

InterCap: Joint Markerless 3D Tracking of Humans and Objects in Interaction

Yinghao Huang¹, Omid Taheri¹, Michael J. Black¹, and Dimitrios Tzionas²

¹ Max Planck Institute for Intelligent Systems, Tübingen, Germany

² University of Amsterdam, Amsterdam, The Netherlands

{yhuang2, otaheri, black}@tue.mpg.de, d.tzionas@uva.nl

Abstract. Humans constantly interact with daily objects to accomplish tasks. To understand such interactions, computers need to reconstruct these from cameras observing whole-body interaction with scenes. This is challenging due to occlusion between the body and objects, motion blur, depth/scale ambiguities, and the low image resolution of hands and graspable object parts. To make the problem tractable, the community focuses either on interacting hands, ignoring the body, or on interacting bodies, ignoring hands. The GRAB dataset addresses dexterous whole-body interaction but uses marker-based MoCap and lacks images, while BEHAVE captures video of body-object interaction but lacks hand detail. We address the limitations of prior work with InterCap, a novel method that reconstructs interacting whole-bodies and objects from multi-view RGB-D data, using the parametric whole-body model SMPL-X and known object meshes. To tackle the above challenges, InterCap uses two key observations: (i) Contact between the hand and object can be used to improve the pose estimation of both. (ii) Azure Kinect sensors allow us to set up a simple multi-view RGB-D capture system that minimizes the effect of occlusion while providing reasonable inter-camera synchronization. With this method we capture the InterCap dataset, which contains 10 subjects (5 males and 5 females) interacting with 10 objects of various sizes and affordances, including contact with the hands or feet. In total, InterCap has 223 RGB-D videos, resulting in 67,357 multi-view frames, each containing 6 RGB-D images. Our method provides pseudo ground-truth body meshes and objects for each video frame. Our InterCap method and dataset fill an important gap in the literature and support many research directions. Our data and code are available for research purposes at <https://intercap.is.tue.mpg.de>.

1 Introduction

A long-standing goal of Computer Vision is to understand human actions from videos. Given a video people effortlessly figure out what objects exist in it, the spatial layout of objects, and the pose of humans. Moreover, they deeply understand the depicted action. What is the subject doing? Why are they doing this? What is their goal? How do they achieve this? To empower computers with the ability to infer such abstract concepts from pixels, we need to capture rich datasets and to devise appropriate algorithms.

Since humans live in a 3D world, their physical actions involve interacting with objects. Think of how many times per day one goes to the kitchen, grabs a cup of

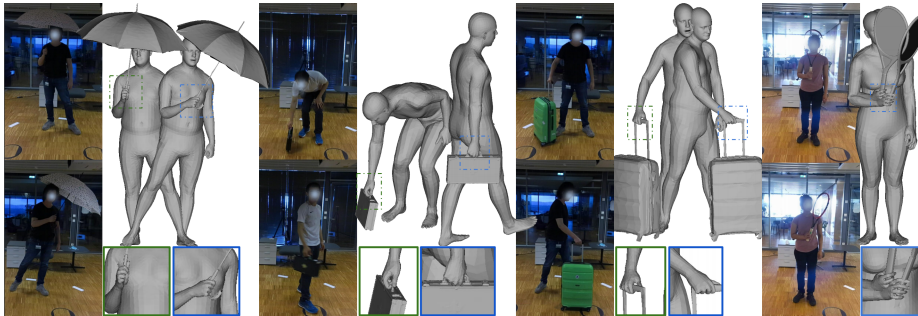


Fig. 1. Humans interact with objects to accomplish tasks. To understand such interactions we need the tools to reconstruct them from whole-body videos in 4D, i.e., as 3D meshes in motion. Existing methods struggle, due to the strong occlusions, motion blur, and low-resolution of hands and object structures in such videos. Moreover, they mostly focus on the main body, ignoring the hands and objects. We develop InterCap, a novel method that reconstructs plausible interacting whole-body and object meshes from multi-view RGB-D videos, using contact constraints to account for strong ambiguities. With this we capture the rich InterCap dataset of 223 RGB-D videos (67,357 multi-view frames, with 6 Azure Kinects) containing 10 subjects (5 fe-/males) interacting with 10 objects of various sizes and affordances; note the hand-object grasps.

water, and drinks from it. This involves contacting the floor with the feet, contacting the cup with the hand, moving the hand and cup together while maintaining contact, and drinking while the mouth contacts the cup. Thus, to understand human actions, it is necessary to reason in 3D about humans and objects *jointly*.

There is significant prior work on estimating 3D humans without taking into account objects [4] and estimating 3D objects without taking into account humans [68]. There is even recent work on inserting bodies into 3D scenes such that their interactions appear realistic [16, 30, 65]. But there is little work on estimating 3D humans interacting with scenes and moving objects, in which the human-scene/object contact is explicitly modeled and exploited. To study this problem, we need a dataset of videos with rich human-object interactions and reliable 3D ground truth.

PROX [15] takes a step in this direction by estimating the 3D body in a known 3D scene. The scene mesh provides information that helps resolve pose ambiguities commonly encountered when a single camera is used. However, PROX involves only coarse interactions of bodies, static scenes with no moving objects, and no dexterous fingers. The recent BEHAVE dataset [3] uses multi-view RGB-D data to capture humans interacting with objects but does not include detailed hand pose or fine hand-object contact. Finally, the GRAB dataset [54] captures the kind of detailed hand-object and whole-body-object interaction that we seek but is captured using marker-based MoCap and, hence, lacks images paired with the ground-truth scene.

We argue that what is needed is a new dataset of RGB videos containing natural human-object interaction in which the whole body is tracked reliably, the hand pose is captured, objects are also tracked, and the hand-object contact is realistic; see Fig. 1. This is challenging and requires technical innovation to create. To that end, we design a system that uses multiple RGB-D sensors that are spatially calibrated and temporally

synchronized. To this data we fit the SMPL-X body model, which has articulated hands, by extending the PROX [15] method to use multi-view data and grasping hand-pose priors. We also track the 3D objects with which the person interacts. The objects used in this work are representative of items one finds in daily life. We obtain accurate 3D models for each object with a handheld Artec scanner. Altogether we collect 223 sequences (67,357 multi-view frames), with 10 subjects interacting with 10 objects.

The problem, however, is that separately estimating the body and objects is not sufficient to ensure accurate 3D body-object contact. Consequently, a key innovation of this work is to estimate these *jointly*, while exploiting information about *contact*. Objects do not move independently, so, when they move, it means the body is in contact. We define likely contact regions on objects and on the body. Then, given frames with known likely contacts, we enforce contact between the body and the object when estimating the body and object poses. The resulting method produces natural body poses, hand poses, and object poses. Uniquely, it provides detailed pseudo ground-truth contact information between the whole body and objects in RGB video.

