

# HOI-Blender: A Unifying Blender Add-on for Standardization and Visualization of Diverse Human-Object Interaction Datasets

Anonymous ICCV submission

Paper ID \*\*\*\*

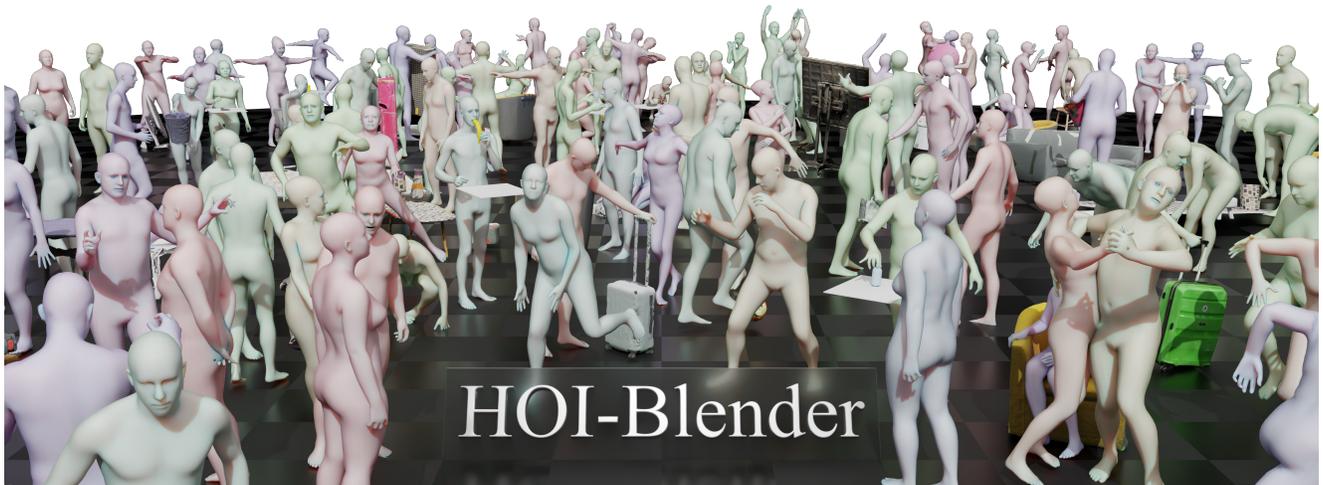


Figure 1. An intersection for interaction: HOI-Blender takes position as a consolidating toolkit, standardizing heterogeneous datasets such that actions may become mutual actors. Coexistence of diverse motions within coherent scenes is facilitated via acquired unified formatting, assisting conjoined visualization, equitable evaluation and aggregated downstream application.

## Abstract

001 Human-Object Interaction (HOI) datasets differ in coordi-  
 002 nate frames, formats, data structures, and directory lay-  
 003 outs, forcing bespoke loaders and brittle dataset-specific  
 004 rendering setups. Thus, we present HOI-Blender, a Blender  
 005 add-on which standardizes HOI data as a unified loader:  
 006 it normalizes coordinates and scale, harmonizes metadata,  
 007 and maps motion parameters to skinned bodies and object  
 008 meshes across congruent scenes. Once standardized, cam-  
 009 era rigs, lighting composition and render settings are reap-  
 010 plicable across datasets, enabling one-click import, anima-  
 011 tion preview, and image rendering. The add-on further de-  
 012 couples motion and appearance, allowing the same motion  
 013 to drive different human identities—or the same identity to  
 014 enact different motions. Additionally, to streamline anno-  
 015 tation, HOI-Blender includes an auto-captioning module  
 016 which forwards rendered frames to Vision-Language Mod-  
 017 els to produce action-object descriptions, supporting rapid  
 018 weak labeling and dataset curation. We demonstrate sup-

port for an initial 15 HOI public datasets and report a  
 comprehensive qualitative evaluation spanning motion fi-  
 delity, mesh integrity, hand-object interactions, and inter-  
 mesh collisions. By adapting heterogeneous HOI resources  
 into coherent scenes, HOI-Blender reduces implementation  
 overhead and facilitates efficient visualization, evaluation  
 and culminating application in human-centric research.

019  
020  
021  
022  
023  
024  
025

## 1. Introduction

Human-Object Interaction (HOI) describes the interplay be-  
 tween human and object motion. Improvement upon its un-  
 derstanding is core to advancing human-centric vision, as  
 it ties into an extensive range of applications, including as-  
 sistive robotics [2], scene understanding [9], synthetic data  
 generation [35], and behavioral modeling [12]. More distin-  
 ctly, character design and animation [7], augmented and  
 virtual reality [22], sports analytics [21], and human-agent  
 interaction [23]. Despite the rapid growth in count of HOI  
 datasets, working with motion data from various sources

026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036

037 within established communal tools remains a significant  
038 bottleneck. Issues consistently stem from the heterogeneity  
039 of dataset publications, where differences include variabil-  
040 ity in conventions, data formatting, distribution structures  
041 and a multitude of individual oddities. Hence, researchers  
042 are often forced to implement dataset-specific importers,  
043 manually assemble scenes with unique setup and properties,  
044 and deal with undesired complications, which slow progress  
045 and restrict accessibility of datasets for downstream usage.

046 To bridge the gap between HOI datasets and their practi-  
047 cality, we introduce HOI-Blender, a Blender add-on which  
048 transforms the aforementioned fragmented, repetitive per-  
049 dataset workflow into one seamless, single-click and stan-  
050 dardized. HOI-Blender applies dataset-specific preprocess-  
051 ing and affords a multi-purpose unified interface such that  
052 all dataset handling is made equal, eliminating any need for  
053 one-off scripts. Human and object motion data is standard-  
054 ized into an internal universal structure, enabling ease of ap-  
055 plication onto skinned bodies and object meshes, together  
056 assembled across coherent scenes. This underpins shared  
057 camera, lighting and rendering setups across datasets, elimi-  
058 nating need for manual per-dataset adjustments. Further, vi-  
059 sualization flexibility is enhanced via decoupled human mo-  
060 tion and appearance, allowing cross-combination from di-  
061 verse sources without additional configuration. Lastly, inte-  
062 gration for automatic captioning forwards rendered frames  
063 to OpenAI-compatible API for action-object descriptions,  
064 facilitating rapid weak labeling and dataset curation.

065 Beyond support for standardized HOI dataset import, as  
066 well as integrated visualization, rendering and motion cap-  
067 tioning capability, across an initial 15 HOI public datasets,  
068 we present a comprehensive qualitative analysis of key us-  
069 ability metrics, covering motion quality, mesh fidelity, inter-  
070 action plausibility, and collision robustness. Our evaluation  
071 underscores a type of efficiency and broad accessibility pro-  
072 vided by HOI-Blender over prior ad-hoc approaches, iden-  
073 tifiable in further cumulative applications.

