

Interaction with Virtual Objects using Human Pose and Shape Estimation

Hong Son Nguyen
Korea University
nguyenhongson303@gmail.com

DaEun Cheong
Korea University
wjdekdm001@korea.ac.kr

Andrew Chalmers
Victoria University of Wellington
andrew.chalmers@vuw.ac.nz

Myoung Gon Kim
Korea University
m_gon_kim@korea.ac.kr

Taehyun Rhee
The University of Melbourne
taehyun.rhee@unimelb.edu.au

JungHyun Han
Korea University
jhan@korea.ac.kr

Abstract

In this paper, we propose an AR system that facilitates a user's natural interaction with virtual objects in an augmented reality environment. The system consists of three modules: human pose and shape estimation, camera-space calibration, and physics simulation. The first module estimates a user's 3D pose and shape from a single RGB video stream, thereby reducing the system setup cost and broadening potential applications. The camera-space calibration module estimates the user's camera-space position to align the user with the input RGB image. The physics simulation enables seamless and physically natural interaction with virtual objects. Two prototyping applications built upon the system prove an enhancement in the quality of interaction, fostering a more immersive and intuitive user experience.

Keywords: Augmented reality, Full-body interaction, Pose and shape estimation.

1. Introduction

For various augmented reality (AR) applications, such as virtual try-on and virtual training, it is essential to provide natural seamless interactions between real humans and virtual objects. Enabling such interactions usually necessitates capturing the humans' motions at real time. In an AR volleyball game, for example, actions such as spiking or blocking a virtual ball must be convincingly mirrored in the AR context.

For capturing the full-body joints of a human, various motion sensors and motion capture suits have been used [31, 8, 14, 76]. However, the intrusive nature of marker-based systems and the encumbrance of motion capture suits pose significant barriers to widespread adoption in everyday AR use. An alternative approach involves using RGB-D sensors to compute 3D poses and approximate 3D shapes based on color and depth information [53, 16, 62, 32]. Recently, deep learning-based approaches further the field by estimating 3D pose information using only RGB sensors [39, 40, 61], leading to

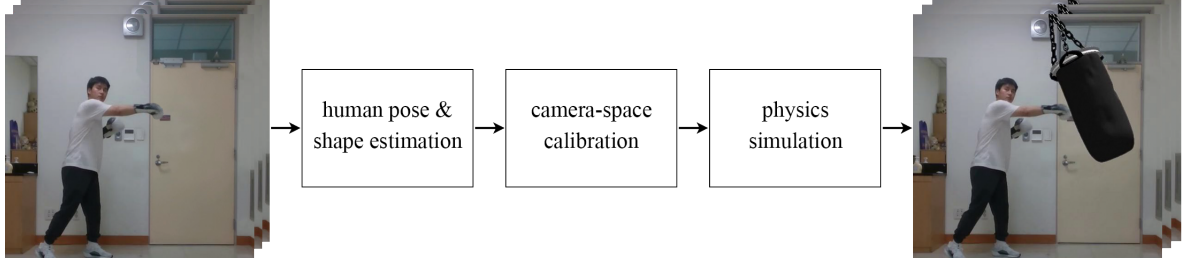


Figure 1: Our AR system is composed of three modules that facilitate interactions between a real human and virtual objects. In this AR boxing game, the user is interacting with the virtual punching bag.

the development of a variety of virtual interaction scenarios [73, 47, 26]. The deep learning-based approaches have been extended to provide 3D shape information, which is crucial for handling collision and occlusion in AR environments [5, 25, 49, 29].

In this paper, we propose an AR system that takes an RGB video as input, estimates a user’s 3D pose and shape, and facilitates the user’s interaction with 3D virtual objects. As shown in Figure 1, the system comprises three modules for (1) human pose and shape estimation, (2) camera-space calibration, and (3) physics simulation.

The human pose and shape (HPS) estimated by the first module are represented in the parametric human model, SMPL-X [49]. As the SMPL-X model is a geometric entity, handling its interactions with virtual objects may appear to be straightforward. Unfortunately, the estimated model is *not* defined in the world or camera spaces. To enable the human model to move around and interact with dynamic virtual objects in the AR environment, we need to *calibrate* the human model every frame so that it is correctly transformed into the world or camera spaces. This paper’s focus is on *camera-space calibration*. Finally, the third module for physics simulation is implemented using a game engine. Our contributions are summarized as follows:

- We present a working AR solution that enables interactions of a real human with 3D virtual objects. Taking only a single RGB video feed as input, it supports 3D shape-based natural interactions.
- Built upon deep learning techniques that are well balanced between accuracy and efficiency, the solution reconstructs and cali-

brates the human model in a robust manner.

- The solution works in an easy-to-setup AR environment, and therefore it can be used across a variety of AR applications.

2. Related Work

2.1. Human Pose and Shape Estimation

A significant challenge in full-body interaction lies in the accurate tracking of human motion. Traditional methods in VR, AR, and graphics often rely on motion capture suits and sensors [59, 9], while others explore marker-based setups [31, 8]. Beyond full-body suits, alternative systems use sparse markers or wearable sensors combined with inverse kinematics or SLAM. Roth *et al.* [57] demonstrated that a reduced rigid-body marker set with inverse kinematics lowers latency and task load in VR without sacrificing body ownership. Yi *et al.* [76] proposed EgoLocate, which integrates inertial mocap and monocular SLAM to achieve real-time human motion capture and localization using six IMUs and a monocular camera. Ghorbani *et al.* [21] deployed deep learning models to estimate human motion and shape from the raw input of marker-based motion capture systems. Chen *et al.* [14] employed a Bi-LSTM model to reconstruct human motion from only four markers. Kim *et al.* [28] proposed DAMO, a deep solver that generalizes across marker configurations via end-to-end optimization. Despite these advances, marker- and sensor-based systems still suffer from restricted mobility [11] and high setup costs due to specialized hardware and calibration.