In summary, our major contributions are as follows: (1) We develop a novel Motion Capture method utilizing multiple RGB-D cameras. It is relatively lightweight and flexible, yet accurate enough, thus suitable for data capture of daily scenarios. (2) We extend previous work on fitting SMPL-X to images to fit it to multi-view RGB-D data while taking into account body-object contact. (3) We capture a novel dataset that contains whole-body human motions and interaction with objects, as well as multi-view RGB-D imagery. Our data and code are available at <https://intercap.is.tue.mpg.de>.

2 Related Work

There is a large literature on estimating 3D human pose and shape from images or videos [4, 7, 25, 29, 37, 41, 44, 57]. Here we focus on the work most closely related to ours, particularly as it concerns, or enables, capturing human-object interaction.

MoCap from Multi-view Videos and IMUs. Markerless MoCap from multi-view videos [8, 22, 31] is widely studied and commercial solutions exist (e.g., Theia Markerless). Compared with traditional marker-based MoCap, markerless offers advantages of convenience, applicability in outdoor environments, non-intrusiveness, and greater flexibility. However, traditional MoCap methods, both marker-based and markerless, focus on extracting a 3D skeleton. This is useful for biomechanics but our goal is to reason about body-scene contact. To enable that, we need to capture the body surface.

Various 3D human representations have been proposed, with recent work focused on learning a parametric mesh-based model of body shape from large-scale collections of 3D scans [2, 33, 42–44, 50, 59]. Here we use the SMPL-X model [44] because it contains fully articulated hands, which are critical for reasoning about object manipulation. The body parameters are often estimated by fitting the 3D generative model to various 2D cues like landmarks detected by Convolutional Neural Networks [6, 39, 58] or silhouettes [1, 47, 60]. Though effective, these monocular video-based methods suffer from depth ambiguity and occlusions. To address this issue, researchers have proposed to combine IMUs with videos to obtain better and more robust results [36, 45].

Many methods estimate 3D bodies from multi-view images but focus on skeletons and not 3D bodies [9, 10, 19, 24, 46, 55, 66]. Recent work addresses 3D body shape estimation from multiple views [11, 22, 67]. Most related to our work are two recent datasets. The RICH dataset [21], fits SMPL-X bodies to multi-view RGB videos taken both indoors and outdoors. The method uses a detailed 3D scan of the scene and models the contact between the body and the world. RICH does not include any object motion; the scenes are completely rigid. In contrast, BEHAVE [3] contains SMPL bodies interacting with 3D objects that move. We go beyond that work, however, to integrate novel contact constraints and to capture hand pose, which is critical for human-object interaction. Additionally, BEHAVE focuses on large objects like boxes and chairs, whereas we have a wider range of object sizes, including smaller objects like cups.

Human-Object Interaction. There has been a lot of work on modeling or analyzing human-object interactions [3, 13, 14, 18, 26, 40, 48, 56, 61]. A detailed discussion is out of the scope of this work. Here, we focus on modeling and analyzing human-object interaction in 3D space. Most existing work, however, only focuses on estimating hand pose [14, 17, 18, 49], ignoring the strong relationship between body motion, hand motion, and object motion. Recent work considers whole-body motion. For example, the GRAB [54] dataset provides detailed object motion and whole-body motion in a parametric body format (SMPL-X). Unfortunately, it is based on MoCap and does not include video. Here our focus is on tracking the whole-body motion, object motion, and the detailed hand-object contact to provide ground-truth 3D information in RGB video.

Joint Modeling of Humans and Scenes. There is some prior work addressing human-object contact in both static images and video. For example, PHOSA estimates a 3D body and a 3D object with plausible interaction from a single RGB image [63]. Our focus here, however, is on dynamic scenes. Motivated by the observation that natural human motions always happen inside 3D scenes, researchers have proposed to model human motion jointly with the surrounding environment [5, 15, 51, 62]. In PROX [15] the contact between humans and scenes is explicitly used to resolve ambiguities in pose estimation. The approach avoids bodies interpenetrating scenes while encouraging contact between the scene and nearby body parts. Prior work also tries to infer the most plausible position and pose of humans given the 3D scene [16, 30, 65]. Most recently, MOVER [62] estimates the 3D scene and the 3D human directly from a static monocular video in which a person interacts with the scene. While the 3D scene is ambiguous and the human motion is ambiguous, by exploiting contact, the method resolves many ambiguities, improving the estimates of both the scene and the person. Unfortunately, this assumes a static scene and does not model hand-object manipulation.

Datasets. Traditionally, MoCap is performed using marker-based systems inside lab environments. To capture object interaction and contact, one approach uses MoSh [32] to fit a SMPL or SMPL-X body to the markers [35]. An advanced version of this is used for GRAB [54]. Such approaches lack synchronized RGB video. The HumanEva [52] and Human3.6M [23] datasets combine multi-camera RGB video capture with synchronized ground-truth 3D skeletons from marker-based MoCap. These datasets lack ground-truth 3D body meshes, are captured in a lab setting, and do not contain human-object manipulation. 3DPW [36] is the first in-the-wild dataset that jointly features natural human appearance in video and accurate 3D pose. This dataset does not

track objects or label human-object interaction. PiGraphs [51] and PROX [15] provide both 3D scenes and human motions but are relatively inaccurate, relying on a single RGB-D camera. This makes these datasets ill-suited as evaluation benchmarks. The recent RICH dataset [21] addresses many of these issues with indoor and outdoor scenes, accurate multi-view capture of SMPL-X, 3D scene scans, and human-scene contact. It is not appropriate for our task, however, as it does not include object manipulation. An alternative approach is the one of GTA-IM [5] and SAIL-VOS [20], which generate human-scene interaction data using either 3D graphics or 2D videos. They feature high-accuracy ground truth but lack visual realism. In summary, we believe that a 3D human-object interaction dataset needs to have accurate hand poses to be useful, since hands are how people most often interact with objects. We compare our InterCap dataset with other ones in Tab. 1.

Name	# of Seq.	Natural Appear.	Moving Objects	Accurate Motion	With Image	Artic. Hands
HumanEva [52]	56	✓	✗	✓	✓	✗
Human3.6M [23]	165	✓	✗	✓	✓	✗
AMASS [35]	11265	✓	✗	✓	✗	✗
GRAB [54]	1334	✓	✓	✓	✗	✓
3DPW [36]	60	✓	✗	✓	✓	✗
GTA-IM [5]	119	✗	✗	✓	✓	✗
SAIL-VOS [20]	201	✗	✗	✗	✗	✗
PiGraphs [51]	63	✓	✗	✓	✓	✗
PROX [15]	20	✓	✗	✗	✓	✗
RICH [21]	142	✓	✗	✓	✓	✗
BEHAVE [3]	321	✓	✓	✓	✓	✗
InterCap (ours)	223	✓	✓	✓	✓	✓

Table 1. Dataset statistics. Comparison of our InterCap dataset to existing datasets.