074 By lowering the technical barriers of accessibility to  
075 HOI assets, HOI-Blender supports many diverse avenues  
076 of HOI research and workflows. Some notable examples  
077 include (i) rapid visualization for inspection, correction,  
078 and qualitative evaluation across HOI datasets for bench-  
079 marking and reproducible comparison, (ii) integration into  
080 Blender, a widely adopted, general-purpose 3D creation en-  
081 vironment, enabling its full functionality to be leveraged for  
082 interactive assembly, manipulation, and refinement of HOI  
083 scenes for development of 3D HOI datasets [33], (iii) ex-  
084 tending of prior synthetic rendering pipelines from human-  
085 only datasets [4, 26] to high-quality rendering of generated  
086 HOI scenes for construction of image-to-HOI datasets and  
087 further various image-based synthetic training data, (iv) in-  
088 corporation with physics-based simulation for development  
089 of physically grounded HOI data or robotic policy learning.

Our work culminates upon the following contributions: 090

- A Blender add-on which standardizes HOI data from an 091  
initial 15 datasets into a common universal motion param- 092  
eter format, enabling one-click import for skinned bodies 093  
and object meshes within coherent scenes. High-quality 094  
visualization capabilities are enhanced via decoupled mo- 095  
tion and appearance for human meshes. 096
- A streamlined workflow for consistent animation preview, 097  
evaluation, and rendering with lightweight scene assem- 098  
bly and cross-dataset reapplicable preset spanning camera 099  
rigging, lighting setup and render settings; where frames 100  
may be forwarded for automatic action-object captioning. 101
- A comprehensive qualitative evaluation of standardized 102  
HOI datasets for motion fidelity, mesh integrity, inter- 103  
action plausibility, and collision accuracy, emphasizing 104  
manner of support for future HOI workflows and diverse 105  
downstream applications. 106

## 2. Related Works 107

### 2.1. HOI Datasets 108

HOI datasets demonstrate heterogeneity over a manner of 109  
characteristics, one of which being scope of action, rang- 110  
ing from isolated body motion to full-body, human-object 111  
or multi-human interaction. Their uniting factor (of those 112  
in our interest) is a common foundation on SMPL-Family 113  
body models (SMPL, SMPL-H, SMPL-X), however, over- 114  
lap does not pertain further across conventions, data for- 115  
mats, hierarchical structures or otherwise. Our standardiza- 116  
tion considers grouping of representative datasets by inter- 117  
action type, which is delineated alongside supported HOI 118  
datasets below, with high-level factors outlined in Table 1. 119

**Isolated Human Motion.** AMASS [20] and AIST++ 120  
[15] provide large-scale full-body motion sequences with- 121  
out paired objects. AMASS unifies motion capture of multi- 122  
ple sources onto common SMPL-Family parameters, while 123  
AIST++ focuses on multi-view dance recordings with syn- 124  
chronized music. These datasets are widely used as motion 125  
priors for retargeting, generation, and animation control. 126

**Human-Human Interaction.** Hi4D [37] and Duolando 127  
[30] capture dyadic motion via purposing of registered 128  
SMPL-Family bodies. Hi4D emphasizes close-contact mo- 129  
tion with dense 4D segmentation, while Duolando con- 130  
tributes paired leader-follower dance sequences highlight- 131  
ing social coordination. Together they support learning so- 132  
cial motion priors, role assignment, and contact reasoning. 133

**Full-Body Human-Object Interaction.** BEHAVE [3], 134  
InterCap [8], COUCH [40], OMOMO [14], NeuralDome 135  
[38], and IMHD<sup>2</sup> [41] provide motion parameters for 136  
aligned SMPL-Family human bodies and mid-sized object 137  
meshes. They differ in address: BEHAVE and NeuralDome 138  
offer multi-object coverage with detailed meshes and con- 139  
tacts, InterCap and IMHD<sup>2</sup> explore monocular or inertial 140

Table 1. High-Level Comparison of SMPL-Family HOI Motion Datasets.

Motion Type	Dataset	# Objects	Object Category	Object Texture	Body Model	Published at
Isolated Human Motion	AIST++ [15]	-	-	-	SMPL	ICCV 2021
	AMASS [20]	-	-	-	SMPL-H/-X	ICCV 2019
Human-Human Interaction	Hi4D [37]	-	-	-	SMPL	CVPR 2023
	Duolando [30]	-	-	-	SMPL-X	ICLR 2024
Full-Body Human-Object Interaction	Behave [3]	20	Mid-Size	✓	SMPL-H	CVPR 2022
	InterCap [8]	10	Mid-Size	✓	SMPL-X	GCPR 2022
	COUCH [40]	4	Chairs	✗	SMPL-H	ECCV 2022
	OMOMO [14]	15	Mid-Size	✗	SMPL-X	SGAsia 2023
	NeuralDome [38]	23	Mid-Size	✓	SMPL-X	CVPR 2023
	IMHD <sup>2</sup> [41]	10	Mid-Size	✓	SMPL-H	CVPR 2024
Hand-Object Interaction	GRAB [31]	51	Small-Size	✗	SMPL-X	ECCV 2020
	Arctic [5]	11	Small-Size	✓	SMPL-X	CVPR 2023
	HIMO[19]	53	Small-Size	✗	SMPL-X	ECCV 2024
Multi-Human Object Interaction	HOI-M3 [39]	90	Large-Size	✓	SMPL	CVPR 2024
	CORE4D[17]	6	Large-Size	✗	SMPL-X	CVPR 2025

141 capture for tracking, COUCH specializes in seated interactions, and OMOMO emphasizes object-conditioned motion  
 142 interactions, and OMOMO emphasizes object-conditioned motion  
 143 synthesis. These datasets enable studies of contact, affordances, and perception in realistic interaction settings.  
 144

145 **Hand-Object Interaction.** GRAB [31], ARCTIC [5],  
 146 and HIMO [19] focus on dexterous manipulation with detailed hand motions upon small objects. GRAB emphasizes  
 147 full-body grasps across a wide object set, ARCTIC captures  
 148 bimanual articulated object manipulation, and HIMO extends its regard of objects across diverse scenarios. In conjunction,  
 149 these datasets are well suited for grasp synthesis, contact modeling, and fine-grained hand pose estimation.  
 150  
 151  
 152

153 **Multi-Human Object Interaction.** HOI-M3 [39] and  
 154 CORE4D [17] extend to collaborative multi-actor interaction scenes with large object coordination in shared environments.  
 155 They promote analysis of role assignment, group dynamics, and cooperative manipulation.  
 156  
 157

158 Despite shared SMPL-Family foundation, divergence in  
 159 dispense of HOI datasets is prominent across many facets. Structural variation span file packaging (per-frame vs. bundled  
 160 sequences; differed containers: .npz, .pk1, .json, .p), path conventions (absolute vs. relative references to objects,  
 161 textures), metadata allocation (lack-of, co-located with pose data, scattered into auxiliary files), dictionary key consistency  
 162 (trans vs. translation, poses vs. body\_pose), as well as parameter grouping (full-body vectors vs. assorted sets  
 163 of orientation, body, hands, face). Objects provide irregularities via varied mesh formats (.obj, .fbx, .g1b, .ply) and  
 164 ranging motion schemas from static meshes with implicit alignment to per-frame transforms at partial attributes; generally,  
 165 common denomination is deficient. Representational discrepancies further compound upon structural issues. Coordinate  
 166 frames differ across up-axis and handedness (Y-up vs. Z-up; left-handed vs. right-handed bases), units in scale  
 167  
 168  
 169  
 170  
 171  
 172  
 173  
 174

