

4D Human Body Capture from Egocentric Video via 3D Scene Grounding

Miao Liu¹, Dexin Yang², Yan Zhang², Zhaopeng Cui³, James M. Rehg¹, Siyu Tang²

¹ Georgia Institute of Technology, Atlanta, United States

² ETH Zürich, Switzerland

³ Zhejiang University, Hangzhou, China

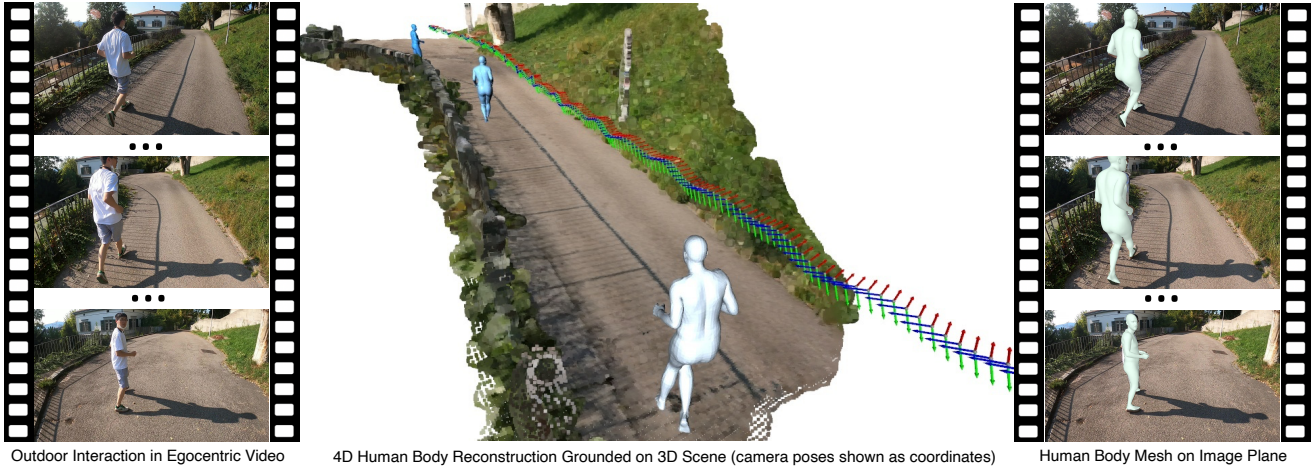


Figure 1. As shown in the middle figure, we seek to reconstruct 4D second-person human body meshes that are grounded on the 3D scene captured in an egocentric view. Our method exploits 2D observations from the entire video sequence and the 3D scene context to optimize human body models over time, and thereby leads to more accurate human motion capture and more realistic human-scene interaction.

Abstract

To understand human daily social interaction from egocentric perspective, we introduce a novel task of reconstructing a time series of second-person¹ 3D human body meshes from monocular egocentric videos. The unique viewpoint and rapid embodied camera motion of egocentric videos raise additional technical barriers for human body capture. To address those challenges, we propose a novel optimization-based approach that leverages 2D observations of the entire video sequence and human-scene interaction constraint to estimate second-person human poses, shapes and global motion that are grounded on the 3D environment captured from the egocentric view. We conduct detailed ablation studies to validate our design choice. Moreover, we compare our method with previous state-of-the-art method on human motion capture from monocular video, and show that our method estimates more accurate human-

body poses and shapes under the challenging egocentric setting. In addition, we demonstrate that our approach produces more realistic human-scene interaction. Our project page is available at: <https://aptx4869lm.github.io/4DEgocentricBodyCapture/>

1. Introduction

Continuous advancements in the capabilities of Augmented Reality (AR) headsets promise new trends of entertainment, communication, healthcare, and productivity, and point towards a revolution in how we interact with the world and communicate with each other. Egocentric vision is a key building block for these emerging capabilities, as AR experiences can benefit from an accurate understanding of the user’s perception, attention, and actions. Substantial progress has been made in understanding human-object interaction [41, 11, 9, 30, 28, 13, 29, 32, 36] from egocentric videos. Additional works investigated social interactions by leveraging egocentric videos to reason about social

¹We refer the social partner of the camera-wearer as “second-person” throughout the paper. The same notation was also adopted in [38, 59]

signals of the second-person [7, 58, 12, 48, 59, 57, 8]. However, these works are largely limited to the analysis of head pose, gaze behavior, and simple gestures. Future intelligent AR headsets should also have the capacity of capturing the subtle nuances of second-person body pose or even generating plausible interactive 3D avatar that grounded on the 3D scene captured from egocentric point of view. To this end, we introduce a novel task of 4D second-person full body capture from egocentric videos. As shown in Fig. 1, we seek to reconstruct time series of motion plausible 3D second-person body meshes that are grounded on 3D scene captured from egocentric perspective.

3D human body capture from videos is a key challenge in computer vision, which has received substantial attention over the years [22, 25, 19, 52]. However, none of previous works considered the challenging setting of reconstructing 3D second-person human body from egocentric perspective². The unique viewpoints and embodied camera motions that arise in egocentric video create formidable technical obstacles to 3D body estimation, causing previous SOTA methods for video-based motion capture to fail. For example, the close interpersonal distances that characterize social interactions result in partial observation of the second-person as body parts move in and out of frame. The drastic camera motion also leads to additional barrier of human kinematic estimation, as the second-person motion is entangled with the embodied movement of the camera wearer.

To address the challenging artifacts of egocentric videos, we propose to a novel optimization-based method that jointly considers time series of 2D observations and 3D scene information. Our key insight is that combining the 2D observations from the entire video sequence provides additional evidence for estimating human body models from frames with only partial observation, and 3D scene also constrains the human body pose and motion. Our approach begins with the use of Structure-from-Motion (SfM) to estimate the camera trajectory and to reconstruct the 3D environment. Note that the 3D scene and body reconstruction from monocular videos is up to a scale. Therefore, directly projecting the 3D body meshes into the reconstructed 3D scene and enforcing human-scene contact will result in unrealistic human-scene interaction. To overcome this challenge, we carefully design the optimization method so that it can not only encourage human-scene contact, but also estimate scale difference between 3D human body and scene reconstruction. We further enforce temporal coherency by uniting time series of body model with temporal prior to recover more plausible *global human motion* even when the second-person body captured by the egocentric view is only partially observable.

