematic chain and stability during optimization. The detailed steps to recover the kinematic chain are explained in the supplementary material.

As we have the fixed associations between anchors and the initial kinematic chain built in the canonical space as α_i , β_i , and \mathbf{G}_i by Eq. (3). We then use Eq. (4) to recover each anchor’s revised position $\tilde{\mathbf{a}}_i$ and thus apply Eq. (6) to infer a revised anchor transformation $\hat{\mathbf{T}}_i^t$ that satisfies the association constraints between the anchors and the kinematic chain, while not breaking the kinematic chain at the same time. During the optimization, we introduce a novel loss that minimizes the differences between the unconstrained anchors and the revised anchors:

$$\mathcal{L}_{\text{anchors}} = \sum_{i=1}^{n_a} \left\| \tilde{\mathbf{a}}_i - \hat{\mathbf{T}}_i^t \mathbf{a}_i \right\|_2^2. \quad (11)$$

This loss aims to make the optimized deformation anchors and the MLP network \mathcal{F}_A kinematic chain aware.

Additionally, we optimize the kinematic chain by introducing a learnable *additive residual* term for each kinematic chain link. The length of each link in the kinematic chain is updated using the learnable residual during optimization such that the kinematic chain can better adapt to the object’s shape and the learned deformation anchors. The details on how the kinematic chain is updated without breaking itself are included in the supplementary material.

In the kinematic chain aware optimization stage, all the MLP networks are also optimized jointly to finetune the object’s shape and appearance with the loss functions in the initial optimization stage. At the end of this stage, a fully controllable articulated model for the object is built, and animations of the object can be easily done by directly manipulating the final kinematic chain.

4. Experiments

4.1. Datasets

Synthetic Datasets. We validate our approach on two synthetic datasets. The Eagle [75] dataset consists of 5 single-view videos (900 frames in total) of an animated eagle model from Turbosquid. We also created a new dataset (iiwa dataset) that contains 4 single-view videos (1200 frames in total) of an animated iiwa robot arm model from BlenderKit. Please find the details and examples of our iiwa dataset in the supplementary material.

AMA Dataset. We also evaluate our approach on the Articulated Mesh Animation (AMA) Dataset [59]. It contains 10 sets of multi-view videos of 3 different clothed humans using synchronized cameras. We use 2 challenging sets of videos and treat the multi-view videos as a collection of monocular videos by discarding the time synchronization information, which yields 2600 monocular frames in total.

In-the-wild Datasets. In addition, we validate our approach on in-the-wild monocular captures of various deformable and articulated objects datasets. In the main paper, we focus on experimental results on the synthetic and

AMA datasets. Please see the supplementary material for implementation details and results on in-the-wild datasets.

4.2. Implementation Details

For synthetic and AMA datasets, we use ground-truth foreground masks during the optimization. For root pose initializations, we use ground-truth root poses for synthetic datasets, while a pre-trained PoseNet [75] is used to initialize the root poses for AMA Dataset. The root poses are also updated during optimization with learned residual terms for AMA Dataset. The optical flows for all datasets are computed by VCN [72] at frame intervals of $\{1, 2, 4, 8, 16, 32\}$ in both forward and backward time directions. We use the keys of the last transformer layer from a pre-trained ViT-S/8 model [3] as supervisions for the canonical features. We perform PCA on the 2D image features to reduce the dimension from 384 to 16, to match our canonical feature’s dimension. To generate mesh from our implicit representations of the object shape, we run marching cubes on $256 \times 256 \times 256$ grids to extract the zero-level set of the learned SDF. We use 10, 25, and 36 anchors for iiwa, Eagle, and AMA Dataset, respectively. Our code and iiwa dataset will be made publicly available.

4.3. Pose Manipulation Results

The advantage of having a kinematic chain is that users can manipulate the 3D object to poses that had never occurred in training videos. Although some surface reconstruction methods [74, 75] allow users to move their learned “floating” control points around for custom animations, these control points are not intuitive to manipulate because it’s unclear how each “floating” control point contributes to the overall deformation. Moreover, these “floating” control points are optimized for only memorizing the poses and movements in training videos. In our approach, the pose manipulations are done by directly applying rotations at kinematic chain joints to rotate and twist corresponding kinematic chain links, which is much more straightforward and intuitive for animation purposes. We do not impose any constraint on the kinematic chain configuration as long as the order of joint connections and the length of each link are preserved. Users are expected to apply valid transformations to the kinematic chain to obtain feasible shapes without mesh degeneration. We show the learned canonical shapes, optimized kinematic chains, and the re-posed objects in novel poses that do not occur in training videos in Figure 4. Please see our supplementary material for additional re-posing results.

4.4. Reconstruction Results

There is no prior work on building category-agnostic and animatable 3D models from monocular videos, and reconstruction quality is an important factor for realistic novel

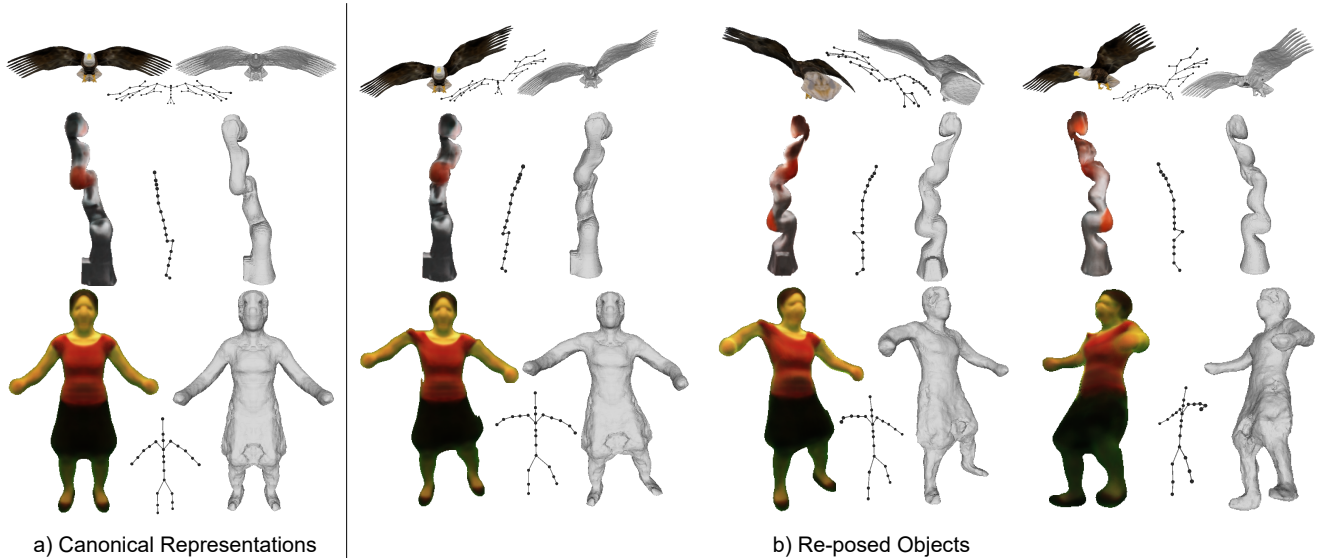


Figure 4. **Pose manipulation examples.** In a) we show the learned canonical representations of the object for each dataset. We perform pose manipulations using optimized kinematic chain and show the transformed kinematic chain, re-posed mesh as well as the rendered object in three different views in b).