3 InterCap Method

Our core goal is to accurately estimate the human and object motion throughout a video. Our markerless motion capture method is built on top of the PROX-D method of Hassan et al. [15]. To improve the body tracking accuracy we extend this method to use multiple RGB-D cameras; here we use the latest Azure Kinect cameras. The motivation is that multiple cameras observing the body from different angles give more information about the human and object motion. Moreover, commodity RGB-D cameras are much more flexible to deploy out of controlled lab scenarios than more specialized devices.

The key technical challenge lies in accurately estimating the 3D pose and translation of the objects while a person interacts with them. In this work we focus on 10 variously sized rigid objects common in daily life, such as cups and chairs. Being rigid does not make the tracking of the objects trivial because of the occlusion by the body and hands. While there is a rich literature on 6 DoF object pose estimation, much of it ignores hand-object interaction. Recent work in this direction is promising but still focuses on scenarios that are significantly simpler than ours, cf. [53].

Similar to previous work on hand and object pose estimation [14] from RGB-D videos, in this work we assume that the 3D meshes of the objects are known in advance. To this end, we first gather the 3D models of these objects from the Internet whenever possible and scan the remaining objects ourselves. To fit the known object models to image data, we first perform semantic segmentation, find the corresponding object regions in all camera views, and fit the 3D mesh to the segmented object contours via differentiable rendering. Since heavy occlusion between humans and objects in some views may make the segmentation results unreliable, aggregating segmentation from all views boosts the object tracking performance.

In the steps above, both the subject and object are treated separately and processing is per frame, with no temporal smoothness or contact constraint applied. This produces jittery motions and heavy penetration between objects and the body. Making matters worse, our human pose estimation exploits OpenPose for 2D keypoint detection, which struggles when the object occludes the body or the hands interact with it. To mitigate this issue and still get reasonable body, hand and object pose in these challenging cases, we manually annotate the frames where the body or the hand is in contact with the object, as well as the body, hand and object vertices that are most likely to be in contact. This manual annotation can be tedious; automatic detection of contact is an open problem and is left for future work. We then explicitly encourage the labeled body and hand vertices to be in contact with the labeled object vertices. We find that this straightforward idea works well in practice. More details are described in the following.

3.1 Multi-Kinect Setup

We use 6 Azure Kinects to track the human and object together, deployed in a “ring” layout in an office; see Sup. Mat. Multiple RGB-D cameras provide a good balance between body tracking accuracy and applicability to real scenarios, compared with costly professional MoCap systems like Vicon, or cheap and convenient but not-so-accurate monocular RGB cameras. Moreover, this approach does not require applying any markers, making the images natural. Intrinsic camera parameters are provided by the manufacturer. Extrinsic camera parameters are obtained via camera calibration with Azure Kinect’s API [38]. However, these can be a bit noisy, as non-neighbouring cameras in a sparse “ring” layout don’t observe the calibration board well at the same time. Thus, we manually refine in MeshLab the extrinsics by comparing the point clouds for neighbouring cameras for several iterations. The hardware synchronization of Azure Kinects is empirically reasonable. Given the calibration information, we choose one camera’s 3D coordinate frame as the global frame and transform the point clouds from the other frames into the global frame, which is where we fit the SMPL-X and object models.

3.2 Sequential Object-Only Tracking

Object Segmentation. To track an object during interaction, we need reliable visual cues about it to compare with the 3D object model. To this end, we perform semantic segmentation by applying PointRend [28] to the whole image. We then extract the object instances that correspond to the categories of our objects; for examples see Sup. Mat. We assume that the subject interacts with a single object. Note that, in contrast to previous approaches where the objects occupy a large portion of the image [14, 15, 40, 56], in our case the entire body is visible, thus, the object takes up a small part of the image and is often occluded by the body and hands; our setting is much more challenging. We observe that PointRend works reasonably well for large objects like chairs, even with heavy occlusion between the object and the human, while for small objects, like a bottle or a cup, it struggles significantly due to occlusion.

In extreme cases, it is possible for the object to not be detected in most of the views. But even when the segmentation is good, the class label for the objects may be wrong. To resolve this, we take two steps: (1) For every frame, we detect all possible object

segmentation candidates and their labels. This step takes place offline and only once. (2) During the object tracking phase, for each view, we compare the rendering of the tracked object from the i^{th} frame with all the detected segmentation candidates for the $(i + 1)^{\text{th}}$ frame, and preserve only the candidate with the largest overlap ratio. This render-compare-and-preserve operation takes place iteratively during tracking.

Object Tracking. Given object masks via semantic segmentation over the whole sequence, we track the object by fitting its model to observations via differentiable rendering [27, 34]. This is similar to past work for hand-object tracking [14]. We assume that the object is rigid and its mesh is given. The configuration of the rigid object in the t^{th} frame is specified via a 6D rotation and translation vector ξ . For initialization, we manually obtain the configuration of the object for the first frame by matching the object mesh into the measured point clouds. Let R_S and R_D be functions that render a synthetic mask and depth image for the tracked 3D object mesh, M . Let also $S = \{S_\nu\}$ be the ‘‘observed’’ object masks and $D = \{D_\nu\}$ be corresponding depth values for the current frame, where ν is the camera view. Then, we minimize:

$$E_O(\xi; S, D) = \sum_{\text{view } \nu} \lambda_{\text{segm}} \|(R_S(\xi, M, \nu) - S_\nu) * S_\nu\|_F^2 + \lambda_{\text{depth}} \|(R_D(\xi, M, \nu) - D_\nu) * S_\nu\|_F^2, \quad (1)$$

where the two terms compute how well the rendered object mask and depth image match the detected mask and observed depth; the $*$ is an element-wise multiplication, and $\|\cdot\|_F$ the Frobenius norm; λ_{segm} and λ_{depth} are steering weights set empirically. For simplicity, we assume that transformations from the master to other camera frames are encoded in the rendering functions R_S, R_D ; we do not denote these explicitly here.

3.3 Sequential Human-Only Tracking

We estimate body shape and pose over the whole sequence from multi-view RGB-D videos in a frame-wise manner. This is similar in spirit with the PROX-D method [15], but, in our case, there is no 3D scene constraint and multiple cameras are used. The human pose and shape are optimized independently in each frame. We use the SMPL-X [44] model to represent the 3D human body. SMPL-X is a function that returns a water-tight mesh given parameters for shape, β , pose, θ , facial expression, ψ , and translation, γ . We follow the common practice of using a 10-dimensional space for shape, β , and a 32-dimensional latent space in VPoser [44] to represent body pose, θ .

