

IAAO: Interactive Affordance Learning for Articulated Objects in 3D Environments

Supplementary Material

6. Appendix

6.1. 3D-2D Correspondence Matching

We start with computing the feature similarity matrix α_{ip} between each pixel in part mask o and the sampled Gaussian point p^t . We then normalize similarity matrix α_{ip} using a softmax across the entire mask to obtain the weight β_{ip} . Finally, we identify the 2D point $s_{p \rightarrow o}^{t'}(I_n^{t'})$ corresponding to the 3D point using a weighted sum. The computation steps are as follows:

1. Compute the feature distance:

$$\alpha_{ip} = \|F_o^{t'}(I_n^{t'})[u_i] - F_{3D,o}^t(p)\|_2,$$

between the i -th pixel u_i of $I_n^{t'}$ and sampled Gaussian point $g_p^t(o)$.

2. Normalize α_{ip} using a softmax across the entire image to obtain the weight:

$$\beta_{ip} = \frac{\exp(-s\alpha_{ip})}{\sum_{i=1}^{|M_o^{t'}(n)|} \exp(-s\alpha_{ip})}.$$

3. Identify the 2D point:

$$s_{p \rightarrow o}^{t'}(I_n^{t'}) = \sum_{i=1}^{|M_o^{t'}(n)|} \beta_{ip} u_i,$$

corresponding to the 3D point using a weighted sum. $F_o^{t'}(I_n^{t'})$ represents the DINOv2 features extracted from $I_n^{t'}$, and s is a hyperparameter that adjusts the smoothness of the heatmap β_{ij} .

6.2. Scene State Fusion

After obtaining the transformation $\xi^t = (s^t, R^t, T^t)$ (or its inverse function $\xi^{t'} = (s^{t'}, R^{t'}, T^{t'})$) for each part in scene state t , the next step is merging the two Gaussian Splatting (GS) models in the two states. We adopt the Gaussian splatting fusion and filtering strategy from [1]. To transform the Gaussians from the coordinate system of $G^{t'} = \{g_p^{t'}\}_{p=1}^{P^{t'}}$ to $G^t = \{g_p^t\}_{p=1}^{P^t}$, the position of each 3D Gaussian $g_p^{t'}$ is transformed as follows:

$$(x_p^{t' \rightarrow t}, y_p^{t' \rightarrow t}, z_p^{t' \rightarrow t})^\top = s^{t'} R^{t'}(x_p^{t'}, y_p^{t'}, z_p^{t'})^\top + T^{t'}.$$

The opacity remains unchanged during this transformation, i.e., $\alpha_p^{t' \rightarrow t} = \alpha_p^{t'}$. The rotation matrix $R_p^{t' \rightarrow t} \in \mathbb{R}^{3 \times 3}$ and scale $S_p^{t' \rightarrow t} \in \mathbb{R}^3$ are computed as:

$$R_p^{t' \rightarrow t} = R^{t'} R_p^{t'}, \quad S_p^{t' \rightarrow t} = s^{t'} S_p^{t'}.$$

Spherical harmonics (SH) coefficients undergo a linear transformation based on their rotation, which can be handled independently for each order. To this end, for any i -th order of SH coefficients, the following steps are performed:

1. Choose $2i + 1$ unit vectors u_0, \dots, u_{2i+1} and compute their corresponding SH coefficients as $Q = (\text{SH}(u_0), \dots, \text{SH}(u_{2i+1}))$.
2. Apply the transformation $\xi^{t'} = (s^{t'}, R^{t'}, T^{t'})$ to the vectors u_0, \dots, u_{2i+1} to obtain transformed vectors $\hat{u}_0, \dots, \hat{u}_{2i+1}$.
3. Compute the transformation matrix for SH coefficients as:

$$(\text{SH}(\hat{u}_0), \dots, \text{SH}(\hat{u}_{2i+1}))Q^{-1}.$$

Finally, the 3D Gaussians in $G^t = \{g_p^t\}_{p=1}^{P^t}$ closer to the center of scene t are merged with those in $G^{t'} = \{g_p^{t'}\}_{p=1}^{P^{t'}}$ near the center of scene t' , producing $G^{t+t'}$.

