

SmartMocap: Joint Estimation of Human and Camera Motion using Uncalibrated RGB Cameras

Nitin Saini^{1,2}, Chun-Hao P. Huang¹, Michael J. Black¹, and Aamir Ahmad^{2,1}

Abstract—Markerless human motion capture (mocap) from multiple RGB cameras is a widely studied problem. Existing methods either need calibrated cameras or calibrate them relative to a static camera, which acts as the reference frame for the mocap system. The calibration step has to be done *a priori* for every capture session, which is a tedious process, and re-calibration is required whenever cameras are intentionally or accidentally moved. In this paper, we propose a mocap method which uses multiple static and moving extrinsically uncalibrated RGB cameras. The key components of our method are as follows. First, since the cameras and the subject can move freely, we select the ground plane as a common reference to represent both the body and the camera motions unlike existing methods which represent bodies in the camera coordinate system. Second, we learn a probability distribution of short human motion sequences (~ 1 sec) relative to the ground plane and leverage it to disambiguate between the camera and human motion. Third, we use this distribution as a motion prior in a novel multi-stage optimization approach to fit the SMPL human body model and the camera poses to the human body keypoints on the images. Finally, we show that our method can work on a variety of datasets ranging from aerial cameras to smartphones. It also gives more accurate results compared to the state-of-the-art on the task of monocular human mocap with a static camera. A video demo and our code are available at <https://tinyurl.com/yeykrb67> and <https://tinyurl.com/2p9rme9y>.

Index Terms—Gesture, Posture and Facial Expressions; Human Detection and Tracking; Deep Learning for Visual Perception

I. INTRODUCTION

MODERN markerless methods use RGB cameras to estimate human motion without the need for markers or sensors on the subject's body [1]. They use sparse 2D keypoints to either fit a 3D body model or train a neural network to output the parameters of the body model. Existing monocular methods [2] take images from a single static camera and estimate the motion of the person relative to it. Since, most applications need human motion relative to the world, the camera should be calibrated relative to the world. Another problem with this setup is that the subject's body parts can often get self-occluded. One solution is to make the camera freely moving such that it can optimize its view for motion estimation [3]. However, a moving camera is much more

difficult to calibrate relative to the world [1]. Estimating the subject's global motion is also difficult using a single camera because the camera motion and the subject's motion cannot be disambiguated using only the sparse 2D keypoints. Multi-view methods employ multiple static cameras to handle self-occlusions. They calibrate the cameras relative to the world in a separate calibration step, which increases the preparation time. The cameras should remain static after the calibration step. In case they are moved intentionally or accidentally, the calibration has to be performed again, making the capture process highly inconvenient and time-consuming.

To address the aforementioned issues, in this paper, we present a system for outdoor human mocap using a set of RGB cameras, where some cameras are static while others are moving. This system is quick to set up, as users can place the cameras and immediately start the capture session. Any camera can be moved during the mocap session to get better visibility of the subject, and each camera, using our method, extrinsically calibrates itself relative to the world using only the sparse 2D keypoints of the human body. Our system does not need a pre-calibration of the extrinsic parameters of the camera (pose of the camera in the 3D space). It, however, needs the camera intrinsics (related to the camera sensor and lens). Since these remain constant for any camera, the intrinsic calibration needs to be done only once and then can be used in multiple mocap sessions.

Our mocap method takes in the synchronized videos from multiple RGB cameras and estimates the camera poses and the subject's motion, defined as the trajectory of human poses (articulated and global), and shape in all the frames. All the estimates are relative to a global reference frame, which is the ground projection of the human root joint onto the ground plane in the first frame. The ground plane is defined as the XY plane. We learn a probability distribution of the human motion relative to a ground plane by training a variational autoencoder using a large human motion dataset (AMASS) [4]. The state-of-the-art human motion prior [5] learns the distribution of pose transitions, which is defined as the difference between two consecutive poses. This is highly sensitive to noise in long-term motion generation/estimation. Contrarily, we learn the distribution of the trajectory of body joint positions and joint angles relative to the above-defined world frame. The length of each motion sequence is fixed as 25 frames at the rate of 30 frames per second. We use this probability distribution to fit the SMPL human body model [6] to the sparse 2D keypoints in all the views. These 2D keypoints are obtained using the openpose 2D keypoint detector [7]. While our motion prior encodes the distribution of human motion relative

Manuscript received: December, 09, 2022; Revised January, 18, 2023; Accepted March, 13, 2023.

This paper was recommended for publication by Editor Gentiane Venture upon evaluation of the Associate Editor and Reviewers' comments.

¹Max Planck Institute for Intelligent Systems, 72076 Tübingen, Germany. firstname.lastname@tuebingen.mpg.de

²Institute of Flight Mechanics and Controls, University of Stuttgart, 70569 Stuttgart, Germany. firstname.lastname@ifr.uni-stuttgart.de
Digital Object Identifier (DOI): see top of this page.

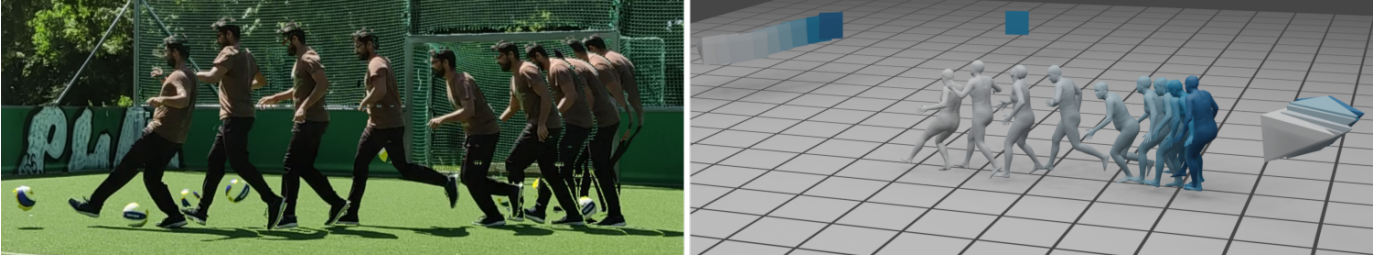


Fig. 1: Multi-exposure image of a person playing football [left] and the reconstructed motion of the person and the cameras using our method [right].