Facing the limitations of marker and sen-

sensor based methods, researchers have leveraged RGB-D sensors, such as Kinect [44], for full-body motion estimation [62, 16, 27, 53]. Shum *et al.* [62] devised a framework for real-time motion tracking using Kinect. Cui *et al.* [16] employed two Kinect sensors, positioned at the front and back, to comprehensively capture human motion. Kim *et al.* [27] combined high-speed RGB and ToF sensors to capture dynamic human poses. Ren *et al.* [53] utilized a convolutional neural network (CNN) along with a depth sensor to capture human motions in large-scale scenarios.

RGB-based deep learning solutions [39, 40, 61] have emerged as effective alternatives for 3D pose estimation without specialized hardware. Existing methods for 3D pose estimation using monocular RGB camera fall into two categories: single-stage vs. two-stage. The single-stage techniques [63, 43] directly localize 3D body keypoints from input images. Sun *et al.* [63] used the soft-argmax operation to obtain joint locations from 3D heatmaps, whereas Moon *et al.* [43] proposed a top-down approach for multi-person pose estimation. Additionally, Apple introduced ARKit [1], a framework that estimates 3D human pose from RGB images for AR applications. In contrast, the two-stage methods [13, 36, 50, 80] leverage accurate 2D pose estimation, lifting keypoints to 3D space. Pavllo *et al.* [50] introduced dilated temporal convolutions, and Zheng *et al.* [80] used a Vision Transformer to capture spatial and temporal dependencies.

The Skinned Multi-Person Linear (SMPL) model [33] has catalyzed a series of studies that aim at estimating 3D HPS using deep learning models. By representing 3D human body as differentiable functions of the HPS parameters, the model is well-suited for end-to-end training of neural networks. Notable contributions in this domain include SMPLify [5], Human Mesh Recovery (HMR) [25], SMPLify-X [49] and VIBE [29]. HPS estimation methods can be performed through optimization or regression techniques. Optimization-based methods [5, 74, 19, 75] align a body model with image cues such as joints, dense vertex correspondences or 2D segmentation masks. In the context of optimization, for example, Luvizon *et al.* [34] extended 3D HPS estimation to

multi-human, scene-aware settings from a single RGB video. On the other hand, regression-based methods [25, 29, 64, 65, 52, 66] train a network using a loss akin to the optimization objective function, enabling it to predict body model parameters.

2.2. Full-body Interaction

In the past, in order to study full-body interaction in virtual environments, marker-based motion capture systems were commonly used [7, 51, 18, 77]. Debarba *et al.* [18] enhanced interaction with virtual objects by utilizing a setup involving 4 LED markers and 14 cameras to capture human motions. Young *et al.* [77] implemented high five in virtual environments by employing a motion capture system.

The advent of RGB-D sensors has spurred extensive research into full-body interaction applications. As emphasized by Caserman *et al.* [10], VR applications employing full-body motion data from RGB-D cameras span a broad spectrum of activities. These include virtual climbing of ladders [67], cycling-based exergames [60, 6, 69], collaborative manufacturing tasks involving human-robot interaction [37], and the creation of role-playing games [15].

Beyond VR applications, commercial products like Kinect and RealSense have enabled full-body interaction in AR. Despite the availability of these commercial solutions, there has been relatively little scholarly exploration into full-body interaction in AR [20]. Moreover, although these systems provide more flexibility compared to marker-based systems, they often necessitate the use of multiple RGB-D sensors, leading to intricate setups and higher costs.

Recent advancements in deep learning have enabled a number of approaches to tackle the challenge of human-virtual interaction. Hwang *et al.* [23] proposed a lightweight 3D human pose estimation model applied to virtual avatar reconstruction, with potential extensions to human-virtual interaction. Nguyen *et al.* [47] developed an AR fitness application that supports full-body interaction by integrating 3D pose estimation into an exercise interface. Wu *et al.* [73] created an AR-based martial arts training system enabling real-time full-body in-

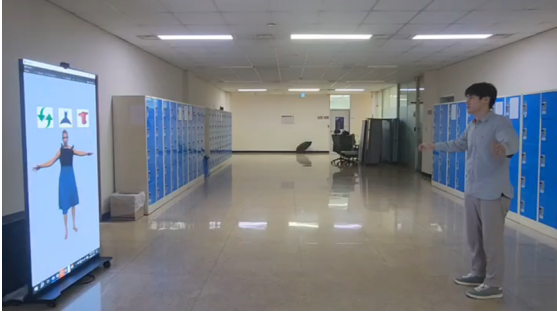


Figure 2: In this AR environment, the user can see his full body and the virtual objects in the large screen.

teraction through pose-guided motion evaluation. Recent studies [55, 56, 12] have expanded avatar-based interactions into immersive environments.

Alongside skeleton-based methods, human body shape has garnered significant attention for full-body interaction, particularly the SMPL [33] based approaches. De *et al.* [17] used HMR to estimate HPS and created a virtual agent for interactive VR. Additionally, Šarić *et al.* [58] introduced an extended reality-based telemedicine collaboration system that presents a 3D avatar of the patient to a remote clinician, leveraging 3D HPS estimation derived from a deep learning model. Recent studies have expanded full-body interaction into immersive environments [54] and AR environments [45, 46].

Inspired by the learning-based methods, this paper proposes an HPS estimation network, which utilizes the shape information to interact with virtual objects in AR environments. Our system is easy to set up, requiring no additional devices, and is demonstrated to be cost-effective.

3. Camera-space Calibration

Our AR setup is shown in Figure 2: A user moves on the floor and the large screen in front of the user displays the mirrored view of the environment using an RGB webcam, which is mounted on top of the screen and captures the *full body* of the user. The captured environment is augmented with virtual objects so that the user can interact with them.

Among the three modules presented in Figure 1, this paper focuses on the second module,

camera-space calibration, while briefly sketching the first (HPS estimation) and third (physics simulation) modules, for which there exist many off-the-shelf methods and commercial/open-source programs.

