

CARI4D: Category-Agnostic 4D Reconstruction of Human-Object Interaction

Anonymous CVPR submission

Paper ID ***



Figure 1. **Category-agnostic 4D interaction reconstruction.** From a monocular RGB video, CARI4D reconstructs a metric-scale object mesh, tracks human and object motion over time, and predicts coherent hand-object contacts for categories unseen during training.

Abstract

001 Reconstructing 4D human-object interaction from monoc-
 002 ular RGB video would make large-scale interaction cap-
 003 ture possible outside controlled studios, but the task is ill
 004 posed: object geometry is unknown, depth and scale are
 005 ambiguous, and contacts are often heavily occluded. We
 006 present CARI4D, a category-agnostic framework that recon-
 007 structs a metric-scale object, estimates temporally consistent
 008 human and object motion, and reasons about hand-object
 009 contact from a single RGB video. CARI4D combines founda-
 010 tion models for object generation, metric depth, human
 011 pose, and object pose, but explicitly aligns their outputs

through coarse-to-fine scale selection, pose-hypothesis 012
 filtering, learned render-and-compare contact reasoning, 013
 and contact-aware joint optimization. On BEHAVE and zero-shot 014
 InterCap evaluations, CARI4D improves combined recon- 015
 struction Chamfer distance by more than 35% over prior 016
 video-based baselines, while also generalizing to in-the-wild 017
 videos with previously unseen object categories. 018

1. Introduction 019

Human-object interaction capture is central to human 020
 understanding, AR/VR, gaming, and robot learning. Existing 021

022 high-fidelity pipelines require calibrated multi-view studios
023 or RGBD rigs [2, 3, 7], which limits scale. Monocular RGB
024 videos are far easier to obtain, but metric 4D reconstruction
025 from them is difficult because object geometry is unknown,
026 depth is ambiguous, and contact-rich motion creates severe
027 occlusion.

028 Prior RGB-based video methods assume known object
029 templates or fixed training categories [17–19]. Founda-
030 tion models now provide strong object generation, pose
031 estimation, human reconstruction, and metric-depth pri-
032 ors [10, 11, 13, 14], but directly composing them is insuffi-
033 cient: their outputs lie in different coordinate systems and
034 do not satisfy interaction contacts.

035 We introduce CARI4D, a category-agnostic monocular
036 RGB framework for full-body 4D human-object interaction
037 reconstruction. CARI4D reconstructs the object and met-
038 ric scale, initializes human and object poses in a shared
039 space, refines them with the learned contact reasoning net-
040 work CoCoNet, and finally performs contact-aware joint
041 optimization. Our contributions are a template-free metric
042 4D pipeline, robust foundation-model integration, and con-
043 tact reasoning that improves zero-shot reconstruction and
044 contact coherence.

045 2. Related Works

046 **Foundation models for 3D reasoning.** Recent models re-
047 construct object shape [13, 15, 20], estimate object or human
048 pose [11, 14], and predict metric scene depth [9, 10]. They
049 provide strong priors but reason about humans, objects, and
050 scenes separately.

051 **Interaction reconstruction.** Image-based interaction meth-
052 ods estimate hand-object or full-body contacts from a sin-
053 gle frame [4, 5, 16, 21], while video trackers improve tem-
054 poral consistency but typically require object templates or
055 category-specific training [17, 19]. CARI4D instead recon-
056 structs the object, tracks human and object motion, and
057 reasons about contacts from monocular RGB video without
058 assuming a known category.

059 3. Method

060 Given RGB frames $\{\mathbf{I}_i\}_{i=1}^N$, our goal is to recover an object
061 mesh \mathbf{O} with per-frame 6DoF poses $\{\mathcal{O}_i = (\mathbf{R}_i^o, \mathbf{t}_i^o)\}_{i=1}^N$
062 and SMPL-H human parameters $\{\mathcal{H}_i = (\boldsymbol{\theta}_i, \boldsymbol{\beta}, \mathbf{t}_i^h)\}_{i=1}^N$ in
063 a single metric coordinate system. Figure 2 summarizes the
064 pipeline.