to the ground plane, the 2D keypoints contain information about the subject's articulated poses and the camera poses relative to the subject. In the optimization formulation, we directly optimize for the camera poses and the human poses in the world frame jointly, and condition the human motion using our learned probability distribution of human motions. This keeps the reconstructed poses aligned with the ground plane. However, this formulation is highly non-convex and susceptible to converging to a local minimum. Therefore, we first initialize the human position as the mean of human motion in the learned latent space. The articulated human pose and the camera poses are initialized using the estimates of a human pose regressor [8]. Since a motion sequence in our learned latent space is of a fixed length of 25 frames and starts from the origin, we take a multi-resolution optimization approach. After the initialization stage, we run the optimization stage, where we split the full sequence into chunks of 25 frames and run optimization on these chunks independently. This optimization treats the starting of each chunk as the origin. In the next stage (stitching stage), we stitch these chunks together such that the last frame of a sequence aligns with the first frame of the next sequence. In the final stage, we run the optimization again on this stitched sequence to get the final estimates. For longer sequences, stitching the full sequence can accumulate noise in the orientation estimates, which leads to poor initialization for the next optimization stage. To avoid this, we stitch together a smaller number of them, perform optimization, stitch, optimize and iterate. This way, we slowly increase the temporal resolution at each stitching stage and perform the alternate stitching and optimization until the final optimization for the full sequence.

In summary, we have the following novel contributions:

- A human motion prior which encodes the global and articulated human motion relative to a global reference frame (ground plane).
- A multi-resolution optimization method for estimating camera poses along with the human motion and shape relative to the ground plane using single/multiple uncalibrated RGB cameras.

II. RELATED WORK

A. Multi-view methods

Most of the existing markerless human mocap methods utilize videos from multiple calibrated and time-synchronized cameras. They first detect 2D features (keypoints, silhouette etc.) on the image plane and then use the camera calibration parameters to either project them in 3D space [9] or fit the 3D human body (model parameters, body joints etc.) to the 2D

features [10]. Since calibrating the cameras in a separate step is not always possible, [11] and [12] utilize the human body to calibrate the static cameras. While [11] use simple human motion constraints such as constant velocity or acceleration of the human joints, [12] learn a human motion prior and use it to fit the SMPL body model [6] to the 2D keypoints. [13] use pan-tilt-zoom cameras and triangulate the human annotated 2D keypoints to reconstruct the 3d motion of people skiing and the pan, tilt and zoom of the cameras. The static cameras cannot actively change the viewpoint for better mocap. Therefore, [14] and [15] use moving cameras along with the static cameras as their mocap setup. However, all the above methods for uncalibrated cameras estimate the human motion relative to one static camera, which should be calibrated relative to the world such that the estimated human motion can be transformed to the world reference frame. Calibrating cameras is particularly hard if all the cameras are moving. [16] use multiple handheld smartphone cameras and calibrates them relative to a static background using a structure-from-motion (SfM) method. However, such calibration is not reliable and only works for a static background with suitable texture. [1] use cameras mounted on custom-designed multiple micro-aerial vehicles and calibrates them using the onboard IMU and GPS sensors. However, such sensors are not available for ordinary RGB cameras.

B. Monocular methods

Early monocular methods like [17], fit a human body model [18] to the 2D keypoints and body silhouette on the image plane. Later ones use deep neural networks to directly regress the parameters of a human body model directly from the RGB image [19]. All the above, estimate the human pose/motion relative to the camera or relative to some local coordinate system on the human body. Recent methods like [5] try to estimate human motion in a world reference frame using a monocular video. HuMoR [5] learned the distribution of human motion transitions and used it to fit the SMPL model to the 2D keypoints on a monocular video from a static camera. However, encoding only the motion transitions is very sensitive to noisy observations (2D keypoints) which can lead to unrealistic human motion estimates. Methods like [20] estimate human pose relative to the camera and use SLAM to track the camera poses that need textured background. Methods such as [21] train a regressor network which can output the trajectory of global pose given the articulated poses. This makes the final output more sensitive to the noisy articulated poses estimated in the previous stages. They do not utilize the fact that the articulated poses can also be improved

using the global pose information, which we do. Additionally, they represent the global poses trajectory using relative local differences (similar to [5]) which makes the future results sensitive to the noise in past poses.

III. APPROACH

A. Goal and preliminaries

Given synchronized image sequences of length T frames from C cameras looking at a moving person, the goal is to estimate the camera motion, the person's shape, and the person's motion, which is defined as the trajectory of the person's poses (articulated and global). We use the SMPL human body model [6] to represent the human poses. SMPL is parameterized by joint angles ($\theta \in \mathbb{R}^{63}$), body shape parameters ($\beta \in \mathbb{R}^{10}$), root orientation ($\phi \in \mathbb{R}^3$) and root position ($\tau \in \mathbb{R}^3$). We use $N = 22$ body joints from SMPL, which includes 21 body joints and 1 root joint. We exclude the 2 hand joints from the original 24 joints in the SMPL model. Instead of representing the articulated pose as joint angles, we represent the subject's articulated pose at any time t in the latent space of VPoser [2] ($z \in \mathbb{R}^{32}$), which is a learned probability distribution of human poses. It is a variational autoencoder (VAE) with encoder (\mathcal{V}_E) and decoder (\mathcal{V}_D). The full human motion is then represented as $((\tau_1, \phi_1, z_1), \dots, (\tau_T, \phi_T, z_T), \beta)$. The position and orientation of a camera c at any time t is represented as $p_{c,t} \in \mathbb{R}^3$ and $r_{c,t} \in \mathbb{R}^6$ respectively. Unless explicitly stated, we use the 6D representation [22] to represent the rotations in this paper. The camera motion for any camera c is represented as $((r_{c,1}, p_{c,1}), \dots, (r_{c,T}, p_{c,T}))$. Our human motion prior uses a different representation of the body pose. The body pose x_t at any time t is the orientation and position of each body joint relative to the world frame, i.e. $x_t \in \mathbb{R}^{22 \times (6+3)}$. The origin of the world frame is defined as the ground projection of the SMPL root joint in the first frame and the ground plane is defined as the XY plane. The motion prior encodes the fixed length of 25 consecutive poses, thus, the motion sequence for the motion prior is represented as $\mathbf{x} = (x_1, \dots, x_{25})$. We represent the estimated value of any parameter by putting a tilde over it, e.g. $\tilde{\mathbf{x}}$ is the estimated value of \mathbf{x} .

B. Human motion prior

We use a VAE to learn a distribution in the latent space of human motion sequences of a fixed length. Each motion sequence consists of 25 consecutive human body poses at the rate of 30 frames per second. The forward-facing direction of the SMPL root joint in the first frame is aligned with the +Y axis of the origin.