²We note that another branch of prior work addresses the related but quite different task of predicting the 3D body pose of the *camera-wearer* from egocentric video [38, 20, 60, 50].

To study this challenging problem of reconstructing 4D second-person body pose and shape from egocentric videos and to validate our proposed approach, we introduce a new egocentric video dataset – EgoMoCap. This dataset captures various human social behaviors in outdoor environment, which serves as an ideal vehicle to study the problem of second-person human body reconstruction from egocentric perspective. We conduct detailed ablation studies on this dataset to show the benefits of our method. We further compare our approach with previous state-of-the-art method on human motion capture from monocular videos, and show our method can address the challenging cases where second-person human body is partially observable. Besides improving the body reconstruction accuracy, we also demonstrate that our method solves the relative scale difference between 3D scene reconstruction and 3D human body reconstruction from monocular videos, and thereby produces more realistic human-scene interaction.

In summary, our work has the following contributions:

- We introduce a new problem of reconstructing time series of second-person poses and shapes from egocentric videos. To the best of our knowledge, we are also the first to address capturing global human motion grounded on the 3D environment.
- We propose a novel optimization-based approach that jointly considers time series of 2D observation and 3D scene context for accurate 4D human body capture. In addition, our approach seeks to address the scale ambiguity of 3D reconstruction from monocular videos.
- We present a new egocentric dataset – EgoMoCap that captures human social interactions in outdoor environment. And we conduct detailed experiments on EgoMoCap dataset and show that our approach can reconstruct more accurate 4D second-person human body, and encourage more realistic human-scene interaction.

2. Related Work

The most relevant works to ours are those investigations on 4D human body reconstruction and human-scene interaction. Our work is also related to recent efforts on reasoning about social interaction from egocentric perspective. Furthermore, we compare our EgoMoCap dataset with other egocentric human interaction datasets.

4D Human Body Reconstruction. A rich set of literature has covered the topic of human body reconstruction. Previous approaches [4, 39, 26, 21, 3, 34, 49, 42] have demonstrated great success on inferring 3D human pose and shape from a single image. Here, we focus on discussing those works on inferring time series of 3D human body poses and shapes from videos. Alldieck et al. [2] proposed to

use optical flow to estimate temporal coherent human bodies from monocular videos. Tung et al. [51] introduced a self-supervised learning method that uses optical flow, silhouettes, and keypoints to estimate SMPL human body parameters from two consecutive video frames. [23, 40] used fully convolutional network to predict 3D human pose from 2D images sequences. Kocabas et al. [25] proposed an adversarial learning framework to produce realistic and accurate human pose and motion from video sequences. Shimada et al. [46] used physical engine to capture physically plausible and temporally stable global 3D human motion. All those deep learning based methods *assumed a fixed camera view and fully observable human body*. Those assumptions do not hold under egocentric setting. Several optimization-based methods [52, 19, 55] considered the moving camera scenarios. [52] proposed to jointly optimize the camera pose and human body model, yet their method requires additional IMU sensor data. [19] enforced temporal coherence to reconstruct reasonable body pose from monocular videos with moving camera. Wang et al. [55] proposed to utilize multiple cameras for outdoor human motion capture. Those methods only targeted at local human kinematic motion without reasoning the 3D scene context. In contrast, we seek to estimate the global human motion grounded on 3D scene from *only monocular egocentric videos*.

Human-Scene Interaction. Several investigations on human-scene interaction seek to reason about environment affordance [36, 16, 15, 10, 54, 27, 6, 35]. Our work is more relevant to those efforts on using the environment cues to better capture 3D human body. Savva et al. [44] proposed to learn a probabilistic model that captures how human interact with the indoor scene from RGB-D sensors. Li et al. [31] factorized estimating 3D person-object interactions into an optimal control problem, and used contact constraints to recover human motion and contact forces from monocular videos. Zhang et al. [61] proposed an optimization-based framework that incorporates the scale loss to jointly reconstruct the 3D spatial arrangement and shape of humans and objects in the scene from a single image. Hassan et al. [17] made use of the 3D scene context – obtained from 3D scan, to estimate more accurate human pose and shape from single image. Zhang et al. [62, 63] further studied the problem of generating plausible human body grounded on 3D scene prior. Despite those progress on using scene information to estimate 3D human body model parameters, none of them considered the egocentric camera motion, 3D scene context from monocular videos, and global human motion grounded on 3D scene in one-shot as in our proposed approach.

Egocentric Social Interaction. Understanding human social interaction has been the subject of many recent efforts in egocentric vision. Several previous works studied human attention during social interaction. Ye et al. [58] proposed to use pose-dependent appearance model to estimate

the eye contact of children. Chong et al. [7] introduced a novel multi-task learning method to predict gaze directions from various kinds of datasets. Park et al. [48] considered the challenging problem of social saliency prediction. Fathi et al. [12] utilized face, attention, and head motion to recognize social interactions. More recently, a few works considered novel vision problems in egocentric social interaction. Yagi [57] addressed the task of localizing future position of target person from egocentric videos. Yonetani et al. [59] proposed to use features from both the first-person and second-person points-of-view for recognizing micro-actions and reactions during social interaction. Ng et al. [38] proposed to use the second-person body pose as additional cues for predicting the egocentric body pose during human interaction. Those previous works studied various signals during human social interaction, however none of them targeted at second-person full body capture. Our work seeks to bridge this gap and points to new research directions in egocentric social interaction.

Egocentric Human Interaction Datasets. Several egocentric datasets target the analysis of human social behavior during naturalistic interactions. Fathi et al. [12] presented an egocentric dataset for the detection and recognition of fixed categories of conversational interactions within a social group. The NUS Dataset [37] and JPL Dataset [43] support more general human interaction classification tasks. Yonetani et al. [59] collected a paired egocentric human interaction dataset to study human action and reaction. While prior datasets focused on social interaction recognition, Park et al. introduced an RGB-D egocentric dataset – EgoMotion [47], for forecasting a walking trajectory based on interaction with the environment. More recently, the You2Me dataset [38] was proposed to study the problem of egocentric body pose prediction. However, none of those datasets were designed to study the *second-person body pose*, which is the focus and contribution of our work. In prior datasets, the majority of second-person body captures are either largely occluded by objects or frequently truncated by the frustum, which makes their utilization for full body capture infeasible. In contrast, our EgoMoCap dataset focuses on outdoor social interaction scenarios that have less foreground occlusion on second-person body.