6.3. More Visualization

Fig. 7 and Fig. 8 show the qualitative results on the “PARIS Two-Part Object Dataset”. Column 4 shows the results of our IAAO and Column 3 shows the comparison with DigitalTwinArt. The input states at t and t' are shown in Columns 1 and 2. The part segments are shown in different colors and the red arrows indicate the joint prediction. Rotation is around the arrow and translation is along the arrow. Generally, as compared with DigitalTwinArt on both datasets, we can see from all figures that our IAAO produces better shape reconstruction with clearer details, more precise part segmentation with lesser erroneous labels, and more accurate joint predictions with correct motion directions. Fig. 9, 10 and 11 show snapshots of the generated motions on several examples at object and scene levels in a simulator. It shows IAAO can produce smooth interpolations of articulated objects with good part geometry modeling. Fig 12 illustrates the qualitative results on a sample scene from the Indoor Scene OmniSim dataset. Unlike existing baselines, which lack semantic meaning in their reconstructed neural fields, our model demonstrates the ability to perform object-level and fine-grained part localization based on prompts within complex indoor environments.

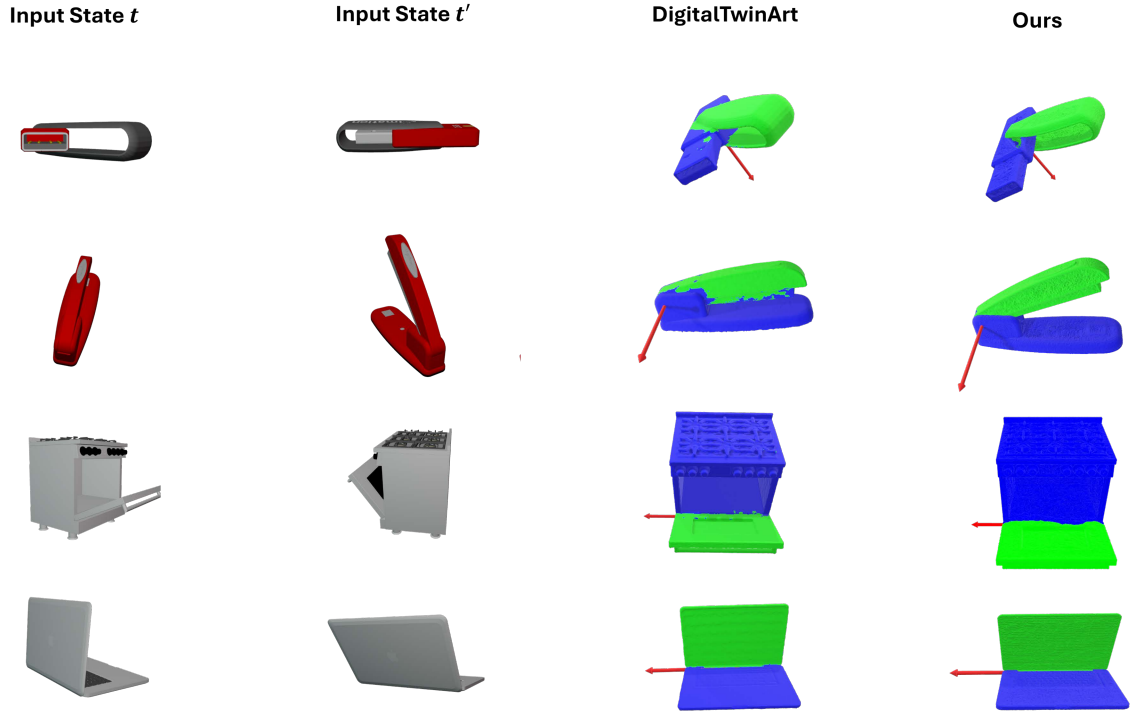


Figure 7. Qualitative analysis of shape reconstruction, part segmentation, and joint prediction results on the PARIS dataset. Our method shows better shape reconstructions and precise part segmentations on the thumbdrive and stapler shown in row 1 and 2.