1) *Training data*: We use AMASS dataset [4] to train our network. AMASS is a collection of multiple human mocap datasets, unifying them with the SMPL body representation. We follow the preprocessing steps in HuMoR [5] which removes the motion sequences where the person's feet are skating or sliding over the static ground plane. In such motions, the person doesn't interact with a stationary ground plane but with an object such as a treadmill or skates. The preprocessing step also changes the frame rate to 30 FPS and gives out a total of 11893 motion sequences. We randomly

select 25 consecutive frames from any of these sequences and canonicalize them such that the origin is the ground projection of the root joint at the first frame and the person's forward direction is aligned with the origin's +Y axis. Even though the shape of the subjects in AMASS varies, we follow [5] and keep the body shape constant to the mean shape. This reduces the complexity of the model by ignoring the body shape variations at the expense of some possible artefacts such as foot-skating.

2) *Model architecture and training*: We use convolutional architecture for both the encoder and the decoder networks. The encoder (\mathcal{M}_E) consists of a 1D convolutional (conv) layer at the input and 4 identical residual blocks. We modify the ResNet [23] residual blocks to create these blocks. We replace the 2D convolutions with 1D convolutions. We further replace the ReLU units with the GELU units within the blocks. We also place a GELU unit after the input 1D conv layer and each of the residual blocks. The output dimension of the first conv layer is 1024. The input and output dimension of each residual block is 1024. Furthermore, two linear layers transform the output of the last residual block to the mean ($\mu \in \mathbb{R}^{1024}$) and log of variance ($\log(\sigma^2) \in \mathbb{R}^{1024}$) of the gaussian distribution in the latent space, from which the latent value ($m \in \mathbb{R}^{1024}$) is sampled using the reparametrization trick [24]. The decoder (\mathcal{M}_D) architecture also consists of 4 consecutive residual blocks, similar to the ones in the encoder. The first residual block acts as the input layer, and there is a deconvolutional layer at the output of the decoder.

We train our motion VAE network using a combination of reconstruction (\mathcal{L}_{rec}) and KL divergence loss (\mathcal{L}_{KL}).

$$\mathcal{L} = \mathcal{L}_{rec} + w_{kl} \mathcal{L}_{KL}, \quad (1)$$

where

$$\mathcal{L}_{rec} = \|\mathbf{x} - \tilde{\mathbf{x}}\|^2 \quad \text{and} \quad \mathcal{L}_{KL} = -\frac{0.5}{1024} \sum_i^{1024} (1 + \log(\sigma_i^2) - \sigma_i^2 - \mu_i^2). \quad (2)$$

We employ a 20 epochs cyclic annealing scheme for the parameter w_{kl} [25]. Initially, the value of w_{kl} starts at 0 and increases linearly with the training epochs. After 10 epochs, the value reaches 1 and stays constant for another 10 epochs. The value again drops to 0 and the cycle continues.

C. Camera and human pose estimation

First, we use openpose [7] to detect 2D keypoints of the subject in each image. Then we use them in our method which consists of the following steps, 1) Initialize, 2) Optimize, 3) Stitch and 4) Optimize-final (see fig. 2).

• **Step 1: Initialize** We initialize the SMPL and camera poses for each frame using the results from PARE [8]. PARE gives the camera pose relative to the person and the person's articulated pose for each image. We take the mean of the articulated poses in all the views (θ_{init}), project it to the VPoser latent space (z_{init}) and use it as the initial articulated pose of the subject. For the initialization of SMPL position (τ_{init}) and orientation (ϕ_{init}) relative to the ground plane, we use the decoded output of the mean value in the motion prior latent space. We use the initialized pose of the person to calculate the position (r_{init}) and orientation (p_{init}) of the cameras relative

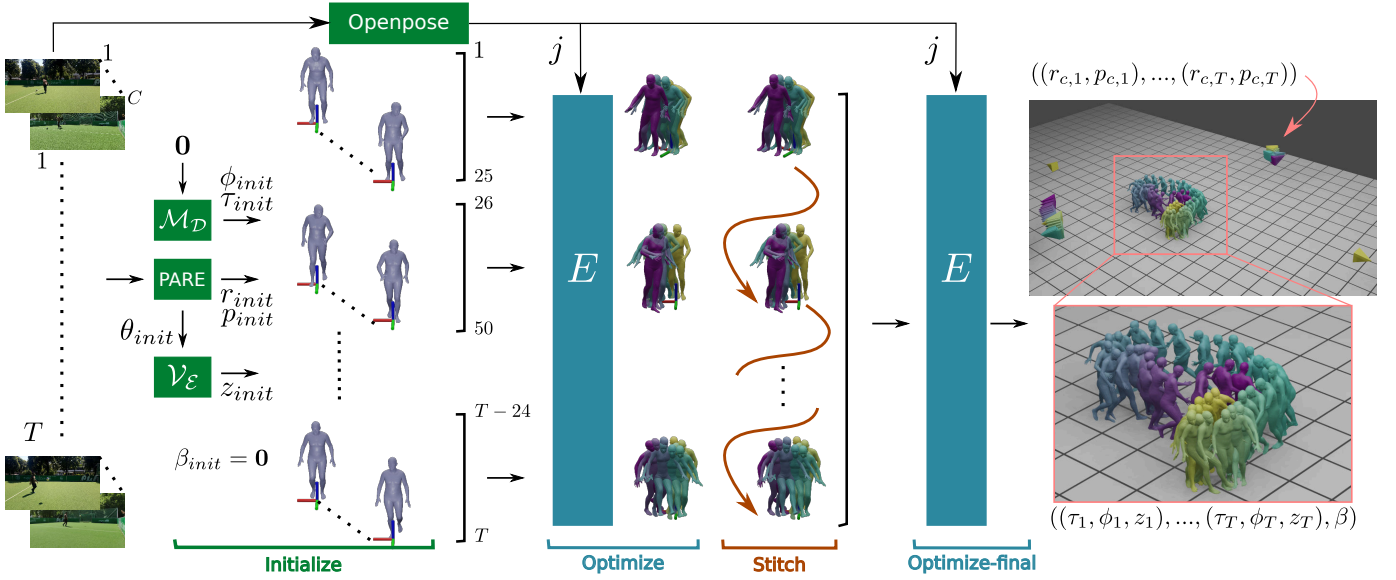


Fig. 2: Our method takes in synchronized images of a moving person from multiple intrinsically uncalibrated cameras, and processes them in 4 steps 1) Initialize, 2) Optimize, 3) Stitch and 4) Optimize-final to give the motion of the person and the cameras in the world frame.

to the ground plane. Since PARE assumes an orthographic camera with a focal length of 5000, we use the method in [26] to transform the SMPL position estimate in the actual camera. The SMPL shape (β_{init}) is initialized with a vector of zero values.