We minimize the loss defined below. For each frame we essentially extend the major loss terms used in PROX [15] to multiple views:

$$E_B(\beta, \theta, \psi, \gamma; K, J_{\text{est}}) = E_J + \lambda_D E_D + \lambda_{\theta_b} E_{\theta_b} + \lambda_{\theta_h} E_{\theta_h} + \lambda_{\theta_f} E_{\theta_f} + \lambda_\alpha E_\alpha + \lambda_\beta E_\beta + \lambda_\mathcal{E} E_\mathcal{E} + \lambda_\mathcal{P} E_\mathcal{P}, \quad (2)$$

where $E_\beta, E_{\theta_b}, E_{\theta_h}, E_{\theta_f}, E_\mathcal{E}$ are prior loss terms for body shape, body pose, hand pose, facial pose and expressions. Also, E_α is a prior for extreme elbow and knee bending. For detailed definitions of these terms see [15]. E_J is a 2D keypoint re-projection loss:

$$E_J(\beta, \theta, \gamma; K, J_{\text{est}}) = \sum_{\text{view } \nu} \sum_{\text{joint } i} k_i^\nu w_i^\nu \rho_J(\Pi_K^\nu(R_{\theta_\gamma}(J(\beta)_i)) - J_{\text{est},i}^\nu), \quad (3)$$



Fig. 2. The objects of our InterCap dataset. **Left:** Color photos. **Right:** Annotations for object areas that are likely to come in contact during interaction, shown in red.

where $\theta = \{\theta_b, \theta_h, \theta_f\}$, ν and i iterate through views and joints, k_i^ν and w_i^ν are the per-joint weight and detection confidence, ρ_J is a robust Geman-McClure error function [12], Π_K^ν is the projection function with K camera parameters, $R_{\theta\gamma}(J(\beta)_i)$ are the posed 3D joints of SMPL-X, and $J_{est,i}^\nu$ the detected 2D joints. The term E_D is:

$$E_D(\beta, \theta, \gamma; K) = \sum_{view \nu} \sum_{p \in P^\nu} \min_{v \in V_b^\nu} \|v - p\|, \quad (4)$$

where P^ν is Azure Kinect’s segmented point cloud for the ν^{th} view, and V_b^ν are SMPL-X vertices that are visible in this view. This term measures how far the estimated body mesh is from the combined point clouds, so that we minimize this discrepancy. Note that, unlike PROX, we have multiple point clouds from all views, i.e., our E_D is a multi-view extension of PROX’s [15] loss. For each view we dynamically compute the visible body vertices, and “compare” them against the segmented point cloud for that view.

Finally, the term $E_{\mathcal{P}}$ penalizes self-interpenetration of the SMPL-X body mesh; see PROX [15] for a more detailed and formal definition of this:

$$E_{\mathcal{P}}(\theta, \beta, \gamma) = E_{\mathcal{P}_{self}}(\theta, \beta). \quad (5)$$

3.4 Joint Human-Object Tracking Over All Frames

We treat the result of the above optimization as initialization for refinement via *joint* optimization of the body and the object *over all frames*, subject to *contact* constraints.

For this we fix the body shape parameters, β , as the mean body shape computed over all frames from the first stage, as done in [22]. Then, we jointly optimize the object pose and translation, ξ , body pose, θ , and body translation, γ , over all frames. We add a temporal smoothness loss to reduce jitter for both the human and the object. We also penalize the body-object interpenetration, as done in PROX [15]. A key difference is that in PROX the scene is static, while here the object is free to move.

To enforce contact, we annotate the body areas that are most likely to be in contact with the objects and, for each object, we label vertices most likely to be contacted. These annotations are shown in Fig. 3 and Fig. 2-right, respectively, in red. We also annotate frame sub-sequences where the body is in contact with objects, and enforce contact

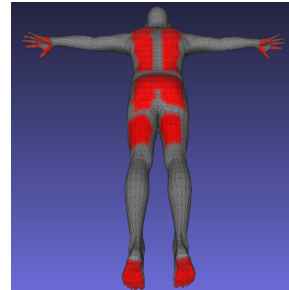


Fig. 3. Annotation of likely body contact areas (red color).

between them explicitly to get reasonable tracking even when there is heavy interaction and occlusion between hands and objects. Such interactions prove to be challenging for state-of-the-art 2D joint detectors, e.g., OpenPose, especially for hands.

Formally, we perform global optimization over all T frames, and minimize a loss, E , that is composed of an object fitting loss, E_O , a body fitting loss, E_B , a motion smoothness prior [64] loss, E_S , and a loss penalizing object acceleration, E_A . We also use a ground support loss, E_G , that encourages the human and the object to be above the ground plane, i.e., to not penetrate it. Last, we use a body-object contact loss, E_C , that attaches the body to the object for frames with contact. The loss E is defined as:

$$\begin{aligned}
E = & \frac{1}{T} \sum_{\text{frame } t} \left[E_O(\Xi_t; \mathcal{S}_t, \mathcal{D}_t) + E_B(\beta^*, \Theta_t, \Psi_t, \Gamma_t; \mathcal{J}_{est}) \right] + \\
& \frac{1}{T} \sum_{\text{frame } t} \left[E_{\mathcal{P}}(\Theta_t, \beta^*, \Gamma_t) + E_C(\beta^*, \Theta_t, \Psi_t, \Gamma_t, \Xi_t, M) \right] + \\
& \frac{\lambda_G}{T} \sum_{\text{frame } t} \left[E_G(\beta^*, \Theta_t, \Psi_t, \Gamma_t) + E_{G'}(\Xi_t, M) \right] + \\
& \frac{\lambda_Q}{T} \sum_{\text{frame } t} \left[Q_t * E_C(\beta^*, \Theta_t, \Psi_t, M', \Xi_t) \right] + \\
& \lambda_S E_S(\Theta, \Psi, \Gamma, A; \beta^*, T) + \\
& \lambda_A E_A(\Xi, T, M),
\end{aligned} \tag{6}$$

where E_O comes from Eq. 1 and E_B from Eq. 2, and both go through all views ν , while $E_{\mathcal{P}}$ comes from Eq. 5. For all frames $t = \{1, \dots, T\}$ of a sequence, $\Theta = \{\theta_t\}$, $\Psi = \{\psi_t\}$, $\Gamma = \{\gamma_t\}$, are the body poses, facial expressions and translations, $\Xi = \{\xi_t\}$ is the object rotations and translations, $\mathcal{S} = \{S_t\}$ and $\mathcal{D} = \{D_t\}$ are masks and depth patches, and $\mathcal{J}_{est} = \{J_{est,t}\}$ are detected 2D keypoints. M is the object mesh, and β^* the mean body shape. E_C encourages body-object contact for frames in contact, which are indicated by the manually annotated binary vectors $Q = \{Q_t\}$, $t = \{1, \dots, T\}$; Q_t is set to 1 if in the t^{th} frame any body part (e.g., hand, foot, thighs) is in contact with the object, and set to 0 otherwise. The motion smoothness loss E_S penalizes abrupt position changes for body vertices, and the vertex acceleration loss E_A encourages smooth object trajectories. We estimate the ground plane surface by fitting a plane to chosen floor points in the observed point clouds. The terms E_G and $E_{G'}$ measure whether the body and object penetrate the ground, respectively. For more details on the above loss terms, please see Sup. Mat. Finally, the parameters λ_G , λ_Q , λ_S , and λ_A are steering weights that are set empirically.