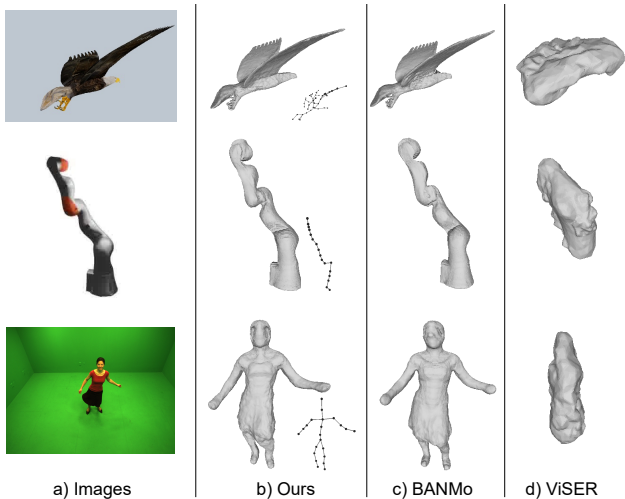


Figure 5. **Qualitative comparison of our method with BANMo [75] and ViSER [74].**

view synthesis and animations. Therefore we evaluate the performance of our method on 3D surface reconstruction using ground-truth 3D meshes provided by each dataset.

Metrics. We align the reconstructed 3D mesh to the ground-truth mesh using Iterative Closest Point (ICP) if they are in different scales and orientations before evaluation. Following [9, 11, 53, 75], we report 3D Chamfer Distances between reconstructed mesh vertices and the ground-truth mesh vertices measured in centimeters. In addition, we follow [56] to report the F-scores at a distance threshold of 2%. All the reported numbers are averaged across all frames for each dataset. In general, these metrics measure

Table 1. 3D surface reconstruction results evaluated in 3D Chamfer Distances (\downarrow) and F-scores (\uparrow) at a distance threshold of 2% averaged across all frames.

Method	Eagle		iiwa		AMA-swing	
	CD	F@2%	CD	F@2%	CD	F@2%
ViSER	32.6	9.4	20.6	18.1	17.9	39.2
BANMo	4.44	82.72	5.55	54.58	9.28	56.24
Ours	4.21	83.38	5.52	56.53	9.69	53.29

the following two aspects of the reconstruction quality for each frame in the training videos: the overall quality of the reconstructed 3D mesh surface, and whether the model is able to predict an accurate shape corresponding to the correct object body pose.

Baselines. We compare our approach with two baselines in 3D surface reconstruction for deformable objects: BANMo [75] and ViSER [74]. We provide the same root pose initializations to all the methods for fair comparison on each dataset. Note that BANMo and ViSER *cannot* perform direct pose manipulations, and BANMo utilizes DensePose CSE feature embeddings [38] designed for humans and quadruped animals, which makes it a category-specific method. We show qualitative comparison of our method with BANMo [75] and ViSER [74] in Figure 5. Our approach can achieve similar 3D surface reconstruction quality using a category-agnostic approach while enabling direct pose manipulations. We also show the quantitative comparison with the same baselines: BANMo [75] and ViSER [74] on each dataset in Table 1.

Table 2. Ablation on kinematic chain aware optimization evaluated in 3D Chamfer Distances (\downarrow) and F-scores (\uparrow) at a distance threshold of 2% averaged across all frames.

Method	Eagle		iiwa		AMA-swing	
	CD	F@2%	CD	F@2%	CD	F@2%
w/o opt.	7.44	57.87	5.69	52.44	11.13	48.88
with opt.	4.21	83.38	5.52	56.53	9.69	53.29
Improv.	-3.23	+25.51	-0.17	+4.09	-1.44	+4.41

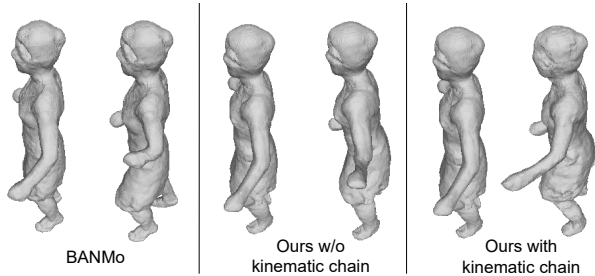


Figure 6. **Ablation on kinematic chain.** We show canonical shapes (left) and the predicted shapes (right) in the same time instance for each method. The length of the left arm is preserved with the help of kinematic chain.

4.5. Ablations

We run ablation studies on important components of our method to evaluate their effects. Please refer to our supplementary materials for more comprehensive results.

Kinematic chain. A kinematic chain enables explicit pose manipulations and effectively reduces the unrealistic surface deformations that often occur in 3D surface reconstruction methods. We show the improvements of having a kinematic chain in Figure 6. The kinematic chain correctly re-poses the person’s arm without volume loss, while BANMo [75] and our method with unconstrained anchors fail to model the correct deformations. We also compare the reconstruction results of directly using kinematic chain initialization from [71] and using optimized kinematic chain in Table 2. We are able to gain considerable improvements with the kinematic chain aware optimization stage (Section 3.4.2) because the optimized kinematic chain aligns better with the underlying object body pose in each frame.

Canonical features. DensePose CSE feature embeddings [38] provide pre-trained 2D-3D correspondences for humans and are used by BANMo [75] on the AMA dataset. We also test our approach by replacing the DINO-ViT [3] features with the same DensePose CSE feature embeddings to supervise our canonical features. As shown in Figure 7, the arm is merged with the human body if we disable the canonical feature learning in our method. This is because the model cannot differentiate between the human body and the arm at this pose due to the lack of the object parts level visual cues. With DINO-ViT features, our method can separate the arm from the human body. However, it does not

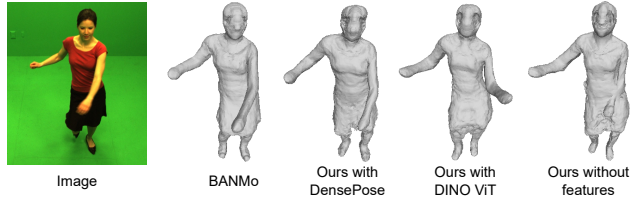


Figure 7. **Ablation on canonical features.** All methods recover high-fidelity surface, but only BANMo and our approach with DensePose supervision correctly predict the left arm’s pose.

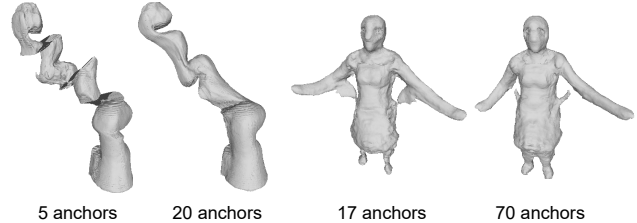


Figure 8. **Ablation on number of anchors.** We show the results of using half and double amounts of anchors on iiwa and AMA.

accurately predict the pose of the arm. The small gap between our approach with DINO-ViT and BANMo in Table 1 mainly comes from these inaccurate pose predictions because the quality of the reconstructed mesh surface is similar, as shown in Figure 5.

Number of deformation anchors. The number of deformation anchors is an object-specific hyper-parameter that plays a crucial role in modeling surface deformations. We show the results of using half and double amounts of anchors on AMA and iiwa datasets in Figure 8. Having a small number of anchors limits our method’s ability to model deformations at the robot arm’s joints and the human’s elbows, while having too many anchors makes the optimization process harder due to the over-parameterization of the deformations at the object’s surface.