3. Method

We denote an input monocular egocentric video as $x = (x^1, \dots, x^t)$ with its frame x^t indexed by time t . We estimate the human body pose and shape at each time step from input x . Due to the unique viewpoint of egocentric video, the captured second-person body is partially observable within a time window. In addition, the second-person body motion is entangled with the camera motion, and therefore incurs additional barrier to enforce temporal coherency. To address those challenges, we propose a novel optimization

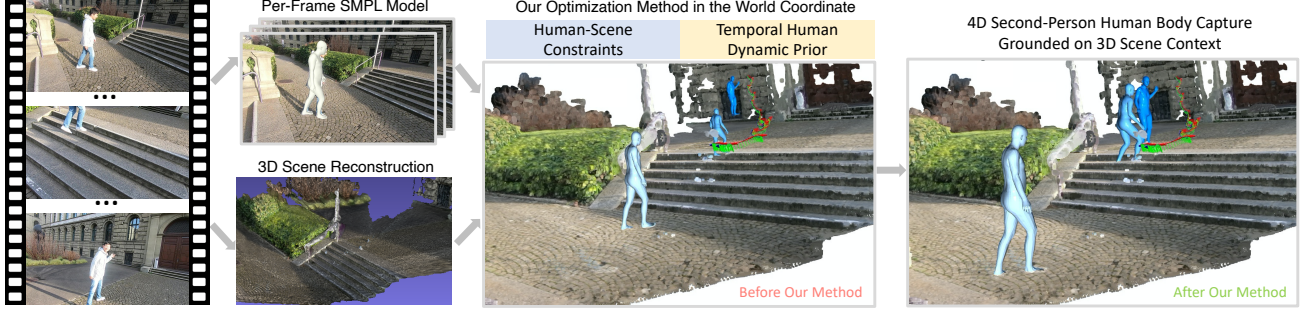


Figure 2. Overview of our method. We introduce an optimization-based method that makes use of human-scene constraints and temporal human dynamic prior to reconstruct time series of 4D human body poses and shapes that grounded on the 3D environment. Our method thereby addresses challenging cases where human body is partially observable (middle figure on the left) and encourages more realistic human-scene interaction (figure on the right).

method that jointly considers the 2D observation of the entire video sequence and 3D scene for more accurate 4D human body reconstruction. We illustrate our method in Fig. 2. Specifically, we first recover the 3D human body at each time instant from the 2D observation of x^t . We then use Structure from Motion (SfM) to project a sequence of 3D body meshes into the 3D world coordinate, and further adopt a contact term to encourage human-scene interaction. In addition, we combine the 2D cues from entire video sequences for reconstructing temporal coherent time series of body poses using human dynamic prior. In following sections, we introduce each component of our method.

3.1. Human Body Model

To better understand various signals during social interaction, we use the differentiable body model SMPL-X [39] to jointly capture human body, hands, and facial expression. SMPL-X produces a body mesh of a fixed topology with 10,475 vertices, using a compact set of body configuration parameters. Specifically, the shape parameter β represents how individuals vary in height, weight, and body proportions, θ encodes the 3D body pose, hand pose and facial expression information, and γ denotes the body translation. Formally, the SMPL-X function is defined as $M_b(\beta, \theta, \gamma)$. It outputs a 3D body mesh as $M_b = (V_b, F_b)$, where $V_b \in \mathbb{R}^{N_b \times 3}$ and F_b denote the body vertices and triangular faces, respectively.

Similar to [39, 4], we factorize fitting the SMPL-X model to each video frame as an optimization problem. Formally, we optimize (β, θ, γ) by minimizing:

$$E_M(\beta, \theta, \gamma, K, J_{est}) = E_J(\beta, \theta, \gamma, K, J_{est}) + \lambda_\beta E_\beta(\beta) + \lambda_\theta E_\theta(\theta), \quad (1)$$

where K is the intrinsic camera parameters; the shape prior term $E_\beta(\beta)$ is learned from SMPL-X model body shape training data and the pose prior term $E_\theta(\theta)$ is learned from CMU MoCap dataset [1]; λ_β and λ_θ denote the weights of $E_\beta(\beta)$ and $E_\theta(\theta)$; E_J refers to the energy function that

minimizes the weighted robust distance between the 2D projection of the body joints, hand joints and face landmarks, and the corresponding 2D joints estimation from OpenPose [5, 56]. E_J is given by:

$$E_J(\beta, \theta, \gamma, L, J_{est}) = \sum_{joint\ i} k_i w_i \rho_J(\Pi_K(R_{\theta_\gamma}(J^i(\beta)) - J_{est}^i), \quad (2)$$

where $J(\cdot)$ returns 3D joints location based on embedded shape parameters β , and $R_{\theta_\gamma}(\cdot)$ transforms the joints along the kinematic tree according to the pose θ and body translation γ ; Π_K is the 3D to 2D projection function based on intrinsic parameters K ; J_{est} refers to the 2D joints estimation from OpenPose; w_i is the 2D joints detection confident score which accounts for the noises of 2D joints estimation; k_i is the per-joint weights for annealed optimization as in [39]; ρ_J denotes a robust Geman-McClure error function [14] that downweights outliers, which is given by:

$$\rho_J(e) = \frac{e^2}{\sigma_J^2 + e^2} \quad (3)$$

where e is the residual error, and σ_J is the robustness constant chosen empirically.