• **Step 2: Optimize** We estimate the motion of the person in small intervals (25 frames). We split the full sequence into chunks of length 25 and run the optimization for each chunk independently in three phases to avoid local minimum, similar to [15]. In the first phase, we optimize the camera poses only. In the second phase, we optimize camera poses along with SMPL position and orientation. In the final phase, we optimize all the parameters, except the SMPL position and orientation in the first frame. The initial SMPL position and orientation in the first frame act as the pivot for all the optimizing parameters.

• **Step 3: Stitch** We stitch together the estimated motions. Since the origin for each chunk is defined as the ground projection of the root joint, we stitch consecutive sequences together such that the root ground projection of the last frame of a chunk is aligned with the first frame of the next chunk. For very long sequences, we stitch together fewer sequences, optimize, stitch and repeat until all the sequences are stitched.

• **Step 4: Optimize-final** In the final optimization step (optimize-final), we again optimize all the parameters for the fully stitched sequences in three phases, the same as in the previous optimization stage. This step is the final optimization step if all the sequences are stitched. For very long sequences, we stitch together fewer chunks instead of all. Then we optimize and repeat the stitching and optimization cycle until the whole sequence is done.

In all the optimization stages, we minimize the same loss function, which is a weighted combination of multiple loss terms. It is given as

$$E = w_{2D}E_{2D} + w_m E_m + w_{3DS}E_{3DS} + w_{COS}E_{COS} + w_{CPS}E_{CPS} + w_\beta E_\beta + w_z E_z + w_{GHP}E_{HGP} + w_{CGP}E_{CGP}. \quad (3)$$

The component E_{2D} is the 2D reprojection loss given as

$$E_{2D} = \frac{1}{NT} \sum_{n,c,t} w_n \|\Pi(r_{c,t}, p_{c,t}, \mathcal{J}_n(\mathcal{V}_D(z_t), \tau_t, \phi_t, \beta)) - j_{c,t}^n\|^2, \quad (4)$$

where, \mathcal{J} is the SMPL 3D joint regressor function [6], \mathcal{V} is the VPoser decoder [2], Π is the camera projection function, $j_{c,t}^n$ is the 2D keypoint corresponding to the joint n in camera c at time instant t , and w_n is the confidence score given by the 2D detector for the joint n . The loss component E_m is the motion prior loss given as

$$E_m = \sum_t^{T-25} \|\mathcal{M}_{\mathcal{E}_\mu}(\mathcal{V}_D(z_{t:t+25}), \tau_{t:t+25}, \phi_{t:t+25}, \beta)\|^2, \quad (5)$$

where $\mathcal{M}_{\mathcal{E}_\mu}$ is the μ part of the motion prior encoder. The loss component E_{3DS} is a temporal smoothing term for the 3D joint positions. It is given as

$$E_{3DS} = \sum_t \|\mathcal{J}(\mathcal{V}_D(z_t), \tau_t, \phi_t, \beta) - \mathcal{J}(\mathcal{V}_D(z_{t-1}), \tau_{t-1}, \phi_{t-1}, \beta)\|^2. \quad (6)$$

E_{COS} and E_{CPS} are the camera motion smoothing terms. Following [15], we use L2 loss on the positions and the 6D representation of the camera orientations, given as

$$E_{COS} = \frac{1}{CT} \sum_{c,t} \|r_{c,t} - r_{c,t-1}\|^2 ; \quad E_{CPS} = \frac{1}{CT} \sum_{c,t} \|p_{c,t} - p_{c,t-1}\|^2. \quad (7)$$

E_β and E_z are the SMPL shape and VPoser regularization terms [15], given as

$$E_\beta = \|\beta\|^2 \quad \text{and} \quad E_z = \|z\|^2. \quad (8)$$

E_{HGP} and E_{CGP} are the ground penetration terms for both the human and the cameras. These terms avoid the scenarios where the cameras or the human goes below the ground plane. These are given as

$$E_{HGP} = \frac{1}{T} \sum_t \max(0, \mathcal{J}^z(\mathcal{V}_D(z_t), \tau_t, \phi_t, \beta)) \quad \text{and} \quad (9)$$

$$E_{CGP} = \frac{1}{CT} \sum_{c,t} \max(0, p_{c,t}^z), \quad (10)$$

where, \mathcal{J}^z and p^z are the vertical (z) component of the 3D joint positions and the camera positions, respectively. In all the above equations, T is replaced with 25 in the **Step 2**.



Fig. 3: Camera setup used to collect our smartphone dataset.

IV. EXPERIMENTS AND RESULTS

A. Datasets

We evaluate our method using a sequence taken from each of the following datasets

1) RICH dataset [27]: it is collected using 7 static and one moving IOI cameras. It has the ground truth poses of the person and the camera poses for the static cameras. The GT poses of the moving camera are not available.

2) AirPose real-world dataset [15]: It is collected using RGB cameras mounted on two DJI unmanned micro aerial vehicles (UAVs). One UAV is kept hovering and other is encircling the subject while looking at him.

3) Our smartphone dataset: it is collected using 4 smartphone cameras, where two of the cameras are static while the other two are moving. The data is collected when the subject is playing with a football in a small field. The camera setup can be seen in Fig. 3. We show the rest of the three cameras and the subject in the frame of the camera *Cam1*. *Cam1* and *Cam3* are static, and *Cam2* and *Cam4* are moved along the boundary walls of the field. We used OpenCamera app [28] to collect the video data on the smartphones. Each smartphone records the videos at 30 FPS. Since the frame rates are the same and constant, we manually synchronize the videos by synchronizing just one frame. Similar to AirPose dataset, the smartphone dataset also does not contain any GT.

