

GHOST: Fast Category-agnostic Hand-Object Interaction Reconstruction from RGB Videos using Gaussian Splatting

Supplementary Material

In our supplementary material, we give more illustrations of different steps in our preprocessing pipeline in Figs. 8 and 9. Furthermore, we compare different design choices and their qualitative gains in Fig. 10. Fig. 11 shows quantitative evaluation of applying \mathcal{L}_{geo} on 5 sequences from the HO3D dataset [15]. We observe an improvement in the interaction distance relative to the right hand root (CD_r) on 3 sequences. Fig. 12 shows limitations of $\mathcal{L}_{bkg,h}$ and \mathcal{L}_{geo} . In addition, Appendix A and B discuss different hyperparameters in our approach. Finally, we show examples of our animatable hand avatar visualization inherited from GaussianAvatars [42] viewer in Fig. 13.

A. Postprocessing HaMeR Predictions

For image I_t with area A_{img} , RTMPose [22, 23] provides hand keypoints for left and right hands with confidence for each keypoint producing right and left hand bounding boxes (B_t^r, B_t^l) with areas (A_t^r, A_t^l) and averaged keypoint confidence (c_t^r and c_t^l). Using the hand bounding boxes, HaMeR [39] generates initial hand reconstructions as stated in Section 3.1.2. For timestep t , our algorithm decides based on the combined rejection rule in Eq. 20 for each hand $h \in \{r, l\}$ if the predictions of the frames should be discarded and interpolated or not. The individual conditions are defined as follows:

1. Pose jitter condition:

$$C_p = (\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}\|_2 > \tau_p \wedge \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t+1}\|_2 > \tau_p), \quad (13)$$

2. Orientation jitter condition:

$$C_o = (\|\mathbf{R}_t^h - \mathbf{R}_{t-1}^h\|_2 > \tau_o \wedge \|\mathbf{R}_t^h - \mathbf{R}_{t+1}^h\|_2 > \tau_o), \quad (14)$$

3. Translation jitter condition (x-y plane):

$$C_t = (\|\mathbf{T}_t^h - \mathbf{T}_{t-1}^h\|_2 > \tau_t \wedge \|\mathbf{T}_t^h - \mathbf{T}_{t+1}^h\|_2 > \tau_t), \quad (15)$$

4. Shape deviation condition:

$$C_s = \frac{|\beta_t - \text{median}(\beta)|}{\text{std}(\beta) + \varepsilon} > \tau_{shape}. \quad (16)$$

5. Confidence threshold:

$$C_c = c_t^h < \tau_{conf}, \quad (17)$$

6. Bounding-box area constraint:

$$C_a = (A_t^h < A_{\min} \vee A_t^h > A_{\max}), \quad (18)$$

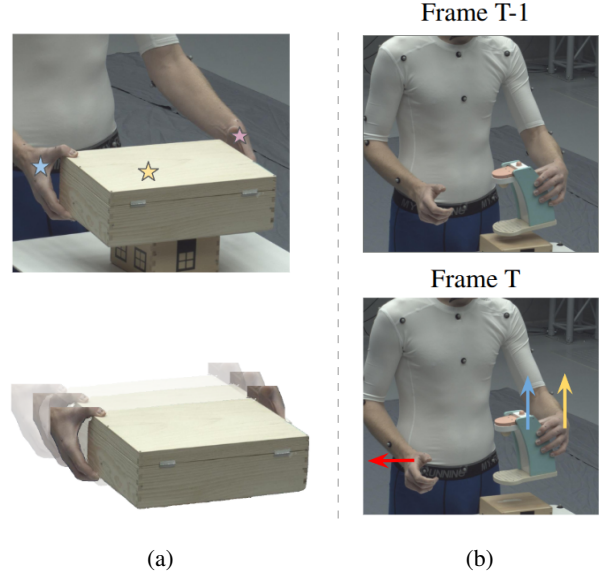


Figure 8. (a) SAM2 [43] is initialized with 3 seed pixels to segment and track the hands (\mathcal{M}_t^h) and the object (\mathcal{M}_t^o) in the scene. (b) During the grasping detection, the object’s motion vector \hat{T}_{xy}^o (blue arrow) is compared with left hand’s motion vector \hat{T}_{xy}^l (orange arrow) and right hand’s motion vector \hat{T}_{xy}^r (red arrow). The example shows a left hand grasp based on the similarity between motion vectors.