3.2. Egocentric Camera Representation

To capture 4D second-person bodies that are grounded on the 3D scene from egocentric videos, we need to take the embodied camera motion into consideration. Here we elaborate the egocentric camera representation adopted in our method. Formally, we denote $T_{cb} \in \mathbb{R}^{4 \times 4}$ as the transformation from the human body coordinate to the egocentric camera coordinate, and T_{wc} as the transformation from the egocentric camera coordinate to the world coordinate. Note that $T_{cb} \in \mathbb{R}^{4 \times 4}$ is derived from the translation parameter γ of SMPL-X model fitting introduced aforementioned section, while T_{wc} is returned from COLMAP Structure from Motion (SfM) [45]. In order to utilize the 3D scene context and enforce the temporal coherency on reconstructed human body meshes, we project the 3D second-person body vertices V_b into world coordinate using human

body to world transformation T_{wb} , which is given by:

$$\hat{V}_{wb}^t = T_{wb}^t \hat{V}_b^t = T_{wc}^t T_{cb}^t \hat{V}_b^t, \quad (4)$$

where \hat{V}_b^t refers to the body vertices at time step t , represented in homogeneous coordinate.

3.3. Optimization with 3D Scene

3D Scene Representation. The 3D scene conveys useful information of human behavior, and therefore plays an important role in 3D human body recovery. As human-scene interaction is often grounded on the surfaces, we adopt a mesh representation for the 3D scene. Formally, we denote the 3D scene mesh as $M_s = (V_s, F_s)$, where $V_s \in R^{N_s \times 3}$ denotes the vertices of the scene representation, and F_s denotes the corresponding triangular faces. We use the dense environment reconstruction from COLMAP to represent M_s .

Human-Scene Contact. Note that the reconstructed 3D scene from the monocular video is up to a scale. To address such scale ambiguity, we design a novel energy function that not only encourages contact between human body and 3D scene, but also estimates the scale difference between 3D scene mesh M_s and 3D body mesh M_b . Specifically, we make use of the annotation from [18], where a candidate set of SMPL-X mesh vertices $V_c \in V_b$ to contact with the world were provided. We then multiply an optimizable scale parameter $S \in R$ to human body vertices V_s during optimization. Therefore, the energy function for enforcing human-scene contact is given by:

$$E_C(\beta, \theta, \gamma, V_s, S) = \sum_{i=1}^t \sum_{v_c \in V_c^t} \rho_c(\min_{v_s \in V_s} \|T_{wb}^t(Sv_c) - v_s\|), \quad (5)$$

where ρ_c is the robust Geman-McClure error function introduced in Eq. 3, and T_{wb} is human body to world transformation introduced in Eq. 4. Note that the scale factor S is shared across the video sequence. This is because we estimate a consistent 3D shape parameter θ from the entire sequence by taking the median of all the shape parameters obtained from the per-frame SMPL-X model fitting.

3.3.1 Human Dynamics Prior

Fitting SMPL-X human body model to each video frame will incur notable temporal inconsistency. Due to drastic camera motion, this problem is further amplified under egocentric scenarios. Here, we propose to use the empirical human dynamics priors to enforce temporal coherency on human body models in the world coordinate. Formally, we have the following energy function:

$$E_T(\beta, \theta, \gamma) = \sum_{i=2}^t \sum_J (1 - w_J) \rho_T((J_{wb}^{i+1} - J_{wb}^i) - (J_{wb}^i - J_{wb}^{i-1})), \quad (6)$$

where J_{wb}^i is the 3D human body joints position at time step i , transformed in world coordinate as in Eq. 4; ρ_T is another robust Geman-McClure error function that accounts for possible outliers; and w_J is confident score of 2D human keypoints estimation. As shown in Eq. 6, we design this energy function to focus on body parts that do not have reliable 2D observation, due to the unique egocentric view-point. Notably, we assume a zero acceleration motion prior. We show that this naive prior can effectively capture human motion in the outdoor environment.

3.3.2 Optimization

Putting everything together, we have the following energy function for our optimization method:

$$E_{total} = \sum_{i=1}^t E_M^i + \lambda_C E_C + \lambda_T E_T, \quad (7)$$

where E_M^i denotes the SMPL-X model fitting energy function for video frame x^i ; λ_C and λ_T represent the weights for human-scene contact term and human dynamic prior term, respectively. We optimize Eq. 7 using a gradient-based optimizer Adam [24] w.r.t. SMPL-X body parameters β, θ, γ , scale parameter S , and camera to world transformation T_{wc} . Note that the SfM already provides a initialization of T_{wc} , making T_{wc} optimizable can further smooth the global second-person human motion.

Note that E_M performs model fitting at each time step, while E_C and E_T optimize time series of body models. In addition, both E_C and E_T seek to optimize human body parameters in world coordinate, the scale ambiguity will cause the gradients of the contact term shift the body global position in wrong direction. Therefore, we carefully design a multi-stage optimization strategy. Specifically, we set λ_C and λ_T to be zero, so that the optimizer will only look at the 2D observation at stage one. We then set λ_C to be 0.1, keep λ_T as zero, and freeze the T_{wc} , so that the optimizer will focus on recovering the scale parameter S . At the final stage, we set λ_T to 0.1 and enable the gradients of T_{wc} to enforce temporal coherency. Our method is implemented in PyTorch and will be made publicly available.

4. Experiments

In this section, we discuss our experiments and results. To begin with, we introduce our dataset and evaluation metrics. We then present detailed ablation studies to validate our model design, and compare our approach with state-of-the-art on 3D body recovery from monocular videos. Finally, we provide a discussion of our method.

4.1. Dataset and Metrics

Datasets. To study the problem of second-person human body reconstruction, we present a new egocentric social in-

teraction dataset – EgoMoCap. This dataset consists of 36 video sequences from 4 participants. Each recording scenario incorporates two participants interacting in the wild. The camera wearer is equipped with head-mounted GoPro camera, and the other participant is asked to interact with the camera wearer in a natural manner. This dataset captures 4 types of outdoor human social interactions: *Greeting*, *Touring*, *Jogging Together*, and *Throw and Catch*.

Evaluation Metrics. For our experiments, we evaluate the human body reconstruction accuracy, motion smoothness, and the plausibility of human-scene interaction.

