# Online Adaptation for Consistent Mesh Reconstruction in the Wild

**Xueting Li**[1], **Sifei Liu**[2], **Shalini De Mello**[2], **Kihwan Kim**[2],
**Xiaolong Wang**[3], **Ming-Hsuan Yang**[1], **Jan Kautz**[2]
[1]University of California, Merced,  [2]NVIDIA,  [3]University of California, San Diego

## Abstract

This paper presents an algorithm to reconstruct temporally consistent 3D meshes of deformable object instances from videos in the wild. Without requiring annotations of 3D mesh, 2D keypoints, or camera pose for each video frame, we pose video-based reconstruction as a self-supervised online adaptation problem applied to any incoming test video. We first learn a category-specific 3D reconstruction model from a collection of single-view images of the same category that jointly predicts the shape, texture, and camera pose of an image. Then, at inference time, we adapt the model to a test video over time using self-supervised regularization terms that exploit temporal consistency of an object instance to enforce that all reconstructed meshes share a common texture map, a base shape, as well as parts. We demonstrate that our algorithm recovers temporally consistent and reliable 3D structures from videos of non-rigid objects including those of animals captured in the wild – an extremely challenging task rarely addressed before. Codes and other resources will be maintained at https://sites.google.com/nvidia.com/vmr-2020.

## 1 Introduction

When we humans try to understand the object shown in Fig. 1(a), we instantly recognize it as a "duck". We also instantly perceive and imagine its shape in the 3D world, its viewpoint, and its appearance from other views. Furthermore, when we see it in a video, its 3D structure and deformation become even more apparent to us. Our ability to perceive the 3D structure of objects contributes vitally to our rich understanding of them.

While 3D perception is easy for humans, 3D reconstruction of deformable objects remains a very challenging problem in computer vision, especially for objects in the wild. For learning-based algorithms, the key bottleneck is the lack of supervision. It is extremely challenging to collect 3D annotations such as 3D shape and camera pose [4, 18]. Consequently, existing research mostly focuses on limited domains (e.g., rigid objects [26], human bodies [15, 53] and faces [50]) for which 3D annotations can be captured in constrained environments. However, these approaches do not generalize well to non-rigid objects captured in naturalistic environments (e.g., animals). In non-rigid structure from motion methods [2, 30], the 3D structure can be partially recovered from correspondences between multiple viewpoints, which are also hard to label. Due to constrained environments and limited annotations, it is nearly impossible to generalize these approaches to the 3D reconstruction of non-rigid objects (e.g., animals) from images and videos captured in the wild.

Instead of relying on 3D supervision, weakly supervised or self-supervised approaches have been proposed for 3D mesh reconstruction. They use annotated 2D object keypoints [14], category-level templates [23, 22] or silhouettes [24]. However, to scale up learning with 2D annotations to hundreds of thousands of images is still non-trivial. This limits the generalization ability of current models to new domains. For example, a 3D reconstruction model trained on single-view images, e.g., [14], produces unstable and erratic predictions for video data. This is unsurprising, due to perturbations

---

Preprint. Under review.

(a) Learning temporally invariant shape, texture and parts

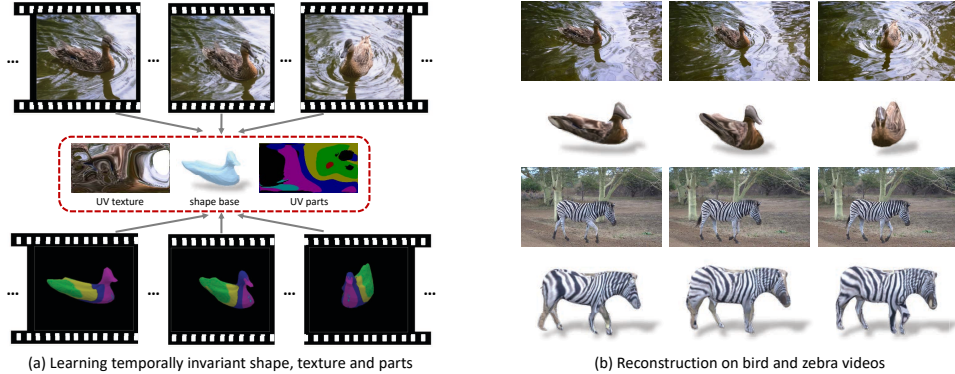(b) Reconstruction on bird and zebra videos

Figure 1: By utilizing the consistency of texture, shape and object parts correspondences in videos (red box) as self-supervision signals in (a), we learn a model that reconstructs temporally consistent meshes of deformable object instances in videos in (b).

over time. However, the temporal signal in videos should provide us an advantage instead of a disadvantage, as recently shown on the task of optimizing a 3D *rigid* object mesh w.r.t. a particular video [56, 26]. The question is, can we also take advantage of the redundancy in temporal sequences as a form of self-supervision in order to improve the reconstruction of dynamic non-rigid objects?

In this work, we address this problem with two important innovations. First, we strike a balance between model generalization and specialization. That is, we train an image-based network on a set of images, while at test time we adapt it online to an input video of a particular instance. Test-time training [40] is non-trivial since no labels are provided for the video. The key is to introduce self-supervised objectives that can continuously improve the model. To do so, we exploit the UV texture space, which provides a parameterization that is invariant to object deformation. We encourage the sampled texture, as well as a group of object parts, to be consistent among all the individual frames in the UV space, as shown in Fig. 1(a). Using this constraint of temporal consistency, the recovered shape and camera pose are stabilized considerably and are adapted to the current video.

One bottleneck of existing image-based 3D mesh reconstruction methods [14, 24] is that the predicted shapes are assumed to be symmetric. This assumption does not hold for most non-rigid animals, e.g., birds tilting their heads, or walking horses, etc. Our second innovation is to remove this assumption and to allow the reconstructed meshes to fit more complex, non-rigid poses via an as-rigid-as-possible (ARAP) constraint. As another constraint that does not require any labels, we enforce ARAP during test-time training as well, to substantially improve shape prediction. We use two image-based 3D reconstruction models for training (i) a weakly supervised one (i.e., with object silhouettes and 2D keypoints provided), and (ii) a self-supervised one where only object silhouettes are available. The image-based models are then adapted to in-the-wild bird and zebra videos collected from the internet. We show that for both models, our innovations lead to an effective and robust approach to deformable, dynamic 3D object reconstruction of non-rigid objects captured in the wild.

## 2 Related Work

**3D object reconstruction from images.** A triangular mesh has long been used for object reconstruction [14, 18, 27, 16, 47, 31, 48]. It is a memory-efficient representation with vertices and faces, and is amenable to differentiable rendering techniques [18, 27]. The task of 3D reconstruction entails the simultaneous recovery of the 3D shape, texture, and camera pose of objects from 2D images. It is highly ill-posed due to the inherent ambiguity of correctly estimating both the shape and camera pose together. A major trend of recent works is to gradually reduce supervision from 3D vertices [4, 48, 47], shading [10], or multi-view images [52, 18, 49, 36, 26] and move towards weakly supervised methods that instead use 2D semantic keypoints [14], or a category-level 3D template [23]. This progress makes the reconstruction of objects, e.g., birds, captured in the wild possible. More recently, self-supervised methods [24, 50, 17] have been developed to further remove the need for annotations. Our method exploits different levels of supervision: weak supervision (i.e., using 2D semantic keypoints) and self-supervision to learn an image-based 3D reconstruction network from a collection of images of a category (Sec. 3.1).

**Non-rigid structure from motion (NR-SFM).** NR-SFM aims to recover the pose and 3D structure of a non-rigid object, or object deforming non-rigidly over time, solely from 2D landmarks without 3D supervision [2]. It is a highly ill-posed problem and needs to be regularized by additional shape priors [2, 57]. Recently, deep networks [21, 30] have been developed that serve as more powerful priors than the traditional approaches. However, obtaining reliable landmarks or correspondences for videos is still a bottleneck. Our method bears resemblances to deep NR-SFM [30], which jointly predicts camera pose and shape deformation. Differently from them, we reconstruct dense meshes instead of sparse keypoints, without requiring labeled correspondences from videos.

**3D object reconstruction from videos.** Existing video-based object reconstruction methods mostly focus on specific domains, e.g., videos of faces [6, 41] or human bodies [42, 1, 5, 15, 53], where dense labelling is possible [43]. To augment video labels, [15] formulates dynamic human mesh reconstruction as an omni-supervision task, where a combination of labeled images and videos with pseudo-ground truth are used for training. For human video-based 3D pose estimation, [33] introduces semi-supervised learning to leverage unlabeled videos with a self-supervised component. Dealing with specific application domains, all the aforementioned works rely on predefined shape priors, such as a parametric body model (e.g., SMPL [28]) or a morphable face model. While our work also exploits unlabeled videos, we do not assume any predefined shape prior, which, practically, is hard to obtain for the majority of objects captured in the wild.

**Optimization-based methods.** Optimization-based methods have also been extensively explored for scene or object reconstruction from videos. Several works [55, 46, 37, 34] first obtain a single-view 3D reconstruction and then optimize the mesh and skeletal parameters. Another line of methods is developed to optimize the weights of deep models instead, to render more robust results for a video of a particular instance [42, 26, 29, 58]. Our method falls into this category. While [42] enforces consistency between observed 2D and a re-projection from 3D, [26, 29] take a further step and encourage consistency between frames via a network that inherently encodes an entire video into an invariant representation. In this work, instead of limiting to rigid objects as [26], or depth estimation as [29], we recover dynamic meshes from videos captured in the wild – a much more challenging problem that is rarely explored.

## 3 Approach

Our goal is to recover coherent sequences of mesh shapes, texture maps and camera poses from unlabeled videos, with a two-stage learning approach: (i) first, we learn a 3D mesh reconstruction model on a collection of single-view images of a category, described in Sec. 3.1; (ii) at inference time, we adapt the model to fit the sequence via temporal consistency constraints, as described in Sec. 3.2. We focus on the weakly-supervised setting in Sec. 3.1 and 3.2, where both silhouettes and keypoints are annotated in the image dataset. We then describe how to generalize the approach to a self-supervised setting, where only silhouettes are available in the image dataset in Sec. 3.3.