(m, cm, arbitrary), and rotation in encoding (axis-angle, Euler  
 175 triples, rotation matrices, quaternions, 6D; potentially varied  
 176 in-dataset). Body root pivots are inconsistent (pelvis-centered,  
 177 floor-aligned, world origin), as are object pivots (centroid,  
 178 mesh bottom, CAD default). Object representations range from  
 179 rigid single-part meshes to multi-part articulated assemblies  
 180 with divergent parent-child conventions, and texturization  
 181 fluctuates (vertex color, UV-mapped textures, external references).  
 182 HOI-Blender precisely bridges all aforementioned gaps between  
 183 datasets, deterministically resolving incongruities, and yielding  
 184 an internal standardized structure primed for purpose throughout contexts.  
 185  
 186

## 2.2. 3D Human-Centric Visualization Tools 187

188 While our focus lies in unifying HOI datasets, this challenge  
 189 sits within a broader landscape of human-centric, interactive  
 190 visualization tools, where prior efforts have addressed related  
 191 but narrower aspects of body model motion.

192 **Viser** [36] is a general-purpose 3D visualization library,  
 193 which provides an imperative Python API and a web-based viewer,  
 194 exposing scene primitives and GUI controls. SMPL mesh  
 195 visualization is supported via pose and shape parameter sliders,  
 196 enabling rapid debugging and interactive exploration. Its design  
 197 simplifies embedding in research pipelines and sharing browser-based  
 198 viewers.

199 **AIT-Viewer** [11] spans to broader coverage, supporting  
 200 SMPL, SMPL-H, SMPL-X, MANO, FLAME, and STAR. Both GUI  
 201 and headless rendering modes are available, allowing either  
 202 manual posing or sequence loading. Its applicability for inspection  
 203 and visualization thus particularly stands out when working with  
 204 diverse body model types.

205 **Meshcapade’s SMPL Blender Add-on** [25] enables the  
 206 import of SMPL, SMPL-H, SMPL-X, and SUPR models, material  
 207 application, shape and pose parameter adjustment,

208 as well as loading of AMASS animation sequences. While  
209 it brings parametric body manipulation into Blender, its  
210 scope remains focused on the body models themselves.

211 These established toolkits provide essential accessibil-  
212 ity for visualization, inspection, and interaction with para-  
213 metric human bodies, yet their functionality is largely con-  
214 fined to sliders and rudimentary sequence playback. Our  
215 work is complementary: HOI-Blender integrates standard-  
216 ized dataset import, which yields coherent scene assembly  
217 alongside consistent preferences, cameras, lighting, materi-  
218 als and further prospective desired arrangement. Such con-  
219 solidation transforms various fragmented resources into in-  
220 terchangeable foundations for diverse downstream applica-  
221 tions, acting as a bridge across datasets, eliminating ad-hoc  
222 setup and enabling analyses, augmentations, and simula-  
223 tions to be repeated seamlessly upon a one-click workflow.

### 224 3. Background

225 Human-Object Interaction (HOI) can be viewed as the time-  
226 varying spatial configuration of articulated human bodies  
227 and rigid objects, both represented as 3D meshes. We func-  
228 tion upon the SMPL-Family [16, 18, 27, 29] of parametric  
229 bodies to model humans, which efficiently encodes identity  
230 via shape parameters and motion via joint rotations, yield-  
231 ing skinned, watertight meshes suitable for animation and  
232 rendering. Objects are ordinarily provided as mesh assets  
233 comprised of geometry and material, occasionally with co-  
234 hesive part-level kinematics. Such mesh-centric representa-  
235 tion aligns both with common datasets and with established  
236 frameworks, enabling direct motion control via joint trans-  
237 forms for humans and rigid-body transforms for objects.

#### 238 3.1. Parametric Body Model

239 The SMPL-Family is a set of widely adapted and applied  
240 parametric body models for representing human shapes and  
241 poses. Its included model types begin with a canonical tem-  
242 plate mesh  $\bar{\mathbf{T}} \in \mathbb{R}^{3N}$  at  $N$  vertices, and subsequently in-  
243 corporate both shape blend shapes  $B_S : \mathbb{R}^{|\beta|} \mapsto \mathbb{R}^{3N}$  (con-  
244 trolled by shape parameters  $\beta \in \mathbb{R}^{|\beta|}$ ) as well as pose blend  
245 shapes  $B_P : \mathbb{R}^{|\theta|} \mapsto \mathbb{R}^{3N}$  (controlled by pose parameters  
246  $\theta \in \mathbb{R}^{|\theta|}$ , at  $K$  joints) as additional vertex offsets:

$$247 \begin{aligned} T_C(\beta, \theta) &= \bar{\mathbf{T}} + B_S(\beta) + B_P(\theta), \\ B_S(\beta; \mathcal{S}) &= \sum_{n=1}^{|\beta|} \beta_n \mathbf{S}_n, \\ B_P(\theta; \mathcal{P}) &= \sum_{n=1}^{9K} R_n(\theta) \mathbf{P}_n, \end{aligned} \quad (1)$$

248 where  $\mathbf{S}_n \in \mathbb{R}^{3N}$  and  $\mathbf{P}_n \in \mathbb{R}^{3N}$  are the respective shape  
249 and pose blend shape deformation matrices, and  $R_n(\theta) \in \mathbb{R}$   
250 is the matching rotation matrix element of the pose  $\theta$ . The  
251 resulting additive mesh  $T_C \in \mathbb{R}^{3N}$  is in canonical space,

and is to be transformed into deformation space  $T_P \in \mathbb{R}^{3N}$   
via Linear Blend Skinning (LBS) with joint rotation:

$$T_P = \left( \sum_{i=1}^K w_i \mathbf{G}_i(\theta) \right) T_C, \quad (2)$$

where  $\mathbf{G}_i(\theta) \in \mathbb{R}^{4 \times 4}$  is the transformation matrix of joint  $i$   
and  $w_i \in \mathbb{R}^N$  are the respective skinning weights. Please  
refer to Sec. 3 and Fig. 3 of SMPL [18] for more details on  
the parametric blend shapes and skinning.