4 InterCap Dataset

We use the proposed InterCap algorithm (Sec. 3) to capture the InterCap dataset, which uniquely features whole-body interactions with objects in multi-view RGB-D videos.

Data-capture Protocol. We use 10 everyday objects, shown in Fig. 2-left, that vary in size and “afford” different interactions with the body, hands or feet; we focus mainly on hand-object interactions. We recruit 10 subjects (5 males and 5 females) that are



Fig. 4. Samples from our InterCap dataset, drawn from four sequences with different subjects and objects. The estimated 3D object and SMPL-X human meshes have plausible contacts that agree with the input images. Best viewed zoomed in.

between 25-40 years old. The subjects are recorded while interacting with 7 or more objects, according to their time availability. Subjects are instructed to interact with objects as naturally as possible. However, they are asked to avoid very fast interactions that cause severe motion blur (Azure Kinect supports only up to 30 FPS), or misalignment between the RGB and depth images for each Kinect (due to technicalities of RGB-D sensors). We capture up to 3 sequences per object depending on object shape and functionality, and by picking an interaction intent from the list below, as in GRAB [54]:

- **“Pass”**: The subject passes the object on to another imaginary person standing on their left/right side; a graspable area needs to be free for the other person to grasp.
- **“Check”**: The subject inspects visually the object from several viewpoints by first picking it up and then manipulating it with their hands to see several sides of it.
- **“Use”**: The subject uses the object in a natural way that “agrees” with the object’s affordances and functionality for everyday tasks.

We also capture each subject performing a freestyle interaction of their choice. All subjects gave informed written consent to publicly share their data for research.

4D Reconstruction. Our InterCap method (Sec. 3) takes as input multi-view RGB-D videos and outputs 4D meshes for the human and object, i.e., 3D meshes over time. Humans are represented as SMPL-X meshes [44], while object meshes are acquired with an Artec hand-held scanner. Some dataset frames along with the reconstructed meshes are shown in Fig. 1 and Fig. 4; see also the video on our website. Reconstructions look natural, with plausible contact between the human and the object.

Dataset Statistics. InterCap has 223 RGB-D videos with a total of 67,357 multi-view frames (6 RGB-D images each). For a comparison with other datasets, see Tab. 1.

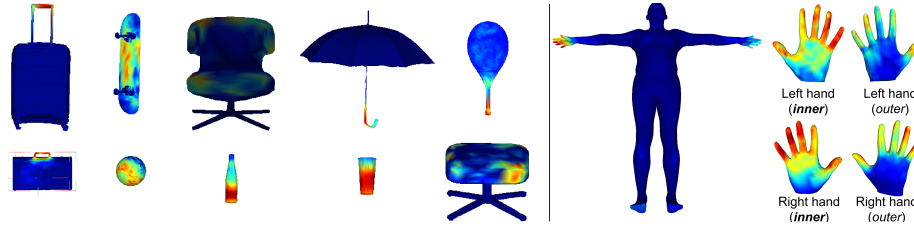


Fig. 5. Contact heatmaps for each object (across all subjects) and the human body (across all objects and subjects). Contact likelihood is color-coded; high likelihood is shown with red, and low with blue. Color-coding is normalized separately for each object, the body, and each hand.

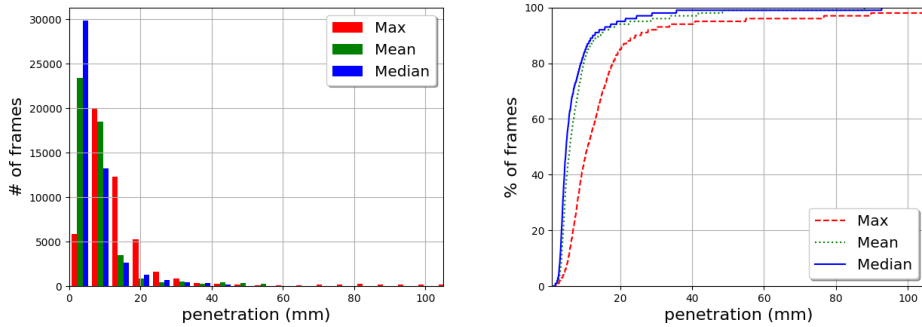


Fig. 6. Statistics of human-object mesh penetration for all InterCap sequences. **Left:** The number of frames (Y-axis) with a certain penetration depth (X-axis). **Right:** The percentage of frames (Y-axis) with a penetration depth below a threshold (X-axis). In the legend, “Max”, “Mean” and “Median” refer to three ways of reporting the penetration for each frame, i.e., taking the maximum, mean and median value of the penetration depth of all vertices, respectively.

5 Experiments

Contact Heatmaps. Figure 5-left shows contact heatmaps on each object, across all subjects. We follow the protocol of GRAB [54], which uses a proximity metric on reconstructed human and object meshes. First, we compute per-frame binary contact maps by thresholding (at 4.5mm) the distances from each body vertex to the closest object surface point. Then, we integrate these maps over time (and subjects) to get “heatmaps” encoding contact likelihood. InterCap reconstructs human and object meshes accurately enough so that contact heatmaps agree with object affordances, e.g., the handle of the suitcase, umbrella and tennis racquet are likely to be grasped, the upper skateboard surface is likely to be contacted by the foot, and the upper stool surface by the buttocks.

Figure 5-right shows heatmaps on the body, computed across all subjects and objects. Heatmaps show that most of InterCap’s interactions involve mainly the right hand. Contact on the palm looks realistic, and is concentrated on the fingers and MCP joints. The “false” contact on the dorsal side is attributed to our challenging camera setup and interaction scenarios, as well as some reconstruction jitter.

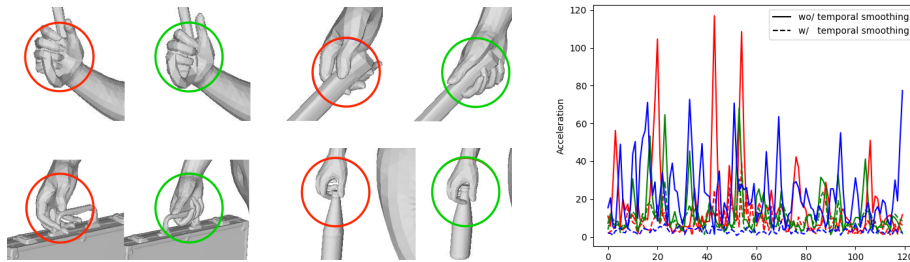


Fig. 7. **Left:** Qualitative ablation of our contact term. Each pair of images shows results w/o (red) and w/ (green) the contact term. Encouraging contact results in more natural hand poses and hand-object grasps. **Right:** Acceleration of a random vertex w/ (dashed line) and w/o (solid line) temporal smoothing for 3 sequences (shown with different color) over the first 120 frames. Dashed lines (w/ temporal smoothing) correspond to lower acceleration, i.e., less jitter.

