

REACTO: Reconstructing Articulated Objects from a Single Video



Chaoyue Song¹, Jiacheng Wei¹, Chuan-Sheng Foo², Guosheng Lin¹, Fayao Liu²

¹Nanyang Technological University, ²A*STAR



NANYANG
TECHNOLOGICAL
UNIVERSITY
SINGAPORE



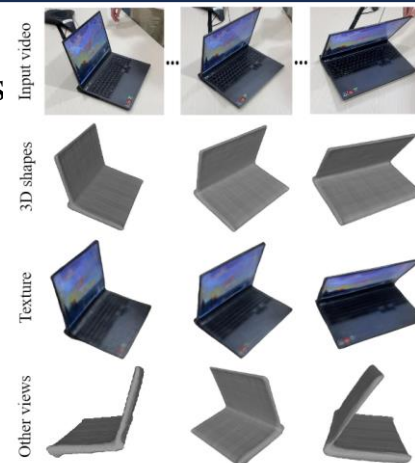
Introduction

Motivation

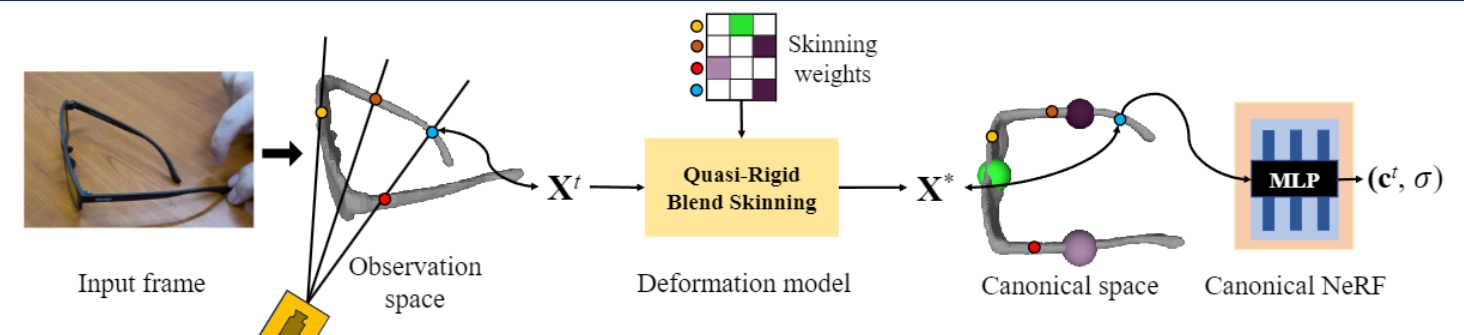
- Current research focuses on modeling humans or animals it's essential to give greater consideration to general articulated objects, which are prevalent in our daily lives.
- Part-rigidity nature of general articulated objects sets them apart from humans or animals, need to design rigging strategy and deformation model to describe their motion.

Task Description

- Given a single casual video capturing a general articulated object (e.g., stapler, laptop, etc.), we model the geometry, texture, and motions of the object.



Overview



We model an articulated 3D object using a shape and appearance model based on a canonical NeRF and a deformation model (Quasi-Rigid Blend Skinning) for transforming 3D points between the observation space and the canonical space. The colors in skinning weights signify the assigned bone for each point.

Method

Canonical NeRF for shape and appearance

$$\mathbf{c}^t = \text{MLP}_{\text{color}}(\mathbf{X}^*, \mathbf{D}^t, \psi_a^t),$$

$$\sigma = \Phi_\beta(\text{MLP}_{\text{SDF}}(\mathbf{X}^*)),$$

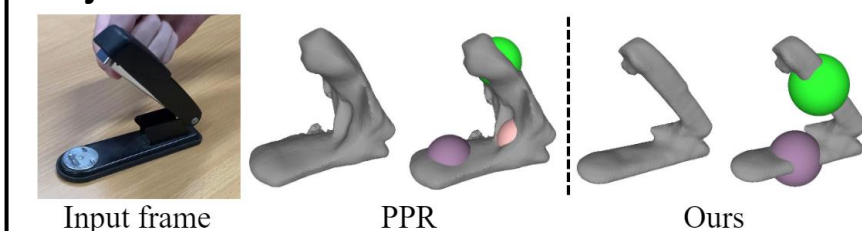
Quasi-rigid blend skinning for deformation

1. Motion representation.

$$\mathbf{X}^t = \mathcal{D}^{t,c \rightarrow o}(\mathbf{X}^*) = \mathbf{T}_{\text{global}}^t \mathbf{T}_{\text{obj}}^{t,c \rightarrow o} \mathbf{X}^*,$$

$$\mathbf{X}^* = \mathcal{D}^{t,o \rightarrow c}(\mathbf{X}^t) = \mathbf{T}_{\text{obj}}^{t,o \rightarrow c} (\mathbf{T}_{\text{global}}^t)^{-1} \mathbf{X}^t,$$

2. Bone definition: rig on bones (ours) v.s. rig on joints.



3. Quasi-sparse skinning weights.

$$\mathbf{W}^s = \text{softmax}\left(\frac{d_M(\mathbf{X}) + \mathbf{W}_\Delta}{\gamma}\right).$$

4. Geodesic point assignment.

Algorithm 1 Geodesic point assignment

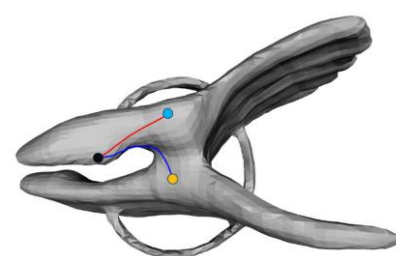
Input: Point assignment $\mathbf{M} = \mathbf{0} \in \mathbb{R}^B$, Mahalanobis distance d_M^i and d_M^j , geodesic distance d_G^i and d_G^j , bone index i and j , hyperparameters η, ζ .