5. Conclusion

In this paper, we present a novel reconstruction pipeline that builds an animatable 3D model for any articulated objects from a collection of monocular videos. We leverage the foreground masks, optical flow, and pre-trained image features to iteratively optimize the kinematic chain along with the object’s shape, appearance and deformation parameters. Our method can be generalized to any articulated object as it does not rely on category-specific priors. It allows users to render the object in arbitrary viewpoints and to perform pose manipulations in 3D directly while achieving state-of-the-art 3D surface reconstruction results. However, our kinematic chain does not check for unreachable states, such as chain collisions and foldings. We leave the task of learning physically plausible kinematic chains to future

work. In addition, our deformation anchors are optimized based on the observed deformations from training videos. Therefore, our model does not generalize well to the novel poses that involve unseen object surface deformations.

A. Appendix Summary

In Section B, we present additional details on the method. Section C provides additional details on the loss functions and optimization process. In Section D, we show more experimental results on each dataset. In Section E, we show examples of learned anchors and their corresponding associations. Section F provides additional discussions on the proposed two-stage optimization pipeline. Section G and Section H include the discussions on our method’s limitations and potential negative impact. In Section I, we show examples of failure cases. In Section J, we provide additional details on our iiwa dataset. Finally, we provide a list of the important variables used in the paper along with their state space and descriptions in Section K.

B. Method Details

B.1. Canonical representation

We use the same canonical appearance, shape and feature representations as in BANMo [75], where the color $\mathbf{c} \in \mathbb{R}^3$, the Signed Distance Function (SDF) value $\mathbf{v} \in \mathbb{R}$, and the canonical feature $\phi \in \mathbb{R}^{16}$ of a point $\mathbf{x} \in \mathbb{R}^3$ in the canonical space are given by separate Multi-Layer Perception (MLP) Networks \mathcal{F}_C , \mathcal{F}_S , and \mathcal{F}_ϕ , respectively:

$$\begin{aligned} \mathbf{c} &= \mathcal{F}_C(\mathbf{x}, \mathbf{d}, \psi_l) \\ \mathbf{v} &= \mathcal{F}_S(\mathbf{x}) \\ \phi &= \mathcal{F}_\phi(\mathbf{x}), \end{aligned} \quad (12)$$

where $\mathbf{d} \in \mathbb{R}^2$ is the viewing direction and $\psi_l \in \mathbb{R}^{64}$ is a learnable latent code that captures the environment illumination conditions following [35]. Similar to [62, 77], we compute the density $\sigma \in \mathbb{R}$ in the canonical space as:

$$\sigma = \Psi_\beta(\mathbf{v}), \quad (13)$$

where Ψ_β is cumulative of a Laplacian distribution with zero mean and learnable scale β , and \mathbf{v} is the value of the SDF at the point \mathbf{x} .

To perform volumetric rendering for the pixel of interest $\mathbf{x}_I^t \in \mathbb{R}^2$ at time t , we first sample N points along the camera ray corresponding to the pixel in deformed space. As our implicit representations of the shape and appearance are defined in canonical space, we first transform the sampled points back to the canonical space [41, 75] using the learned backward kinematics. Let the i -th point along the canonical space camera ray be \mathbf{x}_i , and its corresponding color and density queried from the MLP networks be \mathbf{c}_i and σ_i , we

then follow the volumetric rendering process in NeRF [37] to get the rendered color at the pixel of interest \mathbf{c}_r as:

$$\begin{aligned} \mathbf{c}_r &= \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i \\ T_i &= \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right), \end{aligned} \quad (14)$$

where T_i is the accumulated transmittance between the camera center to the i -th sampled point, δ_i is the interval between consecutive sampled points.

B.2. Recovering proper kinematic chain

In this section, we provide additional details on how to recover a proper and connected kinematic chain after all the joints in the canonical space are transformed into the deformed space using unconstrained transformations.

Given the canonical kinematic chain as $\mathbf{P} = \{\mathbf{p}_i \mid i = 1, \dots, n_j\}$, and the initial optimization stage’s learned unconstrained anchor transformations $\hat{\mathbf{T}}^t$ from the canonical space to the deformed space, we can then apply $\hat{\mathbf{T}}^t$ to transform the canonical kinematic chain joints \mathbf{P} to the deformed space as:

$$\hat{\mathbf{P}}^t = \mathbf{C}^t \sum_{i=1}^{n_a} (\mathbf{w}_{a_i} \hat{\mathbf{T}}_i^t) \mathbf{P}, \quad (15)$$

where $\hat{\mathbf{P}}^t = \{\hat{\mathbf{p}}_i \mid i = 1, \dots, n_p\}$ is the set of unconstrained kinematic chain joints in the deformed space at time t . These unconstrained joints cannot form a proper kinematic chain that preserves the hierarchical order of connections and consistent kinematic chain length, because the unconstrained transformations $\hat{\mathbf{T}}^t$ are learned without any knowledge of the kinematic chain. To recover the proper kinematic chain, we follow Algorithm 1 to obtain the revised location of the joints in hierarchical order denoted by $\tilde{\mathbf{P}} = \{\tilde{\mathbf{p}}_i \mid i = 1, \dots, n_j\}$. Given the joints’ revised locations, we can infer the anchors’ revised transformations $\tilde{\mathbf{T}}^t$ by enforcing pre-defined associations between anchors and the kinematic chain links in the deformed space.

B.3. Additive residuals to kinematic chain

We introduce a learnable additive residual term for each kinematic chain link during the kinematic chain optimization stage. Given the canonical kinematic chain as $\mathbf{P} = \{\mathbf{p}_i \mid i = 1, \dots, n_j\}$ and its corresponding links as $\mathbf{L} = \{\ell_{jk} = (\mathbf{p}_j, \mathbf{p}_k)\}^{(n_j-1)}$, let $\mathbf{R} = \{\mathbf{r}_i \mid i = 1, \dots, (n_j - 1)\}$ be the set of learnable residuals. We clip the residuals to avoid drastic changes on the kinematic chain initialization for stable training as:

$$\tilde{\mathbf{R}} = \gamma \tanh(\mathbf{R}), \quad (16)$$

Algorithm 1 Recovering proper kinematic chain

```
for  $j = 1 : n_j$  in hierarchical order do
  for  $\mathbf{p}_k \in \text{Children}(\mathbf{p}_j)$  do
     $\boldsymbol{\mu} \leftarrow \frac{\|\mathbf{p}_k - \mathbf{p}_j\|_2}{\|\hat{\mathbf{p}}_k - \hat{\mathbf{p}}_j\|_2}$ 
     $\tilde{\mathbf{p}}_k = \tilde{\mathbf{p}}_j + \boldsymbol{\mu}(\hat{\mathbf{p}}_k - \hat{\mathbf{p}}_j)$ 
     $\mathbf{t}_k \leftarrow \tilde{\mathbf{p}}_k - \hat{\mathbf{p}}_k$ 
    for  $\mathbf{p}_d \in \text{Descendants}(\mathbf{p}_k)$  do
       $\hat{\mathbf{p}}_d \leftarrow \hat{\mathbf{p}}_d + \mathbf{t}_k$ 
    end for
  end for
end for
 $j \leftarrow j + 1$ 
end for
```

where $\tilde{\mathbf{R}} = \{\tilde{\mathbf{r}}_i \mid i = 1, \dots, (n_j - 1)\}$ are the set of clipped residuals and γ is a hyper-parameter that ensures the change of each kinematic chain link’s length is within $(-\gamma, \gamma)$.