Penetration. We evaluate the penetration between human and object meshes for all sequences of our dataset. We follow the protocol of GRAB et al. [54]; we first find the “contact frames” for which there is at least minimal human-object contact, and then report statistics for these. In Fig. 6-left we show the distribution of penetrations, i.e., the number of “contact frames” (Y axis) with a certain mesh penetration depth (X axis). In Fig. 6-right we show the cumulative distribution of penetration, i.e., the percentage of “contact frames” (Y axis) for which mesh penetration is below a threshold (X axis). Roughly 60% of “contact frames” have $\leq 5\text{mm}$, 80% $\leq 7\text{mm}$, and 98% $\leq 20\text{mm}$ mean penetration. The average penetration depth over all “contact frames” is 7.2 mm.

Fitting Accuracy. For every frame, we compute the distance from each mesh vertex to the closest point-cloud (PCL) point; for each human or object mesh we take into account only the respective PCL area obtained with PointRend [28] segmentation. The mean vertex-to-PCL distance is 20.29 mm for the body, and 18.50 mm for objects. In comparison, PROX-D [15], our base method, achieves an error of 13.02 mm for the body. This is expected since PROX-D is free to change the body shape to fit each individual frame, while our method estimates a single body shape for the whole sequence. SMPLify-X [44] achieves an mean error of 79.54 mm, for VIBE the mean error is 55.59 mm, while ExPose gets an mean error of 71.78 mm. These numbers validate the effectiveness of our method for body tracking. Note that these methods are based on monocular RGB images only, so there is not enough information for them to accurately estimate the global position of the 3D body meshes. Thus we first align the output meshes with the point clouds, then compute the error. Note that the error is bounded from below for two reasons: (1) it is influenced by factory-design imperfections in the synchronization of Azure Kinects, and (2) some vertices reflect body/object areas that are occluded during interaction and their closest PCL point is a wrong correspondence. Despite this, InterCap empirically estimates reasonable bodies, hands and objects in interaction, as reflected in the contact heatmaps and penetration metrics discussed above.

Ablation of Contact Term. Figure 7-left shows results with-/out our term that encourages body-object contact; visualization “zooms” into hand-object grasps. We see

that encouraging contact yields more natural hand poses and fewer interpenetrations. This is backed up by the contact heatmaps and penetration metrics discussed above.

Ablation of Temporal Smoothing Term. Figure 7-right shows results with-/out our temporal smoothing term. Each solid line shows the acceleration of a randomly chosen vertex without the temporal smoothness term; we show 3 different motions. The dashed lines of the same color show the same motions with the smoothness term; these are clearly smoother. We empirically find that the learned motion prior of Zhang et al. [64] produces a more natural motion than handcrafted ones [22].

Discussion on Jitter. Despite the smoothing, some jitter is still inevitable. We attribute this to two factors: (1) OpenPose and Mask-RCNN are empirically relatively sensitive to occlusions and illumination (e.g., reflections, shadows, poor lighting); the data terms for fitting 3D models depend on these. (2) Azure Kinects have a reasonable synchronization, yet, there is still a small delay among cameras to avoid depth-camera interference; the point cloud “gathered” across views is a bit “patchy” as information pieces have a small time difference. The jitter is more intense for hands relatively to the body, due to their low image resolution, motion blur, and coarse point clouds. Despite these challenges, InterCap is a good step towards capturing everyday whole-body interactions with commodity hardware. Future work will study advanced motion priors.

6 Discussion

Here we focus on whole-body human interaction with everyday rigid objects. We present a novel method, called InterCap, that reconstructs such interactions from multi-view full-body videos, including natural hand poses and contact with objects. With this method, we capture the novel InterCap dataset, with a variety of people interacting with several common objects. The dataset contains reconstructed 3D meshes for the whole body and the object over time (i.e., 4D meshes), as well as plausible contacts between them. In contrast to most previous work, our method uses no special devices like optical markers or IMUs, but only several consumer-level RGB-D cameras. Our setup is lightweight and has the potential to be used in daily scenarios. Our method estimates reasonable hand poses even when there is heavy occlusion between hands and the object. In future work, we plan to study interactions with smaller objects and dexterous manipulation. Our data and code are available at <https://intercap.is.tue.mpg.de>.

7 Acknowledgements

We thank Chun-Hao Paul Huang, Hongwei Yi, Jiayang Shang, as well as Mohamed Hassan for helpful discussion about technical details. We thank Taylor McConnell, Galina Henz, Marku Höschle, Senya Polikovsky, Matvey Safroshkin and Tsvetelina Alexiadis for the data collection and data cleaning. We thank all the participants of our experiments. We also thank Benjamin Pellkofer for the IT and website support.

The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting OT. This work was supported by the German Federal Ministry of Education and Research (BMBF): Tübingen AI Center, FKZ: 01IS18039B.

Conflict of Interest. Disclosure: https://files.is.tue.mpg.de/black/CoI_GCPR_2022.txt.

Bibliography

- [1] Alldieck, T., Magnor, M., Xu, W., Theobalt, C., Pons-Moll, G.: Video based reconstruction of 3D people models. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 8387–8397 (2018) [3](#)
- [2] Angelov, D., Srinivasan, P., Koller, D., Thrun, S., Rodgers, J., Davis, J.: SCAPE: Shape completion and animation of people. *Transactions on Graphics (TOG)* **24**(3), 408–416 (2005) [3](#)
- [3] Bhatnagar, B.L., Xie, X., Petrov, I.A., Sminchisescu, C., Theobalt, C., Pons-Moll, G.: BEHAVE: Dataset and method for tracking human object interactions. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 15935–15946 (2022) [2](#), [4](#), [5](#)
- [4] Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In: *European Conference on Computer Vision (ECCV)*. vol. 9909, pp. 561–578 (2016) [2](#), [3](#)
- [5] Cao, Z., Gao, H., Mangalam, K., Cai, Q., Vo, M., Malik, J.: Long-term human motion prediction with scene context. In: *European Conference on Computer Vision (ECCV)*. vol. 12346, pp. 387–404 (2020) [4](#), [5](#)
- [6] Cao, Z., Hidalgo, G., Simon, T., Wei, S.E., Sheikh, Y.: OpenPose: Realtime multi-person 2D pose estimation using part affinity fields. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **43**(1), 172–186 (2019) [3](#)
- [7] Choutas, V., Pavlakos, G., Bolkart, T., Tzionas, D., Black, M.J.: Monocular expressive body regression through body-driven attention. In: *European Conference on Computer Vision (ECCV)*. vol. 12355, pp. 20–40 (2020) [3](#)
- [8] De Aguiar, E., Stoll, C., Theobalt, C., Ahmed, N., Seidel, H.P., Thrun, S.: Performance capture from sparse multi-view video. *Transactions on Graphics (TOG)* **27**(3), 1–10 (2008) [3](#)
- [9] Dong, J., Fang, Q., Jiang, W., Yang, Y., Huang, Q., Bao, H., Zhou, X.: Fast and robust multi-person 3D pose estimation and tracking from multiple views. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **14**(8), 1–12 (2021) [4](#)
- [10] Dong, J., Jiang, W., Huang, Q., Bao, H., Zhou, X.: Fast and robust multi-person 3D pose estimation from multiple views. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 7792–7801 (2019) [4](#)
- [11] Dong, Z., Song, J., Chen, X., Guo, C., Hilliges, O.: Shape-aware multi-person pose estimation from multi-view images. In: *International Conference on Computer Vision (ICCV)*. pp. 11158–11168 (2021) [4](#)
- [12] Geman, S., McClure, D.E.: Statistical methods for tomographic image reconstruction. In: *Proceedings of the 46th Session of the International Statistical Institute, Bulletin of the ISI*. vol. 52 (1987) [8](#)
- [13] Hamer, H., Schindler, K., Koller-Meier, E., Van Gool, L.: Tracking a hand manipulating an object. In: *International Conference on Computer Vision (ICCV)*. pp. 1475–1482 (2009) [4](#)