In general, existing methods for estimating a human’s motion compute the HPS parameters for SMPL-X [49], which are denoted as θ and β , respectively. Using θ and β , the human’s 3D triangular mesh is reconstructed, which we denote as \mathcal{V} . Then, we derive the joints’ 3D positions from \mathcal{V} . What we call the *human model* consists of the 3D mesh \mathcal{V} and the 3D joints denoted as \mathcal{J} . In order to estimate the human model, we use the Human Mesh Recovery (HMR) framework [25]. For more details, readers are referred to Appendix A.

In our AR setup, the human model is invisible to the user, but its mesh is used, for example, to detect *collisions* with the virtual objects. In the current implementation, the physics simulation module is implemented using Unity engine. However, it can be seamlessly replaced by other commercial or open-source engines.

3.1. Calibration Network

HMR employs the *weak-perspective camera model*, which is implemented as an *orthographic projection* plus a *scaling*. Consequently, the reconstructed human model is not correctly located in the 3D coordinate system of the real camera, e.g., its camera-space depth may be over- or under-estimated. Then, not only would we miss the *collision* between the human model and virtual objects, but we would also handle the *occlusion* between them incorrectly. Therefore, the human model must be *calibrated* so that it is transformed into the *full-perspective camera space*.

Figure 3 illustrates the calibration process. The essential step for calibration is to estimate the camera-space coordinates of the *root joint*. Once they are estimated correctly, all the other joints’ coordinates can be determined immediately using their positions relative to the root, which are stored in the 3D pose \mathcal{J} .

For this, we employed a deep learning network proposed by Pavllo *et al.* [50], which processes “a sequence of 2D poses” to output the camera-space positions of the root. It was mod-

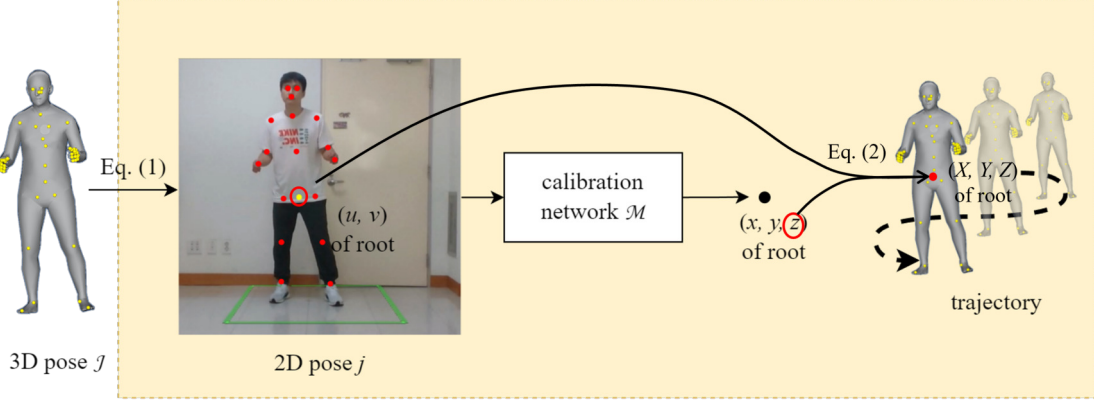


Figure 3: For calibration, we use the 3D pose \mathcal{J} computed by the HPS estimation module. The sequence of the calibrated root nodes makes up the trajectory.

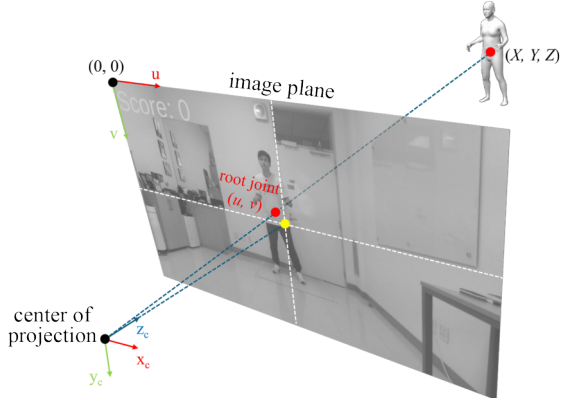


Figure 4: The root’s coordinates, (u, v) , in the image plane are transformed by \mathbf{K}^{-1} , and then are scaled by z , which is estimated by \mathcal{M} , to define the camera-space coordinates, (X, Y, Z) .

ified and retrained to take “a single frame of 2D pose” as input and estimate the root’s camera-space position. Let \mathcal{M} denote the network and j denote the 2D pose input to \mathcal{M} . By projecting the 3D pose \mathcal{J} onto the image plane, we can obtain j . It is simple to implement because HMR returns not only θ and β but also the parameters of the weak-perspective camera model, i.e., the scale factor $s \in \mathbb{R}$ and the 2D translation vector $t \in \mathbb{R}^2$. The 2D pose j is obtained using the 3D pose \mathcal{J} as follows:

$$j = s\Pi(\mathcal{J}) + t \quad (1)$$

where $\Pi(\cdot)$ denotes orthographic projection. $\mathcal{M}(j)$ is the root’s coordinates, (x, y, z) , in the camera space.

The calibration network, \mathcal{M} , is trained using the Human3.6M dataset [24]. Note that the Human3.6M dataset contains only “square” images

with aspect ratio one. In contrast, our system presented in Figure 1 is designed to take general “rectangular” images. Consequently, the x - and y -coordinates estimated by \mathcal{M} are not reliable.

Therefore, we take only the z -coordinate, which is reliable. Let (X, Y, Z) denote the *calibrated* coordinates of the root. Whereas $Z = z$, X and Y are computed using the root’s pixel coordinates, (u, v) , in the image plane. See Figure 4; (u, v) are transformed by the inverse of the camera intrinsic matrix \mathbf{K} and then are scaled by z :

$$\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = z\mathbf{K}^{-1} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} \quad (2)$$

For more details, readers are referred to Appendix B.

3.2. Discussion

In general, the HPS methods based on the weak-perspective model are inaccurate at predicting 3D pose, due to its near-orthographic projection. For example, high levels of perspective distortion in selfies often lead to unacceptable 3D pose. In contrast, the predicted pose would be acceptable if the user is sufficiently far away from the camera. It is the case in our AR setup shown in Figure 2, where a fairly long distance between the camera and the user should be maintained so as to capture the user’s full-body motion.