#### 259 3.2. Human Appearance Model

260 As the base SMPL-Family model types provide plain mesh,  
261 they well represent a naked human equivalent by specified  
262 influence via pose  $\theta$  and shape  $\beta$  parameters, however, hu-  
263 man appearance such as coloration, hair or clothing are not  
264 taken into consideration. Following ClothCap [28], we rep-  
265 resent clothed humans by applying adjunct displacement.  
266 Notably, cloth blend shapes  $B_C \in \mathbb{R}^{3N}$  are appended in  
267 canonical space, adding to the canonical naked-body mesh:

$$T_C(\beta, \theta) = \bar{\mathbf{T}} + B_S(\beta) + B_P(\theta) + B_C. \quad (3)$$

269 To enable coloration, artist-designed UV-spaces  $\mathbf{U} \in \mathbb{R}^{2N}$   
270 are employed for texturization. Despite possibility of arti-  
271 facts arising from applying naked-body skinning weights on  
272 displacement—as often evident under loose clothing akin to  
273 skirts—this low-cost, user-friendly representation well inte-  
274 grates into HOI-Blender, contrasting other potential choices  
275 and their attributes.

#### 276 4. HOI-Blender: A Unifying Visualization Tool

277 We introduce HOI-Blender (see Figure 2), an add-on which  
278 upon one-click (i) standardizes diverse HOI datasets via re-  
279 solving of incidental incompatibilities (file containers, co-  
280 ordinate conventions, rotation encodings, missing channels,  
281 scale ambiguities, and inconsistent object metadata); (ii) as-  
282 sembles attained parameters into Blender-canonical compo-  
283 nents (armatures, shape keys, color attributes), enabling in-  
284 teractive animation viewing, parallel to appearance compo-  
285 sitioning; (iii) provides render-ready setups with established  
286 camera placement, lighting arrangement and base material  
287 across datasets; and (iv) arranges automatic captioning on  
288 rendered frames, with outlook towards downstream applica-  
289 tions, extended via Blender-native options, such as auxil-  
290 iary render passes. The HOI-Blender User Interface is de-  
291 lined in Figure 3, with demonstrated workflow as part of  
292 our Supp. Video.

#### 293 4.1. Standardization of Heterogeneous Datasets

294 To reconcile inconsistencies, HOI-Blender performs auto-  
295 mated characteristic-based dataset identification and applies

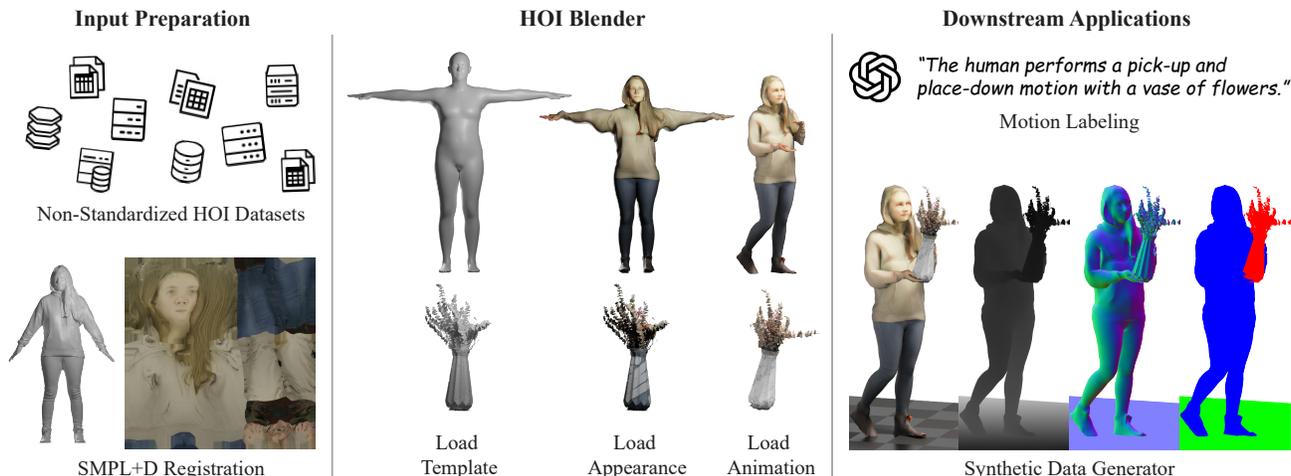


Figure 2. **HOI-Blender Overview.** Our add-on adapts heterogeneous HOI datasets into unified Blender-canonical components, facilitated over three key stages: (1) **Data Standardization:** ingests diverse HOI datasets and normalizes them into consistent representation; (2) **Scene Assembly:** maps SMPL-Family parameters to rigged human bodies, with decoupled motion and appearance, and assembles object meshes within a coherent 3D scene; (3) **Rendering & Applications:** enables scalable rendering with auxiliary render passes, such as depth, normals, and segmentation masks. The unified pipeline supports downstream applications encompassing synthetic supervision generation, VLM-based motion captioning, and cross-dataset evaluation.

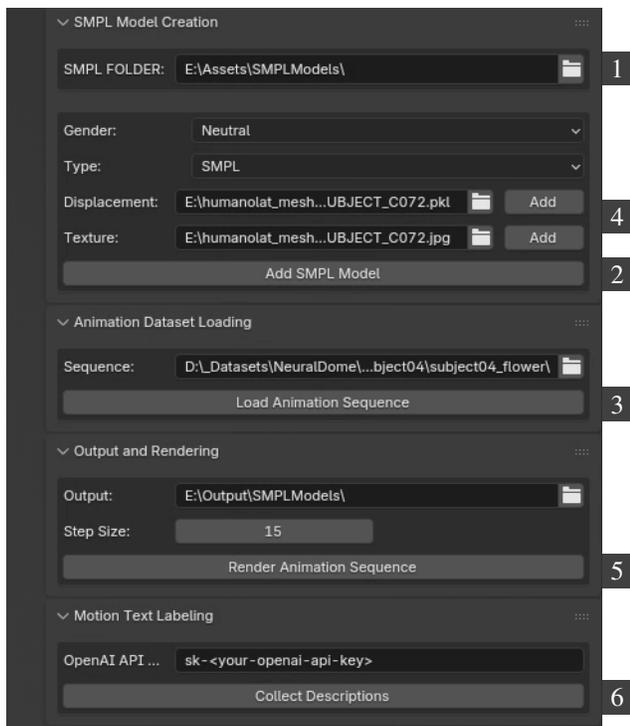


Figure 3. **HOI-Blender User Interface.** After imparting an asset path to SMPL-Family models (1), users can either import a specified model type (2) or load individual sequences directly (3). Both texture and displacement may be optionally provided (4). The select sequence can be directly rendered at specified step-size (5), with output frames electively passed for automatic captioning (6).

tailored loaders, such that a compact per-sequence dictionary of standardized parameters is acquired. Blender compatibility is attained via resolved sequence and asset structures, key mapping, parsed containers, canonicalized conventions, aligned coordinate frames, unified rotations, and accounting for dataset-specific idiosyncrasies and exceptions. Invisibility of dataset dissensions to all later processing modules via standardization is hence achieved.

#### 4.2. 4D HOI Representation

As generating individual meshes per-frame is highly expensive and inconsistent, we instead leverage the native rigging tools provided by our chosen framework for advantage.