- [14] Hampali, S., Rad, M., Oberweger, M., Lepetit, V.: HONotate: A method for 3D annotation of hand and object poses. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3193–3203 (2020) [4](#), [5](#), [6](#), [7](#)
- [15] Hassan, M., Choutas, V., Tzionas, D., Black, M.J.: Resolving 3D human pose ambiguities with 3D scene constrains. In: *International Conference on Computer Vision (ICCV)*. pp. 2282–2292 (2019) [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [12](#)
- [16] Hassan, M., Ghosh, P., Tesch, J., Tzionas, D., Black, M.J.: Populating 3D scenes by learning human-scene interaction. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 14708–14718 (2021) [2](#), [4](#)
- [17] Hasson, Y., Tekin, B., Bogo, F., Laptev, I., Pollefeys, M., Schmid, C.: Leveraging photometric consistency over time for sparsely supervised hand-object reconstruction. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 568–577 (2020) [4](#)
- [18] Hasson, Y., Varol, G., Tzionas, D., Kalevatykh, I., Black, M.J., Laptev, I., Schmid, C.: Learning joint reconstruction of hands and manipulated objects. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 11807–11816 (2019) [4](#)
- [19] He, Y., Yan, R., Fragkiadaki, K., Yu, S.I.: Epipolar transformers. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 7776–7785 (2020) [4](#)
- [20] Hu, Y.T., Chen, H.S., Hui, K., Huang, J.B., Schwing, A.G.: SAIL-VOS: Semantic amodal instance level video object segmentation - a synthetic dataset and baselines. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3105–3115 (2019) [5](#)
- [21] Huang, C.H.P., Yi, H., Höschle, M., Safroshkin, M., Alexiadis, T., Polikovsky, S., Scharstein, D., Black, M.J.: Capturing and inferring dense full-body human-scene contact. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 13274–13285 (2022) [4](#), [5](#)
- [22] Huang, Y., Bogo, F., Lassner, C., Kanazawa, A., Gehler, P.V., Romero, J., Akhter, I., Black, M.J.: Towards accurate marker-less human shape and pose estimation over time. In: *International Conference on 3D Vision (3DV)*. pp. 421–430 (2017) [3](#), [4](#), [8](#), [13](#)
- [23] Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **36**(7), 1325–1339 (2014) [4](#), [5](#)
- [24] Isakov, K., Burkov, E., Lempitsky, V., Malkov, Y.: Learnable triangulation of human pose. In: *International Conference on Computer Vision (ICCV)*. pp. 7717–7726 (2019) [4](#)
- [25] Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 7122–7131 (2018) [3](#)
- [26] Karunratanakul, K., Yang, J., Zhang, Y., Black, M.J., Muandet, K., Tang, S.: Grasping Field: Learning implicit representations for human grasps. In: *International Conference on 3D Vision (3DV)*. pp. 333–344 (2020) [4](#)
- [27] Kato, H., Ushiku, Y., Harada, T.: Neural 3D mesh renderer. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3907–3916 (2018) [7](#)

- [28] Kirillov, A., Wu, Y., He, K., Girshick, R.: PointRend: Image segmentation as rendering. In: Computer Vision and Pattern Recognition (CVPR). pp. 9799–9808 (2020) [6](#), [12](#)
- [29] Kocabas, M., Athanasiou, N., Black, M.J.: VIBE: Video inference for human body pose and shape estimation. In: Computer Vision and Pattern Recognition (CVPR). pp. 5252–5262 (2020) [3](#)
- [30] Li, X., Liu, S., Kim, K., Wang, X., Yang, M., Kautz, J.: Putting humans in a scene: Learning affordance in 3D indoor environments. In: Computer Vision and Pattern Recognition (CVPR). pp. 12368–12376 (2019) [2](#), [4](#)
- [31] Liu, Y., Stoll, C., Gall, J., Seidel, H.P., Theobalt, C.: Markerless motion capture of interacting characters using multi-view image segmentation. In: Computer Vision and Pattern Recognition (CVPR). pp. 1249–1256 (2011) [3](#)
- [32] Loper, M., Mahmood, N., Black, M.J.: MoSh: Motion and shape capture from sparse markers. *Transactions on Graphics (TOG)* **33**(6), 1–13 (2014) [4](#)
- [33] Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. *Transactions on Graphics (TOG)* **34**(6), 248:1–248:16 (2015) [3](#)
- [34] Loper, M.M., Black, M.J.: OpenDR: An approximate differentiable renderer. In: European Conference on Computer Vision (ECCV). vol. 8695, pp. 154–169 (2014) [7](#)
- [35] Mahmood, N., Ghorbani, N., F. Troje, N., Pons-Moll, G., Black, M.J.: AMASS: Archive of motion capture as surface shapes. In: International Conference on Computer Vision (ICCV). pp. 5441–5450 (2019) [4](#), [5](#)
- [36] von Marcard, T., Henschel, R., Black, Michael J. and Rosenhahn, B., Pons-Moll, G.: Recovering accurate 3D human pose in the wild using IMUs and a moving camera. In: European Conference on Computer Vision (ECCV). vol. 11214, pp. 614–631 (2018) [3](#), [4](#), [5](#)
- [37] Mehta, D., Sridhar, S., Sotnychenko, O., Rhodin, H., Shafiei, M., Seidel, H.P., Xu, W., Casas, D., Theobalt, C.: VNect: Real-time 3D human pose estimation with a single RGB camera. *Transactions on Graphics (TOG)* **36**(4), 44:1–44:14 (2017) [3](#)
- [38] Microsoft: Azure Kinect SDK (K4A). <https://github.com/microsoft/Azure-Kinect-Sensor-SDK> (2022) [6](#)
- [39] Newell, A., Yang, K., Deng, J.: Stacked hourglass networks for human pose estimation. In: European Conference on Computer Vision (ECCV). vol. 9912, pp. 483–499 (2016) [3](#)
- [40] Oikonomidis, I., Kyriazis, N., Argyros, A.A.: Full DOF tracking of a hand interacting with an object by modeling occlusions and physical constraints. In: International Conference on Computer Vision (ICCV). pp. 2088–2095 (2011) [4](#), [6](#)
- [41] Omran, M., Lassner, C., Pons-Moll, G., Gehler, P., Schiele, B.: Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In: International Conference on 3D Vision (3DV). pp. 484–494 (2018) [3](#)
- [42] Osman, A.A.A., Bolkart, T., Tzionas, D., Black, M.J.: SUPR: A sparse unified part-based human body model. In: European Conference on Computer Vision (ECCV) (2022) [3](#)
- [43] Osman, A.A., Bolkart, T., Black, M.J.: STAR: Sparse trained articulated human body regressor. In: European Conference on Computer Vision (ECCV). vol. 12351, pp. 598–613 (2020)