Recently, the problem of the weak-perspective model has been directly tackled by

a few efforts that estimate the camera’s intrinsic parameters [48, 30] or the user’s depth [71]. While these methods improve the 3D accuracy, they introduce additional processing and training complexities. Our strategy of adopting HMR and integrating it with the calibration network can be considered as well-balanced between accuracy and efficiency. (We train HMR with the BEDLAM dataset [4], which provides diverse and realistic synthetic data.)

The root joint’s camera-space position keeps changing over the frames, and it is generally called *trajectory*. Existing methods for calculating the trajectory involve complex processes, ranging from optimization-based frameworks [79, 39, 38] to estimating relative scene depths by exploiting scene constraints [72, 78]. Whereas these methods do not estimate the 3D shape, they estimate not only the trajectory but also the 3D pose, which is already computed by our “human pose and shape estimation” module. Instead of adopting such methods, which consume significant amount of computing resources to produce redundant information, i.e., the 3D pose \mathcal{J} , we compute the 2D pose j from \mathcal{J} “with little cost” and provide it for the modification of the network proposed by Pavllo *et al.* [50] to estimate the trajectory.

The original network proposed by Pavllo *et al.* [50] is based on dilated temporal convolutions, predicting the trajectory using both past and future data. On the other hand, they proposed *causal convolutions* for real-time applications, which have access only to past frames and predict the trajectory of the current frame. Readers are referred to Appendix C for visualization of the two types of convolutions.

To evaluate accuracy across the input sequences of different lengths, we made an experiment, where the causal convolutions were trained respectively with sequence lengths of 1, 27 and 81 frames. The case with a single-frame sequence corresponds to our calibration network. Training was made with Human3.6M

dataset [24] in a *semi-supervised manner*, where Subjects 1, 5 and 6 were labeled, i.e., both (u, v) and (x, y, z) coordinates were utilized, whereas Subjects 7 and 8 were treated as unlabeled, i.e., only (u, v) coordinates were used.

We evaluated the three trajectory estimation models with Subjects 9 and 11 of the Human3.6M dataset. Obviously, our calibration network requires fewer parameters and FLOPs than the others, as shown in Table 1, making it more suitable for real-time applications. It is interesting to find that with respect to the *trajectory error*, which is defined as the average difference of the ground-truth and predicted positions of the root joint, our network is the most accurate, and the model with 27 frames is more accurate than that with 81 frames, i.e., the longer input sequence leads to the larger error. With these findings, we speculate that in the causal convolutions, the current frame holds the most critical information for predicting the root joint position, and less relevant information of the past frames is accumulated over the frame sequence to hamper precise prediction.

4. Applications

The equipment of our AR setup shown in Figure 2 consists of a webcam, a PC and a large screen that can be replaced by a TV. These are easily found in everyday environments, providing an easy-to-setup gaming platform. The setup can also be used for many other AR applications such as virtual try-on.

4.1. Application 1: Boxing Game

Human shape estimation opens the door to a wide range of AR applications. We have developed a boxing game prototype, where the user punches a virtual object. In Figure 5, the user moves around a punching bag and keeps punching it. To convincingly composite the virtual object into the real-world environment, we use the estimated shape of the user. Figure 5 clearly shows that our system is able to handle *occlusion* between the real and virtual objects. For enhancing the realism, the shadow of the virtual punching bag is generated and composited into the scene using a virtual floor. The virtual floor

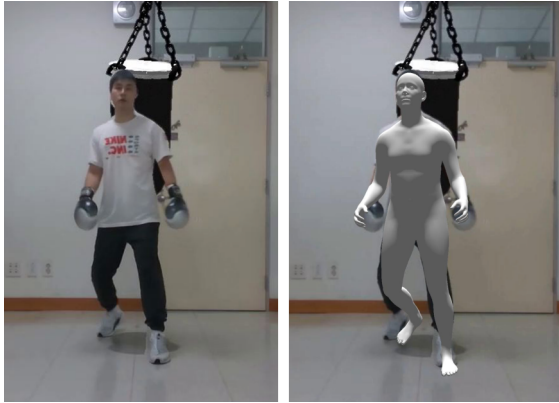
Table 1: Calibration model evaluation.

# input frames	parameters	FLOPs	trajectory errors
1 frame	4.24M	8.46M	147.0mm
27 frames	8.51M	16.99M	147.7mm
81 frames	12.70M	25.38M	181.3mm

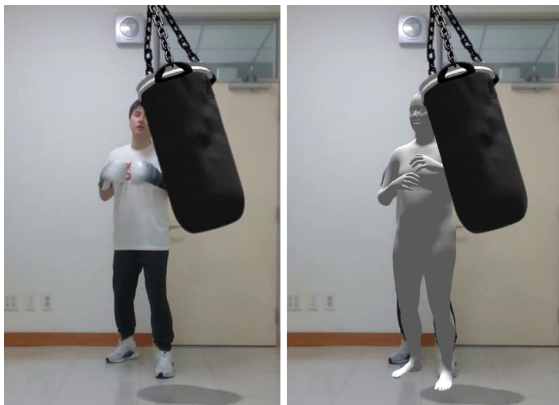
can also be used, for example, to make a virtual ball bouncing in the real scene.

4.2. Application 2: Virtual Try-on

Our AR setup presented in Figure 2 can be used for many other AR applications such as virtual try-on. Figure 6 shows a sequence of snapshots that visualize complex real-time interactions between a virtual garment and the human model reconstructed by the method presented in this paper. In this proof-of-concept implementation, the garment is composed of 8K vertices. It is simulated with XPBD (extended position-based dynamics) [35] implemented in Taichi [22], a high-level language for GPU programming, and the proximity query for colliding triangles is accelerated via spatial hashing [68].



(a)



(b)

Figure 5: Occlusion: (a) The real human occludes the virtual object. (b) The virtual Object occludes the real human.

Table 2: Runtime performances of HPS estimation and calibration.