**Notations.** We represent a textured mesh with $|V|$ vertices ($V \in \mathbb{R}^{|V| \times 3}$), $|F|$ faces ($F \in \mathbb{R}^{|F| \times 3}$) and a UV texture image ($I_{\mathrm{uv}} \in \mathbb{R}^{H_{\mathrm{uv}} \times W_{\mathrm{uv}} \times 3}$) of height $H_{\mathrm{uv}}$ and width $W_{\mathrm{uv}}$. Similarly to [14], we use a weak perspective transformation to represent the camera pose $\theta \in \mathbb{R}^7$ of an input image. We denote $\mathcal{R}(\cdot)$ as a general projection, which can represent (i) a differentiable renderer [27, 18] to render a mesh to a 2D silhouette as $\mathcal{R}(V, \theta)$, or a textured mesh to an RGB image as $\mathcal{R}(V, \theta, I_{\mathrm{uv}})$ (we omit mesh faces $F$ for conciseness); (ii) or a projection of a 3D point $v$ to the image space as $\mathcal{R}(v, \theta)$. The Soft Rasterizer [27] is used as the differentiable renderer in this work.

### 3.1 Single-view Mesh Reconstruction

In the first stage, we train a network with a collection of category-specific images that jointly estimates the shape, texture, and camera pose of an input image. Similarly to [14], we predict a texture flow $I_{\mathrm{flow}} \in \mathbb{R}^{H_{\mathrm{uv}} \times W_{\mathrm{uv}} \times 2}$ that maps pixels from the input image to the UV space. A predefined UV mapping function $\Phi$ [11, 14] is then used to map these pixels from the UV space to the mesh surface. With a differentiable renderer [27], we train the network with supervision from object silhouettes, texture, and the Laplacian objectives as in [14, 24]. More details of learning texture and camera pose can be found in [14]. An overview of our reconstruction framework is shown in Fig. 2(a).
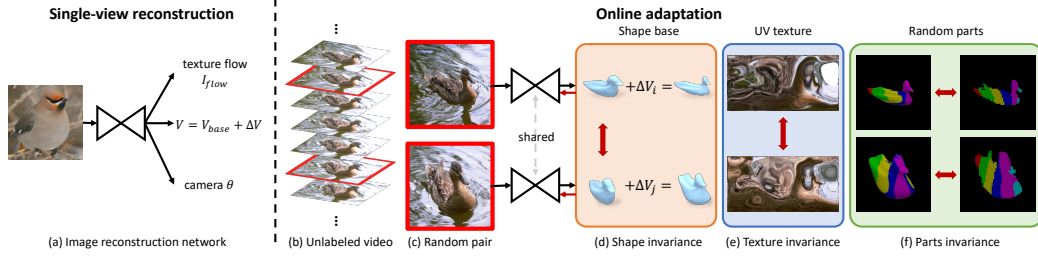
Figure 2: Overview. We show the single-view image reconstruction network on the left and the test-time training procedure to adapt it to a video on the right. Bold red arrows indicate invariance constraints in Sec. 3.2.

**Recovering asymmetric shapes.** We propose a novel shape reconstruction module as shown in Fig. 2(a). The key idea is to remove the symmetry requirement of object shapes, which is employed by many prior works [14, 24]. This is particularly important for recovering dynamic meshes in sequences, e.g., when a bird rotates its head as shown in Fig. 4, its mesh is no longer mirror-symmetric. Prior works [14, 24] model object shape by predicting vertex offsets from a jointly learned 3D template. Simply removing the symmetry assumption for the predicted vertex offsets leads to excessive freedom in shape deformation, e.g., see Fig. 4(f). To resolve this, we learn a group of $N_b$ shape bases $\{V_i\}_{i=1}^{N_b}$, and replace the template by a weighted combination of them, denoted as the base shape $V_{\text{base}}$. Compared to a single mesh template, the base shape $V_{\text{base}}$ is more powerful in capturing the object's identity and saves the model from predicting large motion deformation, e.g., of deforming a standing bird template to a flying bird. The full shape reconstruction can be obtained by:

$$V = V_{\text{base}} + \Delta V, \qquad V_{\text{base}} = \sum_{i=1}^{N_b} \beta_i V_i, \tag{1}$$

where the $\Delta V$ encodes the object's asymmetric non-rigid motion and $\{\beta_i\}_{i=1}^{N_b}$ are learned coefficients.

The computation of our shape bases is inspired by parametric models [28, 59, 60], where the basis shapes are extracted from an existing mesh dataset [38] or toy scans [60]. However, we make our model completely free of 3D supervision and obtain the bases by applying K-Means clustering to all meshes reconstructed by CMR [14]. We use each cluster center as a basis shape in our model.

**Keypoint re-projection.** In the weakly-supervised setting, the 2D keypoints are provided that semantically associate different instances. When projected onto the mesh surface, the same semantic keypoint (e.g., the tail keypoint in the orange circle in Fig. 3(a)) for different object instances should be matched to the same face on the mesh (the tail keypoint in the orange circle in Fig. 3(d)). To model the mapping between the 3D mesh surface and the 2D keypoints, prior work [14] learns an affinity matrix that describes the probability of each 2D keypoint mapping to each vertex on the mesh. The affinity matrix is shared among all instances and is independent of individual shape variations. However,
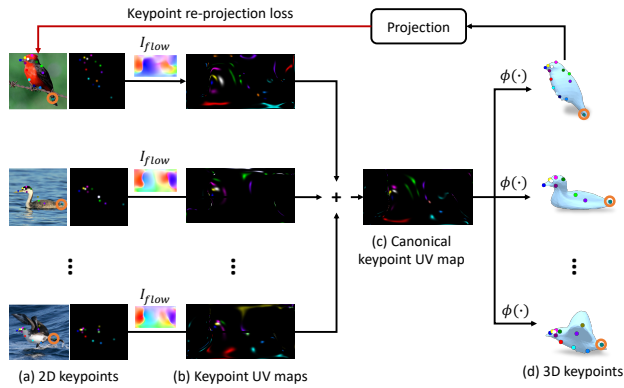


Figure 3: 3D canonical keypoints computation: (a) annotated 2D keypoints and their location heatmaps; (b) keypoint heatmaps mapped to the UV space using learned texture flows; (c) aggregated canonical keypoint heatmaps in the UV space; (d) canonical keypoints on different instance mesh surface. $\Phi(\cdot)$ is the UV mapping function discussed in Sec. 3.1.

this approach is sub-optimal because: (i) Mesh vertices are a subset of discrete points on a continuous mesh surface and so their weighted combination defined by the affinity matrix may not lie on it, leading to inaccurate mappings of 2D keypoints. (ii) The mapping from the image space to the mesh surface described by the affinity matrix, in our case, however, is already modeled by the texture flow. Hence, it is potentially redundant to learn both of them independently.

4

In this work, we re-utilize texture flow to map 2D keypoints from each image to the mesh surface. We first map each 2D keypoint to the UV space that is independent of shape deformation (Fig. 3(b)). Ideally, each semantic keypoint from different instances should map to the same point in the UV space as discussed above. In practice, this does not hold due to inaccurate texture flow prediction. To accurately map each keypoint to the UV space, we compute a canonical keypoint UV map as shown in Fig. 3(c) by: (i) mapping the keypoint heat map in Fig. 3(a) for each instance to the UV space via its predicted texture flow, and (ii) aggregating these keypoint UV maps in Fig. 3(b) across all instances to eliminate outliers caused by incorrect texture flow prediction.

We further utilize the pre-defined UV mapping function $\Phi$ discussed above to map each semantic keypoint from the UV space to the mesh surface. Given the 3D correspondence (denoted as $K_{3D}^i$) of each 2D semantic keypoint $K_{2D}^i$, the keypoint re-projection loss enforces the projection of the former to be consistent with the latter by:

$$L_{kp} = \frac{1}{N_k} \sum_{i=1}^{N_k} \left\| \mathcal{R}(K_{3D}^i, \theta) - K_{2D}^i \right\|, \tag{2}$$

where $N_k$ is the number of keypoints.

**As-rigid-as-possible (ARAP) constraint.** Without any pose-related regularization, the predicted motion deformation $\Delta V$ often leads to erroneous random deformations and spikes as shown in Fig. 4(f), which do not faithfully describe the motion of a non-rigid object. Therefore, we introduce an as-rigid-as-possible (ARAP) constraint [39, 8] to encourage rigidity of local transformations and the preservation of the local mesh structure. Instead of solving the optimization in [39, 8], we reformulate it as an objective that ensures that the predicted shape $V$ is a locally rigid transformation from the predicted base shape $V_{\text{base}}$ by:

$$L_{\text{arap}}(V_{\text{base}}, V) = \sum_{i=1}^{|V|} \sum_{j \in \mathcal{N}(i)} w_{ij} \left\| (V^i - V^j) - R_i(V_{\text{base}}^i - V_{\text{base}}^j) \right\|, \tag{3}$$

where $\mathcal{N}(i)$ represents the neighboring vertices of a vertex $i$, $w_{ij}$ and $R_i$ are the cotangent weight and the best approximating rotation matrix, respectively, as described in [39].

### 3.2 Online Adaptation for Videos

Applying the image-based model developed in Sec. 3.1 independently to each frame of an unseen video usually results in inconsistent mesh reconstruction (Fig. 5(a)), mainly due to the domain differences in video quality, lighting conditions, etc. In this section, we propose to perform online adaptation to fit the model to individual test video, which contains a single object instance that moves over time. Inspired by the keypoint re-projection constraint described in Sec. 3.1, we resort to the UV space, where the (i) RGB texture, and (ii) object parts of an instance should be constant when mapped from 2D via the predicted texture flow, and invariant to shape deformation. By enforcing the predicted values for (i) and (ii) to be consistent in the UV space across different frames, the adapted network is regularized to generate coherent reconstructions over time. In the following, we describe how to exploit the aforementioned temporal invariances as self-supervisory signals to tune the model.

**Part correspondence constraint.** We propose a part correspondence constraint that utilizes corresponding parts of each video frame to facilitate camera pose learning. The idea bears resemblance to NR-SFM methods [30, 21], but in contrast, we do not know the ground truth correspondence between frames. Instead, we resort to an unsupervised video correspondence (UVC) method [25]. The UVC model learns an affinity matrix that captures pixel-level correspondences among video frames. It can be used to propagate any annotation (e.g, segmentation labels, keypoints, part labels, etc.), from an annotated keyframe to the other unannotated frames. In this work, we generate part correspondence within a clip: we "paint" a group of random parts on the object, e.g., the vertical stripes in Fig. 2(f), on the first frame and propagate them to the rest of the video using the UVC model.