We update the canonical kinematic chain with the latest residuals $\tilde{\mathbf{R}}$ in each iteration to obtain the updated kinematic chain $\tilde{\mathbf{P}} = \{\tilde{\mathbf{p}}_i \mid i = 1, \dots, n_j\}$ following Algorithm 2. Specifically, we start from the root joint of the kinematic chain and update each link’s length by changing the corresponding joint’s position in a hierarchical manner to ensure that we do not break the kinematic chain. After the residual update, we compute the new set of association parameters for each anchor with respect to the updated kinematic chain for kinematic chain driven deformations.

Algorithm 2 Kinematic chain update with residuals

```
for  $j = 1 : n_j$  in hierarchical order do
  for  $\mathbf{p}_k \in \text{Children}(\mathbf{p}_j)$  do
     $\boldsymbol{\mu} \leftarrow \frac{\|\mathbf{p}_k - \mathbf{p}_j\|_2 + \tilde{\mathbf{r}}_k}{\|\mathbf{p}_k - \mathbf{p}_j\|_2}$ 
     $\tilde{\mathbf{p}}_k = \tilde{\mathbf{p}}_j + \boldsymbol{\mu}(\mathbf{p}_k - \mathbf{p}_j)$ 
     $\mathbf{t}_k \leftarrow \tilde{\mathbf{p}}_k - \mathbf{p}_k$ 
    for  $\mathbf{p}_d \in \text{Descendants}(\mathbf{p}_k)$  do
       $\mathbf{p}_d \leftarrow \mathbf{p}_d + \mathbf{t}_k$ 
    end for
  end for
end for
 $j \leftarrow j + 1$ 
end for
```

C. Additional Implementation Details

In this section, we provide additional details on the loss functions and optimizations.

C.1. Loss functions

Similar to [75], we impose reconstruction losses on the 2D observations, including 2D images, foreground masks, and optical flow. We follow the volumetric rendering process described in B.1 to obtain predicted color and foreground mask values. To render optical flow, we follow [75]

to transform canonical space points to consecutive frames’ deformed spaces, and use the camera projection model to find their corresponding positions on the 2D image. We then compute their difference as the optical flow at the point of interest. The reconstruction losses are formulated in the same way as in [37, 73–75, 78], where we compute the difference between the rendered and the actual observations:

$$\mathcal{L}_{\text{recon}} = \sum_{\mathbf{x}_I} \|\mathbf{o}_r - \mathbf{o}_{gt}\|^2, \quad (17)$$

where \mathbf{o}_r and \mathbf{o}_{gt} are the pairs of rendered and ground-truth 2D observations (2D images, foreground masks, optical flow) at pixels of interest $\mathbf{x}_I^t \in \mathbb{R}^2$ at time t .

We apply soft argmax descriptor matching [16, 33, 74, 75] at the pixel of interest \mathbf{x}_I^t , to find the most probable corresponding surface point $\mathbf{x}_m \in \mathbb{R}^3$ in the canonical space by matching learned canonical feature embeddings at the object’s surface with the pre-trained DINO ViT [3] features. Let $\mathbf{x}_c \in \mathbb{R}^3$ be the point that corresponds to the pixel of interest \mathbf{x}_I^t based on backward kinematics. To supervise the canonical feature embedding, we follow [75] to minimize the difference between these two points’ positions as:

$$\mathcal{L}_{3d\text{-match}} = \sum_{\mathbf{x}_I} \|\mathbf{x}_m - \mathbf{x}_c\|_2^2. \quad (18)$$

Let π^t be the camera project model at time t , we transform \mathbf{x}_m from the canonical space to the deformed space using learned forward kinematics, and use π^t to project it to image coordinates. We follow [22, 74, 75] to enforce the consistency at the image coordinate level by minimizing the difference between the obtained point and \mathbf{x}_I^t :

$$\mathcal{L}_{2d\text{-match}} = \sum_{\mathbf{x}_I} \left\| \pi^t \left(\mathbf{C}^t \sum_{i=1}^{n_a} (\mathbf{w}_{a_i} \mathbf{T}_i^t) \mathbf{x}_m \right) - \mathbf{x}_I^t \right\|_2^2. \quad (19)$$

In addition, we follow [27, 75] to encourage consistency in terms of learned forward kinematics and backward kinematics. Specifically, we transform each point \mathbf{x}_i^t along the sampled camera ray in the deformed space to the canonical space using backward kinematics, and then transform each point back to the deformed space. We encourage the resulting point $\mathbf{x}_i^{t'}$ to be at the same position as \mathbf{x}_i^t :

$$\mathcal{L}_{\text{transform}} = \sum_i T_i \left\| \mathbf{x}_i^{t'} - \mathbf{x}_i^t \right\|_2^2, \quad (20)$$

where \mathbf{x}_i^t is the original sampled point along the camera ray, $\mathbf{x}_i^{t'}$ is the sampled point that undergoes a backward and a forward transform, and T_i the accumulated transmittance between the camera center to the i -th sampled point, which we use as the weight for individual sampled points.

During the kinematic chain optimization stage, we also introduce regularization on the anchors and learned transformations by minimizing the differences between the revised anchors based on the kinematic chain and the unconstrained anchors from learned unconstrained anchor transformations $\hat{\mathbf{T}}_i^t$ predicted by \mathcal{F}_A :

$$\mathcal{L}_{\text{anchors}} = \sum_{i=1}^{n_a} \left\| \tilde{\mathbf{a}}_i - \hat{\mathbf{T}}_i^t \mathbf{a}_i \right\|_2^2, \quad (21)$$

where $\tilde{\mathbf{a}}_i$ is the revised anchor based on kinematic chain constraints and $\hat{\mathbf{T}}_i^t$ is the unconstrained rigid transformation from the canonical space to the deformed space given by the MLP network \mathcal{F}_A for the canonical space anchor \mathbf{a}_i .

C.2. Optimizations

Our model is trained using the Adam [17] optimizer with weight decay on a single RTX 3090 GPU. In the initial optimization stage, we use an initial learning rate of 0.0005 and update it using one-cycle policy [50] and cosine annealing [32]. Inspired by [75], we introduce a separate MLP network after the initial optimization stage that takes in a 3D point, and outputs a skinning weight residual with respect to each anchor that adds to the Euclidean distance based skinning weights before normalization, to obtain smoother deformed surface. All the MLP networks and learnable parameters are optimized jointly during training. For the iiwa dataset, we do not update the kinematic chain link lengths with additive residuals because iiwa is considered a rigid robot arm without deformable surface, and its kinematic chain initialization is generally a single chain that already captures its articulation modes.

D. Additional Results

D.1. More qualitative results

We test on in-the-wild monocular captures of various deformable and articulated objects [75]. We use off-the-shelf pre-trained models to obtain foreground masks [18], initial root pose estimation [75] and optical flow [72] for each frame of the in-the-wild dataset videos. We show manual pose manipulation results in Figure 9, and 3D surface reconstruction results in Figure 10 for in-the-wild datasets.

We also include video examples of manual pose manipulations by re-posing the kinematic chain, and the 3D reconstruction results for poses in training videos for each dataset. Please refer to the videos in the supplementary material for better visualization.