Model	Time (ms)
Yolo v8n	3.04
HMR	2.22
Calibration Network	0.24
total	5.5

Table 3: Runtime performances of applications.

Application	Time (ms)	FPS
Boxing Game	6.6	151
Virtual Try-on	36.2	27

4.3. Evaluation and Discussion

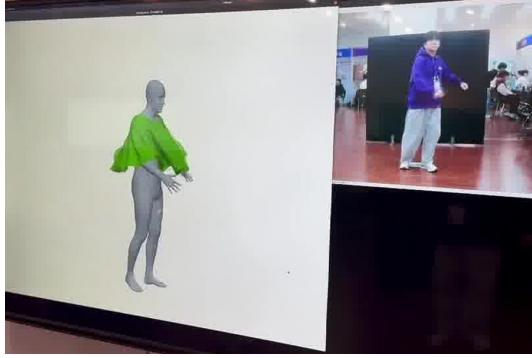
Table 2 shows the times consumed by the first module for HPS estimation (i.e., by Yolo v8n + HMR) and the second one for camera calibration (i.e., by the calibration network). The total time is 5.5ms. When they are used for the boxing game, the time is increased to 6.6ms, as shown in Table 3. It implies that Unity consumes 1.1ms for rigid-body simulation and rendering. In contrast, the time is increased to 36.2ms for the virtual try-on application, which implies that our XPBD implementation consumes approximately 30ms.

Using a simplistic cloth model, Figure 6 shows a proof-of-concept implementation. More realistic garments can be simulated with XPBD, but the design of realistic garments introduces a separate challenge that must be carefully managed in practical systems. For example, they have to be designed manually so that the clothes do not interpenetrate with the human mesh at initialization (i.e., at T pose). We believe that such an effort is orthogonal to the current contribution of this paper. We envision that the matter can be addressed in a future work targeted at serious clothing applications.

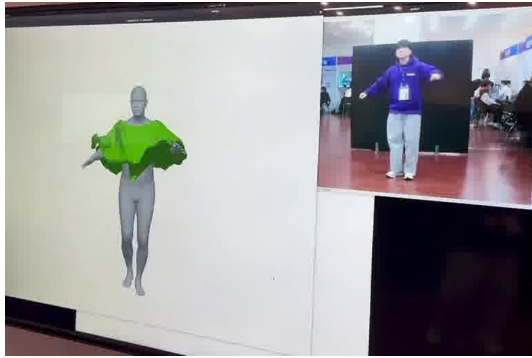
Our AR setup has its own disadvantages. Unlike AR headsets such as HoloLens [42], Meta Quest 3 [41] and Apple Vision Pro [2], the large screen cannot provide binocular depth perception, often failing to create immersive experiences for users. Despite the drawbacks of lacking depth perception and immersion, however, we believe that the applications that can be built upon our AR setup are different from those on AR headsets, and the virtual try-on is a good ex-



(a)



(b)



(c)

Figure 6: In this virtual try-on demo, a real human physically interacts with a virtual garment in real time.

ample for the first category of applications.

5. Conclusion and Future Work

This paper presents a practical solution that enables natural interactions of a human with 3D virtual objects in AR environments using only a conventional RGB video camera. Built upon deep learning techniques that are well balanced between accuracy and efficiency, the solution reconstructs and calibrates the human model in a robust manner so as to handle collision and oc-

clusion successfully.

While our approach may be able to facilitate significantly full-body interaction in AR, there is potential for further refinement. Our AR system relies on the HPS estimation technique, yet it does not offer precise 3D localization due to the limitations of the SMPL model. Our solution addresses this limitation through a trajectory model. However, the trajectory model is limited within the cases relevant to training dataset. This can fail with challenging poses, e.g., when the user is lying down. This would limit the potential applications of the current method to games and entertainment content. To extend the application range, an advanced technique needs to be developed, thus improving the quality of interactions.

A. Reconstructing Full-Body Pose and Shape

In our system, the bounding-box image returned by an object detector is scaled to 224×224 and then fed into the CNN encoder. The encoder-generated features are iteratively processed by a regression module to produce (i) the pose parameters, $\theta \in \mathbb{R}^{3K}$, where K denotes the number of skeletal joints ($K = 22$ in the current implementation), and (ii) the shape parameters, $\beta \in \mathbb{R}^{10}$.

The pose parameters, θ , which represent the relative rotations of $K = 22$ joints, are converted into the axis-angle representations for integration with SMPL-X model. The shape parameters, β , are characterized by the first 10 principal components of the PCA shape space. SMPL-X is a differentiable function, W , which takes θ and β as input and produces a triangular mesh, \mathcal{V} , of 10,475 vertices:

$$\mathcal{V} = W(\theta, \beta) \quad (3)$$

The 3D joint positions, denoted by \mathcal{J} , are derived from \mathcal{V} via a linear regression function R :

$$\mathcal{J} = R(\mathcal{V}) \quad (4)$$

The 3D human model consists of the 3D mesh \mathcal{V} and the 3D joints \mathcal{J} .

B. Calibration

Given the focal lengths, f_x and f_y , and the principal point, (c_x, c_y) , the camera intrinsic matrix \mathbf{K} is defined as follows:

$$\mathbf{K} = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix}$$

The root's homogeneous coordinates, $(u, v, 1)$, are transformed by the inverse of the intrinsic matrix \mathbf{K} and scaled by z estimated by \mathcal{M} :

$$\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} = z \mathbf{K}^{-1} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = z \begin{pmatrix} \frac{1}{f_x} & 0 & -\frac{c_x}{f_x} \\ 0 & \frac{1}{f_y} & -\frac{c_y}{f_y} \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} z \left(\frac{u - c_x}{f_x} \right) \\ z \left(\frac{v - c_y}{f_y} \right) \\ z \end{pmatrix} \quad (5)$$

Note that Z is identical to z estimated by \mathcal{M} .

C. Symmetric convolutions vs. causal convolutions

Figure 7 is a copy of Fig. 6 from Pavllo et al. [50] and compares symmetric convolutions and casual ones.