- [44] Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3D hands, face, and body from a single image. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 10975–10985 (2019) [3](#), [7](#), [10](#), [12](#)
- [45] Pons-Moll, G., Baak, A., Helten, T., Müller, M., Seidel, H.P., Rosenhahn, B.: Multisensor-fusion for 3D full-body human motion capture. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 663–670 (2010) [3](#)
- [46] Qiu, H., Wang, C., Wang, J., Wang, N., Zeng, W.: Cross view fusion for 3D human pose estimation. In: *International Conference on Computer Vision (ICCV)*. pp. 4341–4350 (2019) [4](#)
- [47] Rhodin, H., Robertini, N., Casas, D., Richardt, C., Seidel, H.P., Theobalt, C.: General automatic human shape and motion capture using volumetric contour cues. In: *European Conference on Computer Vision (ECCV)*. vol. 9909, pp. 509–526 (2016) [3](#)
- [48] Rogez, G., III, J.S.S., Ramanan, D.: Understanding everyday hands in action from RGB-D images. In: *International Conference on Computer Vision (ICCV)*. pp. 3889–3897 (2015) [4](#)
- [49] Romero, J., Kjellström, H., Kragic, D.: Hands in action: Real-time 3D reconstruction of hands in interaction with objects. In: *International Conference on Robotics and Automation (ICRA)*. pp. 458–463 (2010) [4](#)
- [50] Romero, J., Tzionas, D., Black, M.J.: Embodied hands: Modeling and capturing hands and bodies together. *Transactions on Graphics (TOG)* **36**(6), 245:1–245:17 (2017) [3](#)
- [51] Savva, M., Chang, A.X., Hanrahan, P., Fisher, M., Nießner, M.: PiGraphs: Learning interaction snapshots from observations. *Transactions on Graphics (TOG)* **35**(4), 139:1–139:12 (2016) [4](#), [5](#)
- [52] Sigal, L., Balan, A., Black, M.J.: HumanEva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion. *International Journal of Computer Vision (IJCV)* **87**(1-2), 4–27 (2010) [4](#), [5](#)
- [53] Sun, J., Wang, Z., Zhang, S., He, X., Zhao, H., Zhang, G., Zhou, X.: OnePose: One-shot object pose estimation without CAD models. In: *CVPR*. pp. 6825–6834 (2022) [5](#)
- [54] Taheri, O., Ghorbani, N., Black, M.J., Tzionas, D.: GRAB: A dataset of whole-body human grasping of objects. In: *European Conference on Computer Vision (ECCV)*. vol. 12349, pp. 581–600 (2020) [2](#), [4](#), [5](#), [10](#), [11](#), [12](#)
- [55] Tu, H., Wang, C., Zeng, W.: VoxelPose: Towards multi-camera 3D human pose estimation in wild environment. In: *European Conference on Computer Vision (ECCV)*. vol. 12346, pp. 197–212 (2020) [4](#)
- [56] Tzionas, D., Ballan, L., Srikantha, A., Aponte, P., Pollefeys, M., Gall, J.: Capturing hands in action using discriminative salient points and physics simulation. *International Journal of Computer Vision (IJCV)* **118**(2), 172–193 (2016) [4](#), [6](#)
- [57] Varol, G., Laptev, I., Schmid, C.: Long-term temporal convolutions for action recognition. *Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* **40**(6), 1510–1517 (2017) [3](#)
- [58] Wei, S.E., Ramakrishna, V., Kanade, T., Sheikh, Y.: Convolutional pose machines. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 4724–4732 (2016) [3](#)

- [59] Xu, H., Bazavan, E.G., Zanfir, A., Freeman, W.T., Sukthankar, R., Sminchisescu, C.: GHUM & GHUML: Generative 3D human shape and articulated pose models. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 6183–6192 (2020) [3](#)
- [60] Xu, W., Chatterjee, A., Zollhöfer, M., Rhodin, H., Mehta, D., Seidel, H.P., Theobalt, C.: MonoPerfCap: Human performance capture from monocular video. *Transactions on Graphics (TOG)* **37**(2), 1–15 (2018) [3](#)
- [61] Yao, B., Fei-Fei, L.: Modeling mutual context of object and human pose in human-object interaction activities. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 17–24 (2010) [4](#)
- [62] Yi, H., Huang, C.H.P., Tzionas, D., Kocabas, M., Hassan, M., Tang, S., Thies, J., Black, M.J.: Human-aware object placement for visual environment reconstruction. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 3959–3970 (2022) [4](#)
- [63] Zhang, J.Y., Pepose, S., Joo, H., Ramanan, D., Malik, J., Kanazawa, A.: Perceiving 3D human-object spatial arrangements from a single image in the wild. In: *European Conference on Computer Vision (ECCV)* (2020) [4](#)
- [64] Zhang, S., Zhang, Y., Bogo, F., Pollefeys, M., Tang, S.: Learning motion priors for 4D human body capture in 3D scenes. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 11323–11333 (2021) [9](#), [13](#)
- [65] Zhang, Y., Hassan, M., Neumann, H., Black, M.J., Tang, S.: Generating 3D people in scenes without people. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 6193–6203 (2020) [2](#), [4](#)
- [66] Zhang, Y., An, L., Yu, T., Li, X., Li, K., Liu, Y.: 4D association graph for realtime multi-person motion capture using multiple video cameras. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 1321–1330 (2020) [4](#)
- [67] Zhang, Y., Li, Z., An, L., Li, M., Yu, T., Liu, Y.: Light-weight multi-person total capture using sparse multi-view cameras. In: *International Conference on Computer Vision (ICCV)*. pp. 5560–5569 (2021) [4](#)
- [68] Zollhöfer, M., Stotko, P., Görlitz, A., Theobalt, C., Nießner, M., Klein, R., Kolb, A.: State of the art on 3D reconstruction with RGB-D cameras. *Computer Graphics Forum (CGF)* **37**(2), 625–652 (2018) [2](#)