D.2. More quantitative results

We report quantitative 3D surface reconstruction results for iiwa (Table 3), Eagle (Table 4), AMA-swing (Table 5), and AMA-samba (Table 6) evaluated in 3D Chamfer Distances in centimeters and F-scores at distance thresholds of

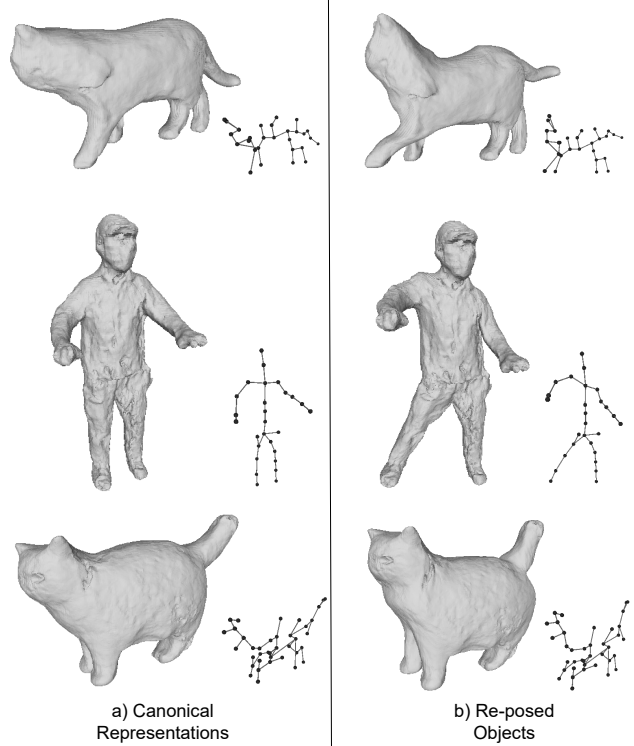


Figure 9. **Pose manipulation examples.** In a) we show the learned canonical representations of the objects for in-the-wild datasets. We perform pose manipulations using optimized kinematic chain and show the reposed kinematic chain and meshes in b).

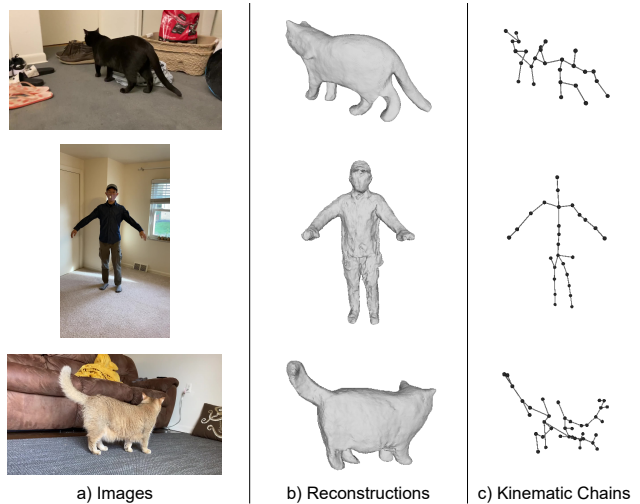


Figure 10. **3D surface reconstruction results.** In a) we show some RGB images from the in-the-wild datasets. In b) we show the reconstructed 3D meshes and in c) we show the corresponding kinematic chains for the same frame.

1%, 2%, and 5%. Note that for the AMA dataset, we train each model using all the sequences (both swing and samba), and report results on swing and samba separately.

Table 3. 3D surface reconstruction results on **iiwa** dataset evaluated in 3D Chamfer Distances (\downarrow) and F-scores (\uparrow) at distance thresholds of 1%, 2%, and 5% averaged across all frames.

Method	CD	F@1%	F@2%	F@5%
5 anchors	7.40	26.76	51.16	79.84
10 anchors	5.52	28.69	56.53	85.61
20 anchors	5.95	28.56	55.17	83.51
with DensePose feat.	5.64	28.85	55.67	85.79
w/o feat.	7.26	25.16	47.33	73.82
w/o kine. chain opt.	5.69	26.90	52.44	81.81

Table 4. 3D surface reconstruction results on **Eagle** dataset evaluated in 3D Chamfer Distances (\downarrow) and F-scores (\uparrow) at distance thresholds of 1%, 2%, and 5% averaged across all frames.

Method	CD	F@1%	F@2%	F@5%
12 anchors	4.51	43.20	81.86	98.88
25 anchors	4.21	43.41	83.38	99.22
50 anchors	4.38	42.21	82.52	98.91
with DensePose feat.	4.31	43.07	82.45	99.11
w/o feat.	4.31	43.51	82.81	98.98
w/o kine. chain opt.	7.44	25.30	57.87	94.40

Table 5. 3D surface reconstruction results on **AMA-swing** evaluated in 3D Chamfer Distances (\downarrow) and F-scores (\uparrow) at distance thresholds of 1%, 2%, and 5% averaged across all frames.

Method	CD	F@1%	F@2%	F@5%
17 anchors	12.41	24.37	45.69	78.10
35 anchors	9.69	29.17	53.29	85.20
70 anchors	12.49	25.62	46.72	77.45
with DensePose feat.	8.96	33.68	58.46	86.49
w/o feat.	9.88	27.60	52.15	84.42
w/o kine. chain opt.	11.13	26.54	48.88	80.69

Table 6. 3D surface reconstruction results on **AMA-samba** evaluated in 3D Chamfer Distances (\downarrow) and F-scores (\uparrow) at distance thresholds of 1%, 2%, and 5% averaged across all frames.

Method	CD	F@1%	F@2%	F@5%
17 anchors	11.76	25.51	49.09	81.52
35 anchors	9.22	29.81	54.20	86.91
70 anchors	11.72	25.74	48.72	80.69
with DensePose feat.	8.34	33.89	60.16	88.66
w/o feat.	9.24	27.06	52.88	87.37
w/o kine. chain opt.	9.80	31.86	55.34	84.69

E. Anchors and Associations

The anchors are optimized to model the modes of the deformations, and the skinning weights are optimized w.r.t.

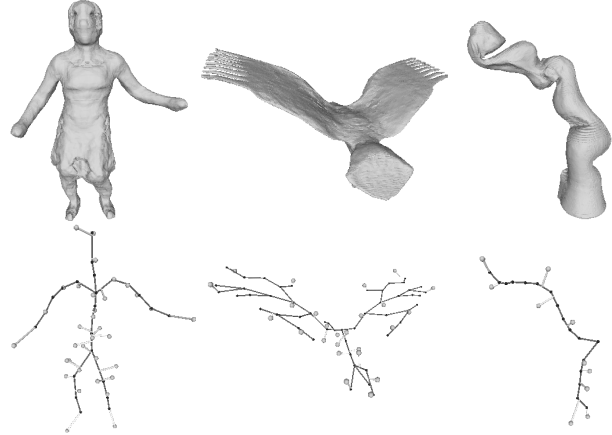


Figure 11. **Anchors and associations.** Anchors (white dots around the kinematic chain) and their corresponding associations (line segments between anchors and kinematic chain links) for the synthetic and AMA datasets.

anchors. By enforcing deterministic associations between the anchors and the kinematic chain, the transformations of anchors can be directly inferred from user-defined kinematic chain transformations. Therefore, we can make good use of optimized anchors and their corresponding skinning weights, instead of optimizing the skinning weights w.r.t. kinematic chain joints from scratch. In Figure 11, we show the learned deformation anchors and their associations to the kinematic chain links for each dataset. The anchors move along with the kinematic chain links to keep the association parameters constant at all times to enable re-positings directly driven by the kinematic chain.

Note that our proposed kinematic chain driven deformations formulation is capable of modeling twists of a kinematic chain link about its axial direction even though our kinematic chain is defined as connected line segments. An additional rotation matrix that represents the rotation about the axial axis of the link is pre-multiplied to rotate the anchor around its associated link for twist effects.