- **Human Body Reconstruction Accuracy:** We acknowledge that the 3D ground truth of human bodies can be obtained from RGB-D data [17], or Motion Capture Systems [53, 33]. However, all those systems adopt constrained capture environments and may result in unnatural social interactions. Our work focuses on outdoor social interaction, where the 3D human body ground truth is extremely difficult to capture. To evaluate the accuracy of human body reconstruction, we annotate our datasets with 2D human keypoints and evaluate the reconstruction quality using per-joint 2D projection error (PJE) on the image plane as in [60]. We report the PJE on both uniformly sampled frames (PJE-U), and frames where second-person body is partially observable (PJE-P). Note that we focus on evaluating human body poses, even though our method has the capability of reconstructing 3D hands and faces. This is because the primary goal of this work is to explore how environment factor affects 4D human body capture, while 3D scene context has minor influence on facial expression and hand pose for outdoor social interaction.

- **Motion Smoothness:** We adopt a physics-based metric [60] that uses average magnitude of joint accelerations to measure the smoothness of the estimated pose sequence. Thus, a lower value indicates that the times series of body meshes have more consistent human motion. Note that the motion smoothness is evaluated on 3D human joints projected in world coordinate. For fair comparison, we normalize the scale factor when reporting the results.

- **Plausibility of Human-Scene Interaction:** To evaluate whether our method leads to more realistic human-scene interaction, we transform the human body meshes into 3D world coordinate, render the results as video sequences, and further upload them to Amazon Mechanical Turk (AMT) for a user study. Specifically, we put the rendered results of all compared methods and our method side-by-side, and ask the AMT worker to choose the instance has the most realistic human-scene interaction.

4.2. Quantitative Results

We now introduce our quantitative experiment results. We first present detailed ablation studies, and then compare our method with state-of-the-art for 3D human body recon-

Method	PJE-U / PJE-P ↓	Smoothness ↓	User Study ↑
E_M	22.19 / 73.14	5.33	7.4
$E_M + E_C$	30.09 / 87.74	5.72	23.2
$E_M + E_T$	23.93 / 75.14	2.23	13.7
$E_M + E_C + E_T$ (Ours)	24.03 / 66.03	1.82	55.7

Table 1. Ablation study for our proposed method. We compare our method with baseline method that uses only 2D observation, and further analyze the role of human dynamic prior and human-scene interaction term. Our approach can not only improve motion smoothness and encourage realistic human-scene interaction, but also recover human body poses and shapes of partial observable second-person human body. (↑/↓ indicates higher/lower is better)

Method	PJE-U / PJE-P ↓	Smoothness ↓	User Study ↑
VIBE [25]	22.45 / 75.91	4.79	17.2
Ours	24.03 / 66.03	1.85	82.8

Table 2. Experiment results comparison with previous state-of-the-art method on human motion capture from monocular videos. (↑/↓ indicates higher/lower is better)

struction from monocular videos.

Ablation Study. Here we analyze the functionality of the terms in Eq. 7. The results are summarized in Table 1. E_M refers to the baseline method that performs per-frame fitting with 2D observation as in SMPLify-X [39]. E_M achieves 22.19 in PJE, yet has undesirable performance on motion smoothness and human-scene interaction user study. In the second row ($E_M + E_C$), we report the method that makes use of both human scene contact term and 2D observations. Though adding the contact term alone leads to more realistic human-scene interaction, it compromises the performance on 2D projection error and motion smoothness by a notable margin. $E_M + E_T$ in the third row refers to the method that optimizes the 2D observations together with the human dynamic prior term E_T . Not surprisingly, E_T can significantly improve the motion smoothness. In the last row, we present the results of our full optimization approach. Our method achieves the best performance on motion smoothness and plausibility of human-scene interaction. An interesting observation is that ours outperforms $E_M + E_T$ by a notable margin on motion smoothness. We speculate that this is because the physical human scene constraints narrows do the solution space of model fitting, and thereby leads to more optimal performance on temporal coherency. We note that our model performs slightly more worse on PJE-U. This is because PJE is a 2D metric, and therefore favors the method that adopts only 2D projection error as objective function during optimization. However, when the 2D observation can not be robustly estimated due to partial observation, our method outperforms other baselines by a significant margin (66.03 vs. 73.14 in PJE-P). Those results support our claim that our method can address the challenge of partially observable human body, and estimate plausible global human motion grounded on the 3D scene.

Comparison to SOTA Method. In Table 2, we compare our approach with SOTA method of 3D body recovery from monocular videos – VIBE [25]. Since VIBE does