065 3.1. Metric-scale Object Reconstruction

066 We assume the object is mostly visible in the first frame, a
067 common setup for interaction videos. We segment it with
068 f-BRS [12] and reconstruct a normalized object mesh using
069 Hunyuan3D-2 [13]. To place the mesh in metric scale, we
070 estimate depth and intrinsics with UniDepth [10] and run

FoundationPose [14] over candidate object scales. For each
071 scale, the posed mesh is compared with the segmented depth
072 point cloud using a one-way Chamfer distance. We perform
073 a coarse search, refine around the top candidates, and keep
074 the scale with the lowest distance. 075

076 3.2. Human and Object Pose Initialization

077 With the metric object mesh fixed, FoundationPose can pro-
078 vide object pose candidates per frame, but its top-ranked
079 prediction is unreliable under monocular depth noise, im-
080 perfect generated geometry, and hand/body occlusion. We
081 therefore select among its K candidates using two filters.
082 First, each candidate is rendered and scored by mask IoU
083 against the observed object mask after removing human oc-
084 clusion. Second, candidates whose rotation changes too
085 abruptly from the previous accepted frame are rejected. If
086 no candidate survives during an occluded interval, we jump
087 forward to a reliable frame and track backward, applying the
088 same filters. For the human, we run NLF [11] and align its
089 global translation and scale to UniDepth point clouds so the
090 human and object share metric depth.

091 3.3. Contact Reasoning and Refinement

092 The initialization treats human and object separately, so
093 contacts can float, penetrate, or drift over time. We train
094 CoCoNet, a category-agnostic contact reasoning model, to
095 refine poses and predict binary left/right hand contact la-
096 bels. For a short frame window, we render the initialized
097 human and object into RGB, depth, and mask channels, then
098 compare these renderings against the observed RGB, depth,
099 and segmentation cues. A DINOv2-based encoder [8] ex-
100 tracts visual and geometric features, spatiotemporal attention
101 aggregates them across frames, and lightweight heads pre-
102 dict pose updates and contacts. During training, estimated
103 depth is aligned to ground-truth depth before generating
104 initializations, which prevents the network from learning
105 dataset-specific depth biases instead of interaction correc-
106 tions.

107 3.4. Contact-based Joint Optimization

108 The final stage optimizes the full sequence by minimizing a
109 weighted objective

$$L = \lambda_c L_c + \lambda_j L_{j2d} + \lambda_m L_m + \lambda_p L_{pen} + \lambda_a L_{acc}. \quad (1) \quad 110$$

111 L_c attracts contacted hand joints to the object when Co-
112 CoNet predicts contact, L_{j2d} and L_m preserve image evi-
113 dence, L_{pen} discourages human-object penetration, and L_{acc}
114 smooths motion.

115 4. Experiments

116 **Experimental setup and metrics.** We train on BEHAVE [1]
117 and HODome [22], then evaluate on BEHAVE and the un-
118 seen InterCap dataset [6]. Following prior work [4, 17], we