F. Necessity of Two-Stage Optimization

Since we do not have access to any shape prior or template, building both object’s shape and kinematic chain from scratch simultaneously is a highly ill-conditioned optimization problem given only a collection of monocular videos. Our proposed two-stage optimization simplifies this problem by taking advantage of the techniques in the existing template-free surface reconstruction method [75] to extract an initial estimate of the 3D shape, which allows us to use RigNet [71] to extract an initial estimate of the kinematic chain. The two-stage scheme significantly reduces the search space for the underlying rigid structure of the object and makes the optimization much more stable.

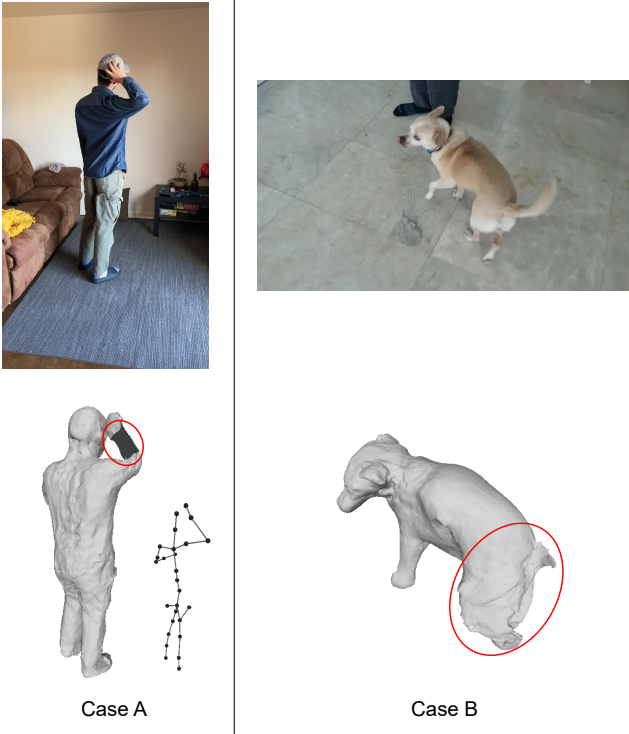


Figure 12. **Failure case examples.** In Case A, the mesh of the person’s right arm collapses despite the kinematic chain being able to represent the correct pose. In Case B, the dog’s tail and back legs are missing due to bad reconstruction.

As RigNet [71] operates on 3D mesh, it’s possible to directly apply RigNet to the initial optimization stage results to obtain an animtable model driven by RigNet’s predicted skinning weights. However, it would yield three problems. The first problem is that the skinning weights predicted by RigNet are based on a static mesh without any knowledge from the videos (dynamic scenes). These weights are inaccurate and not optimized for the deformation modes that occurred in the videos. The second problem is that RigNet suffers from significantly large memory consumption and runtime. The author of RigNet recommends users to only predict skinning weights for kinematic chain joints on low-resolution mesh ($\leq 5,000$ vertices) on its code repository [71] to prevent running out of hardware’s memory and extremely long runtime. Our approach does not have this bottleneck as we do not rely on the skinning weights predicted by RigNet. We build the associations between anchors and the kinematic chain to animate object meshes regardless of the resolution. The third problem is that directly using the RigNet initialization combined with the anchors would lead to poor results. As reported in the Tables in the main paper and the supplementary material, we achieve consistent improvements on all datasets with our proposed kinematic chain optimization stage.

G. Limitations

Input dependencies. In our pipeline, we leverage an off-the-shelf model [72] to extract optical flow from monocular videos, and use ground-truth foreground masks from the datasets for optimizations. For the Eagle and iiwa datasets, we use the ground-truth root poses, while for the AMA dataset, we leverage a pre-trained PoseNet [75] for root pose initializations and jointly optimize them during training. Therefore, our pipeline’s performance is affected by the quality of these inputs, and it requires a generic root pose estimator to enable our pipeline to work on any object of interest because PoseNet [75] is only suitable for humans and quadruped animals.

Kinematic chain initialization. We use a category-agnostic skeleton estimator, RigNet [71], on the initial 3D mesh estimate of the object to obtain kinematic chain initialization. However, the underlying structures and articulation modes vary across different object categories, requiring users to tune a RigNet’s test-time hyper-parameter that controls how dense the joints are distributed within the object’s mesh. We encourage the users to run RigNet multiple times on the same initial estimate of the mesh with different test-time hyper-parameters, and select the kinematic chain that agrees with the object’s underlying structure and articulation modes the most before optimizing it.

H. Potential Negative Impact

As our approach effectively builds a 3D animatable model of an object using monocular videos, it can be potentially used for malicious purposes, such as generating fake videos or other formats of illegal content. Since monocular videos are often easy to acquire in the real world, it is important to ensure our method is used with prior consent.

I. Failure Cases

We show some example failure cases in Figure 12. In Case A, the mesh of the person’s right arm collapses at the pose shown in the figure despite the correct kinematic chain pose. This is because the deformation anchors are not optimized well to handle the desired pose. In the cases of training with fast-moving objects and large self-occlusions in the videos, the reconstruction results sometimes suffer from significant degradation, as shown in Case B. The dog’s tail is moving very fast in the training video and often causes self-occlusions at the back of the dog, making the overall optimization much less stable and sometimes converge to poor results. In such cases, we recommend tuning the hyper-parameters in training or gathering more videos with 360 degrees captures of the object to make the training more stable. In general, objects that are moving slowly with 360 degrees captures are usually easier to train.

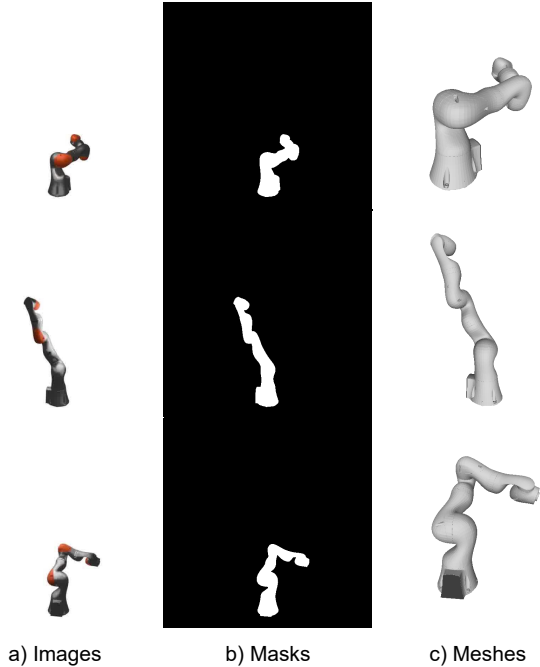


Figure 13. **Examples of our iiwa dataset.** We show examples of a) RGB images, b) corresponding foreground masks, and c) corresponding ground truth meshes for three randomly picked frames in the dataset. Note that the meshes are re-scaled for better visualization in the figure.

J. Additional Details on iiwa Dataset

We created a new dataset that contains four video sequences of an animated KUKA LBR iiwa robot arm. The synthetic model of the robot arm is obtained from BlenderKit. We animate the robot arm in Blender and render the videos of the robot arm’s motion. Each video contains 400 frames in total. Note that for each frame in the video sequences, the camera pose, foreground mask, and 3D mesh of the robot arm are also included in the dataset for quantitative evaluation purposes. We show examples of our dataset in Figure 13.