Topology of canonical template mesh  $\bar{T}$  is explicitly represented in Blender, with blendshapes  $B_S$ ,  $B_P$  exposed as shape keys (additive per-vertex offsets), and joint transformations mapped to an armature (hierarchical skeleton) via per-joint vertex groups  $i$  representing per-vertex weights  $w_i$ , facilitating Blender’s native skinning for vertex position updates. Keyframes for global mesh location, shape key values and bone rotations are inserted respectively as read from standardized translation, shape and pose data—enabling efficient and consistent human animation playback.

Appearance is decoupled from motion, and may be realized in two manners (see Figure 4). Available texture assets are integrated into a principled material node network, with UV-mapping  $U$  carried over from the SMPL-Family artist-designed template. Provided per-vertex displacement  $B_C$  are exposed as shape keys, whose inherent additive nature allows arbitrary layered blending. This separation bears



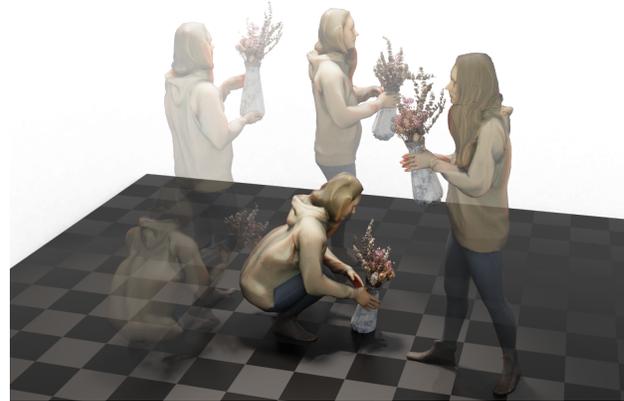
Figure 4. **Motion Appearance Decoupling.** Top row: same identity enacting various motion sequences. Bottom row: same motion enacted by various identities.

325 flexibility for dynamic asset swapping and enables diverse  
326 visual variations without altering the underlying animation.  
327 Objects are typically considered rigid per-sequence assets.  
328 After pivot configuration, keyframes for standardized  
329 translation, rotation, and scale are applied. Material assign-  
330 ment adheres to perdataset convention, supporting texture  
331 and vertexcolor basis.

332 Our standardization facilitates flexible import of HOI ani-  
333 mations across heterogeneous datasets, laying groundwork  
334 for scalable rendering, controlled experimentation, and di-  
335 verse downstream applications.

### 336 4.3. Automated Sequence Rendering

337 Rendering in HOI-Blender is taken care of by dataset-aware  
338 framing, automated instantiation of cameras and lighting,  
339 and preference configuration. Sequence mesh translation is  
340 aggregated over to compute a scene bounding box, guiding  
341 placement of tracked cameras and relative lighting. Frames  
342 are configurably rendered and optionally post-processed to  
343 achieve compact, centered crops. In addition, object-centric  
344 stills may also be produced. Collectively, our standardized  
345 cross-dataset basis enables consistent and customizable vi-  
346 sual output—potentially adorned by auxiliary render passes  
347 or further modifications—serving a multiplicity of founda-  
348 tions for cumulative downstream applications.



 "The human performs a pick-up and place-down motion with a vase of flowers: crouching to grab the vase, lifting and carrying it, then crouching again to set it down."

Figure 5. **VLM-Based Motion Captioning.** HOI-Blender renders image sequences of standardized HOI dataset animations and forwards them to VLMs for automated captioning.

### 4.4. Downstream Applications

**Motion Text Description.** HOI-Blender carries a captioning module, that collects HOI-centric crops and forwards to VLMs under a structured prompt (see Figure 5). Achieved are concise yet semantically rich motion-text pairs describing motion phases, contact patterns, and interaction semantics at scale. Our pipeline ensures standardized visual input with potential for systematic study of visual factors on generated text. Acquired captions not only facilitate rapid weak labeling and dataset curation, but they also serve as training data for downstream tasks such as motion captioning, interaction retrieval, and instruction-following applications.

**Human-Centric World Models** aim to learn how humans act and interact with objects to predict, simulate, and plan actions in realistic context settings [24]. Their training requires large, coherent HOI datasets with precise, temporally aligned camera parameters, interaction poses, and contact information. In practiced RGBD captures, signals are often degraded by sync errors, rollingshutter artifacts, calibration drift, incomplete object geometry, and missing contact labels, making frametoframe correspondence and contact dynamics hard to obtain.

HOI-Blender’s standardized render-ready format enables generating largescale synthetic HOI data with perfect correspondences, metric depth, calibrated cameras, and dense labels (see Figure 6). This approach mirrors a broader trend in worldmodeling, where synthetic data—offering dense, noise-free ground truth—boosts generalization and is thus predominantly trained upon [6, 10, 13, 32].

**3D & 4D HOI Reconstruction** provides the structured, temporally-coherent (geometry, pose, contact) supervision needed for humanoid robot training [1]. Accuracy is hinged for related methods at aforementioned issues, as practiced

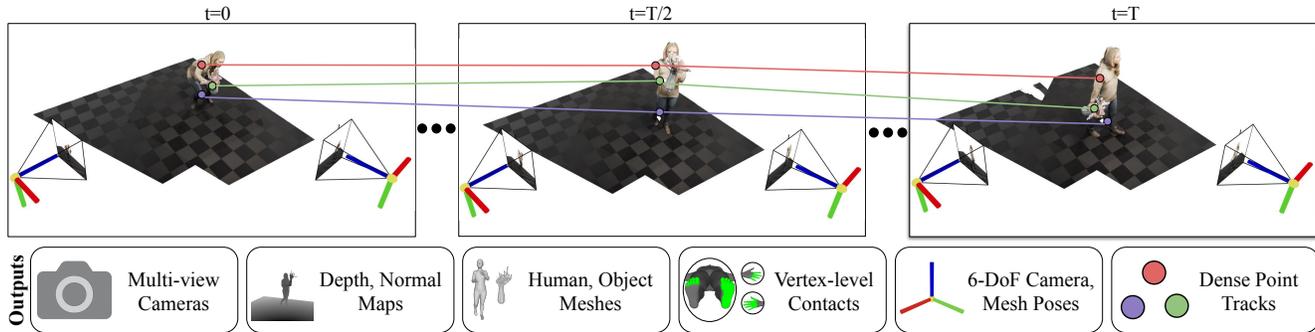


Figure 6. **Generated Data Modalities.** HOI-Blender may assist in generating multi-view synthetic HOI renders with auxiliary render passes and annotations—including depth, normal and segmentation maps, human and object mesh exports, as well as interaction contact points. Additional attainable annotations include 6-DoF camera poses, motion parameters, and temporally consistent dense point tracks; defining comprehensive resources for training models on HOI dynamics and temporal reasoning.