Output: Updated assignment \mathbf{M}

```

1: if  $d_M^i/d_M^j < 1 - \eta$  then
2:    $\mathbf{M}[i] \leftarrow 1$ 
3: else if  $\frac{|d_G^i - d_G^j|}{\min(d_G^i, d_G^j)} < \zeta$  then
4:    $\mathbf{M}[i], \mathbf{M}[j] \leftarrow 1$            ▷ Assigning to joints
5: else
6:    $\mathbf{M}[\text{argmin}(d_G)] \leftarrow 1$ 
7: end if
    
```

$$\mathcal{L}_{\text{sparse}} = \frac{\sum \|\mathbf{W}^s \odot \bar{\mathbf{M}}\|^2}{\sum \bar{\mathbf{M}}},$$

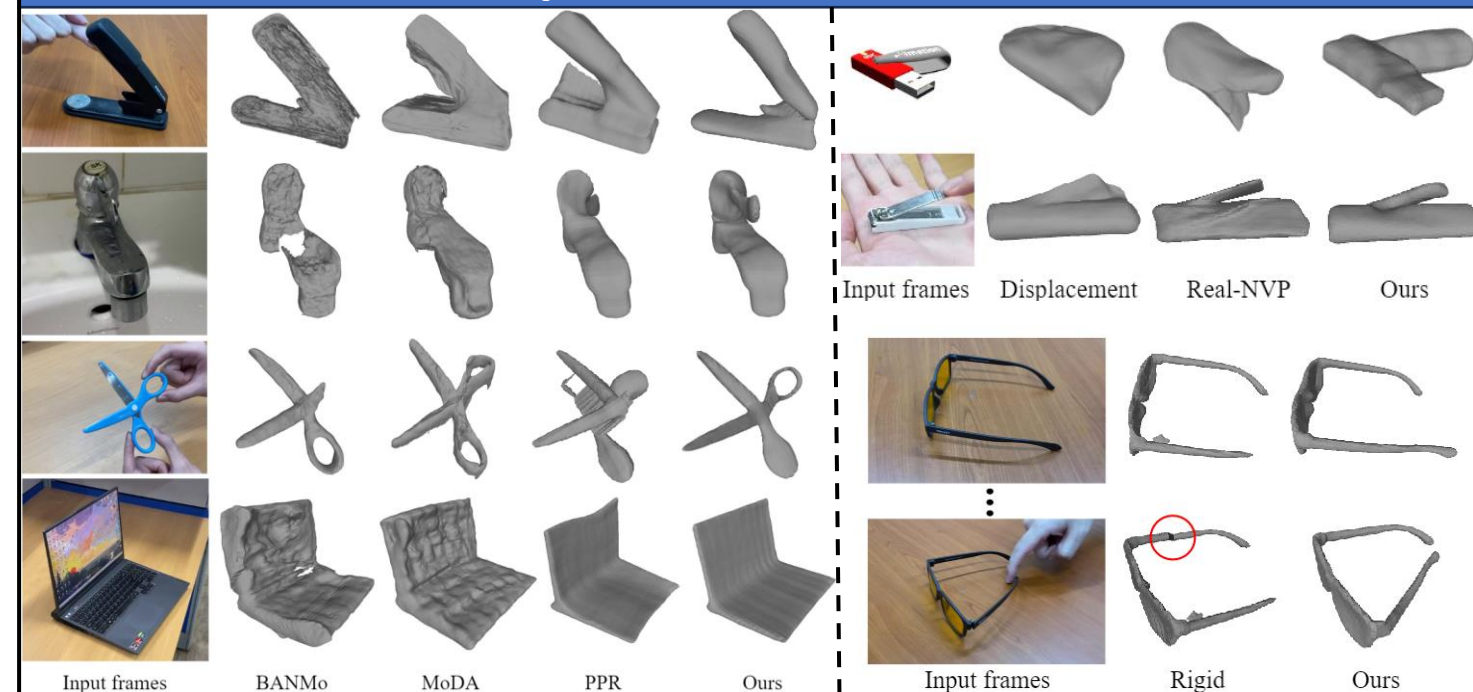


Geodesic distances between 3D point (black) and bones (blue and yellow). Shorter distances indicate stronger associations.

5. Quasi-rigid blend skinning.

$$\mathbf{X}(\psi_p) = \mathbf{T}_{\text{obj}} \mathbf{X} = \left(\sum_{b=0}^{B-1} w_b^S \mathbf{T}_b \right) \mathbf{X},$$

Experimental Results



Method	USB			stapler			scissors		
	CD(↓)	F(10%, ↑)	F(5%, ↑)	CD(↓)	F(10%, ↑)	F(5%, ↑)	CD(↓)	F(10%, ↑)	F(5%, ↑)
BANMo	20.3	65.1	45.0	19.1	57.8	32.8	19.9	66.8	41.4
MoDA	17.1	74.9	49.5	18.8	64.2	40.3	14.8	77.7	42.3
PPR	20.7	65.7	38.9	16.8	67.5	40.0	16.1	71.4	39.9
Ours	15.3	78.6	51.5	14.3	75.5	42.7	14.0	78.2	43.9

We introduce REACTO, a groundbreaking method for reconstructing general articulated 3D objects from single casual videos, achieving enhanced modeling and precision by redefining rigging structures and employing Quasi-Rigid Blend Skinning. QRBS ensures the rigidity of each component while retaining smooth deformation on the joints by utilizing quasi-sparse skinning weights and geodesic point assignment.