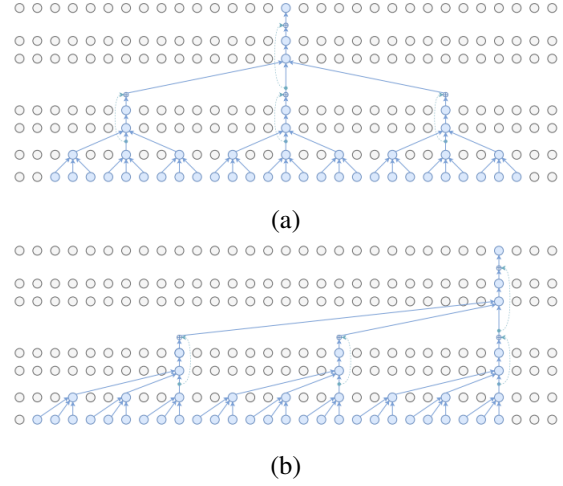


Figure 7: Two types of convolutions: (a) Symmetric convolutions (b) Causal convolutions.

D. Integration with Unity

After the calibration network (presented in Section 3.1) is trained, it is converted into the Open Neural Network Exchange (ONNX) [3] format. This conversion ensures compatibility with Unity via a lightweight cross-platform neural network inference library designed for Unity, Sentic [70], which supports deep learning network execution on both GPU and CPU.

References

- [1] Apple. Arkit. <https://developer.apple.com/documentation/arkit/>.
- [2] Apple. Vision pro. <https://www.apple.com/apple-vision-pro/>, 2024.
- [3] J. Bai, F. Lu, K. Zhang, et al. Onnx: Open neural network exchange. <https://github.com/onnx/onnx>, 2019.
- [4] M. J. Black, P. Patel, J. Tesch, and J. Yang. Bedlam: A synthetic dataset of bodies exhibiting detailed lifelike animated motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8726–8737, 2023.
- [5] F. Bogo, A. Kanazawa, C. Lassner, P. Gehler, J. Romero, and M. J. Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part V 14*, pages 561–578. Springer, 2016.
- [6] J. Bolton, M. Lambert, D. Lirette, and B. Unsworth. Paperdude: a virtual reality cycling exergame. In *CHI’14 Extended Abstracts on Human Factors in Computing Systems*, pages 475–478. 2014.
- [7] P. Bourdin, J. M. T. Sanahuja, C. C. Moya, P. Haggard, and M. Slater. Persuading people in a remote destination to sing by beaming there. In *Proceedings of the 19th ACM Symposium on Virtual Reality Software and Technology*, pages 123–132, 2013.
- [8] A. Cannavò, F. G. Praticò, A. Bruno, and F. Lamberti. Ar-mocap: Using augmented reality to support motion capture acting. In *2023 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*, pages 318–327. IEEE, 2023.
- [9] C. Canton-Ferrer, J. R. Casas, and M. Pardas. Marker-based human motion capture in multiview sequences. *EURASIP Journal on Advances in Signal Processing*, 2010:1–11, 2010.
- [10] P. Caserman, A. Garcia-Agundez, and S. Göbel. A survey of full-body motion reconstruction in immersive virtual reality applications. *IEEE transactions on visualization and computer graphics*, 26(10):3089–3108, 2019.
- [11] S. Chagué and C. Charbonnier. Real virtuality: a multi-user immersive platform connecting real and virtual worlds. In *Proceedings of the 2016 virtual reality international conference*, pages 1–3, 2016.
- [12] A. Chalmers, F. Zaman, and T. Rhee. Avatar360: Emulating 6-dof perception in 360° panoramas through avatar-assisted navigation. In *2024 IEEE Conference Virtual Reality and 3D User Interfaces (VR)*, pages 630–638. IEEE, 2024.
- [13] C.-H. Chen and D. Ramanan. 3d human pose estimation= 2d pose estimation+ matching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7035–7043, 2017.
- [14] X. Chen, X. Jiang, L. Zhan, S. Guo, Q. Ruan, G. Luo, M. Liao, and Y. Qin. Full-body human motion reconstruction with sparse joint tracking using flexible sensors. *ACM Transactions on Multimedia Computing, Communications and Applications*, 20(2):1–19, 2023.
- [15] K.-T. Chou, M.-C. Hsiu, and C. Wang. Fighting gulliver: An experiment with cross-platform players fighting a body-controlled giant. In *Proceedings of the 33rd Annual ACM Conference Extended Abstracts on Human Factors in Computing Systems*, pages 65–68, 2015.
- [16] Y. Cui, W. Chang, and D. Stricker. Fully automatic body scanning and motion capture using two kinects. In *SIGGRAPH Asia 2013 Posters*, pages 1–1. 2013.
- [17] G. D. de Dinechin and A. Paljic. Virtual agents from 360 video for interactive virtual reality. In *Proceedings of the 32nd International Conference on Computer An-*