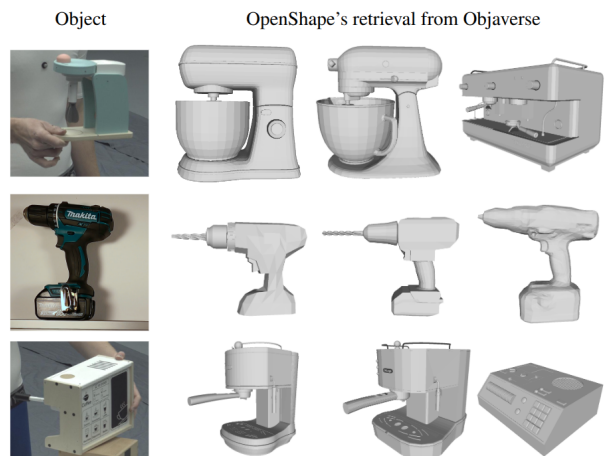
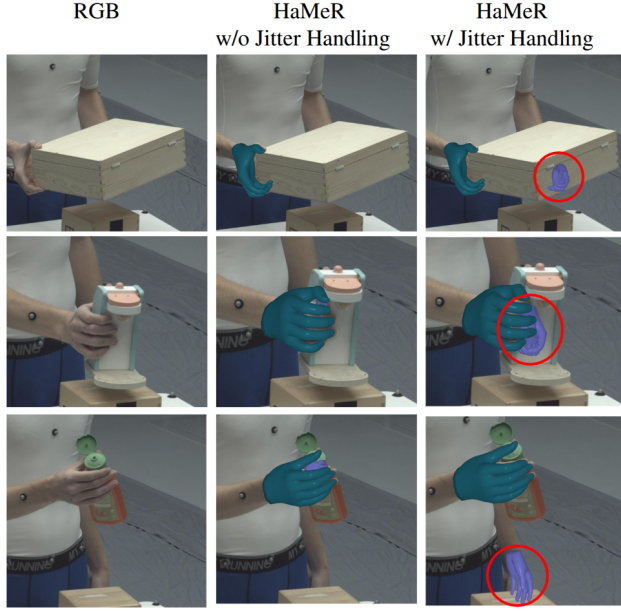
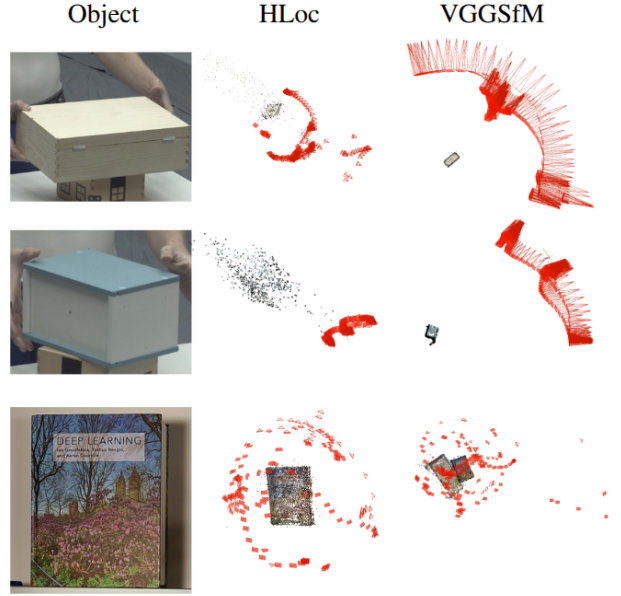


Figure 9. OpenShape [31] retrieves 3D models from Objaverse [10, 11], however, the retrieved 3D models do not always match the geometry of the desired object. Therefore, the final geometric prior \mathcal{O} can be suboptimal.