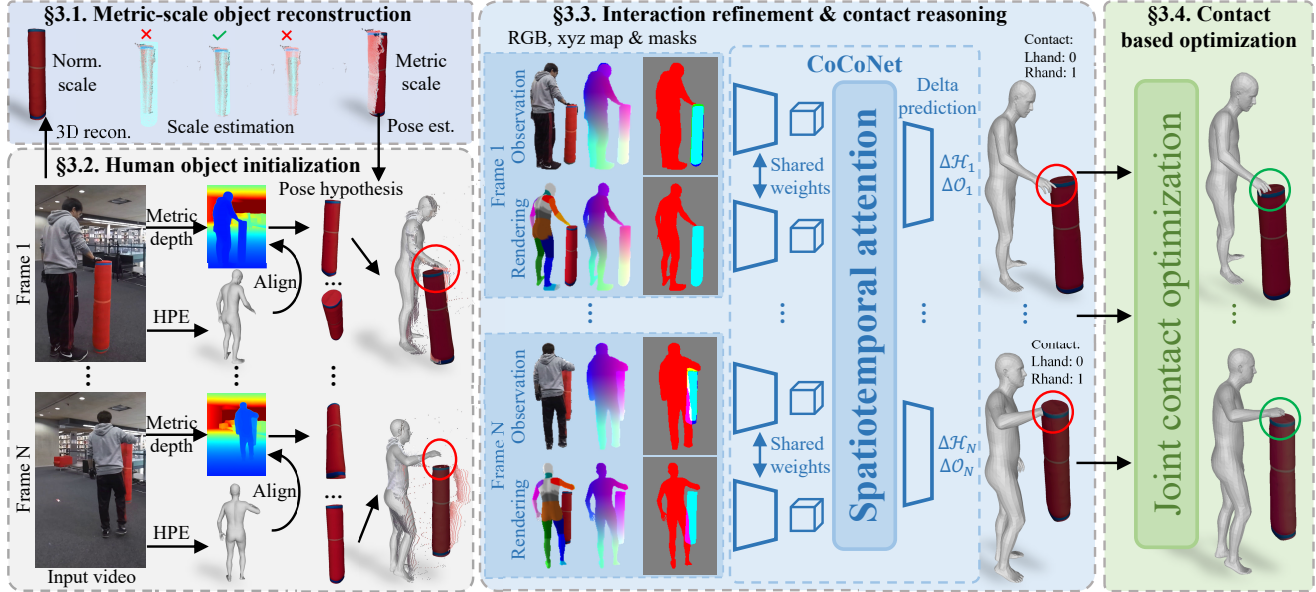


Figure 2. **CARI4D overview.** We reconstruct a metric object mesh, initialize human and object poses in a shared coordinate system, refine them with CoCoNet using render-and-compare contact reasoning, and perform contact-aware joint optimization.

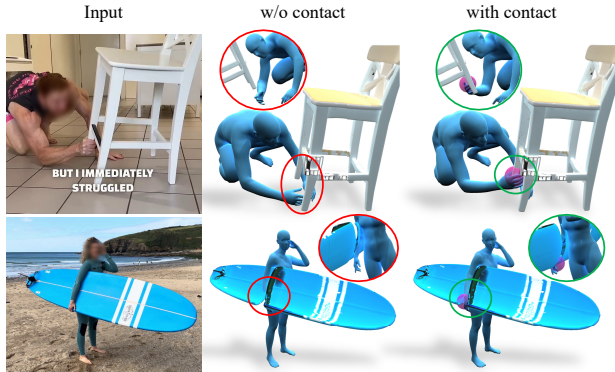


Figure 3. **Contact-aware refinement.** Predicted contacts guide optimization away from floating objects and penetration artifacts. report Chamfer distance in centimeters for the human (CD-h), object (CD-o), and combined mesh (CD-c), along with human and object acceleration errors. To preserve the metric 4D setting, we align only the first frame to ground truth and apply that single transformation to the full video.

119
120
121
122
123

124 4.1. Baseline Comparison

125 Table 1 compares against InterTrack and VisTracker. Inter-
126 Track reconstructs category-specific point clouds but pro-
127 duces inconsistent global translation, yielding high CD-c.
128 VisTracker has temporally smoother motion but requires
129 known object templates and struggles when adapted with
130 reconstructed meshes. CARI4D reconstructs the object di-
131 rectly from video, recovers coherent motion, and reduces

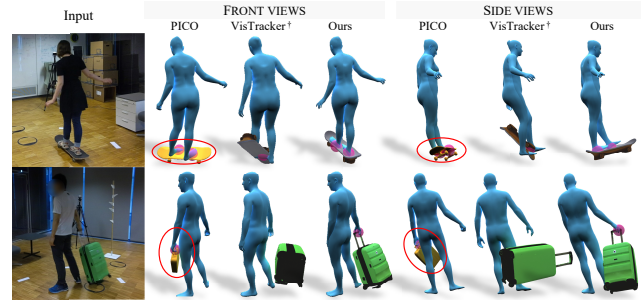


Figure 4. **Zero-shot InterCap examples.** CARI4D reconstructs unseen object geometry and temporally consistent poses without category-specific training or a provided object template.