382 signals become rarely available at high fidelity [34, 35], and  
 383 datasets commonly lack reliable contact or object geometry  
 384 due to emphasis on humans over objects, thus hampering  
 385 interaction modeling [4, 26]. In contrast, HOI-Blender fa-  
 386 cilitates scalable synthetic supervision via multiview RGB,  
 387 ground-truth cameras, auxiliary render passes, and tempo-  
 388 rally consistent contact meshes. Decoupled motion and ap-  
 389 pearance enables controlled ablations and randomized ren-  
 390 ders, yielding dense, perfectly labeled data for reconstruc-  
 391 tion and policy learning—without fragile initialization or  
 392 heavy optimization.

393 **Dataset Evaluation** is accommodated by HOI-Blender  
 394 via direct, controlled comparison of HOI datasets facilitated  
 395 through standardized, scalable import and one-click render-  
 396 ing at shared camera views. This streamlines verification of  
 397 alignment, contact plausibility, geometry gaps, mesh pen-  
 398 etrations, and reprojection accuracy, as well as disentangling  
 399 of content from capture factors via decoupled motion from  
 400 appearance under specified lighting. Stratified sampling and  
 401 gallery renders help identify outliers, guiding dataset splits  
 402 and cross-dataset ablations.

403 We further demonstrate application of HOI-Blender for  
 404 dataset evaluation by applying its provided tools for quali-  
 405 tative analysis of our initial 15 HOI supported datasets.

## 406 5. Results: Dataset Evaluation

407 HOI-Blender’s standardized import and one-click rendering  
 408 enable controlled, equitable comparisons across heteroge-  
 409 neous HOI datasets. To demonstrate use, we render motion  
 410 clips across all supported datasets with automated sequence  
 411 import, camera framing, and lighting—then utilize acquired  
 412 renders as basis for dataset-level analysis.

413 We qualitatively evaluate 50 homogenized motion clips  
 414 per-dataset via following axes across a four-point scale (4 =  
 415 strong performance, 1 = weak performance, see Figure 7):

- 416 • **Animation Smoothness:** Motion continuity, void of tem-  
 417 poral jitter and abrupt frame-to-frame transitions;

- **Motion Realism:** Perceived naturalness of motion; 418
- **Mesh Quality:** Absence of deformation artifacts as well 419  
as topological failures during animation playback; 420
- **HOI Contact Quality:** Tenability of hand-object contact; 421
- **Hand Grip Quality:** Plausibility of grip-action and fin- 422  
ger placement during grasping; 423
- **Inter-Mesh Collision:** Lack of irregular interpenetration. 424

425 Both Human-only and SMPL-type datasets—or rather, 425  
those generally in lack of hand parameters—may not facili- 426  
tate impression for occasional evaluation criteria, hence are 427  
marked via ‘-’. 428

429 **Discussion.** Table 2 enables dataset comparison across 429  
qualitative metrics, and thus the observation of occasional 430  
trends. Isolated Human Motion datasets score high for kine- 431  
matic metrics, likely aided by a lack of dispersed focus be- 432  
tween human and object during capture; AMASS outscores 433  
AIST++ as dance may perform actions of considerable flex- 434  
ibility and complexity (AIST++: AS 3.52, MR 3.78, MQ 3.54; 435  
AMASS: AS 3.94, MR 3.98, MQ 3.96). Human-Human Inter- 436  
action is characterized by natural motion and clean geom- 437  
etry (Hi4D: MR 4.00, MQ 3.88, IMC 3.96; Duolando : MR 3.98, 438  
MQ 3.92, IMC 3.98). Full-Body Human-Object Interaction of- 439  
ten sees deformation fault at occlusions (Behave: MQ 3.10; 440  
NeuralDome : MQ 2.92; IMHD<sup>2</sup>: MQ 3.10) or animation jit- 441  
ter (InterCap: AS 2.00; COUCH: AS 2.86; IMHD<sup>2</sup>: AS 3.02), 442  
however, cases of particularly well captured interaction are 443  
also distinguishable (OMOMO: CQ 3.98; InterCap: GQ 3.68; 444  
NeuralDome: GQ 3.74; IMHD<sup>2</sup>: GQ 3.88). Hand-Object In- 445  
teraction datasets majorly present restricted interaction sce- 446  
narios, thus offer absence of deformation (GRAB: MQ 3.98; 447  
Arctic: MQ 4.00; HIMO: MQ 4.00) with high grip plausibility 448  
(GRAB: GQ 3.94; Arctic: GQ 3.90; HIMO: GQ 3.53). Multi- 449  
Human Object Interaction proffers contrasting perspectives 450  
on deformations and interpenetration (HOI-M3: MQ 4.00; 451  
CORE4D: MQ 3.24; vs. HOI-M3: IMC 2.96; CORE4D: IMC 452  
3.90), with CORE4D at higher overall scoring impression 453  
(CORE4D: AS 3.26; MR 3.41; GQ 3.69). 454

Table 2. Evaluation of Human-Centric Motion Datasets with Automatic Visualization via HOI-Blender.

Dataset	Animation Smoothness	Motion Realism	Mesh Quality	HOI Contact Quality	Hand Grip Quality	Inter-Mesh Collision
AIST++ [15]	3.52	3.78	3.54	-	-	-
AMASS [20]	3.94	3.98	3.96	-	-	-
Hi4D [37]	3.24	4.00	3.88	-	-	3.96
Duolando [30]	3.34	3.98	3.92	-	-	3.98
Behave [3]	2.62	3.61	3.10	2.96	-	2.92
InterCap [8]	2.00	2.56	3.94	2.90	3.68	3.38
COUCH [40]	2.86	2.96	3.67	2.73	-	3.25
OMOMO [14]	3.94	3.94	3.92	3.98	-	3.36
NeuralDome [38]	3.70	3.40	2.92	3.60	3.74	3.29
IMHD <sup>2</sup> [41]	3.02	3.12	3.10	3.31	3.88	3.44
GRAB [31]	3.91	3.52	3.98	3.43	3.94	3.69
Arctic [5]	3.98	3.52	4.00	3.34	3.90	3.46
HIMO[19]	3.82	2.84	4.00	2.90	3.53	3.47
HOI-M3 [39]	3.04	3.26	4.00	3.25	-	2.96
CORE4D[17]	3.26	3.41	3.24	3.71	3.69	3.90

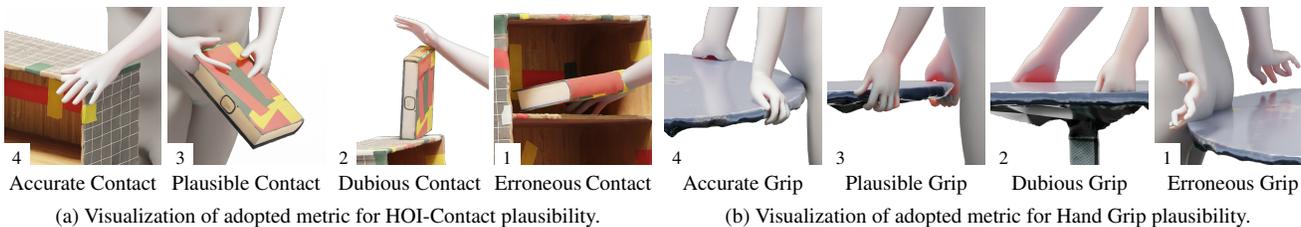


Figure 7. Demonstration of Plausibility Ordering for both Contact and Grip Quality.