Given the propagated part correspondences in all the frames, we map them to the UV space via the texture flow, similar to our approach for the canonical keypoint map ( Sec. 3.1). We then average all part UV maps to obtain a video-level part UV map ("UV parts" in Fig. 1(a)) for the object depicted in the video. We map the part UV map to each individually reconstructed mesh, and render it via the predicted camera pose of each frame (see Fig. 1(a), bottom). Finally, we penalize the discrepancy

between the parts being rendered back to the 2D space, and the propagated part maps, for each frame. As the propagated part maps are usually temporally smooth and continuous, this loss implicitly regularizes the network to predict coherent camera pose and shape over time. In practice, instead of minimizing the discrepancy between the rendered part map and the propagation part map of a frame, we found that it is more robust to penalize the geometric distance between the projections of vertices assigned to each part with 2D points sampled from the corresponding part as:

$$L_c = \sum_{j=1}^{N_f} \sum_{i=1}^{N_p} \frac{1}{|V_i^j|} \text{Chamfer}(\mathcal{R}(V_i^j, \theta^j), Y_i^j), \tag{4}$$

where $N_f$ is the number of frames in the video, $N_p = 6$ is the number of parts and $V_i^j$ are vertices assigned to part $i$. Here we utilize the Chamfer distance because the vertex projections $\mathcal{R}(V_i^j, \theta^j)$ do not strictly correspond one-to-one to the sampled 2D points $Y_i^j$.

**Texture invariance constraint.**    Based on the observation that object texture mapped to the UV space should be invariant to shape deformation and stay constant over time, we propose a texture invariance constraint to encourage consistent texture reconstruction from all frames. However, naively aggregating the UV texture maps from all the frames via a scheme similar to the one described for keypoints and parts, leads to a blurry video-level texture map. We instead enforce texture consistency between random pairs of frames, via a swap loss. Given two randomly sampled frames $I^i$ and $I^j$, we swap their texture maps $I_{\text{uv}}^{\mathbf{i}}$ and $I_{\text{uv}}^{\mathbf{j}}$, and combine them with the original mesh reconstructions $V^i$ and $V^j$ as:

$$L_t = \text{dist}(\mathcal{R}(V^i, \theta^i, I_{\text{uv}}^{\mathbf{j}}) \odot S^i, I^i \odot S^i) + \text{dist}(\mathcal{R}(V^j, \theta^j, I_{\text{uv}}^{\mathbf{i}}) \odot S^j, I^j \odot S^j), \tag{5}$$

where $S^i$ and $S^j$ are the silhouettes of frame $i$ and $j$, respectively and $\text{dist}(\cdot, \cdot)$ is the perceptual metric used in [54, 14, 24].

**Base shape invariance constraint.**    As discussed in Sec. 3.1, our shape model is represented by a base shape $V_{\text{base}}$ and a deformation term $\Delta V$, in which the base shape $V_{\text{base}}$ intuitively corresponds to the "identity" of the instance, e.g., a duck, or a flying bird, etc. During online adaptation, we enforce the network to predict consistent $V_{\text{base}}$ to preserve the identity, via a swapping loss function:

$$L_s = \text{niou}(\mathcal{R}(V_{\text{base}}^{\mathbf{j}} + \Delta V^i, \theta^i), S^i) + \text{niou}(\mathcal{R}(V_{\text{base}}^{\mathbf{i}} + \Delta V^j, \theta^j), S^j), \tag{6}$$

where $V_{\text{base}}^{\mathbf{i}}$ and $V_{\text{base}}^{\mathbf{j}}$ are the base shapes for frame $i$ and $j$; $\Delta V^i$ and $\Delta V^j$ are the motion deformations for frame $i$ and $j$; and $\text{niou}(\cdot, \cdot)$ denotes the negative intersection over union (IoU) objective [16, 24]. All other notations are defined in Eq. 5.

**As-rigid-as-possible (ARAP) constraint.**    Besides the consistency constraints, we keep the ARAP objective, as discussed in Sec. 3.1, during online adaptation since it also does not require any form of supervision. We found that the ARAP constraint can obviously improve the qualitative results, as visualized for the online adaptation procedure in Fig. 6.

**Online adaption.**    During inference, we fine-tune the model on a particular video with the invariance constraints discussed above, along with a silhouette and a texture objective, a Laplacian term as in [14, 24], and the ARAP constraint discussed in Sec. 3.1. The foreground masks used for the silhouette and texture objective are obtained by a segmentation model [3] trained with the ground truth foreground masks available for the image collection. More details of the objectives used for online adaptation can be found in the supplementary material.

To obtain accurate part propagation of object parts by the UVC [25] model, we employ two strategies. Firstly, we fine-tune all parameters in the reconstruction model on sliding windows instead of all video frames. Each sliding window includes $N_w = 50$ consecutive frames and the sliding stride is set to $N_s = 10$. We tune the reconstruction model for $N_t = 40$ iterations with frames in each sliding window. Secondly, instead of "painting" random parts onto the first frame and propagating them to the rest of the frames sequentially in a window, we "paint" random parts onto the *middle* frame (i.e. the $\frac{N_w}{2}$ th frame) in the window and propagate the parts backward to the first frame as well as forward to the last frame in the window. This strategy improves the propagation quality by decreasing the propagation range to half of the window size.

6

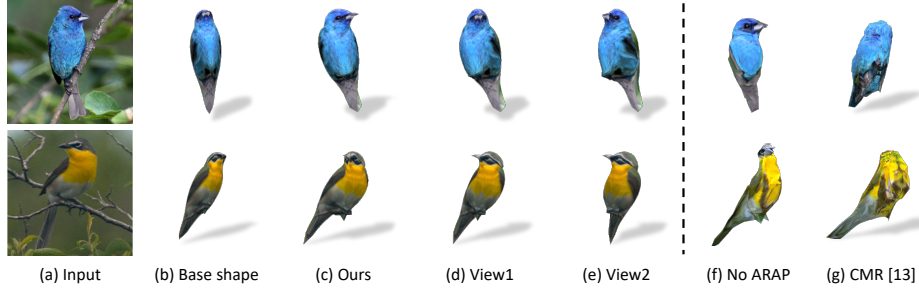|  (a) Input | (b) Base shape | (c) Ours | (d) View1 | (e) View2 | (f) No ARAP | (g) CMR [13] |

Figure 4: Mesh reconstructions from single-view images. All meshes are visualized from the predicted camera pose except for (d) and (e), where the reconstructions in (c) are visualized from two extra views. Meshes in (f) are reconstructed by a model trained without the ARAP constraint.

## 3.3 Self-supervised Setting

Our model can also be easily generalized to a self-supervised setting in which keypoints are not provided for the image datasets. In this setting, the template prior as well as camera poses in the CMR method [14] computed from the keypoints are no longer available. This self-supervised setting is trained differently from the weakly-supervised one in the following: (i) The first stage still assumes shape symmetry to ensure stability when training without keypoints. (ii) It learns a single template from scratch via the progressive training in [24]. (iii) We train this model without the keypoints re-projection and the ARAP constraints in Sec. 3.1. (iv) Without the shape bases, the base shape invariance constraint is thus removed in the online adaptation procedure. Other structure and training settings in the self-supervised model are the same as in the weakly-supervised model discussed in Sec. 3.1 and 3.2.

## 4 Experiments

We conduct experiments on animals, i.e., birds and zebras. We evaluate our contributions in two aspects: (i) the improvement of single-view mesh reconstruction, and (ii) the reconstruction of a sequence of frames via online adaptation. Due to the lack of ground truth meshes for images and videos captured in the wild, we evaluate the reconstruction results via mask and keypoint re-projection accuracy, e.g., we follow, and compare against [14] to evaluate the model trained on the image dataset. We also describe a new bird video dataset that we curate and evaluate the test-time tuned model on it in the following. We focus on evaluations on the bird category in the paper and leave evaluations on the zebra category to the supplementary.

### 4.1 Experimental Settings

**Datasets.** We first train image reconstruction models, discussed in Sec. 3.1, for the CUB bird [45] and the synthetic zebra [59] datasets. For test-time adaptation on videos, we collect a new bird video dataset for quantitative evaluation. Specifically, we collect 19 slow-motion, high-resolution bird videos from YouTtube, and 3 bird videos from the DAVIS dataset [19]. For each slow-motion video collected from the Internet, we apply a segmentation model [3] trained on the CUB bird dataset [45] to obtain its foreground segmentation for online adaptation.

**Evaluation metrics.** We evaluate the image-based model on the testing split of the CUB dataset. Note that for keypoint re-projection, instead of using the keypoint assignment matrix in [14], we apply the canonical keypoint UV map to obtain the 3D keypoints (Sec. 3.1). For the video dataset, we annotate frame-level object masks and keypoints via a semi-automatic procedure. We train a segmentation model and a keypoint detector [7] on the CUB dataset. Then, we manually adjust and filter out inaccurate predictions to ensure the correctness of the ground-truth labels. To evaluate the accuracy of mask re-projection, we compute the Jacaard index $\mathcal{J}$ (IoU) and contour-based accuracy $\mathcal{F}$ proposed in [35], between the rendered masks and the ground truth silhouettes of all annotated frames. Evaluations on keypoint re-projection can be found in the supplementary documents.

In addition, to further quantitatively evaluate shape reconstruction quality, we animate a synthetic 3D bird model and create a video with 520 frames in various poses such as flying, landing, walking etc., as shown in Fig. 7. We then compare the predicted mesh with the ground truth mesh using Chamfer distance every 10 frames.
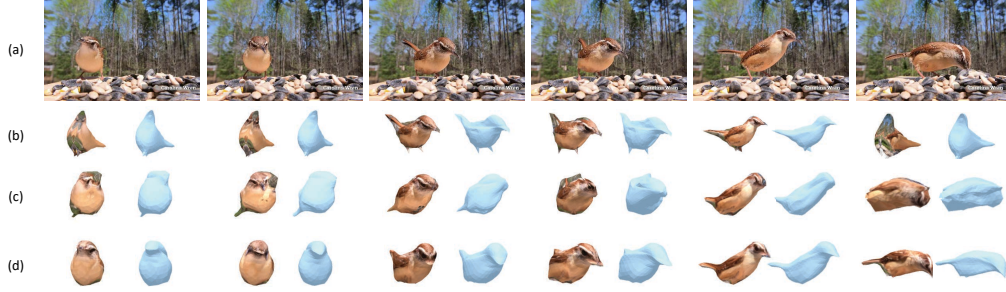
Figure 5: Mesh reconstruction from video frames. (a) Input video frames. (b) Reconstruction from the model trained only on single-view images. (c) Reconstruction from the model test-time trained on the video without the invariance constraints in Sec. 3.2. (d) Reconstruction from the proposed video reconstruction model.