CD-c from 14.22 to 9.23 on BEHAVE. On unseen InterCap
objects, CARI4D improves CD-c from 20.17 to 12.88 over
the strongest video baseline.

132
133
134

4.2. Zero-shot Generalization

135 Compared with image-based PICO [4] on InterCap key
136 frames, CARI4D reduces CD-c from 87.73 to 5.90 when
137 both are evaluated on the same sampled frames. Qualita-
138 tively, CARI4D also handles in-the-wild internet videos with
139 object categories not present in training, as shown in Fig-
140 ure 1.
141

4.3. Ablation Studies

142 Table 2 shows that vanilla NLF plus FoundationPose is
143 not sufficient in our RGB-only setting. Our pose-selection
144 and alignment stage yields a much stronger initialization,
145

142
143
144
145

Method	BEHAVE					InterCap zero-shot				
	CD-h	CD-o	CD-c	Acc-h	Acc-o	CD-h	CD-o	CD-c	Acc-h	Acc-o
InterTrack [19]	25.71	47.66	30.20	5.23	5.64	34.79	40.37	33.53	4.34	5.26
VisTracker [†] [17]	13.52	18.29	14.22	0.54	0.77	16.12	27.41	20.17	0.98	1.25
CARI4D	7.74	12.05	9.23	1.14	0.35	11.06	15.69	12.88	1.25	0.82

Table 1. **Video reconstruction results** (cm, lower is better). [†]VisTracker requires an object template, so we provide our reconstructed object mesh for this comparison. CARI4D improves combined Chamfer distance on both in-distribution BEHAVE and unseen InterCap.



Figure 5. **In-the-wild qualitative comparison.** Prior approaches can produce noisy object shapes, flipped poses, or incorrect contacts on internet videos. CARI4D maintains more coherent object geometry, pose, and contact predictions across diverse categories.

Variant	CD-h	CD-o	CD-c	Acc-h	Acc-o
Raw NLF + FP tracking	11.53	1565.42	405.13	3.06	4.34
Raw NLF + FP pose est.	11.53	40.54	13.02	3.06	9.26
Our initialization	7.81	16.85	10.79	2.21	8.36
+ CoCoNet	7.01	11.59	8.62	1.75	3.78
Full model	8.41	11.57	9.35	1.06	0.38

Table 2. **Ablations on BEHAVE subset.** Pose selection and alignment improve initialization; CoCoNet improves reconstruction; joint optimization improves motion smoothness and contact coherence.

146 CoCoNet further improves Chamfer distance, and the full
 147 contact-aware optimization produces the smoothest object
 148 motion. Figure 3 illustrates how contact reasoning removes
 149 common floating and penetration artifacts.

150 5. Conclusion

151 We presented CARI4D, a category-agnostic framework for
 152 metric 4D reconstruction of human-object interaction from

monocular RGB video. By aligning foundation-model pre-
 153 dictions, reasoning about contacts with CoCoNet, and enforcing
 154 contact-aware sequence optimization, CARI4D improves
 155 reconstruction accuracy on BEHAVE, generalizes zero-shot
 156 to InterCap, and applies to in-the-wild videos without a pro-
 157 vided object template.
 158