455 HOI-Blender streamlined our dataset-level comparisons  
 456 by standardizing import and rendering for equitable contrast  
 457 across datasets. Consequently, we identified isolated  
 458 human motion excel in kinematic fidelity, multi-human pro-  
 459 ficiency in lifelike captures, and hand-object prevail at con-  
 460 strained, refined grasping. Human-object interaction is of-  
 461 ten impeded by occlusion-driven artifacts, however, carries  
 462 datasets of exceptional contact point performance as well.  
 463 Large-scale contrast between datasets can guide the direc-  
 464 tion of upcoming developmental advancement, downstream  
 465 application, and focused intervention, overall leading to en-  
 466 hanced performance across a multitude of HOI tasks.

## 467 6. Outlook and Conclusion

468 We introduced HOI-Blender, a unifying loader Blender add-  
 469 on that turns heterogeneous HOI datasets into a consistent,  
 470 render-ready format. It resolves conflicting conventions and  
 471 metadata, congregates coherent scenes with shared camera  
 472 rigs, lighting composition and render configuration, decou-  
 473 ples motion from appearance for flexible recombination,  
 474 and supports one-click import, animation preview, as well  
 475 as automated rendering. In practice, this makes visualiza-  
 476 tion, cross-dataset comparison, and dataset-driven supervi-  
 477 sion strikingly simpler and more reliable.

478 Our findings from Section 5 guide downstream data gen-  
 479 eration and amendment: alongside HOI-Blenders standard-

480 ized importing and coherent scenes, motion characteristics  
 481 can compose into high-quality sequences (AMASS motion,  
 482 OMOMO contacts, GRAB hands). Smooth fullbody kine-  
 483 matics can be combined with object placements and con-  
 484 tact cues, as well as retargeted hand poses to yield natural  
 485 motion, plausible contacts, and robust grasps. Such compo-  
 486 sitional strategies would enable exploration of large design  
 487 spaces and generation of targeted datasets for both training  
 488 and evaluation.

489 Looking ahead, HOI-Blender provides a practical founda-  
 490 tion for extensions which broaden its significance across  
 491 abundant synthesis, benchmarking, and embodied learning.  
 492 By leveraging standardized scene coherence alongside eq-  
 493 uitable renders, evaluations can grow beyond qualitative in-  
 494 spection into automated, quantitative diagnostics of mesh  
 495 contact, reprojection, and interpenetration. Through mix  
 496 and matching of sequences with differing characteristics,  
 497 those non-interactive can be enriched by interactive compo-  
 498 nents and low-fidelity grasps can be upgraded with higher-  
 499 quality counterparts. Lightweight post-processing—finger  
 500 retargeting, pose priors, and minor optimization passes—  
 501 can reduce hand artifacts and improve contact plausibility,  
 502 while native scripting support can streamline large-scale or  
 503 batched generation of synthetic corpora. Finally, stronger  
 504 physics and contact realism could enable consistent retar-  
 505 geting and contact correction, as well as provide synthetic,  
 506 force-aware supervision for robotics and control.

507

## References

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

- [1] Arthur Allshire, Hongsuk Choi, Junyi Zhang, David McAllister, Anthony Zhang, Chung Min Kim, Trevor Darrell, Pieter Abbeel, Jitendra Malik, and Angjoo Kanazawa. Visual imitation enables contextual humanoid control. *arXiv preprint arXiv:2505.03729*, 2025. 6
- [2] Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robotics and Autonomous Systems*, 57(5):469–483, 2009. 1
- [3] 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)*. IEEE, 2022. 2, 3, 8
- [4] Michael J. Black, Priyanka Patel, Joachim Tesch, and Jinlong Yang. BEDLAM: A Synthetic Dataset of Bodies Exhibiting Detailed Lifelike Animated Motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2, 7
- [5] Zicong Fan, Omid Taheri, Dimitrios Tzionas, Muhammed Kocabas, Manuel Kaufmann, Michael J. Black, and Otmar Hilliges. ARCTIC: A Dataset for Dexterous Bimanual Hand-Object Manipulation. In *Proceedings IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 3, 8
- [6] Adam W. Harley, Yang You, Xinglong Sun, Yang Zheng, Nikhil Raghuraman, Yunqi Gu, Sheldon Liang, Wen-Hsuan Chu, Achal Dave, Pavel Tokmakov, Suya You, Rares Ambrus, Katerina Fragkiadaki, and Leonidas J. Guibas. AllTracker: Efficient Dense Point Tracking at High Resolution, 2025. 6
- [7] Mohamed Hassan, Yunrong Guo, Tingwu Wang, Michael Black, Sanja Fidler, and Xue Bin Peng. Synthesizing Physical Character-Scene Interactions. New York, NY, USA, 2023. Association for Computing Machinery. 1
- [8] 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, 3, 8
- [9] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action Genome: Actions as Composition of Spatio-temporal Scene Graphs, 2019. 1
- [10] Zeren Jiang, Chuanxia Zheng, Iro Laina, Diane Larlus, and Andrea Vedaldi. Geo4D: Leveraging Video Generators for Geometric 4D Scene Reconstruction, 2025. 6
- [11] Manuel Kaufmann, Velko Vechev, and Dario Mylonopoulos. aitviewer, 2022. 3
- [12] Hema S. Koppula and Ashutosh Saxena. Anticipating Human Activities Using Object Affordances for Reactive Robotic Response. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(1):14–29, 2016. 1
- [13] Rosario Leonardi, Antonino Furnari, Francesco Ragusa, and Giovanni Maria Farinella. Are Synthetic Data Useful for Egocentric Hand-Object Interaction Detection? In *ECCV*, 2024. 6
- [14] Jiaman Li, Jiajun Wu, and C Karen Liu. Object Motion Guided Human Motion Synthesis. *ACM Trans. Graph.*, 42(6), 2023. 2, 3, 8
- [15] Ruilong Li, Shan Yang, David A. Ross, and Angjoo Kanazawa. Learn to Dance with AIST++: Music Conditioned 3D Dance Generation, 2021. 2, 3, 8
- [16] Tianye Li, Timo Bolkart, Michael J. Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4D scans. *ACM Trans. Graph.*, 36(6):194:1–194:17, 2017. 4
- [17] Yun Liu, Chengwen Zhang, Ruofan Xing, Bingda Tang, Bowen Yang, and Li Yi. CORE4D: A 4D Human-Object-Human Interaction Dataset for Collaborative Object REarrangement. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 1769–1782, 2025. 3, 8
- [18] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A Skinned Multi-Person Linear Model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, 2015. 4
- [19] Xintao Lv, Liang Xu, Yichao Yan, Xin Jin, Congsheng Xu, Shuwen Wu, Yifan Liu, Lincheng Li, Mengxiao Bi, Wenjun Zeng, and Xiaokang Yang. HIMO: A New Benchmark for Full-Body Human Interacting with Multiple Objects, 2024. 3, 8
- [20] Naureen Mahmood, Nima Ghorbani, Nikolaus F. Troje, Gerard Pons-Moll, and Michael J. Black. AMASS: Archive of Motion Capture as Surface Shapes. In *International Conference on Computer Vision*, pages 5442–5451, 2019. 2, 3, 8
- [21] Fahad Majeed, Maria Nazir, Kamilla Swart, Marco Agus, and Jens Schneider. Real-time analysis of soccer ballplayer interactions using graph convolutional networks for enhanced game insights. *Scientific Reports*, 15(1), 2025. 1
- [22] Madhur Mangalam, Sanjay Oruganti, Gavin Buckingham, and Christoph W. Borst. Enhancing hand-object interactions in virtual reality for precision manual tasks. *Virtual Reality*, 28(4), 2024. 1
- [23] Esteve Valls Mascaro, Daniel Sliwowski, and Dongheui Lee. HOI4ABOT: Human-Object Interaction Anticipation for Human Intention Reading Assistive roBOTS. In *7th Annual Conference on Robot Learning*, 2023. 1
- [24] Russell Mendonca, Shikhar Bahl, and Deepak Pathak. Structured World Models from Human Videos. *arXiv preprint arXiv:2308.10901*, 2023. 6
- [25] Meshcapade. SMPL\_blender\_addon [blender add-on]. [https://github.com/Meshcapade/SMPL\\_blender\\_addon](https://github.com/Meshcapade/SMPL_blender_addon), 2021. Accessed: August 22, 2025. 3
- [26] Priyanka Patel, Chun-Hao P. Huang, Joachim Tesch, David T. Hoffmann, Shashank Tripathi, and Michael J. Black. AGORA: Avatars in Geography Optimized for Regression Analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13468–13478, 2021. 2, 7
- [27] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive Body Capture: 3D Hands, Face, and Body From a Single Image. In *IEEE Conference* 564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620