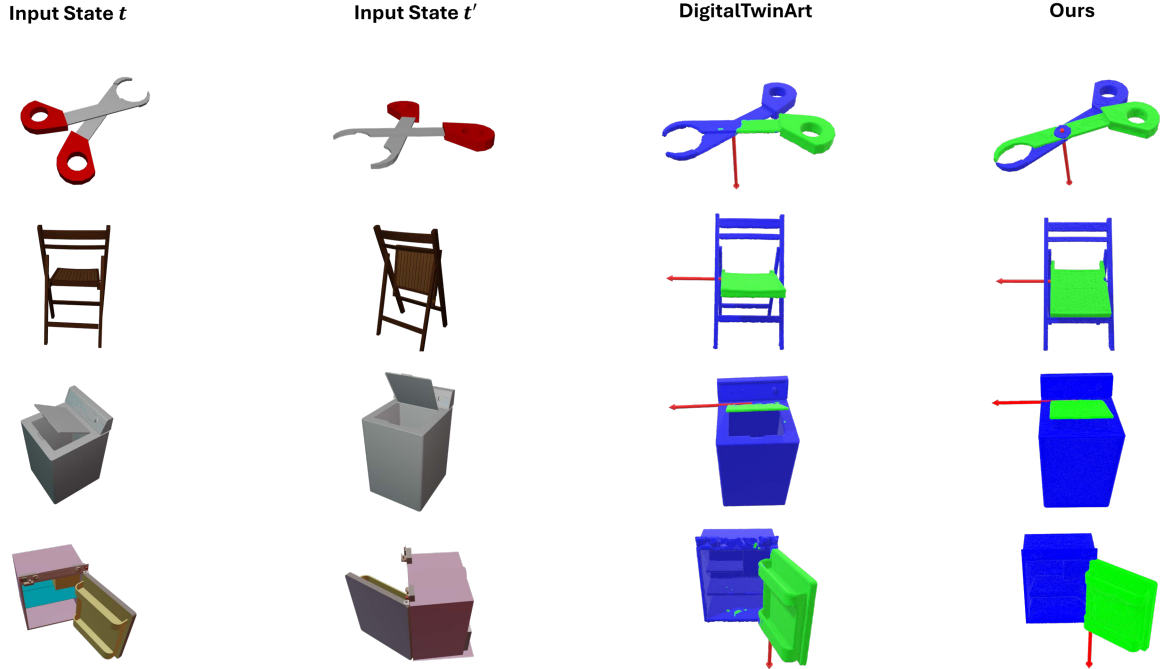


Figure 8. Qualitative analysis of shape reconstruction, part segmentation, and joint prediction results on the PARIS dataset. We can see that our method produces a better part segmentation result compared to DigitalTwinArt on the scissors shown in row 1.