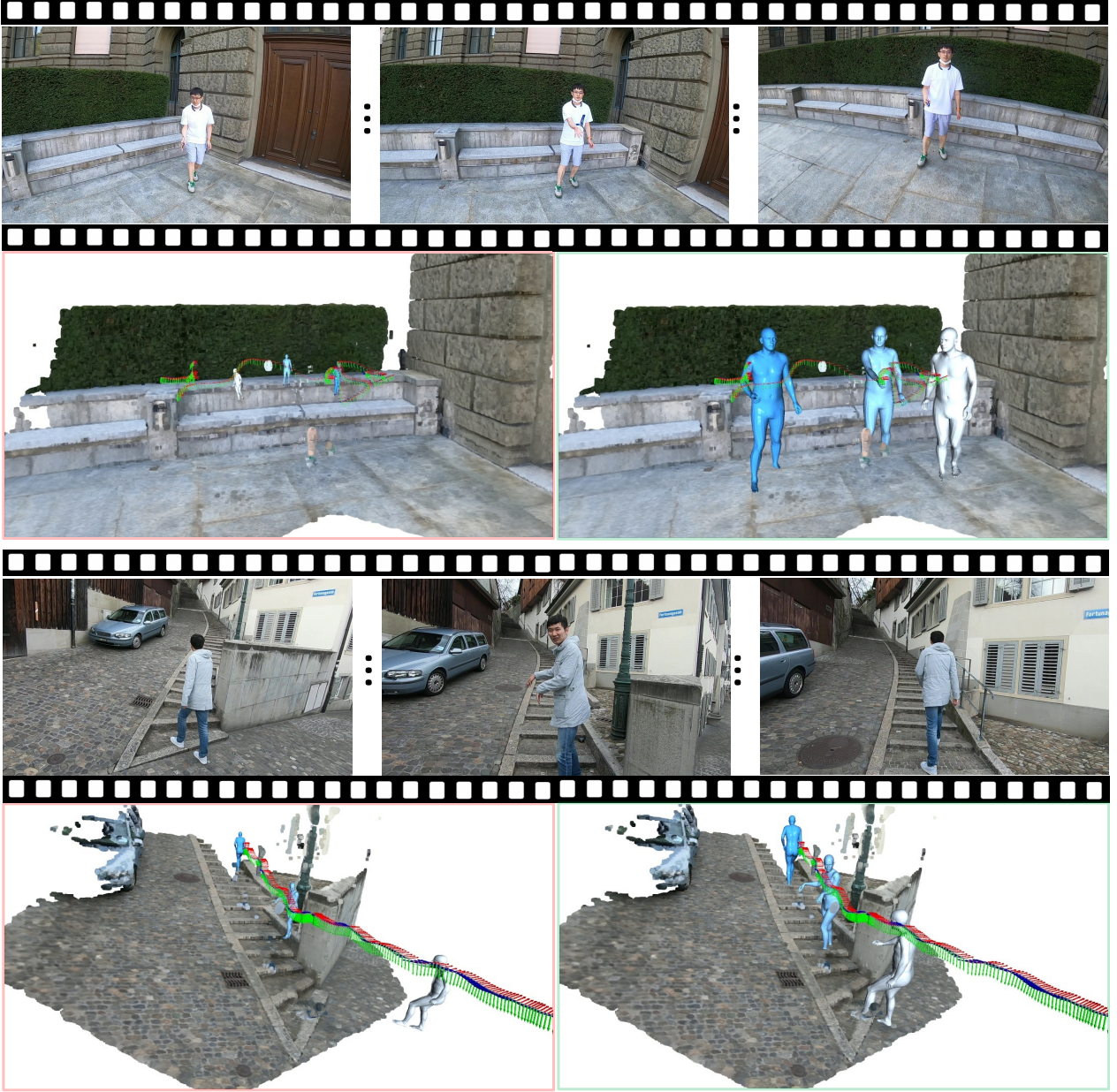


Figure 3. Visualization of time series of human bodies in the world coordinate. We visualize both results of SMPLif-X baseline (Left) and our method (Right) projected into 3D scene reconstruction. Our method recovers the scale ambiguity between 3D scene reconstruction and 3D body reconstruction from monocular video, and therefore leads to more plausible human-scene interaction.

not model the human-scene constraints, simply projecting human body meshes into 3D scene results in unrealistic human-scene interaction. Moreover, the egocentric camera motion causes VIBE failing to capture temporal coherent human bodies. In contrast, our method outperforms VIBE on motion smoothness and human-scene interaction plausibility by a large margin. Though VIBE performs slightly better on PJE-U (22.45 vs. 24.03), it lags far behind of our method on PJE-P (75.91 vs. 66.03). We have to re-emphasize that the 2D projection error can not reflect the true performance improvement of our method. This is be-

cause the 2D keypoints annotation is only available for visible human body parts, and therefore 2D per-joint projection error does not penalize the method that fits wrong 3D body model to partially 2D observation. Take the VIBE result shown in the third row of Fig. 4 for an instance, the 2D projection error may have decent performance, even though the reconstructed 3D human body is completely wrong.

4.3. Qualitative Results

We now present the qualitative results of our method. As shown in Fig. 3, we visualized the results of both E_M

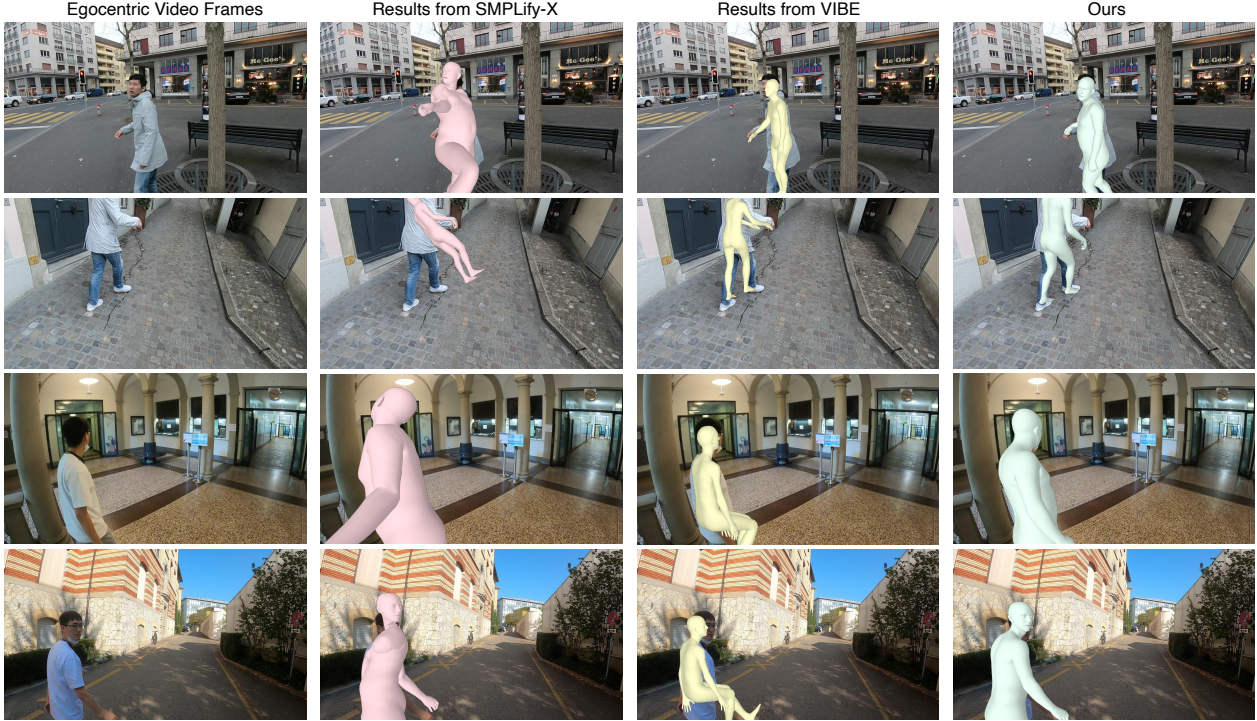


Figure 4. Qualitative comparison between our method and other approaches. The first column is the original video frames; the second column is the results from SMPLify-X, the third column is the results from VIBE, and the last shows our results. Our approach can address the challenging cases when second-person body is partially observable.