(a)



(b)

Figure 10. a) Initial hand reconstructions \mathcal{V}_t^h obtained from HaMeR [39] suffer from Jitter under occlusion. Detecting jitter based on temporal cues, detection confidence, and interpolation results in improving initial hand meshes \mathcal{V}_t^h . b) Structure-from-Motion (SfM) has a large impact on subsequent steps. VGGsFM [51] improved SfM when applied to the arctic data compared to the HLoc+COLMAP [46, 48] pipeline. However, VGGsFM is sensitive to hyperparameter selection and does not always show similar improvements as seen in the last row.

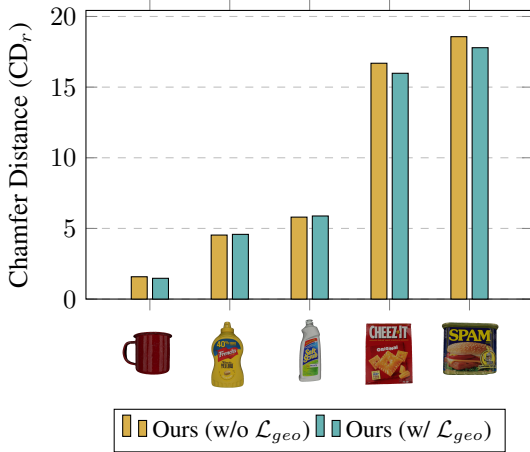


Figure 11. The influence of \mathcal{L}_{geo} (CD_r , lower is better) across five HO3D sequences using distinct YCB objects.

7. Bounding-box overlap (IoU) constraint:

$$\mathcal{C}_{iou} = IoU(B_t^r, B_t^l), \quad (19)$$

where the thresholds are empirically chosen as: $\tau_p = 1.0$, $\tau_o = 1.0$, $\tau_t = 2.0$, $\tau_s = 4.0$, $\tau_c = 0.3$, $\tau_{iou} = 0.3$, $A_{min} = 0.006 A_{img}$, and $A_{max} = 0.2 A_{img}$.



Figure 12. a) Applying $\mathcal{L}_{bkg,h}$ for object reconstruction fails when the hand never changes its contact point with the object. In that case, unwanted gaussians spawn in the hand region. b) In some cases, the retrieved geometric prior \mathcal{O} does not align perfectly with the initial object's point cloud \mathcal{P}_{sfm} (Middle column). The results of applying \mathcal{L}_{geo} in this case will result in moving gaussian centers c_o towards unwanted regions (see blue point cloud).

Hand rejection rule:

$$\mathcal{F}_{reject} = \mathcal{C}_p \vee \mathcal{C}_o \vee \mathcal{C}_t \vee \mathcal{C}_s \vee \mathcal{C}_c \vee \mathcal{C}_a \vee \mathcal{C}_{iou}. \quad (20)$$

Fig 10a shows the importance of this rejection rule on hand meshes \mathcal{V}_t^h .

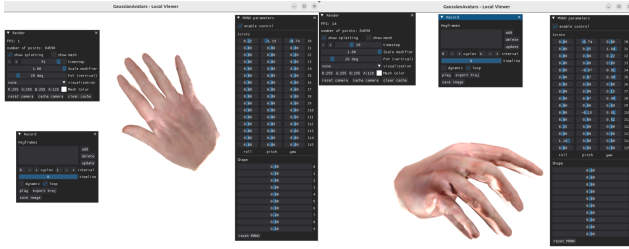


Figure 13. We also provide an interactive 3D viewer for Gaussian hand avatars. The interface visualizes and controls MANO-based [45] hand reconstructions, allowing users to adjust both pose and shape parameters in real time. Imported motion sequences can be played to animate the Gaussian hand avatar.

B. Experimental details

1. Prior Alignment Parameters. Optimizer: AdamW [33], LR: 10^{-2} , Betas: 0.9, 0.99, eps: 10^{-8} , Iterations: 1500.
2. HO Alignment. Optimizer: Adam [4], LR: 0.05, Iterations: 500.
3. Object Gaussian Splatting Optimization. Optimizer: Adam [4]. Iterations: 30000.
4. Hand-Object Gaussian Splatting Optimization. Optimizer: Adam [4]. Iterations: 30000. More details on the Gaussian Splatting hyperparameters are available in the code.