- imation and Social Agents, pages 75–78, 2019.
- [18] H. G. Debarba, S. Perrin, B. Herbelin, and R. Boulic. Embodied interaction using non-planar projections in immersive virtual reality. In *Proceedings of the 21st ACM Symposium on Virtual Reality Software and Technology*, pages 125–128, 2015.
- [19] T. Fan, K. V. Alwala, D. Xiang, W. Xu, T. Murphey, and M. Mukadam. Revitalizing optimization for 3d human pose and shape estimation: A sparse constrained formulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11457–11466, 2021.
- [20] A. Genay, A. Lécuyer, and M. Hachet. Being an avatar “for real”: a survey on virtual embodiment in augmented reality. *IEEE Transactions on Visualization and Computer Graphics*, 28(12):5071–5090, 2021.
- [21] N. Ghorbani and M. J. Black. Soma: Solving optical marker-based mocap automatically. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11117–11126, 2021.
- [22] Y. Hu, T.-M. Li, L. Anderson, J. Ragan-Kelley, and F. Durand. Taichi: a language for high-performance computation on spatially sparse data structures. *ACM Transactions on Graphics (TOG)*, 38(6):1–16, 2019.
- [23] D.-H. Hwang, S. Kim, N. Monet, H. Koike, and S. Bae. Lightweight 3d human pose estimation network training using teacher-student learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 479–488, 2020.
- [24] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1325–1339, 2013.
- [25] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7122–7131, 2018.
- [26] M. Kari, T. Grosse-Puppenthal, L. F. Coelho, A. R. Fender, D. Bethge, R. Schütte, and C. Holz. Transformr: Pose-aware object substitution for composing alternate mixed realities. In *2021 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 69–79. IEEE, 2021.
- [27] J. Kim and M. Kim. Motion capture with high-speed rgb-d cameras. In *2014 International Conference on Information and Communication Technology Convergence (ICTC)*, pages 394–395. IEEE, 2014.
- [28] K. Kim, S. Seo, D. Han, and H. Kang. Damo: A deep solver for arbitrary marker configuration in optical motion capture. *ACM Transactions on Graphics*, 2024.
- [29] M. Kocabas, N. Athanasiou, and M. J. Black. Vibe: Video inference for human body pose and shape estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5253–5263, 2020.
- [30] M. Kocabas, C.-H. P. Huang, J. Tesch, L. Müller, O. Hilliges, and M. J. Black. Spec: Seeing people in the wild with an estimated camera. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11035–11045, 2021.
- [31] Y. Li, D. Weng, D. Li, and Y. Wang. A low-cost drift-free optical-inertial hybrid motion capture system for high-precision human pose detection. In *2019 IEEE International Symposium on Mixed and Augmented Reality Adjunct (ISMAR-Adjunct)*, pages 75–80. IEEE, 2019.
- [32] Z. Liu, L. Zhou, H. Leung, and H. P. Shum. Kinect posture reconstruction based on a local mixture of gaussian process

- models. *IEEE transactions on visualization and computer graphics*, 22(11):2437–2450, 2015.
- [33] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: a skinned multi-person linear model. *ACM Trans. Graph.*, 34(6), Oct. 2015.
- [34] D. C. Luvizon, M. Habermann, V. Golyanik, A. Kortylewski, and C. Theobalt. Scene-aware 3d multi-human motion capture from a single camera. In *Computer Graphics Forum*, volume 42, pages 371–383. Wiley Online Library, 2023.
- [35] M. Macklin, M. Müller, and N. Chentanez. Xpbd: position-based simulation of compliant constrained dynamics. In *Proceedings of the 9th International Conference on Motion in Games*, pages 49–54, 2016.
- [36] J. Martinez, R. Hossain, J. Romero, and J. J. Little. A simple yet effective baseline for 3d human pose estimation. In *Proceedings of the IEEE international conference on computer vision*, pages 2640–2649, 2017.
- [37] E. Matsas, G.-C. Vosniakos, and D. Batras. Modelling simple human-robot collaborative manufacturing tasks in interactive virtual environments. In *Proceedings of the 2016 Virtual Reality International Conference*, pages 1–4, 2016.
- [38] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *2017 international conference on 3D vision (3DV)*, pages 506–516. IEEE, 2017.
- [39] D. Mehta, O. Sotnychenko, F. Mueller, W. Xu, M. Elgharib, P. Fua, H.-P. Seidel, H. Rhodin, G. Pons-Moll, and C. Theobalt. Xnect: Real-time multi-person 3d motion capture with a single rgb camera. *Acm Transactions On Graphics (TOG)*, 39(4):82–1, 2020.
- [40] D. Mehta, S. Sridhar, O. Sotnychenko, H. Rhodin, M. Shafiei, H.-P. Seidel, W. Xu, D. Casas, and C. Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. *Acm transactions on graphics (tog)*, 36(4):1–14, 2017.
- [41] Meta. Quest 3. <https://www.meta.com/kr/quest/quest-3/>, 2023.
- [42] Microsoft. Hololens 2. <https://www.microsoft.com/en-us/hololens/>, 2019.
- [43] G. Moon, J. Y. Chang, and K. M. Lee. Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10133–10142, 2019.
- [44] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE international symposium on mixed and augmented reality*, pages 127–136. Ieee, 2011.
- [45] H. Nguyen, A. Chalmers, D. Cheong, M. Kim, T. Rhee, and J. Han. A simple but effective ar framework for human-object interaction. In *EuroXR 2024: Proceedings of the 21st EuroXR International Conference*, pages 67–71, 2024.
- [46] H. S. Nguyen, A. Chalmers, D. Cheong, M. G. Kim, T. Rhee, and J. Han. Full-Body Interaction in Mixed Reality using 3D Pose and Shape Estimation . In *2025 IEEE Conference on Virtual Reality and 3D User Interfaces Abstracts and Workshops (VRW)*, pages 1306–1307, Los Alamitos, CA, USA, Mar. 2025. IEEE Computer Society.
- [47] H. S. Nguyen, M. Kim, C. Im, S. Han, and J. Han. Convnextpose: A fast accurate method for 3d human pose estimation and its ar fitness application in mobile devices. *IEEE Access*, 2023.