159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214**References**

- [1] Bharat Lal Bhatnagar, Xianghui Xie, Ilya Petrov, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Behave: Dataset and method for tracking human object interactions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [2] Derek Bradley, Tiberiu Popa, Alla Sheffer, Wolfgang Heidrich, and Tamy Boubekeur. Markerless garment capture. *ACM Trans. Graph.*, 27(3):1–9, 2008. 2
- [3] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Trans. Graph.*, 34(4), 2015. 2
- [4] Alpár Cseke, Shashank Tripathi, Sai Kumar Dwivedi, Arjun Lakshminpathy, Agniv Chatterjee, Michael J. Black, and Dimitrios Tzionas. PICO: Reconstructing 3D people in contact with objects. In *CVPR (CVPR)*, pages 1783–1794, 2025. 2, 3
- [5] Yana Hasson, Gül Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. Learning joint reconstruction of hands and manipulated objects. In *CVPR*, 2019. 2
- [6] Yinghao Huang, Omid Taheri, Michael J. Black, and Dimitrios Tzionas. InterCap: Joint markerless 3D tracking of humans and objects in interaction. In *German Conference on Pattern Recognition (GCPR)*, pages 281–299. Springer, 2022. 2
- [7] Hanbyul Joo, Tomas Simon, and Yaser Sheikh. Total capture: A 3d deformation model for tracking faces, hands, and bodies. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8320–8329, 2018. 2
- [8] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. Dinov2: Learning robust visual features without supervision, 2024. 2
- [9] Luigi Piccinelli, Yung-Hsu Yang, Christos Sakaridis, Mattia Segu, Siyuan Li, Luc Van Gool, and Fisher Yu. UniDepth: Universal monocular metric depth estimation. In *CVPR (CVPR)*, 2024. 2
- [10] Luigi Piccinelli, Christos Sakaridis, Yung-Hsu Yang, Mattia Segu, Siyuan Li, Wim Abbeloos, and Luc Van Gool. UniDepthV2: Universal monocular metric depth estimation made simpler, 2025. 2
- [11] István Sáráncsi and Gerard Pons-Moll. Neural localizer fields for continuous 3d human pose and shape estimation. 2024. 2
- [12] Konstantin Sofiiuk, Iliia Petrov, Olga Barinova, and Anton Konushin. f-brs: Rethinking backpropagating refinement for interactive segmentation. In *CVPR*, pages 8623–8632, 2020. 2
- [13] Tencent Hunyuan3D Team. Hunyuan3d 2.0: Scaling diffusion models for high resolution textured 3d assets generation, 2025. 2
- [14] Bowen Wen, Wei Yang, Jan Kautz, and Stan Birchfield. FoundationPose: Unified 6D Pose Estimation and Tracking of Novel Objects. In *CVPR*, 2024. 2
- [15] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen, Xin Tong, and Jiaolong Yang. Structured 3d latents for scalable and versatile 3d generation. In *CVPR (CVPR)*, 2025. Spotlight. 2
- [16] Xianghui Xie, Bharat Lal Bhatnagar, and Gerard Pons-Moll. Chore: Contact, human and object reconstruction from a single rgb image. In *European Conference on Computer Vision (ECCV)*. Springer, 2022. 2
- [17] Xianghui Xie, Bharat Lal Bhatnagar, and Gerard Pons-Moll. Visibility aware human-object interaction tracking from single rgb camera. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 4
- [18] Xianghui Xie, Bharat Lal Bhatnagar, Jan Eric Lenssen, and Gerard Pons-Moll. Template free reconstruction of human-object interaction with procedural interaction generation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 230
- [19] Xianghui Xie, Jan Eric Lenssen, and Gerard Pons-Moll. Intertrack: Tracking human object interaction without object templates. 2024. 2, 4
- [20] Yuxuan Xue, Xianghui Xie, Riccardo Marin, and Gerard Pons-Moll. Human 3diffusion: Realistic avatar creation via explicit 3d consistent diffusion models. In *Arxiv*, 2024. 232
- [21] Yufei Ye, Poorvi Hebbbar, Abhinav Gupta, and Shubham Tulsiani. Diffusion-guided reconstruction of everyday hand-object interaction clips. In *ICCV*, 2023. 233
- [22] Juze Zhang, Haimin Luo, Hongdi Yang, Xinru Xu, Qianyang Wu, Ye Shi, Jingyi Yu, Lan Xu, and Jingya Wang. Neural-dome: A neural modeling pipeline on multi-view human-object interactions. In *CVPR*, 2023. 234