K. Notations

In Table 7, we list the important variables used in this paper, along with their state space and descriptions.

Table 7. **Notations.** A list of the important variables used in the paper.

Symbol	State space	Description
\mathcal{F}_C	MLP	MLP for color in canonical space
\mathcal{F}_S	MLP	MLP for SDF in canonical space
\mathcal{F}_ϕ	MLP	MLP for canonical feature in canonical space
\mathcal{F}_A	MLP	MLP for unconstrained transformations of anchors
\mathbf{p}_i	\mathbb{R}^3	i -th kinematic chain joint in canonical space
$\hat{\mathbf{p}}_i$	\mathbb{R}^3	i -th unconstrained kinematic chain joint
$\tilde{\mathbf{p}}_i$	\mathbb{R}^3	i -th revised kinematic chain joint
\mathbf{a}_i	\mathbb{R}^3	i -th deformation anchor in canonical space
ℓ_{jk}	link	Kinematic chain link between \mathbf{p}_j and \mathbf{p}_k
\mathbf{C}^t	$SE(3)$	Object root pose w.r.t the canonical root pose at time t
\mathbf{T}_i^t	$SE(3)$	Transformation of anchor \mathbf{a}_i at time t
$\hat{\mathbf{T}}_i^t$	$SE(3)$	Unconstrained transformation of anchor \mathbf{a}_i at time t
$\tilde{\mathbf{T}}_i^t$	$SE(3)$	Revised transformation of anchor \mathbf{a}_i at time t
\mathbf{H}_i^t	$SE(3)$	Forward kinematics of kinematic chain joint \mathbf{p}_i at time t
\mathbf{w}_{a_i}	\mathbb{R}	Forward skinning weight w.r.t anchor \mathbf{a}_i
$\mathbf{w}_{a_i}^t$	\mathbb{R}	Backward skinning weight w.r.t anchor \mathbf{a}_i at time t
\mathbf{R}	\mathbb{R}^{n_j-1}	Unclipped additive residuals for kinematic chain
$\tilde{\mathbf{R}}$	\mathbb{R}^{n_j-1}	Clipped additive residuals for kinematic chain
\mathbf{x}_I^t	\mathbb{R}^2	Pixel of interest at time t
ψ_a^t	\mathbb{R}^{128}	Learnable latent code for anchor transformations at time t
ψ_l	\mathbb{R}^{64}	Learnable latent code illumination conditions
π^t	$\mathbb{R}^{3 \times 4}$	Camera projection model at time t

References

- [1] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep vit features as dense visual descriptors. *arXiv preprint arXiv:2112.05814*, 2021. 2, 4
- [2] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *ICCV*, 2021. 1
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 2, 3, 4, 6, 8, 10
- [4] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *ICCV*, 2019. 2
- [5] Xu Chen, Yufeng Zheng, Michael J Black, Otmar Hilliges, and Andreas Geiger. Snarf: Differentiable forward skinning for animating non-rigid neural implicit shapes. In *ICCV*, 2021. 2
- [6] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *ICCV*, 2019. 1
- [7] Boyang Deng, John P Lewis, Timothy Jeruzalski, Gerard Pons-Moll, Geoffrey Hinton, Mohammad Norouzi, and Andrea Tagliasacchi. Nasa neural articulated shape approximation. In *ECCV*, 2020. 1, 2
- [8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 4
- [9] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, 2017. 7
- [10] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *CVPR*, 2022. 1
- [11] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *CVPR*, 2018. 7
- [12] Tong He, Yuanlu Xu, Shunsuke Saito, Stefano Soatto, and Tony Tung. Arch++: Animation-ready clothed human reconstruction revisited. In *ICCV*, 2021. 2
- [13] Zeng Huang, Yuanlu Xu, Christoph Lassner, Hao Li, and Tony Tung. Arch: Animatable reconstruction of clothed humans. In *CVPR*, 2020. 2
- [14] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *ICCV*, 2021. 1
- [15] Angjoo Kanazawa, Shubham Tulsiani, Alexei A Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *ECCV*, 2018. 2
- [16] Alex Kendall, Hayk Martirosyan, Saumitro Dasgupta, Peter Henry, Ryan Kennedy, Abraham Bachrach, and Adam Bry. End-to-end learning of geometry and context for deep stereo regression. In *ICCV*, 2017. 5, 10
- [17] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*, 2014. 11
- [18] Alexander Kirillov, Yuxin Wu, Kaiming He, and Ross Girshick. Pointrend: Image segmentation as rendering. In *CVPR*, 2020. 11
- [19] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 2
- [20] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 2
- [21] Chen Kong and Simon Lucey. Deep non-rigid structure from motion. In *ICCV*, 2019. 1, 2
- [22] Nilesh Kulkarni, Abhinav Gupta, David F Fouhey, and Shubham Tulsiani. Articulation-aware canonical surface mapping. In *CVPR*, 2020. 10
- [23] Suryansh Kumar. Non-rigid structure from motion: Prior-free factorization method revisited. In *WACV*, 2020. 1, 2
- [24] Jogendra Nath Kundu, Mugalodi Rakesh, Varun Jampani, Rahul Mysore Venkatesh, and R Venkatesh Babu. Appearance consensus driven self-supervised human mesh recovery. In *ECCV*, 2020. 1, 2
- [25] Ruilong Li, Julian Tanke, Minh Vo, Michael Zollhofer, Jürgen Gall, Angjoo Kanazawa, and Christoph Lassner. Tava: Template-free animatable volumetric actors. In *ECCV*, 2022. 2
- [26] Xueting Li, Sifei Liu, Shalini De Mello, Kihwan Kim, Xiaolong Wang, Ming-Hsuan Yang, and Jan Kautz. Online adaptation for consistent mesh reconstruction in the wild. In *NeurIPS*, 2020. 2
- [27] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 1, 2, 10
- [28] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *ICCV*, 2021. 1
- [29] Lingjie Liu, Marc Habermann, Viktor Rudnev, Kripasindhu Sarkar, Jiatao Gu, and Christian Theobalt. Neural actor: Neural free-view synthesis of human actors with pose control. In *SIGGRAPH*, 2021. 1, 2
- [30] Wen Liu, Zhixin Piao, Jie Min, Wenhan Luo, Lin Ma, and Shenghua Gao. Liquid warping gan: A unified framework for human motion imitation, appearance transfer and novel view synthesis. In *ICCV*, 2019. 2
- [31] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *SIGGRAPH*, 2015. 1, 2
- [32] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint*, 2016. 11
- [33] Diogo C Luvizon, Hedi Tabia, and David Picard. Human pose regression by combining indirect part detection and contextual information. In *Computers & Graphics*, 2019. 5, 10
- [34] Liqian Ma, Xu Jia, Qianru Sun, Bernt Schiele, Tinne Tuytelaars, and Luc Van Gool. Pose guided person image generation. In *NeurIPS*, 2017. 2