B. Metrics

We use the following metrics to quantitatively evaluate the reconstruction by our method on the RICH dataset. 1) *Mean camera position error (MCPE)*: This is the mean value of the distance between the estimated and the GT position of all the cameras. It is given as

$$MCPE = \frac{1}{C} \sum_c \|p_c - \tilde{p}_c\|. \quad (11)$$

2) *Mean camera orientation error (MCOE)*: This metric is the mean geodesic distance between the estimated camera orientation and the GT on the 3D manifold of rotation matrices [29]. It is given as

$$MCOE = \frac{1}{C} \sum_c \arccos(0.5 * (Tr(\tilde{R}_c R_c^T) - 1)), \quad (12)$$

where R_c is the camera orientation matrix.

3) *Mean position error (MPE)*: This is the mean value of the distance between the estimated SMPL root position parameter and its GT value provided by the RICH dataset.

$$MPE = \frac{1}{T} \sum_t \|\tau_t - \tilde{\tau}_t\|. \quad (13)$$

4) *Mean orientation error (MOE)*: This metric is the mean geodesic distance between the estimated SMPL root orientation and the GT on the 3D manifold of rotation matrices [29]. It is given as

$$MOE = \frac{1}{T} \sum_t \arccos(0.5 * (Tr(\tilde{R}_\phi R_\phi^T) - 1)), \quad (14)$$

where R_ϕ is the SMPL root orientation matrix at time t .

5) *Root-aligned mean per-joint position error (RA-MPJPE)*: This metric is to quantitatively evaluate the articulated pose estimate [1]. It is the distance between the estimated 3D joints and their corresponding values when the SMPL position, orientation of the root joint and shape are aligned with their corresponding GT. For alignment, all these three parameters are set to zero. The error is given as

$$RA-MPJPE = \frac{1}{NT} \sum \|\mathcal{J}_n(\theta_t) - \mathcal{J}_n(\tilde{\theta}_t)\|. \quad (15)$$

6) *Mean per-vertex position error (MPVPE)*: This metric is to evaluate the shape estimate relative to the GT. We do the SMPL forward pass using only the estimated and GT shape parameters and then calculate the mean distance between the corresponding vertices as

$$MPVPE = \frac{1}{V} \sum \|\mathcal{S}_v(\beta) - \mathcal{S}_v(\tilde{\beta})\|, \quad (16)$$

where \mathcal{S} is the SMPL vertices regressor function [6] and V is the number of vertices in the SMPL model.

C. Results and discussion

1) *RICH dataset*: We show the evaluation metrics and their corresponding standard deviation values of our method on RICH in Table I and II. The results of our method using all 8 cameras are shown in the last row of Table I ($C_{1,...,8}$). We also compare the performance of our method with multiple camera configurations. In the first row, we show the results when only the first camera is used (C_1). Next, we add camera 8, which is moving, and the results are shown in the second row ($C_{1,8}$). We keep adding static cameras one at a time and show the results in further rows. Adding an extra view, only gives information about the articulated pose of the subject and the relative poses of the camera and the subject. This is why we don't see improvement in the global pose estimates of the person or in the camera poses with addition of more views. However, adding more views helps in handling occlusions better and we see in Fig. 4 the RA-MPJPE improves with more camera views but gets saturated after 4 views. This shows that 4 views are sufficient to resolve any uncertainty in the person's articulated pose due to occlusions, and adding more views doesn't provide any additional information.

Note that we do not optimize the SMPL position and orientation in the first frame, as it acts as the pivot and bring the person conforming to the first frame using (5) and cameras using (4). Therefore, the estimate in the first frame is noisy and because of (7), a few initial frames become noisy. Hence, we ignore the first 10 frames for evaluation.

In table II, we compare the monocular and the multi-view version of our method (row 3 and 4) with the reference methods GLAMR [21] (row 1) and HuMor [5] (row 2), which are the state-of-the-art monocular human pose and shape

Cameras	MCPE (cm)	MCOE (rad)	MPE (cm)	MOE (rad)	RA-MPJPE (cm)	MPVPE (cm)
C ₁	72.68	0.20	10.89 ± 7.31	0.27 ± 0.10	6.60 ± 4.92	3.07 ± 1.45
C _{1,8}	55.85	0.14	7.61 ± 4.14	0.22 ± 0.06	6.24 ± 5.02	3.13 ± 1.50
C _{1,2,8}	90.46	0.17	12.72 ± 2.93	0.23 ± 0.07	5.97 ± 4.91	3.00 ± 1.38
C _{1,2,3,8}	89.68	0.14	11.67 ± 2.59	0.20 ± 0.07	5.73 ± 4.80	3.02 ± 1.38
C _{1,...,4,8}	95.03	0.16	12.39 ± 2.69	0.22 ± 0.07	5.68 ± 4.68	2.91 ± 1.30
C _{1,...,5,8}	93.74	0.16	12.49 ± 2.74	0.20 ± 0.07	5.66 ± 4.65	2.91 ± 1.32
C _{1,...,6,8}	92.43	0.17	13.42 ± 2.85	0.20 ± 0.06	5.87 ± 4.77	2.31 ± 0.92
C _{1,...,8}	88.13	0.18	13.69 ± 3.05	0.19 ± 0.06	6.03 ± 4.94	1.86 ± 0.79

TABLE I: Evaluation of our method using multiple camera configurations. Camera no. 1-7 are static in the RICH dataset with available GT, and camera no. 8 is moving but GT is not available. Hence, the MCPE and MCOE metrics for rows 2-9 do not include camera 8.

Camera	MCPE (cm)	MCOE (rad)	MPE (cm)	MOE (rad)	RA-MPJPE (cm)	MPVPE (cm)
GLAMR [21]	250.93	0.27	24.07 ± 4.96	0.48 ± 0.28	9.66 ± 9.51	2.84 ± 1.12
HuMor [5]	90.09	0.17	30.32 ± 8.12	0.48 ± 0.39	10.82 ± 8.14	4.2 ± 2.09
C ₁ (ours)	72.68	0.20	10.89 ± 7.31	0.27 ± 0.10	6.6 ± 4.92	3.07 ± 1.45
C _{1,...,8} (ours)	88.13	0.18	13.69 ± 3.05	0.19 ± 0.06	6.03 ± 4.94	1.86 ± 0.79

TABLE II: Comparison of our method (C_{1,...,8}), monocular version of our method (C₁) and state-of-the-art monocular methods HuMor [5] and GLAMR [21] on RICH dataset.

estimation methods. We see that our method significantly outperforms these methods. Both HuMor and GLAMR uses a motion prior which encodes human motion transitions instead of absolute motions. As we discussed in Sec. II-B, reconstructing the motion from the space of motion transitions is very sensitive to noise and can lead to spurious results. In Fig. 5, we show the qualitative results of our method and the reference method and compare them with the GT. The 3D reconstruction of the human and the cameras relative to the ground plane are shown for our method (blue), the GT (green) and the reference method [5] (red). For a clearer illustration, we render each pose by adding a time-dependent offset to the position estimate at that time. Camera pose estimates are unchanged for the rendering. We can see that our estimates are very close to the GT, while the reference method estimates are quite inaccurate. For example, in the left inset box, we see that both the feet are on the same side, giving a physically implausible global pose estimate. In the right inset box, the person's body suddenly rotates more than 90°, again resulting in a physically implausible motion. In Fig. 6, we do a qualitative comparison of our method with GLAMR [21] by showing the resulting mesh overlaid on the original image. While the results from our method are nearly perfectly aligned with the person in the image, the overlaid GLAMR results does not match the person. This is due to the errors in both the estimated person's poses and the estimated camera pose.