**Network architecture.** For fair comparisons to the baseline method [14], we train our model using the same network as [14], i.e., ResNet18 [9] with batch normalization layers [13] as the encoder, we call this model "ACMR" in the following, which is short for "asymmetric CMR". However, ACMR cannot be well adapted to test videos due to the batch normalization layers and the domain gap between images and videos (see Table 2(d)). Thus, we train a variant model where we use ResNet50 [9] as our encoder and replace all the batch normalization layers in the network with group normalization layers [51]. We call this variant model "ACMR-vid". All test-time training is carried out on the ACMR-vid model unless otherwise specified.

**Network training.** Taking the weakly-supervised setting as an example, to train the image reconstruction model, we first warm up the model without the motion deformation branch, the keypoint re-projection objective, or the ARAP constraint for 200 epochs. This warm-up process effectively avoids the trivial solution where the model solely depends on the motion deformation branch for shape deformation while ignoring the base shape branch. We then train the full image reconstruction network with all objectives for another 200 epochs. Other training details, including the self-supervised setting, and the illustration of our network architecture can be found in the supplementary material.

### 4.2 Qualitative Results

**Mesh reconstruction from images.** In Fig. 4, we show visualizations of reconstructed bird meshes from single-view images. Thanks to the "motion deformation branch" discussed in Sec. 3.1, the proposed ACMR model is able to capture asymmetric motion of the bird such as head rotation (Fig. 4(c)), which cannot be modeled by the baseline method [14] (Fig. 4(g)).
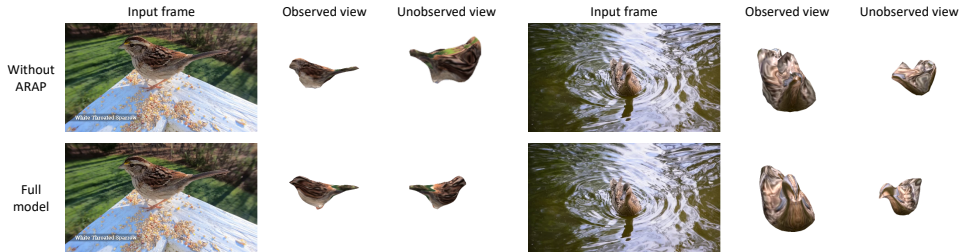


Figure 6: Qualitative comparison of online adaptation with/without the ARAP constraint.

**Online adaptation on a video.** We visualize the reconstructed meshes by our ACMR-vid model for video frames in Fig. 5. Without online adaptation, the ACMR-vid model independently applied to each frame suffers from a domain gap and shows instability over time (Fig. 5(b)). With online adaption as discussed in Sec. 3.2, the ACMR-vid model reconstructs plausible meshes for each video frame as shown in Fig. 5(c) and (d). Specifically, to demonstrate the effectiveness of the proposed invariance constraints, we also show reconstructions of an ACMR-vid model trained without all the invariance constraints in Fig. 5(c), which predicts less reliable shape, camera pose as well as texture compared to our full ACMR-vid model. Finally, we visualize the effectiveness of ARAP for online adaptation in Fig. 6. Without this constraint, the reconstructed meshes are less plausible, especially from unobserved views. More video examples can be found in the supplementary.

8

## 4.3 Quantitative Results

Table 1: Quantitative evaluation of mask IoU and keypoint re-projection (PCK@0.1) on the CUB dataset [45].

| (a) Metric | (b) CMR [14] | (c) ACMR | (d) ACMR, no $\Delta V$ | (e) ACMR, no ARAP | (f) ACMR-vid |
|---|---|---|---|---|---|
| Mask IoU $\uparrow$ | 0.706 | 0.708 | 0.647 | 0.758 | **0.773** |
| PCK@0.1 $\uparrow$ | 0.810 | 0.855 | 0.790 | 0.857 | **0.895** |

**Evaluations on the image dataset.** As shown in Table 1(b) *vs.* (c), our ACMR model achieves comparable mask IoU and higher keypoints re-projection accuracy compared to the baseline model [14] with the same network architecture. This confirms the correctness of both the reconstructed meshes as well as that of the predicted camera poses. In addition, our ACMR-vid model achieves even better performance as shown in Table 1(f). We note that our full ACMR model does not quantitatively outperform the model trained without the ARAP constraint, because the motion deformation $\Delta V$ freely over-fits to the mask and keypoint supervision without any regularization. However, the model without the ARAP constraint visually shows spikes and unnatural deformations as shown in Fig. 4(f) and in the supplementary.

Table 2: Quantitative evaluation of mask re-projection accuracy on the bird video dataset. "(T)" indicates the model is test-time trained on the given video., $L_c$, $L_t$, $L_s$ are defined in Eq. 4, 5, 6 respectively.

| (a) Metric | (b) CMR [14] | (c) ACMR | (d) ACMR (T) | (e) ACMR-vid (T), no $L_c$, $L_t$, $L_s$ | (f) ACMR-vid (T) |
|---|---|---|---|---|---|
| $\mathcal{J}(Mean) \uparrow$ | 0.554 | 0.686 | 0.706 | 0.836 | **0.868** |
| $\mathcal{F}(Mean) \uparrow$ | 0.209 | 0.363 | 0.406 | 0.666 | **0.756** |

**Evaluations on the video dataset.** As shown in Table 2(b) and (c) *vs.* (e), by using the proposed online adaptation method discussed in Sec. 3.2, the model tuned on videos achieves higher $\mathcal{J}$ and $\mathcal{F}$ scores compared to the model trained only on images. This indicates that the test-time trained model successfully adapts to unlabeled videos and can reconstruct meshes that conform well to the frames. The performance of the model is further improved by adding the correspondence, texture, and shape invariance constraints discussed in Sec. 3.2 during online adaptation, as shown in Table 2(f).
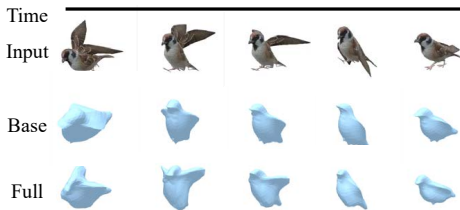


Figure 7: Reconstructions of an animated video clip.

**Evaluation on animated sequences.** We apply the proposed model to an animated video clip and compare the predicted mesh with the ground truth mesh using Chamfer distance every 10 frames. We show the qualitative reconstructions in Fig. 7 and quantitative evaluation results in Table 3. The proposed ACMR method outperforms the baseline CMR [14] model and is further improved via the proposed online adaptation strategy discussed in Sec. 3.2.

**Evaluations for self-supervised setting (Sec. 3.3).** After online adaptation, this model too, achieves both a higher $\mathcal{J}$ score (0.843 *vs.* 0.582) and $\mathcal{F}$ score (0.678 *vs.* 0.256) compared

| Metric | CMR | ACMR | ACMR (T) |
|---|---|---|---|
| Chamfer$\downarrow$ | 0.016 | 0.015 | **0.012** |

Table 3: Evaluation on synthetic data.

to the model pre-trained on the image dataset, i.e., our method reconstructs promising dynamic 3D meshes without annotating any keypoint for training.

## 5 Conclusions

We propose a method to reconstruct temporally consistent 3D meshes of deformable objects from videos captured in the wild. We learn a category-specific 3D mesh reconstruction model that jointly predicts the shape, texture, and camera pose from single-view images, which is capable of capturing asymmetric non-rigid motion deformation of objects. We then adapt this model to any unlabeled video by exploiting self-supervised signals in videos, including those of shape, texture, and part consistency. Experimental results demonstrate the superiority of the proposed method compared to state-of-the-art works, both qualitatively and quantitatively.

## Broader Impact

The developed method will make significant contributions to both the 3D vision and endangered species research. The method provides a way to study animals that can only be captured in the wild as 2D videos, e.g., endangered animal species of birds and zebras. The broader impact includes enhancing our understanding of such endangered animals simply from videos, as they can be reconstructed and viewed in 3D. The method can also be applied to tasks such as bird watching, motion analysis, shape analysis, to name a few. Furthermore, another important application is to simplify an artists workflow, as an initial animated and textured 3D shape can be directly derived from a video.

# Appendix

In the Appendix, we provide additional details, discussions, and experiments to support the original submission. In the following, we first discuss our self-supervised setting in Sec. 6. We then describe evaluation of keypoint re-projection accuracy on videos in Sec. 7. Next, we show more ablation studies in Sec. 8. More qualitative results on both bird and zebra image reconstructions are present in Sec. 9. Details of the network design and implementation are discussed in Sec. 10 and Sec. 11, respectively. Finally, we describe failure cases and limitations in Sec. 12.