- [48] P. Patel and M. J. Black. Camerahmr: Aligning people with perspective. *arXiv preprint arXiv:2411.08128*, 2024.
- [49] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. Osman, D. Tzionas, and M. J. Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10975–10985, 2019.
- [50] D. Pavllo, C. Feichtenhofer, D. Grangier, and M. Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7753–7762, 2019.
- [51] T. C. Peck, S. Seinfeld, S. M. Aglioti, and M. Slater. Putting yourself in the skin of a black avatar reduces implicit racial bias. *Consciousness and cognition*, 22(3):779–787, 2013.
- [52] Z. Qiu, Q. Yang, J. Wang, H. Feng, J. Han, E. Ding, C. Xu, D. Fu, and J. Wang. Psvt: End-to-end multi-person 3d pose and shape estimation with progressive video transformers. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 21254–21263, 2023.
- [53] Y. Ren, C. Zhao, Y. He, P. Cong, H. Liang, J. Yu, L. Xu, and Y. Ma. Lidar-aid inertial poser: Large-scale human motion capture by sparse inertial and lidar sensors. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2337–2347, 2023.
- [54] T. Rhee, A. Chalmers, F. Zaman, A. Stangnes, and V. Roberts. Real-time stage modelling and visual effects for live performances. In *ACM SIGGRAPH 2023 Real-Time Live!*, pages 1–2. 2023.
- [55] T. Rhee, L. Petikam, B. Allen, and A. Chalmers. Mr360: Mixed reality rendering for 360 panoramic videos. *IEEE transactions on visualization and computer graphics*, 23(4):1379–1388, 2017.
- [56] T. Rhee, S. Thompson, D. Medeiros, R. Dos Anjos, and A. Chalmers. Augmented virtual teleportation for high-fidelity telecollaboration. *IEEE transactions on visualization and computer graphics*, 26(5):1923–1933, 2020.
- [57] D. Roth, J.-L. Lugin, J. Büser, G. Bente, A. Fuhrmann, and M. E. Latoschik. A simplified inverse kinematic approach for embodied vr applications. In *2016 IEEE Virtual Reality (VR)*, pages 275–276. IEEE, 2016.
- [58] M. Šarić, M. Russo, L. Kraljević, and D. Meter. Extended reality telemedicine collaboration system using patient avatar based on 3d body pose estimation. *Sensors*, 24(1):27, 2023.
- [59] A. C. Sementille, L. E. Lourenço, J. R. F. Brega, and I. Rodello. A motion capture system using passive markers. In *Proceedings of the 2004 ACM SIGGRAPH international conference on Virtual Reality continuum and its applications in industry*, pages 440–447, 2004.
- [60] L. A. Shaw, B. C. Wünsche, C. Lutteroth, S. Marks, and R. Callies. Challenges in virtual reality exergame design. 2015.
- [61] S. Shimada, V. Golyanik, W. Xu, P. Pérez, and C. Theobalt. Neural monocular 3d human motion capture with physical awareness. *ACM Transactions on Graphics (ToG)*, 40(4):1–15, 2021.
- [62] H. Shum and E. S. Ho. Real-time physical modelling of character movements with microsoft kinect. In *Proceedings of the 18th ACM symposium on Virtual reality software and technology*, pages 17–24, 2012.
- [63] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei. Integral human pose regression. In *Proceedings of the European conference on computer vision (ECCV)*, pages 529–545, 2018.
- [64] Y. Sun, Q. Bao, W. Liu, Y. Fu, M. J. Black, and T. Mei. Monocular, one-stage, regression of multiple 3d people. In *Proceedings*

- of the *IEEE/CVF international conference on computer vision*, pages 11179–11188, 2021.
- [65] Y. Sun, Q. Bao, W. Liu, T. Mei, and M. J. Black. Trace: 5d temporal regression of avatars with dynamic cameras in 3d environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8856–8866, 2023.
- [66] Y. Sun, W. Liu, Q. Bao, Y. Fu, T. Mei, and M. J. Black. Putting people in their place: Monocular regression of 3d people in depth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13243–13252, June 2022.
- [67] T. M. Takala and M. Matveinen. Full body interaction in virtual reality with affordable hardware. In *2014 IEEE Virtual Reality (VR)*, pages 157–157. IEEE, 2014.
- [68] M. Teschner, B. Heidelberger, M. Müller, D. Pomerantes, and M. H. Gross. Optimized spatial hashing for collision detection of deformable objects. In *Vmv*, volume 3, pages 47–54, 2003.
- [69] E. Tuveri, L. Macis, F. Sorrentino, L. D. Spano, and R. Scateni. Fitmersive games: Fitness gamification through immersive vr. In *Proceedings of the International Working Conference on Advanced Visual Interfaces*, pages 212–215, 2016.
- [70] Unity. Sentis. <https://unity.com/products/sentis>. Accessed: 2024-03-08.
- [71] S. Wang, J. Li, T. Li, Y. Yuan, H. Fuchs, K. Nagano, S. De Mello, and M. Stengel. Blade: Single-view body mesh learning through accurate depth estimation. *arXiv preprint arXiv:2412.08640*, 2024.
- [72] Z. Weng and S. Yeung. Holistic 3d human and scene mesh estimation from single view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 334–343, 2021.
- [73] E. Wu and H. Koike. Futurepose-mixed reality martial arts training using real-time 3d human pose forecasting with a rgb camera. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1384–1392. IEEE, 2019.
- [74] D. Xiang, H. Joo, and Y. Sheikh. Monocular total capture: Posing face, body, and hands in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10965–10974, 2019.
- [75] D. Xiang, F. Prada, C. Wu, and J. Hodgins. Monoclothcap: Towards temporally coherent clothing capture from monocular rgb video. In *2020 International Conference on 3D Vision (3DV)*, pages 322–332. IEEE, 2020.
- [76] X. Yi, Y. Zhou, M. Habermann, V. Golyanik, S. Pan, C. Theobalt, and F. Xu. EgoLocate: Real-time motion capture, localization, and mapping with sparse body-mounted sensors. *ACM Transactions on Graphics (TOG)*, 42(4):1–17, 2023.
- [77] M. K. Young, J. J. Rieser, and B. Bodenheimer. Dyadic interactions with avatars in immersive virtual environments: High fiving. In *Proceedings of the ACM SIGGRAPH Symposium on Applied Perception*, pages 119–126, 2015.
- [78] A. Zanfir, E. Marinoiu, and C. Sminchisescu. Monocular 3d pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2148–2157, 2018.
- [79] A. Zanfir, E. Marinoiu, M. Zanfir, A.-I. Popa, and C. Sminchisescu. Deep network for the integrated 3d sensing of multiple people in natural images. *Advances in neural information processing systems*, 31, 2018.
- [80] C. Zheng, S. Zhu, M. Mendieta, T. Yang, C. Chen, and Z. Ding. 3d human pose estimation with spatial and temporal transformers. In *Proceedings of the IEEE/CVF*

International Conference on Computer Vision, pages 11656–11665, 2021.