2) *Airpose dataset:* In Fig. 7, we show qualitative results of our method on the Airpose real-world dataset. We show the cropped version of the original images of the subject, along with the same image with the estimated mesh overlaid on top. Two adjacent columns are the two views at the same time instant. The results show that our method can reconstruct the diverse poses captured from an aerial view. In the bottom-right corner, we show the 3D reconstruction of the subject's poses, shape, and the camera poses for a sub-sequence. The color gradient from yellow to violet is used to show the time transition. We observe that the subject's reconstructed body is not touching the ground, but lies a bit above the ground plane. This is because the actual terrain is not a plane, but a sloped hilly terrain. Even though the motion prior is trained

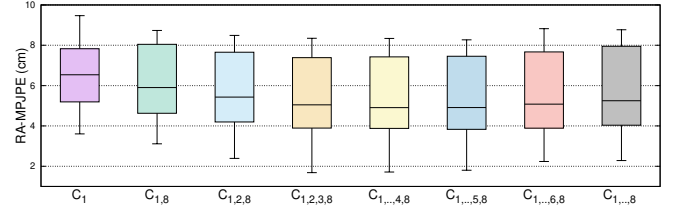


Fig. 4: RA-MPJPE of our method using different camera configurations.

on the motion sequences performed on a plane surface, our method can still recover the global motion on terrain with a small slope.

3) *Smartphone dataset:* We also show the qualitative results of our method on our smartphone dataset. We show the cropped images, the overlaid estimated mesh and the 3D reconstruction of the subject and the cameras in Fig. 8. Our method accurately reconstructed the subject's pose playing football and the camera motion along the wall of the playing arena. We see that the overlays are near-perfectly aligned with the images, showing the accurate reconstruction of the relative pose of the camera and the person. The complete reconstruction is shown in the bottom-right image, and we see that the subject's motion and the camera motion are temporally and spatially coherent.

V. LIMITATIONS

Our method assumes a planar ground surface and human motions which do not involve moving ground (e.g. a treadmill) or sliding motions (e.g. skating, skiing, etc.). However, it can be extended for non-planar ground surfaces by encoding the surface, articulated poses and global poses together in a prior. A major limitation in training such a prior network is the unavailability of human mocap data where the ground surface is also captured. Moreover, most existing datasets are collected with a human subject moving on a planar ground surface.

VI. CONCLUSIONS AND FUTURE WORK

In this paper, we presented a method to reconstruct the 3D human poses, shape, and camera poses relative to a global coordinate frame using synchronized RGB videos from single/multiple extrinsically uncalibrated cameras. We use the ground plane as the reference coordinate system and train a human motion prior using a large amount of human mocap data. We use the latent space of this motion prior to fit the SMPL body model to the 2D keypoints on all the views simultaneously. We show our results on two existing dataset and one new dataset that we collect using smartphones. We show that our method reconstructed the human poses, shape, and camera poses on all the three datasets. We showed the quantitative results on the RICH dataset, demonstrating that our method achieves more accurate results compared to a state-of-the-art method on the task of monocular human motion reconstruction. We also analyzed the effects of multiple views on our method's performance. We show that our method works for diverse types of camera views by showing qualitative results on all the three datasets. The accurate reconstruction by our method on the smartphone data is evidence of the ease of

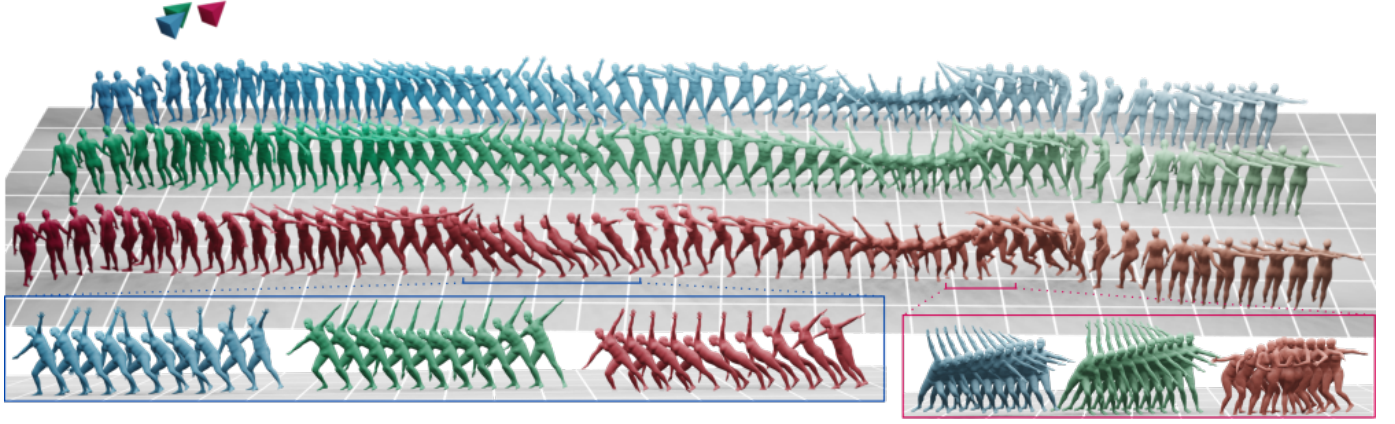


Fig. 5: Comparison: our method (blue), HuMoR [5] (red) and GT (green). For clarity, the human poses are shifted sideways using time-dependent offsets.



Fig. 6: Comparison of GLAMR [21] and our method using a single camera on the RICH dataset. We show images at 4 time instants, each containing the original image, the GLAMR results (green) and results from our method (cyan) using a single camera.