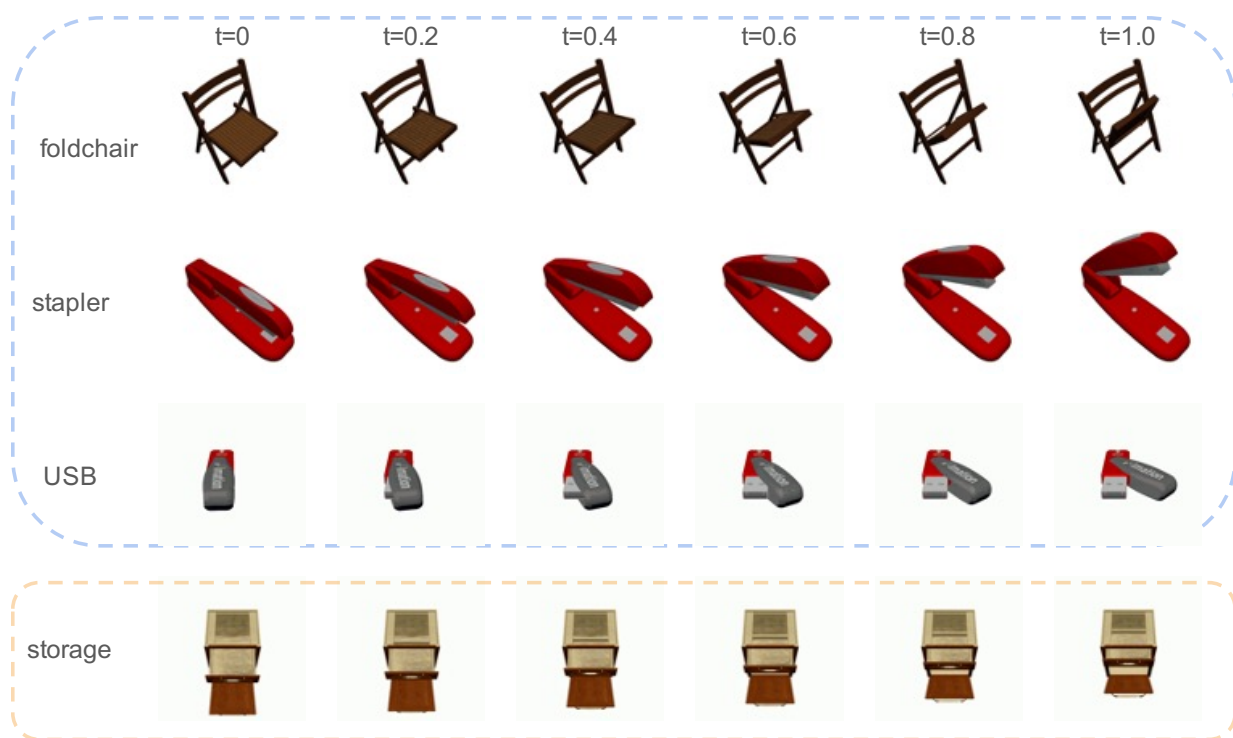


Figure 9. Qualitative analysis of scene interpolation on [PARIS](#) and [PartNet-Mobility](#) datasets.