## 6 Self-supervised Mesh Reconstruction

We train the self-supervised image reconstruction model with only silhouettes and single-view images in a category. To this end, we also learn a template from scratch as in the self-supervised 3D mesh reconstruction approach [24]. Essentially, we do not apply the semantic parts from the SCOPS method [12], which means that no additional modules are required for training. After training the image model, we adapt it to each unlabeled video using the method discussed in Sec.3.2 in the submission. Since neither keypoints annotations, nor additional self-supervised blocks such as SCOPS are adopted, the results of this self-supervised single-view image reconstruction model do not outperform those of existing methods, i.e., [14] and [24], and of the proposed ACMR model (see Sec. 4.3 for the quantitative results). However, the test-time training improves the fidelity and the robustness of the reconstruction results as shown in Fig. 8. The reconstructions are more plausible after online adaptation, especially from unobserved views.
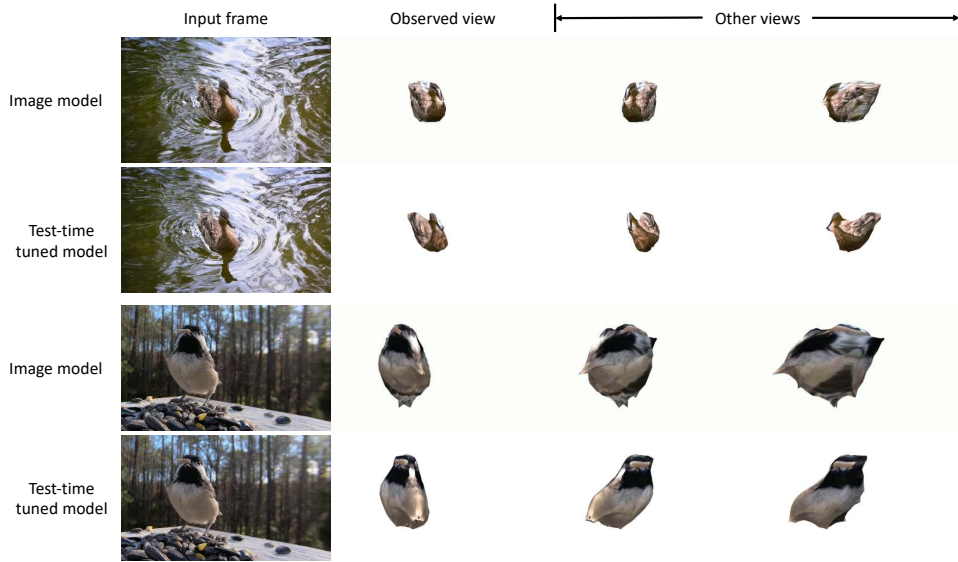


Figure 8: Comparison of the self-supervised image model and the test-time tuned self-supervised model.

## 7 Keypoint Re-projection Accuracy on Videos

**Video Keypoints Annotation.** We evaluate the keypoint re-projection accuracy (see Sec 4.1 in the paper) on the 22 videos we collected. To create the ground truth keypoints, we follow the protocol of the CUB dataset [45] to annotate 15 semantic keypoints every five frames in each video, via the Labelme [44] toolkit (see Fig. 9 for the annotation interface.) Visualizations of the re-projected keypoints by different methods are visualized and compared in Fig. 10.



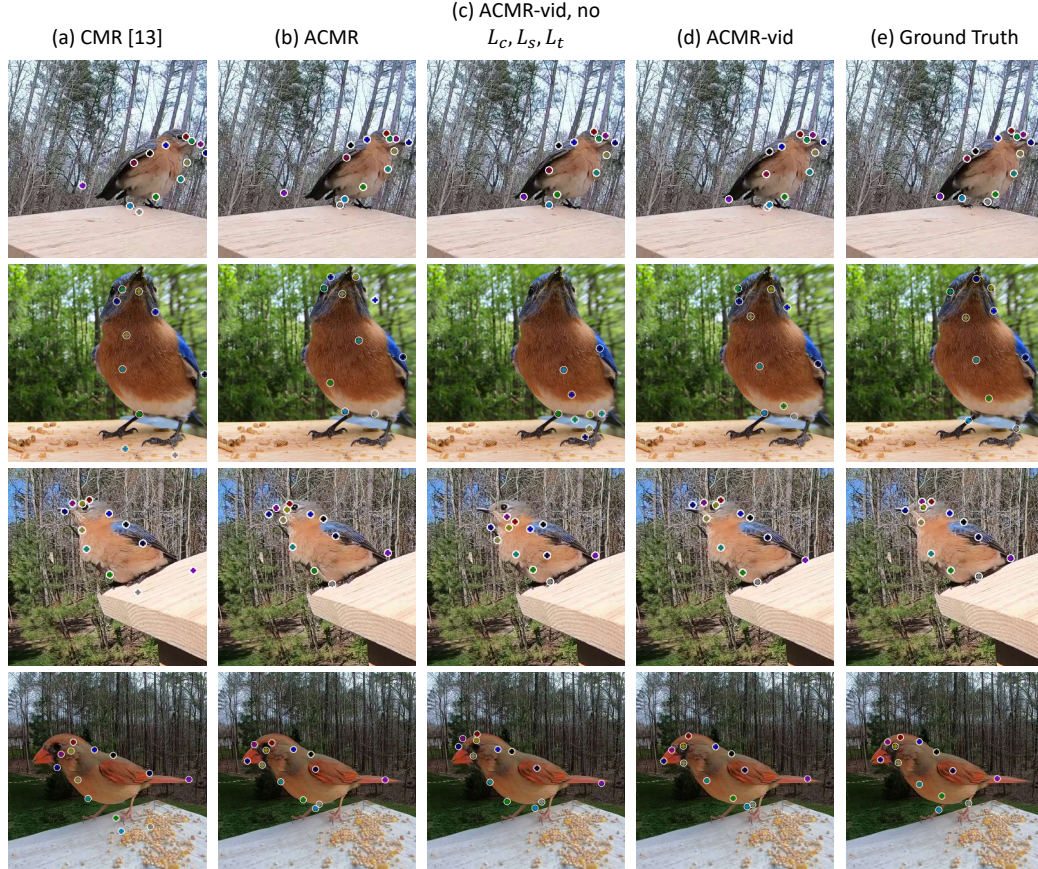Figure 9: Keypoints annotation using the Labelme [44] toolkit.

Figure 10: Visualization of re-projected keypoints on videos. We use white circles to highlight the keypoints for better visualization.

Table 4: Keypoint re-projection evaluation on videos. "(T)" indicates the model is test-time trained on the given video., $L_c$, $L_t$, $L_s$ are defined in Eq.4, 5 and 6 of the main paper, respectively.

| (a) Metric | (b) CMR [14] | (c) ACMR | (d) ACMR (T) | (e) ACMR-vid (T), no $L_c, L_t, L_s$ | (f) ACMR-vid (T) |
|---|---|---|---|---|---|
| $PCK@0.1 \uparrow$ | 0.514 | 0.751 | 0.424 | 0.644 | **0.794** |

**Details of Keypoint Re-projection.** Since we do not have the keypoint assignment matrix proposed in [14], we employ the canonical keypoint UV map to obtain the 3D keypoints (Sec. 3.1 in the paper). The keypoint re-projection is done by (i) warping the canonical keypoint UV map to each individual predicted mesh surface; (ii) projecting the canonical keypoint back to the 2D space via the predicted camera pose; (iii) comparing against the ground truth keypoints in 2D. This evaluation implicitly reveals the correctness of both the predicted shape and camera pose for the mesh reconstruction algorithm, especially for objects that do not have 3D ground truth annotations.

Compared to frame-wisely applying CMR [14] (Table 4 (b)) or ACMR (Table 4 (c)) discussed in Sec. 3.1 in the main paper, the test-time tuned model achieves higher PCK score, as shown in Table 4 (f). It verifies the effectiveness of the proposed test-time training procedure and the invariance constraints. Essentially, as we noted in Sec 4.1, although the original ACMR, i.e., using the ResNet-18 [9] as the image encoder with batch normalization layers [13] in Table 4 (c), achieves relatively promising results, it is hard to adapt this model to new domains like low-quality videos (e.g., when switching from the *.eval()* mode to the *.train()* mode in PyTorch [32]). The performance drops significantly after test-time tuning as shown in Table 4 (d).

# 8 Ablation Studies

## 8.1 The Role of Shape Base

To demonstrate the superiority of using a set of shape bases versus using a single template, we train a baseline model where we replace the shape combination branch with the template obtained by the CMR approach [14]. This setting is equivalent to a using a single shape base (denoted as *single base*). We show quantitative and qualitative comparisons

Table 5: Quantitative comparison of the single base model with the proposed ACMR model.

| (a) Metric | (b) Single base | (c) ACMR |
|---|---|---|
| Mask IoU ↑ | 0.605 | 0.708 |
| PCK@0.1 ↑ | 0.655 | 0.855 |

with the proposed ACMR model in Table 5 and Fig. 11, respectively. As shown in Table 5(b) *vs.* (c), the model trained with a single base template struggles to fit the final shape when the instance is largely different from the given template (also shown in Fig. 11). In contrast, the proposed ACMR model with 8 shape bases performs favorably against the single base model.



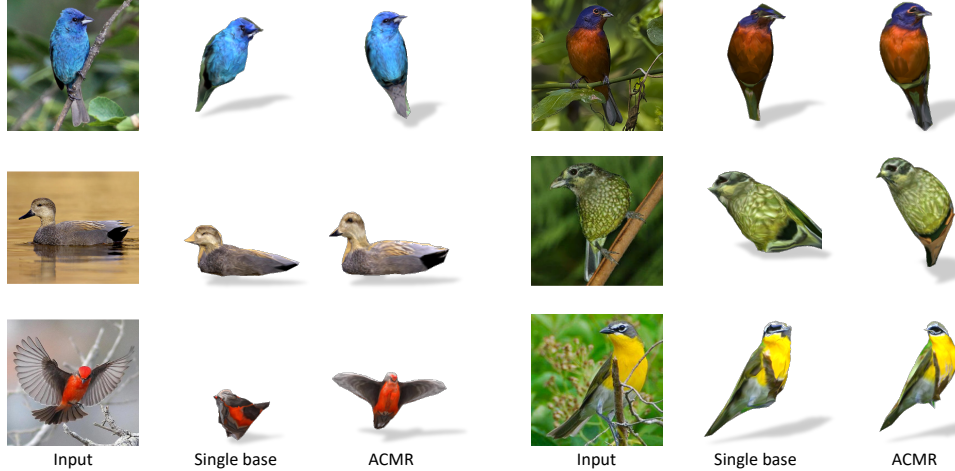| Input | Single base | ACMR | Input | Single base | ACMR |

Figure 11: Qualitative comparison of the single base model with the proposed ACMR model with multiple shape bases. The single base model suffers when the instance is largely different from the template, e.g. flying bird or a duck.

## 8.2 ARAP Constraint in Online Adaptation

To verify the effectiveness of using the ARAP constraint in the online adaptation process, we test-time tune on the videos without this constraint. Although performing online adaptation without the ARAP constraint yields better quantitative evaluations as shown in Table 6, the reconstructed meshes are not plausible from unobserved views, as shown in Fig.6 in the submission.