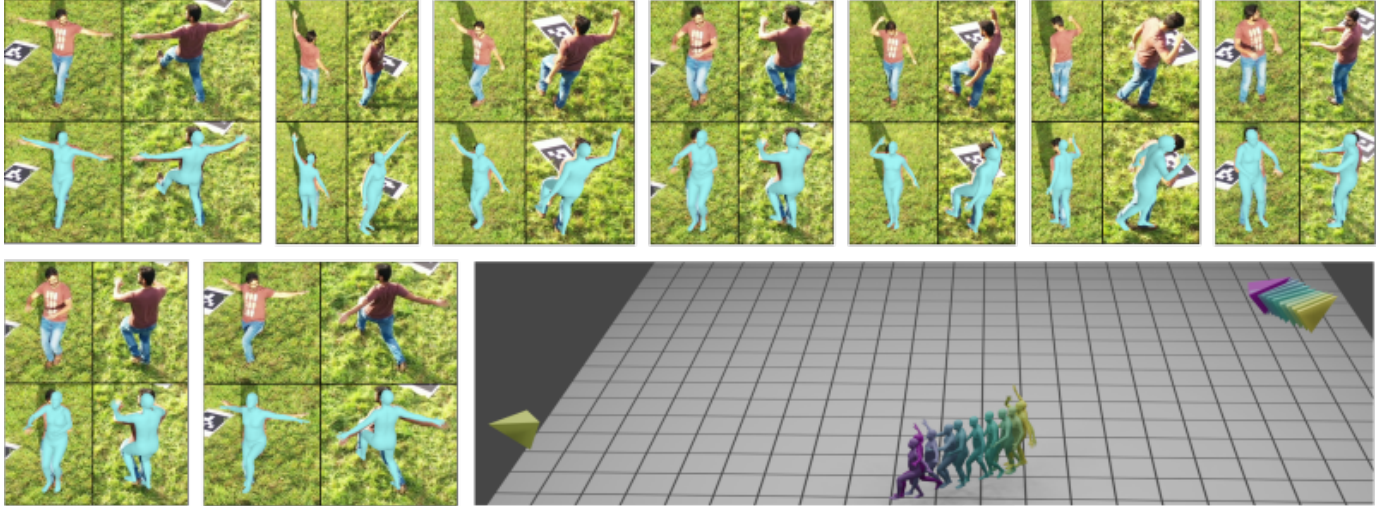


Fig. 7: Results on the Airpose real-world dataset. The result at each time instant is shown using an image grid of 2x2. The top row shows the cropped region of the actual image, while the bottom row shows the estimated mesh overlaid on top of it. Each column corresponds to each camera view. The bottom-right image shows the full 3D reconstruction of the subject's poses, shape, and the camera poses.

use of our method. Our future work includes usage of synthetic data with physics to generate training dataset.

Disclosure: Michael J. Black has received research gift funds from Intel, Nvidia, Adobe, Facebook, and Amazon. While he is a part-time employee of Amazon, his research was performed solely at, and funded solely by, MPI. His has financial interests in Amazon and Meshcapade GmbH.

REFERENCES

- [1] N. Saini, E. Price, R. Tallamraju, R. Enficiaud, R. Ludwig, I. Martinović, A. Ahmad, and M. Black, “Markerless outdoor human motion capture using multiple autonomous micro aerial vehicles,” in *Proceedings 2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, IEEE, Oct. 2019, pp. 823–832.
- [2] G. Pavlakos, V. Choutas, N. Ghorbani, T. Bolkart, A. A. A. Osman, D. Tzionas, and M. J. Black, “Expressive body capture: 3d hands, face, and body from a single image,” in *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [3] R. Tallamraju, N. Saini, E. Bonetto, M. Pabst, Y. T. Liu, M. Black, and A. Ahmad, “Aircaprl: Autonomous aerial human motion capture using deep reinforcement learning,” *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6678 – 6685, Oct. 2020.
- [4] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black, “AMASS: Archive of motion capture as surface shapes,” in *International Conference on Computer Vision*, Oct. 2019, pp. 5442–5451.
- [5] D. Rempel, T. Birdal, A. Hertzmann, J. Yang, S. Sridhar, and L. J. Guibas, “Humor: 3d human motion model for robust pose estimation,” in *International Conference on Computer Vision (ICCV)*, 2021.
- [6] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black, “SMPL: A skinned multi-person linear model,” *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, vol. 34, no. 6, pp. 248:1–248:16, Oct. 2015.
- [7] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh, “Openpose: Realtime multi-person 2d pose estimation using part affinity fields,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [8] M. Kocabas, C.-H. P. Huang, O. Hilliges, and M. J. Black, “PARE: Part attention regressor for 3D human body estimation,” in *Proc. International Conference on Computer Vision (ICCV)*, Oct. 2021, pp. 11 127–11 137.
- [9] K. Isakov, E. Burkov, V. Lempitsky, and Y. Malkov, “Learnable triangulation of human pose,” in *International Conference on Computer Vision (ICCV)*, 2019.