baseline and our method in the world coordinate. By examining the SMPLify-X baseline results, we can observe an obvious mismatched scale between the 3D reconstruction of human body and environment, which results in unrealistic human-scene interaction. In contrast, our method produces more plausible human body motion grounded on 3D scene by resolving the scale ambiguity of 3D reconstruction from monocular videos. In Fig. 4, we visualize our results on 2D image plane. Specifically, we choose instances where the second-person human body is partially observable. Notably, both SMPLify-X and VIBE fail substantially for those challenging cases. Our method, on the other hand, makes use of the 2D cues from entire video sequences and 3D scene for reconstructing temporal coherent time series of body poses, and therefore can successfully reconstruct the human body even when it is partially observable. In the supplementary materials, we provide additional video demos to demonstrate the benefits of our approach.

4.4. Remarks and Discussion

The previous sections have demonstrated, via detailed experimental evaluation and comparisons, that our method can capture more accurate second-person human bodies, and produce more realistic human-scene interaction, compared to prior works. However, our method also has certain limitations. A key issue is the need to retrieve the camera trajectory and 3D scene only from monocular RGB videos

via Structure from Motion (SfM). Therefore, our method has the same bottleneck as SfM: Challenging factors such as dynamic scenes, featureless surfaces, changing illumination, etc., may cause visual feature matching to fail. We note that the camera and environment information can be more robustly estimated using additional sensors (Lidar, Depth Camera, Matterport etc.). Incorporating those sensors into the egocentric capture setting is a very interesting and promising future direction. In addition, our naive human motion prior (zero acceleration), may result in unrealistic motions in some cases. More effort in learning motion priors could potentially address this issue. We believe our efforts constitute an important step forward for a largely unexplored egocentric vision task, and we hope our work can inspire the community to make further investments.

5. Conclusion

In this work, we introduce a novel task of reconstructing a time series of second-person 3D human body meshes that are grounded on the 3D scene information from monocular egocentric videos. We propose a novel optimization-based method to address the challenges of egocentric capture, that exploits the 2D observation of entire video sequence and 3D scene information for second-person human body capture. In addition, we introduce a new egocentric video dataset – EgoMocap, and provide extensive quantitative and qualita-

tive analysis to demonstrate that our method can effectively reconstruct partially-observable second-person human bodies and produce more realistic human-scene interaction.

References

- [1] Cmu mocap dataset. <http://mocap.cs.cmu.edu>. 4
- [2] Thiemo Alldieck, Marc Kassubeck, Bastian Wandt, Bodo Rosenhahn, and Marcus Magnor. Optical flow-based 3d human motion estimation from monocular video. In *GCPR*, 2017. 2
- [3] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*, pages 408–416. 2005. 2
- [4] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J. Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *ECCV*, 2016. 2, 4
- [5] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 4
- [6] Chao-Yeh Chen and Kristen Grauman. Subjects and their objects: Localizing interactees for a person-centric view of importance. *IJCV*, 126(2-4):292–313, 2018. 3
- [7] Eunji Chong, Nataniel Ruiz, Yongxin Wang, Yun Zhang, Agata Rozga, and James M Rehg. Connecting gaze, scene, and attention: Generalized attention estimation via joint modeling of gaze and scene saliency. In *ECCV*, pages 383–398, 2018. 2, 3
- [8] Eunji Chong, Yongxin Wang, Nataniel Ruiz, and James M. Rehg. Detecting attended visual targets in video. In *CVPR*, 2020. 2
- [9] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *ECCV*, 2018. 1
- [10] Vincent Delaitre, David F Fouhey, Ivan Laptev, Josef Sivic, Abhinav Gupta, and Alexei A Efros. Scene semantics from long-term observation of people. In *ECCV*, 2012. 3
- [11] Alireza Fathi, Ali Farhadi, and James M Rehg. Understanding egocentric activities. In *ICCV*, 2011. 1
- [12] Alireza Fathi, Jessica K Hodgins, and James M Rehg. Social interactions: A first-person perspective. In *CVPR. IEEE*, 2012. 2, 3
- [13] Antonino Furnari and Giovanni Maria Farinella. What would you expect? anticipating egocentric actions with rolling-unrolling LSTMs and modality attention. In *ICCV*, 2019. 1
- [14] Stuart Geman. Statistical methods for tomographic image reconstruction. *Bull. Int. Stat. Inst.*, 4:5–21, 1987. 4
- [15] Helmut Grabner, Juergen Gall, and Luc Van Gool. What makes a chair a chair? In *CVPR*, 2011. 3
- [16] Abhinav Gupta, Scott Satkin, Alexei A Efros, and Martial Hebert. From 3d scene geometry to human workspace. In *CVPR*, 2011. 3
- [17] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J. Black. Resolving 3D human pose ambiguities with 3D scene constraints. In *ICCV*, 2019. 3, 6
- [18] Mohamed Hassan, Vasileios Choutas, Dimitrios Tzionas, and Michael J Black. Resolving 3d human pose ambiguities with 3d scene constraints. In *ICCV*, 2019. 5
- [19] Yinghao Huang, Federica Bogo, Christoph Lassner, Angjoo Kanazawa, Peter V. Gehler, Javier Romero, Ijaz Akhter, and Michael J. Black. Towards accurate marker-less human shape and pose estimation over time. In *3DV*, 2017. 2, 3
- [20] Hao Jiang and Kristen Grauman. Seeing invisible poses: Estimating 3d body pose from egocentric video. In *CVPR*, 2017. 2
- [21] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 2
- [22] Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. 2
- [23] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. 3
- [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [25] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 2, 3, 6
- [26] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 2
- [27] Hema S Koppula and Ashutosh Saxena. Anticipating human activities using object affordances for reactive robotic response. *TPAMI*, 38(1):14–29, 2015. 3
- [28] Yin Li, Miao Liu, and James M. Rehg. In the eye of beholder: Joint learning of gaze and actions in first person video. In *ECCV*, 2018. 1
- [29] Yin Li, Miao Liu, and James M Rehg. In the eye of the beholder: Gaze and actions in first person video. *arXiv preprint arXiv:2006.00626*, 2020. 1
- [30] Yin Li, Zhefan Ye, and James M Rehg. Delving into egocentric actions. In *CVPR*, 2015. 1
- [31] Zongmian Li, Jiri Sedlar, Justin Carpentier, Ivan Laptev, Nicolas Mansard, and Josef Sivic. Estimating 3d motion and forces of person-object interactions from monocular video. In *CVPR*, 2019. 3
- [32] Miao Liu, Siyu Tang, Yin Li, and James Rehg. Forecasting human object interaction: Joint prediction of motor attention and actions in first person video. In *ECCV*, 2020. 1
- [33] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *CVPR*, 2019. 6
- [34] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3d human pose estimation. In *ICCV*, 2017. 2
- [35] Tushar Nagarajan, Christoph Feichtenhofer, and Kristen Grauman. Grounded human-object interaction hotspots from video. In *ICCV*, 2019. 3