Table 6: Ablation study on the ARAP constraint in online adaptation.

| (a) Metric | (b) without ARAP | (c) ACMR-vid (T) |
|---|---|---|
| $\mathcal{J}(Mean)$ ↑ | 0.875 | 0.868 |
| $\mathcal{F}(Mean)$ ↑ | 0.782 | 0.756 |
| PCK@0.1 ↑ | 0.815 | 0.794 |

## 9 Qualitative Evaluations

### 9.1 Camera Pose Stability

To visually demonstrate the effectiveness of the test-time training procedure that stabilizes camera pose prediction, we visualize the differences in camera pose predictions between adjacent frames in Fig. 13. Compared to the model that is only trained on images, the proposed method predicts more stable camera poses that change smoothly over time.

### 9.2 Keypoints Re-projection for Image-based Reconstruction

We visualize re-projected keypoints on test images in Fig. 14, where the corresponding quantitative results are presented in Sec. 4.3, Table 1 of the main paper. The proposed ACMR model is able to predict more accurate keypoints compared to the CMR [14] method, especially when the bird performs an asymmetric pose, e.g. first row in Fig. 14.

## 9.3 Single-view Image Reconstructions

In Fig. 12, we show more visualizations of reconstructions from single-view images of the test dataset of CUB birds [45] as well as comparisons with the baseline method [14]. By removing the symmetric assumption, our model is able to reconstruct objects in the input images more faithfully versus the baseline method [14] (Fig. 12(g)).

We also demonstrate the effectiveness of the ARAP constraint, as discussed in Sec.3.1 in the paper. Without this constraint, the reconstructed meshes contain unnatural spikes, as shown in Fig. 12(f).

Finally, we show reconstruction results of zebra images of the test dataset [58] in Fig. 15. Our ACMR model successfully captures motions such as head bending or walking for zebras.
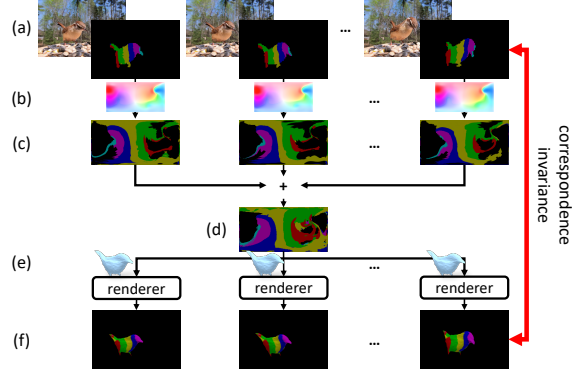


Figure 16: Part correspondence constraint. (a) Input frame and part propagations. (b) Predicted texture flows. (c) Part UV map. (d) Aggregated video-level part UV map. (e) Base shape and differentiable renderer. (f) Part rendering.

## 10 Network Architecture

### 10.1 Bases Visualization

We visualize the eight shape bases obtained by applying KMeans clustering on all reconstructed meshes by the CMR [14] method for birds in Fig. 17(a). We also show the bases obtained by applying PCA to the bottleneck features of the image encoder. Note that the latter fails to discover rare shape modes (e.g., duck and flying bird) in the dataset as shown in Fig. 17(b). Thus we choose to use KMeans to obtain shape bases.

### 10.2 Part Correspondence Constraint

We illustrate the part correspondence constraint in details in Fig. 16. Given the propagated parts in each frame in Fig. 16(a), we map them to the UV space with the predicted texture flow in Fig. 16(b) and obtain part UV maps in Fig. 16(c). By aggregating these part UV maps, i.e., averaging, we minimize noise in each individual part UV map and obtain a video-level part UV map in Fig. 16(d). This video-level part UV map is shared by all frames in the video. Thus, for each frame, we wrap the video-level part UV map onto the base shape prediction and render it under the predicted camera pose as shown in Fig. 16(f). Finally, we encourage consistency between part renderings and part propagations, as shown by the red arrow in Fig. 16. Through the differentiable renderer, the loss implicitly improves the predicted camera pose.

### 10.3 Single-view Reconstruction Network

In Fig. 18(a), we show details of our single-view reconstruction network. Given an input image, the network jointly predicts texture, shape and camera pose. By utilizing a differentiable renderer [27], we are able to utilize 2D supervision, i.e. silhouettes and input images.

### 10.4 Sliding Window Scheme

We show the proposed test-time tuning process in Fig. 18(b). Within each sliding window, we encourage the consistency of UV texture, UV parts as well as base shape of all frames.

## 11 Implementation Details

### 11.1 Objectives for Image Reconstruction Model

We summarize the objectives used in the single-view reconstruction model (Fig. 18(a)) discussed in Sec.3.1 in the submission as follows: (i) *foreground mask loss:* a negative intersection over union objective between rendered and ground truth silhouettes [14, 24, 16]; (ii) *foreground RGB texture loss:* a perceptual metric [14, 24, 54] between rendered and input RGB images; (iii) *mesh smoothness:* a Laplacian constraint [14, 24] to encourage smooth mesh reconstruction; (iv) *keypoint re-projection loss:* as discussed in Sec.3.1 in the paper; and (v) *the ARAP constraint:* described in Sec.3.1 in the paper. The weight for each objective is set to 3.0, 3.0, 0.0008, 5.0 and 10.0.

### 11.2 Objectives for Online Adaptation

We summarize the objectives used in the online adaptation process (Fig. 18(b)) in the following. Since it is feasible to predict a segmentation mask via a pretrained segmentation model, we make use of the predicted foreground mask and compute the (i), (ii), and (iii) losses (mentioned above) similarly to the image-based training. We also adopt the the ARAP constraint described in Sec.3.1 in the paper, and the three invariance constraints as discussed in Sec.3.2 for online adaptation. The weight for each objective is set to 0.1, 0.5, 0.0006, 2.0 and 2.0 (texture invariance), 1.0 (part correspondence), 1.0 (base shape invariance).

### 11.3 Training Details

We implement the proposed method in PyTorch [32] and use the Adam optimizer [20] with a learning rate of 0.0001 for both the image reconstruction model training and online adaptation. The weight of each objective for the image reconstruction model as well as the online adaptation process is discussed in Sec. 11.1 and Sec. 11.2, respectively.

### 11.4 Training on Zebra Images and Videos

We adopt a different scheme to train a single-view reconstruction model on zebras: (i) since natural zebra images labeled with keypoints are not publicly available, we adopt a synthetic dataset [58], (ii) zebras have more complex shapes with large concavities. Therefore, it is not suitable to learn the shape by deforming from a sphere primitive. Instead, we utilize a readily available zebra mesh as a template and learn motion deformation on top of it. We first train an image reconstruction model using the synthetic dataset provided by [58]. Similarly as [58], we utilize the silhouettes, keypoints, texture maps as well as partially available UV texture flow as supervision. For shape reconstruction, instead of the utilizing the SMAL parametric model [60], we use the proposed shape module, i.e. combination of base shapes and motion deformation. Due to the limited motion of zebras, we only use one base shape, which is a readily available zebra mesh with 3889 vertices and 7774 faces. For camera pose prediction, we use the perspective camera pose discussed in Sec.3 in the submission as well as in [14]. Due to the limited capacity of a single UV texture map, we also model the texture map by cutting the UV texture map into four pieces and stitch them together similarly as in [58]. We note that this "cutting and stitching" operation does not influence the mapping and aggregation of the part UV maps discussed in Sec.3.2 in the submission.

## 12 Failure Cases

Our work is the first to explore the challenging task of reconstructing 3D meshes of deformable object instances from videos in the wild. Impressive as the performance is, this challenging task is far from being fully solved. We discuss failure cases and limitations of the proposed method in the following. To begin with, we focus on genus-0 objects such as birds and zebras in this work. Thus our model suffers when it is generalized to objects with large concave holes such as chairs, humans etc. Second, our work struggles to reconstruct meshes from videos with large motion and lighting changes as well as occlusion, (see Fig. 19). This is mainly due to the failure in correctly propagating parts by the self-supervised UVC model [25], which is out of scope of this work. We leave all these failure cases and limitations to future works.

# References

[1] A. Arnab, C. Doersch, and A. Zisserman. Exploiting temporal context for 3d human pose estimation in the wild. In *CVPR*, June 2019. 3

[2] C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3d shape from image streams. In *CVPR*, 2000. 1, 3

[3] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI*, 2017. 6, 7

[4] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *ECCV*, 2016. 1, 2

[5] C. Doersch and A. Zisserman. Sim2real transfer learning for 3d human pose estimation: motion to the rescue. In *NeurIPS*, 2019. 3

[6] Y. Feng, F. Wu, X. Shao, Y. Wang, and X. Zhou. Joint 3d face reconstruction and dense alignment with position map regression network. In *ECCV*, 2018. 3

[7] P. Guo and R. Farrell. Aligned to the object, not to the image: A unified pose-aligned representation for fine-grained recognition. In *WACV*, 2019. 7

[8] M. Habermann, W. Xu, M. Zollhoefer, G. Pons-Moll, and C. Theobalt. Deepcap: Monocular human performance capture using weak supervision. In *CVPR*, 2020. 5

[9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 8, 12

[10] P. Henderson and V. Ferrari. Learning to generate and reconstruct 3d meshes with only 2d supervision. In *BMVC*, 2018. 2

[11] J. F. Hughes, A. Van Dam, J. D. Foley, M. McGuire, S. K. Feiner, and D. F. Sklar. *Computer graphics: principles and practice*. Pearson Education, 2014. 3

[12] W.-C. Hung, V. Jampani, S. Liu, P. Molchanov, M.-H. Yang, and J. Kautz. Scops: Self-supervised co-part segmentation. In *CVPR*, 2019. 11

[13] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 8, 12

[14] A. Kanazawa, S. Tulsiani, A. A. Efros, and J. Malik. Learning category-specific mesh reconstruction from image collections. In *ECCV*, 2018. 1, 2, 3, 4, 6, 7, 8, 9, 11, 12, 13, 14, 15

[15] A. Kanazawa, J. Y. Zhang, P. Felsen, and J. Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. 1, 3

[16] H. Kato and T. Harada. Learning view priors for single-view 3d reconstruction. In *CVPR*, 2019. 2, 6, 15