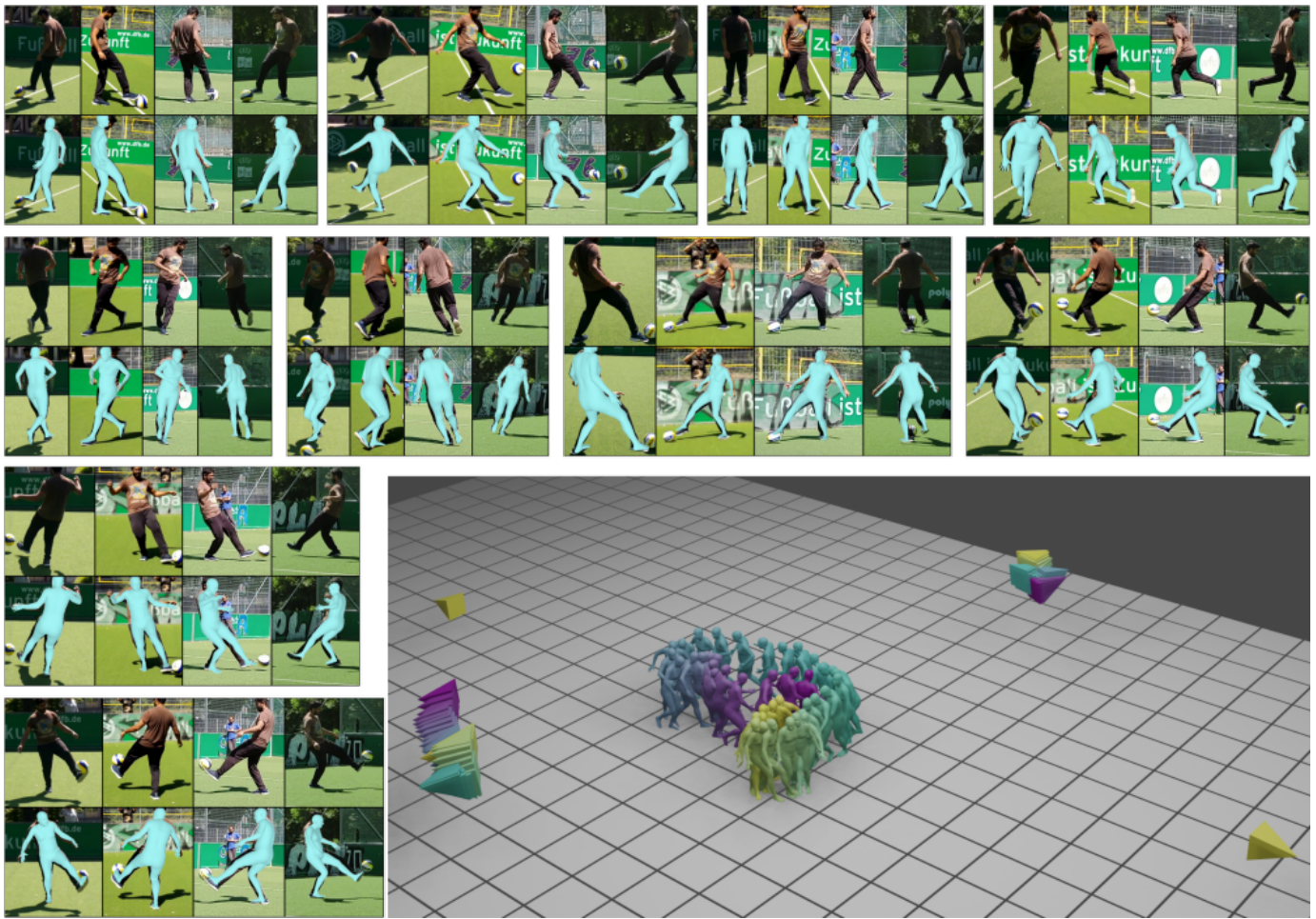


Fig. 8: Results on our smartphone dataset. The result at each time instant is shown using an image grid of 2x4 (check Fig. 7 caption for details).

- [10] Y. Huang, F. Bogo, C. Lassner, A. Kanazawa, P. V. Gehler, J. Romero, I. Akhter, and M. J. Black, “Towards accurate marker-less human shape and pose estimation over time,” in *2017 International Conference on 3D Vision (3DV)*, 2017, pp. 421–430.
- [11] K. Takahashi, D. Mikami, M. Isogawa, and H. Kimata, “Human pose as calibration pattern: 3d human pose estimation with multiple unsynchronized and uncalibrated cameras,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018, pp. 1856–18567.
- [12] B. Huang, Y. Shu, T. Zhang, and Y. Wang, “Dynamic multi-person mesh recovery from uncalibrated multi-view cameras,” in *2021 International Conference on 3D Vision (3DV)*. Los Alamitos, CA, USA: IEEE Computer Society, dec 2021, pp. 710–720.
- [13] R. Bachmann, J. Spörr, P. Fua, and H. Rhodin, “Motion capture from pan-tilt cameras with unknown orientation,” in *2019 International Conference on 3D Vision (3DV)*, 2019, pp. 308–317.
- [14] A. Elhayek, C. Stoll, K. I. Kim, and C. Theobalt, “Outdoor human motion capture by simultaneous optimization of pose and camera parameters,” *Computer Graphics Forum*, vol. 34, no. 6, pp. 86–98, Dec. 2014.
- [15] N. Saini, E. Bonetto, E. Price, A. Ahmad, and M. J. Black, “AirPose: Multi-view fusion network for aerial 3D human pose and shape estimation,” *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4805 – 4812, Apr. 2022.
- [16] N. Hasler, B. Rosenhahn, T. Thormahlen, M. Wand, J. Gall, and H.-P. Seidel, “Markerless motion capture with unsynchronized moving cameras,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 224–231.
- [17] P. Guan, A. Weiss, A. O. Bălan, and M. J. Black, “Estimating human shape and pose from a single image,” in *2009 IEEE 12th International Conference on Computer Vision*, 2009, pp. 1381–1388.
- [18] D. Anguelov, P. Srinivasan, D. Koller, S. Thrun, J. Rodgers, and J. Davis, “SCAPE,” *ACM Transactions on Graphics*, vol. 24, no. 3, pp. 408–416, Jul. 2005.
- [19] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik, “End-to-end recovery of human shape and pose,” in *Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [20] D. F. Henning, T. Laidlow, and S. Leutenegger, “Bodyslam: Joint camera localisation, mapping, and human motion tracking,” *arXiv preprint arXiv: 2205.02301*, 2022.
- [21] Y. Yuan, U. Iqbal, P. Molchanov, K. Kitani, and J. Kautz, “Glamr: Global occlusion-aware human mesh recovery with dynamic cameras,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [22] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li, “On the continuity of rotation representations in neural networks,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5738–5746.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- [24] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” in *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2014.
- [25] H. Fu, C. Li, X. Liu, J. Gao, A. Celikyilmaz, and L. Carin, “Cyclical annealing schedule: A simple approach to mitigating kl vanishing,” 2019. [Online]. Available: <https://arxiv.org/abs/1903.10145>
- [26] I. Kissos, L. Fritz, M. Goldman, O. Meir, E. Oks, and M. Kliger, *Beyond Weak Perspective for Monocular 3D Human Pose Estimation*, 01 2020, pp. 541–554.
- [27] C.-H. P. Huang, H. Yi, M. Höschle, M. Safroshkin, T. Alexiadis, S. Polikovsky, D. Scharstein, and M. J. Black, “Capturing and inferring dense full-body human-scene contact,” in *Proceedings IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, Jun. 2022, pp. 13 274–13 285.
- [28] “Opencamera-sensors android application,” https://f-droid.org/en/packages/com.opencamera_sensors.app/.
- [29] “Distance between rotations,” <http://boris-belousov.net/2016/12/01/quat-dist/>.