- [36] Tushar Nagarajan, Yanghao Li, Christoph Feichtenhofer, and Kristen Grauman. Ego-topo: Environment affordances from egocentric video. In *CVPR*, 2020. 1, 3
- [37] Sanath Narayan, Mohan S Kankanhalli, and Kalpathi R Ramakrishnan. Action and interaction recognition in first-person videos. In *CVPRW*, 2014. 3
- [38] Evonne Ng, Donglai Xiang, Hanbyul Joo, and Kristen Grauman. You2me: Inferring body pose in egocentric video via first and second person interactions. In *CVPR*, 2020. 1, 2, 3
- [39] 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 *CVPR*, 2019. 2, 4, 6
- [40] Dario Pavlo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *CVPR*, 2019. 3
- [41] Yair Poleg, Ariel Ephrat, Shmuel Peleg, and Chetan Arora. Compact CNN for indexing egocentric videos. In *WACV*, 2016. 1
- [42] Javier Romero, Dimitrios Tzionas, and Michael J Black. Embodied hands: Modeling and capturing hands and bodies together. *ACM Transactions on Graphics (ToG)*, 36(6):245, 2017. 2
- [43] Michael S Ryoo and Larry Matthies. First-person activity recognition: What are they doing to me? In *CVPR*, 2013. 3
- [44] Manolis Savva, Angel X. Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. PiGraphs: Learning Interaction Snapshots from Observations. *ACM Transactions on Graphics (TOG)*, 35(4), 2016. 3
- [45] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 4
- [46] Soshi Shimada, Vladislav Golyanik, Weipeng Xu, and Christian Theobalt. Physcap: Physically plausible monocular 3d motion capture in real time. *ACM Transactions on Graphics*, 39(6), dec 2020. 3
- [47] Hyun Soo Park, Jyh-Jing Hwang, Yedong Niu, and Jianbo Shi. Egocentric future localization. In *CVPR*, 2016. 3
- [48] Hyun Soo Park and Jianbo Shi. Social saliency prediction. In *CVPR*, 2015. 2, 3
- [49] Xiao Sun, Jiaxiang Shang, Shuang Liang, and Yichen Wei. Compositional human pose regression. In *ICCV*, 2017. 2
- [50] Denis Tome, Patrick Peluse, Lourdes Agapito, and Hernan Badino. xr-egopose: Egocentric 3d human pose from an hmd camera. In *ICCV*, 2019. 2
- [51] Hsiao-Yu Tung, Hsiao-Wei Tung, Ersin Yumer, and Katerina Fragkiadaki. Self-supervised learning of motion capture. In *NeurIPS*, 2017. 3
- [52] Timo von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 2, 3
- [53] Timo von Marcard, Bodo Rosenhahn, Michael Black, and Gerard Pons-Moll. Sparse inertial poser: Automatic 3d human pose estimation from sparse imus. *Computer Graphics Forum 36(2), Proceedings of the 38th Annual Conference of the European Association for Computer Graphics (Eurographics)*, pages 349–360, 2017. 6
- [54] Xiaolong Wang, Rohit Girdhar, and Abhinav Gupta. Binge watching: Scaling affordance learning from sitcoms. In *CVPR*, 2017. 3
- [55] Yangang Wang, Yebin Liu, Xin Tong, Qionghai Dai, and Ping Tan. Outdoor markerless motion capture with sparse handheld video cameras. *IEEE transactions on visualization and computer graphics*, 24(5):1856–1866, 2017. 3
- [56] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *CVPR*, 2016. 4
- [57] Takuma Yagi, Karttikeya Mangalam, Ryo Yonetani, and Yoichi Sato. Future person localization in first-person videos. In *CVPR*, 2018. 2, 3
- [58] Zhefan Ye, Yin Li, Yun Liu, Chanel Bridges, Agata Rozga, and James M Rehg. Detecting bids for eye contact using a wearable camera. In *2015 11th IEEE International Conference and Workshops on Automatic Face and Gesture Recognition (FG)*, volume 1, pages 1–8. IEEE, 2015. 2, 3
- [59] Ryo Yonetani, Kris M Kitani, and Yoichi Sato. Recognizing micro-actions and reactions from paired egocentric videos. In *CVPR*, 2016. 1, 2, 3
- [60] Ye Yuan and Kris Kitani. 3d ego-pose estimation via imitation learning. In *ECCV*, 2018. 2, 6
- [61] Jason Y. Zhang, Sam Pepose, Hanbyul Joo, Deva Ramanan, Jitendra Malik, and Angjoo Kanazawa. Perceiving 3d human-object spatial arrangements from a single image in the wild. In *ECCV*, 2020. 3
- [62] Siwei Zhang, Yan Zhang, Qianli Ma, Michael J Black, and Siyu Tang. Generating person-scene interactions in 3d scenes. 2020. 3
- [63] Yan Zhang, Mohamed Hassan, Heiko Neumann, Michael J Black, and Siyu Tang. Generating 3d people in scenes without people. In *CVPR*, 2020. 3