[17] H. Kato and T. Harada. Self-supervised learning of 3d objects from natural images. *arXiv preprint arXiv:1911.08850*, 2019. 2

[18] H. Kato, Y. Ushiku, and T. Harada. Neural 3d mesh renderer. In *CVPR*, 2018. 1, 2, 3

[19] A. Khoreva, A. Rohrbach, and B. Schiele. Video object segmentation with language referring expressions. In *ACCV*, 2018. 7

[20] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 15

[21] C. Kong and S. Lucey. Deep non-rigid structure from motion. In *ICCV*, 2019. 3, 5

[22] N. Kulkarni, A. Gupta, D. F. Fouhey, and S. Tulsiani. Articulation-aware canonical surface mapping. In *CVPR*, 2020. 1

[23] N. Kulkarni, A. Gupta, and S. Tulsiani. Canonical surface mapping via geometric cycle consistency. In *ICCV*, 2019. 1, 2

[24] X. Li, S. Liu, K. Kim, S. De Mello, V. Jampani, M.-H. Yang, and J. Kautz. Self-supervised single-view 3d reconstruction via semantic consistency. *arXiv preprint arXiv:2003.06473*, 2020. 1, 2, 3, 4, 6, 7, 11, 15

[25] X. Li, S. Liu, S. D. Mello, X. Wang, J. Kautz, and M.-H. Yang. Joint-task self-supervised learning for temporal correspondence. In *NeurIPS*, 2019. 5, 6, 15

[26] C.-H. Lin, O. Wang, B. C. Russell, E. Shechtman, V. G. Kim, M. Fisher, and S. Lucey. Photometric mesh optimization for video-aligned 3d object reconstruction. In *CVPR*, 2019. 1, 2, 3

[27] S. Liu, T. Li, W. Chen, and H. Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *ICCV*, 2019. 2, 3, 14

[28] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 2015. 3, 4

[29] X. Luo, J. Huang, R. Szeliski, K. Matzen, and J. Kopf. Consistent video depth estimation. *ACM Transactions on Graphics (Proceedings of ACM SIGGRAPH)*, 39(4), 2020. 3

[30] D. Novotny, N. Ravi, B. Graham, N. Neverova, and A. Vedaldi. C3dpo: Canonical 3d pose networks for non-rigid structure from motion. In *ICCV*, 2019. 1, 3, 5

[31] J. Pan, X. Han, W. Chen, J. Tang, and K. Jia. Deep mesh reconstruction from single rgb images via topology modification networks. In *ICCV*, 2019. 2

[32] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *NeurIPS*. 2019. 12, 15

[33] D. Pavllo, C. Feichtenhofer, D. Grangier, and M. Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *CVPR*, 2019. 3

[34] X. B. Peng, A. Kanazawa, J. Malik, P. Abbeel, and S. Levine. Sfv: Reinforcement learning of physical skills from videos. *ACM Trans. Graph.*, 37(6), Nov. 2018. 3

[35] J. Pont-Tuset, F. Perazzi, S. Caelles, P. Arbeláez, A. Sorkine-Hornung, and L. Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 7

[36] D. J. Rezende, S. A. Eslami, S. Mohamed, P. Battaglia, M. Jaderberg, and N. Heess. Unsupervised learning of 3d structure from images. In *NeurIPS*, 2016. 2

[37] H. Rhodin, N. Robertini, D. Casas, C. Richardt, H.-P. Seidel, and C. Theobalt. General automatic human shape and motion capture using volumetric contour cues. In *ECCV*, 2016. 3

[38] K. Robinette, S. Blackwell, H. Daanen, M. Boehmer, S. Fleming, T. Brill, D. Hoeferlin, and D. Burnsides. Civilian american and european surface anthropom- etry resource (caesar) final report. *Tech. Rep. AFRL-HE-WP-TR-2002-0169, US Air Force Research Laboratory*, 2002. 4

[39] O. Sorkine and M. Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry processing*, 2007. 5

[40] Y. Sun, X. Wang, Z. Liu, J. Miller, A. A. Efros, and M. Hardt. Test-time training for out-of-distribution generalization. *arXiv preprint arXiv:1909.13231*, 2019. 2

[41] L. Tran and X. Liu. On learning 3d face morphable model from in-the-wild images. *TPAMI*, 2019. 3

[42] H.-Y. Tung, H.-W. Tung, E. Yumer, and K. Fragkiadaki. Self-supervised learning of motion capture. In *NeurIPS*, 2017. 3

[43] T. von Marcard, R. Henschel, M. J. Black, B. Rosenhahn, and G. Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 3

[44] K. Wada. labelme: Image Polygonal Annotation with Python. https://github.com/wkentaro/labelme, 2016. 11

[45] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset, 2011. 7, 9, 11, 14, 18

[46] B. Wandt, H. Ackermann, and B. Rosenhahn. 3d reconstruction of human motion from monocular image sequences. *TPAMI*, 2016. 3

[47] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *ECCV*, 2018. 2

[48] C. Wen, Y. Zhang, Z. Li, and Y. Fu. Pixel2mesh++: Multi-view 3d mesh generation via deformation. In *ICCV*, 2019. 2

[49] O. Wiles and A. Zisserman. Silnet: Single-and multi-view reconstruction by learning from silhouettes. *arXiv preprint arXiv:1711.07888*, 2017. 2

[50] S. Wu, C. Rupprecht, and A. Vedaldi. Unsupervised learning of probably symmetric deformable 3d objects from images in the wild. In *CVPR*, 2020. 1, 2

[51] Y. Wu and K. He. Group normalization. In *ECCV*, 2018. 8

[52] X. Yan, J. Yang, E. Yumer, Y. Guo, and H. Lee. Perspective transformer nets: Learning single-view 3d object reconstruction without 3d supervision. In *NeurIPS*, 2016. 2

[53] J. Y. Zhang, P. Felsen, A. Kanazawa, and J. Malik. Predicting 3d human dynamics from video. In *ICCV*, 2019. 1, 3

[54] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 6, 15

[55] X. Zhou, M. Zhu, S. Leonardos, K. G. Derpanis, and K. Daniilidis. Sparseness meets deepness: 3d human pose estimation from monocular video. In *CVPR*, 2016. 3

[56] R. Zhu, C. Wang, C.-H. Lin, Z. Wang, and S. Lucey. Object-centric photometric bundle adjustment with deep shape prior. In *WACV*, 2018. 2

[57] Y. Zhu, D. Huang, F. De La Torre, and S. Lucey. Complex non-rigid motion 3d reconstruction by union of subspaces. In *CVPR*, 2014. 3

[58] S. Zuffi, A. Kanazawa, T. Berger-Wolf, and M. J. Black. Three-d safari: Learning to estimate zebra pose, shape, and texture from images in the wild. In *ICCV*, 2019. 3, 14, 15

[59] S. Zuffi, A. Kanazawa, T. Berger-Wolf, and M. J. Black. Three-d safari: Learning to estimate zebra pose, shape, and texture from images "in the wild". In *ICCV*, 2019. 4, 7

[60] S. Zuffi, A. Kanazawa, D. Jacobs, and M. J. Black. 3D menagerie: Modeling the 3D shape and pose of animals. In *CVPR*, 2017. 4, 15
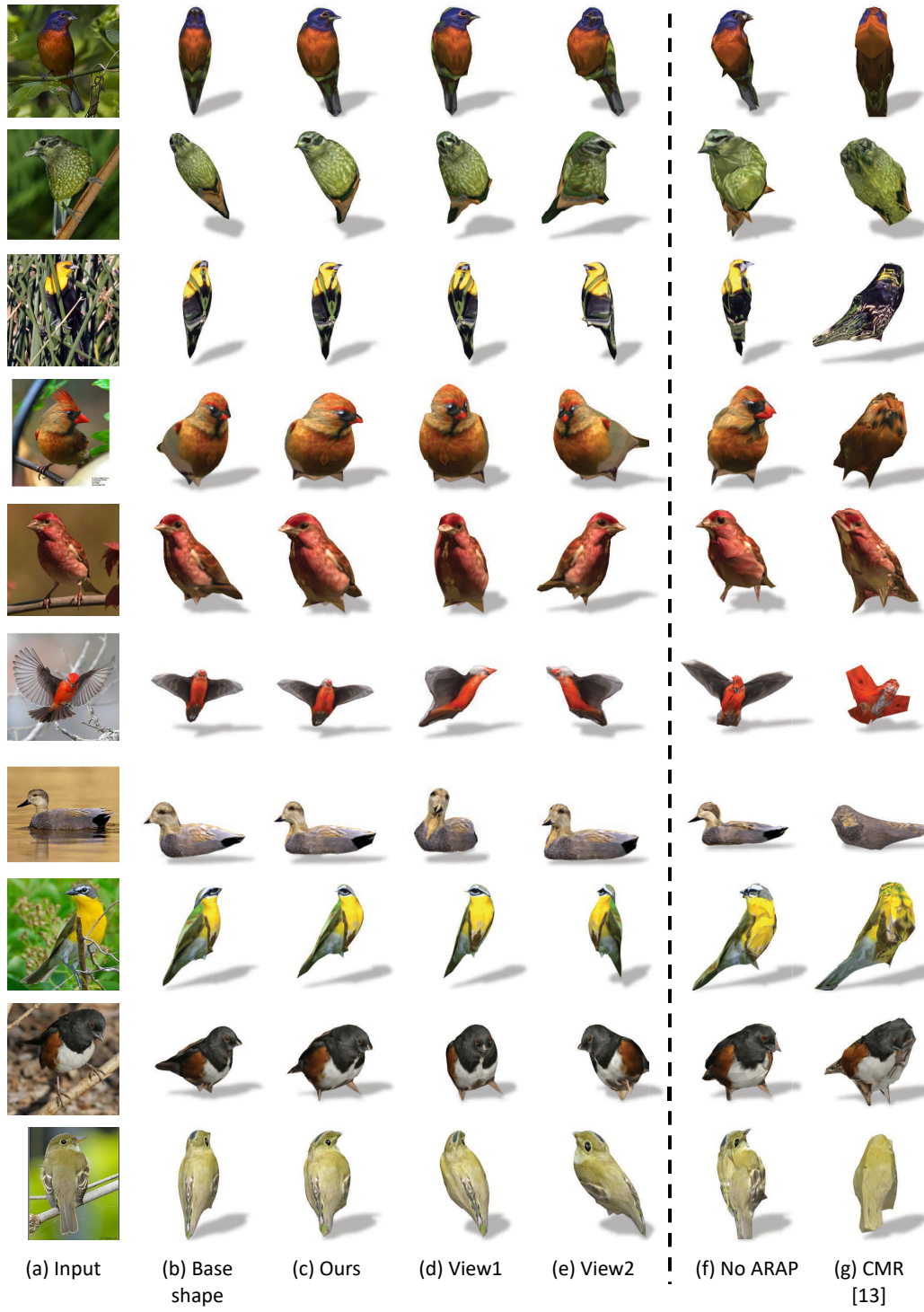
| (a) Input | (b) Base shape | (c) Ours | (d) View1 | (e) View2 | (f) No ARAP | (g) CMR [13] |