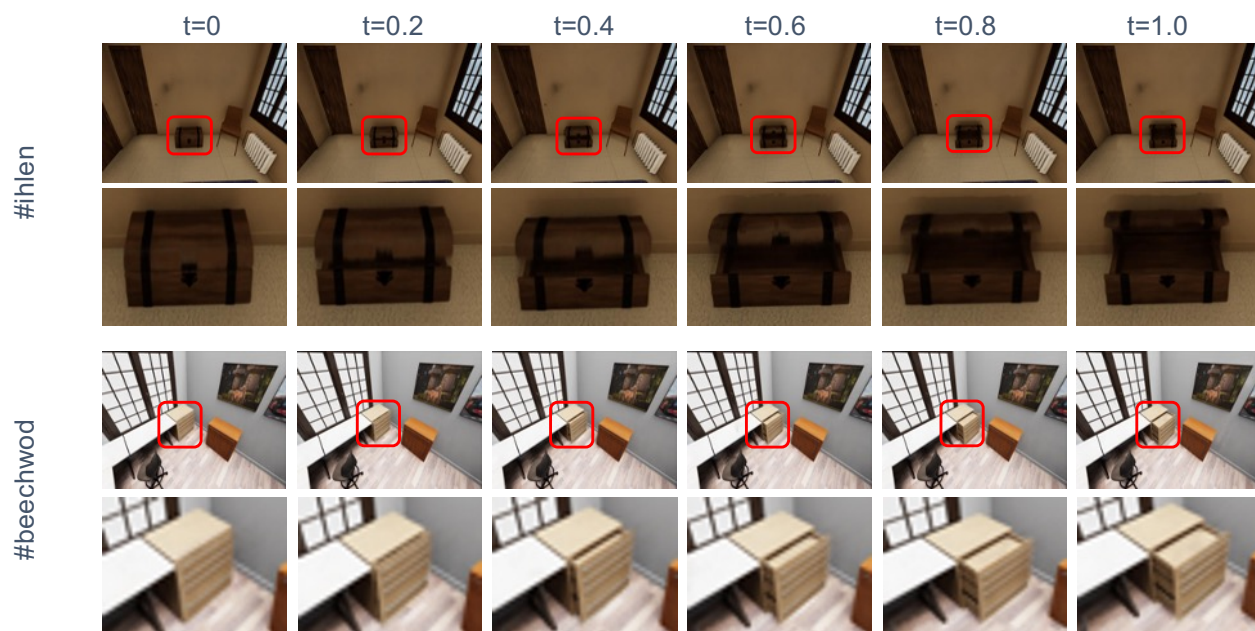


Figure 10. Motion snapshots on [#ihlen](#) and [#beechwod](#) from Indoor scene OmniSim dataset.



Figure 11. Motion snapshots on #merom and #wainscott from Indoor scene OmniSim dataset.

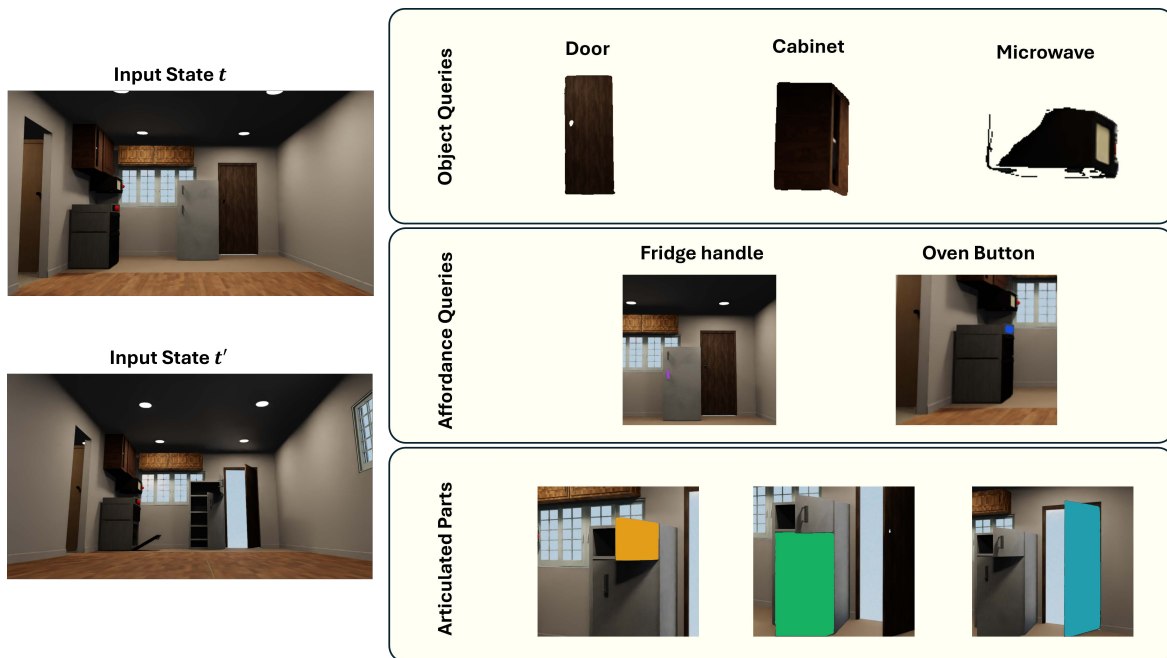


Figure 12. Qualitative analysis of object and affordance retrieval on one example scene from the Indoor scene OmniSim dataset.