621            *on Computer Vision and Pattern Recognition, CVPR 2019,*  
622            *Long Beach, CA, USA, June 16-20, 2019*, pages 10975–  
623            10985. Computer Vision Foundation / IEEE, 2019. 4

624 [28] Gerard Pons-Moll, Sergi Pujades, Sonny Hu, and Michael  
625            Black. ClothCap: Seamless 4D Clothing Capture and Retar-  
626            geting. *ACM Transactions on Graphics, (Proc. SIGGRAPH)*,  
627            36(4), 2017. Two first authors contributed equally. 4

628 [29] Javier Romero, Dimitrios Tzionas, and Michael J. Black.  
629            Embodied hands: modeling and capturing hands and bodies  
630            together. *ACM Trans. Graph.*, 36(6):245:1–245:17, 2017. 4

631 [30] Li Siyao, Tianpei Gu, Zhitao Yang, Zhengyu Lin, Ziwei Liu,  
632            Henghui Ding, Lei Yang, and Chen Change Loy. Duolando:  
633            Follower GPT with Off-Policy Reinforcement Learning for  
634            Dance Accompaniment. In *ICLR*, 2024. 2, 3, 8

635 [31] Omid Taheri, Nima Ghorbani, Michael J. Black, and Dim-  
636            itrios Tzionas. GRAB: A Dataset of Whole-Body Human  
637            Grasping of Objects. In *European Conference on Computer*  
638            *Vision (ECCV)*, 2020. 3, 8

639 [32] Zachary Teed and Jia Deng. RAFT: Recurrent All-Pairs Field  
640            Transforms for Optical Flow, 2020. 6

641 [33] Boran Wen, Dingbang Huang, Zichen Zhang, Jiahong  
642            Zhou, Jianbin Deng, Jingyu Gong, Yulong Chen, Lizhuang  
643            Ma, and Yong-Lu Li. Reconstructing In-the-Wild Open-  
644            Vocabulary Human-Object Interactions. In *Proceedings of*  
645            *the IEEE/CVF Conference on Computer Vision and Pattern*  
646            *Recognition (CVPR)*, 2025. 2

647 [34] Xianghui Xie, Bharat Lal Bhatnagar, and Gerard Pons-Moll.  
648            CHORE: Contact, Human and Object REconstruction from  
649            a single RGB image. In *European Conference on Computer*  
650            *Vision (ECCV)*. Springer, 2022. 7

651 [35] Xianghui Xie, Bharat Lal Bhatnagar, Jan Eric Lenssen, and  
652            Gerard Pons-Moll. Template Free Reconstruction of Human-  
653            object Interaction with Procedural Interaction Generation. In  
654            *IEEE Conference on Computer Vision and Pattern Recogni-*  
655            *tion (CVPR)*, 2024. 1, 7

656 [36] Brent Yi, Chung Min Kim, Justin Kerr, Gina Wu, Rebecca  
657            Feng, Anthony Zhang, Jonas Kulhanek, Hongsuk Choi, Yi  
658            Ma, Matthew Tancik, and Angjoo Kanazawa. Viser: Imper-  
659            ative, Web-based 3D Visualization in Python, 2025. 3

660 [37] Yifei Yin, Chen Guo, Manuel Kaufmann, Juan Zarate, Jie  
661            Song, and Otmar Hilliges. Hi4D: 4D Instance Segmentation  
662            of Close Human Interaction. In *Computer Vision and Pattern*  
663            *Recognition (CVPR)*, 2023. 2, 3, 8

664 [38] Juze Zhang, Haimin Luo, Hongdi Yang, Xinru Xu, Qianyang  
665            Wu, Ye Shi, Jingyi Yu, Lan Xu, and Jingya Wang. Neural-  
666            Dome: A Neural Modeling Pipeline on Multi-View Human-  
667            Object Interactions. In *CVPR*, 2023. 2, 3, 8

668 [39] Juze Zhang, Jingyan Zhang, Zining Song, Zhanhe Shi,  
669            Chengfeng Zhao, Ye Shi, Jingyi Yu, Lan Xu, and Jingya  
670            Wang. HOI-M3: Capture Multiple Humans and Objects  
671            Interaction within Contextual Environment. *arXiv preprint*  
672            *arXiv:2404.00299*, 2024. 3, 8

673 [40] Xiaohan Zhang, Bharat Lal Bhatnagar, Sebastian Starke,  
674            Vladimir Guzov, and Gerard Pons-Moll. COUCH: Towards  
675            Controllable Human-Chair Interactions. 2022. 2, 3, 8

676 [41] Chengfeng Zhao, Juze Zhang, Jiashen Du, Ziwei Shan,  
677            Junye Wang, Jingyi Yu, Jingya Wang, and Lan Xu. I’m hoi:  
678            Inertia-aware monocular capture of 3d human-object interac-  
679            tions. In *Proceedings of the IEEE/CVF Conference on Com-*  
680            *puter Vision and Pattern Recognition*, pages 729–741, 2024.  
681            2, 3, 8