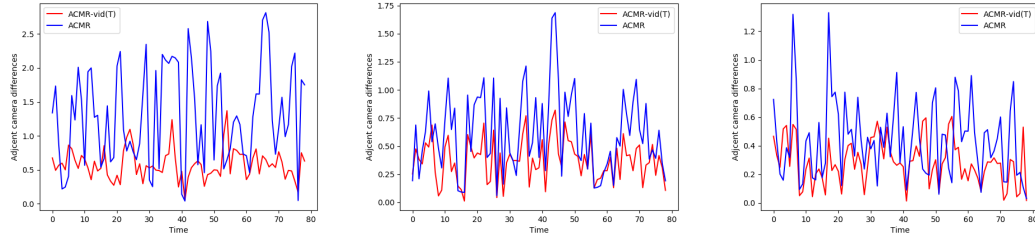Figure 12: More qualitative reconstruction results on CUB birds [45].

18

Figure 13: Camera stability visualization. We visualize differences between adjacent camera pose predictions. The blue and red lines represent the model trained only with images and the test-time tuned model, respectively.
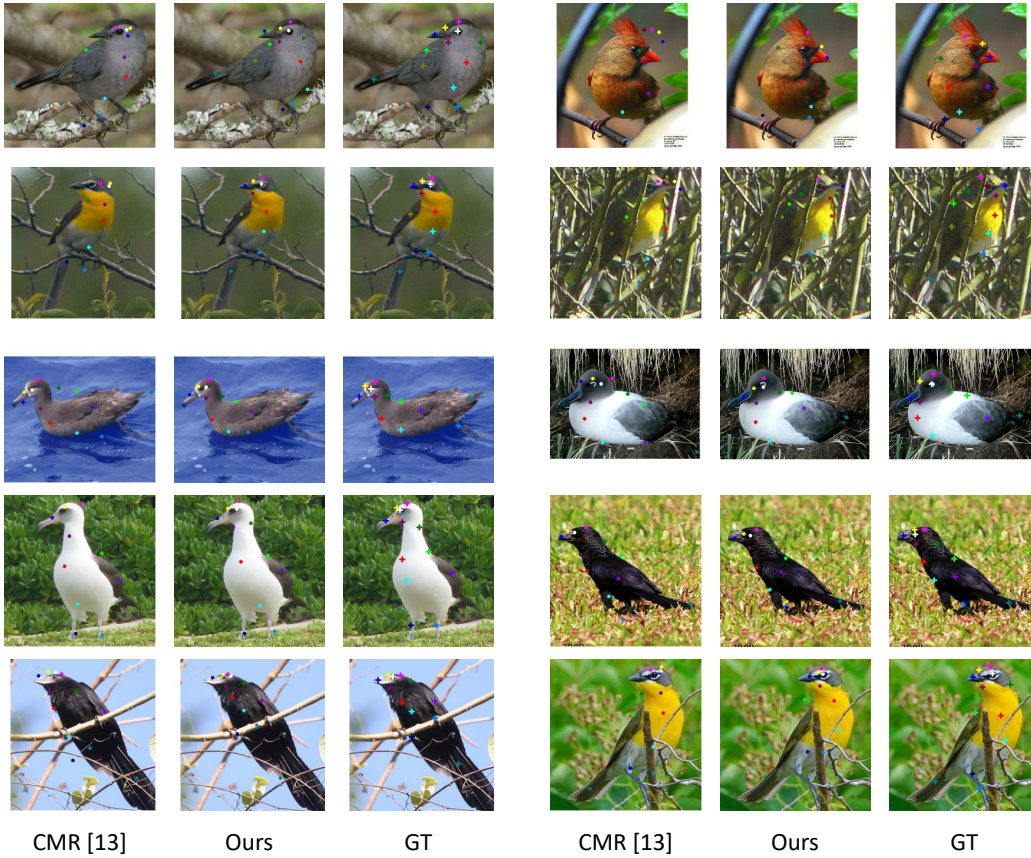


| CMR [13] | Ours | GT | CMR [13] | Ours | GT |

Figure 14: Visualization of re-projected keypoints of the single-view image reconstruction model.

Input       Observed View       Other views

Figure 15: Visualization of reconstructed zebras.

(a) Applying KMeans to reconstructed meshes by CMR [13]



(b) Applying PCA to feature space of CMR [13]

Figure 17: Bases visualization.
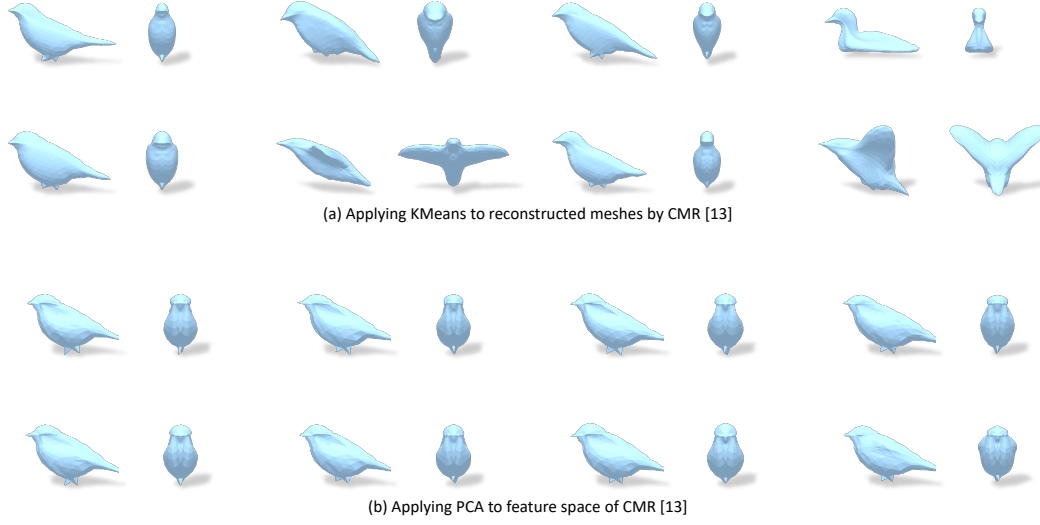


(a) Image reconstruction model
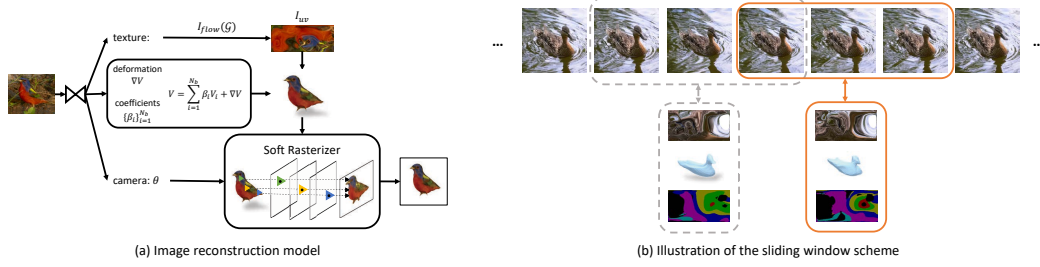


(b) Illustration of the sliding window scheme

Figure 18: Frameworks. For the purposes of illustration only, we show the test-time tuning procedure in (b) with a sliding widow size of 3 and sliding stride of 2. In all our experiments, we use a sliding window size of 50 and stride of 10. The gray dashed box shows the previous sliding window while the orange box shows the current sliding window.
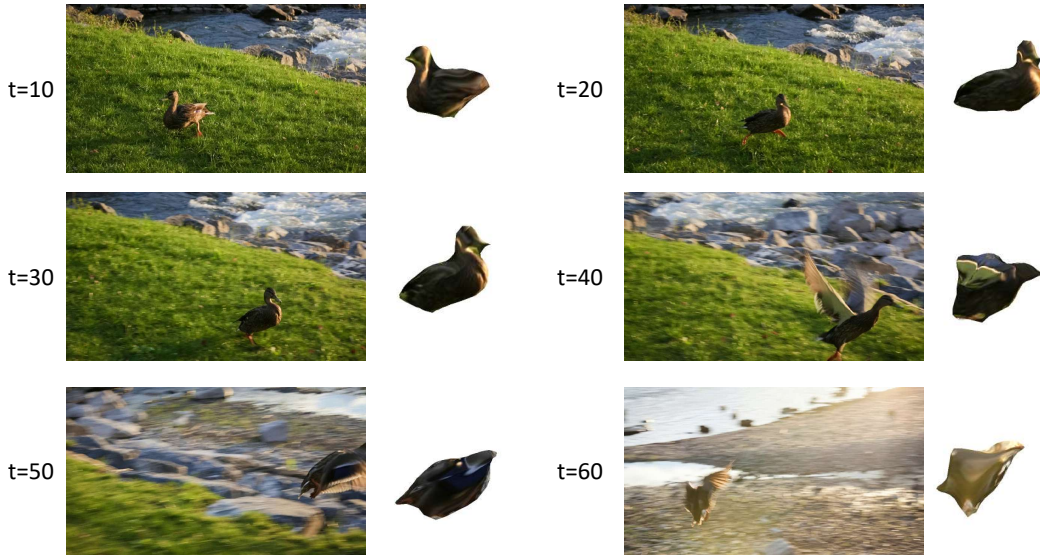


t=10

t=20

t=30

t=40

t=50

t=60

Figure 19: Failure cases. Our model fails when there is occlusion (e.g., $t = 50$, $t$ stands for frame number) or large lighting changes (e.g., $t = 60$).