- [35] 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 *CVPR*, 2021. 1, 9
- [36] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *ICCV*, 2021. 1
- [37] 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. 1, 2, 3, 5, 9, 10
- [38] Natalia Neverova, David Novotny, Marc Szafraniec, Vasil Khalidov, Patrick Labatut, and Andrea Vedaldi. Continuous surface embeddings. In *NeurIPS*, 2020. 4, 7, 8
- [39] Atsuhiko Noguchi, Umar Iqbal, Jonathan Tremblay, Tatsuya Harada, and Orazio Gallo. Watch it move: Unsupervised discovery of 3d joints for re-posing of articulated objects. In *CVPR*, 2022. 1, 2
- [40] Atsuhiko Noguchi, Xiao Sun, Stephen Lin, and Tatsuya Harada. Neural articulated radiance field. In *CVPR*, 2021. 1, 2
- [41] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *ICCV*, 2021. 1, 2, 9
- [42] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *CVPR*, 2019. 1, 2
- [43] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *ICCV*, 2021. 1, 2
- [44] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *CVPR*, 2021. 1, 2
- [45] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 1, 2
- [46] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. In *SIGGRAPH*, 2022. 1
- [47] Shunsuke Saito, Jinlong Yang, Qianli Ma, and Michael J Black. Scanimate: Weakly supervised learning of skinned clothed avatar networks. In *CVPR*, 2021. 2
- [48] Peter Sand and Seth Teller. Particle video: Long-range motion estimation using point trajectories. In *IJCV*, 2008. 4
- [49] Vikramjit Sidhu, Edgar Tretschk, Vladislav Golyanik, Antonio Agudo, and Christian Theobalt. Neural dense non-rigid structure from motion with latent space constraints. In *ECCV*, 2020. 1, 2
- [50] Leslie N Smith. A disciplined approach to neural network hyper-parameters: Part 1—learning rate, batch size, momentum, and weight decay. *arXiv preprint*, 2018. 11
- [51] Shih-Yang Su, Timur Bagautdinov, and Helge Rhodin. Danbo: Disentangled articulated neural body representations via graph neural networks. In *ECCV*, 2022. 2
- [52] Shih-Yang Su, Frank Yu, Michael Zollhöfer, and Helge Rhodin. A-nerf: Articulated neural radiance fields for learning human shape, appearance, and pose. In *NeurIPS*, 2021. 2
- [53] Xingyuan Sun, Jiajun Wu, Xiuming Zhang, Zhoutong Zhang, Chengkai Zhang, Tianfan Xue, Joshua B Tenenbaum, and William T Freeman. Pix3d: Dataset and methods for single-image 3d shape modeling. In *CVPR*, 2018. 7
- [54] Narayanan Sundaram, Thomas Brox, and Kurt Keutzer. Dense point trajectories by gpu-accelerated large displacement optical flow. In *ECCV*, 2010. 4
- [55] Jiapeng Tang, Lev Markhasin, Bi Wang, Justus Thies, and Matthias Nießner. Neural shape deformation priors. In *NeurIPS*, 2022. 2
- [56] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *CVPR*, 2019. 7
- [57] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *ICCV*, 2021. 1, 2
- [58] Narek Tumanyan, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Splicing vit features for semantic appearance transfer. In *CVPR*, 2022. 2, 4
- [59] Daniel Vlastic, Ilya Baran, Wojciech Matusik, and Jovan Popović. Articulated mesh animation from multi-view silhouettes. In *SIGGRAPH*, 2008. 6
- [60] Minh Phuoc Vo, Yaser A Sheikh, and Srinivasa G Narasimhan. Spatiotemporal bundle adjustment for dynamic 3d human reconstruction in the wild. In *TPAMI*, 2020. 1, 2
- [61] Chaoyang Wang, Ben Eckart, Simon Lucey, and Orazio Gallo. Neural trajectory fields for dynamic novel view synthesis. *arXiv preprint arXiv:2105.05994*, 2021. 1, 2
- [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. In *NeurIPS*, 2021. 3, 9
- [63] Shaofei Wang, Katja Schwarz, Andreas Geiger, and Siyu Tang. Arah: animatable volume rendering of articulated human sdf. In *ECCV*, 2022. 1, 2
- [64] Fangyin Wei, Rohan Chabra, Lingni Ma, Christoph Lassner, Michael Zollhöfer, Szymon Rusinkiewicz, Chris Sweeney, Richard Newcombe, and Mira Slavcheva. Self-supervised neural articulated shape and appearance models. In *CVPR*, 2022. 1
- [65] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *CVPR*, 2022. 1, 2
- [66] Shangzhe Wu, Ruining Li, Tomas Jakab, Christian Rupprecht, and Andrea Vedaldi. Magicpony: Learning articulated 3d animals in the wild. *arXiv preprint arXiv:2211.12497*, 2022. 2

- [67] Yuefan Wu, Zeyuan Chen, Shaowei Liu, Zhongzheng Ren, and Shenlong Wang. Casa: Category-agnostic skeletal animal reconstruction. In *NeurIPS*, 2022. 2
- [68] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body, and hands in the wild. In *CVPR*, 2019. 1, 2
- [69] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Ghum & ghuml: Generative 3d human shape and articulated pose models. In *CVPR*, 2020. 1, 2
- [70] Yuanlu Xu, Song-Chun Zhu, and Tony Tung. Denserac: Joint 3d pose and shape estimation by dense render-and-compare. In *ICCV*, 2019. 2
- [71] Zhan Xu, Yang Zhou, Evangelos Kalogerakis, Chris Landreth, and Karan Singh. Rignet: Neural rigging for articulated characters. In *SIGGRAPH*, 2020. 3, 5, 8, 12, 13
- [72] Gengshan Yang and Deva Ramanan. Volumetric correspondence networks for optical flow. In *NeurIPS*, 2019. 6, 11, 13
- [73] 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. 1, 2, 3, 5, 10
- [74] 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. 1, 2, 3, 4, 5, 6, 7, 10
- [75] Gengshan Yang, Minh Vo, Natalia Neverova, Deva Ramanan, Andrea Vedaldi, and Hanbyul Joo. Banmo: Building animatable 3d neural models from many casual videos. In *CVPR*, 2022. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13
- [76] Chun-Han Yao, Wei-Chih Hung, Yuanzhen Li, Michael Rubinstein, Ming-Hsuan Yang, and Varun Jampani. Lassie: Learning articulated shapes from sparse image ensemble via 3d part discovery. In *NeurIPS*, 2022. 1, 2
- [77] Lior Yariv, Jiatao Gu, Yoni Kasten, and Yaron Lipman. Volume rendering of neural implicit surfaces. In *NeurIPS*, 2021. 3, 9
- [78] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *NeurIPS*, 2020. 5, 10
- [79] Yufei Ye, Shubham Tulsiani, and Abhinav Gupta. Shelf-supervised mesh prediction in the wild. In *CVPR*, 2021. 2
- [80] Jae Shin Yoon, Lingjie Liu, Vladislav Golyanik, Kripasindhu Sarkar, Hyun Soo Park, and Christian Theobalt. Pose-guided human animation from a single image in the wild. In *CVPR*, 2021. 2
- [81] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *ICCV*, 2021. 1
- [82] Haitian Zeng, Yuchao Dai, Xin Yu, Xiaohan Wang, and Yi Yang. Pr-rrn: pairwise-regularized residual-recursive networks for non-rigid structure-from-motion. In *ICCV*, 2021. 2
- [83] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. 1
- [84] Zhen Zhu, Tengeng Huang, Baoguang Shi, Miao Yu, Bofei Wang, and Xiang Bai. Progressive pose attention transfer for person image generation. In *CVPR*, 2019. 2
- [85] Silvia Zuffi, Angjoo Kanazawa, and Michael J Black. Lions and tigers and bears: Capturing non-rigid, 3d, articulated shape from images. In *CVPR*, 2018. 1
- [86] Silvia Zuffi, Angjoo Kanazawa, David W Jacobs, and Michael J Black. 3d menagerie: Modeling the 3d shape and pose of animals. In *CVPR*, 2017. 1