

Dynamic Multi-Person Mesh Recovery From Uncalibrated Multi-View Cameras

Buzhen Huang Yuan Shu Tianshu Zhang Yangang Wang*

Southeast University, China

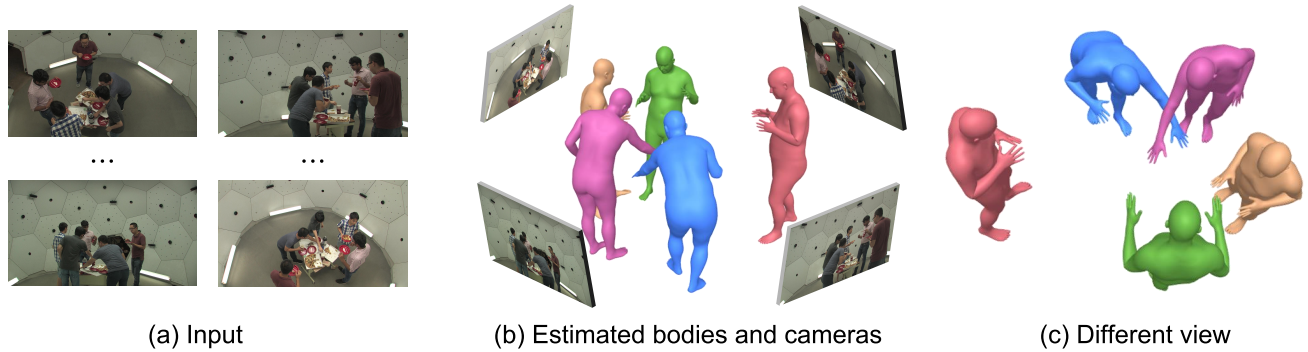


Figure 1: Given multi-person video sequences from sparse uncalibrated cameras, our method simultaneously recovers human motions and extrinsic camera parameters from noisy human semantics.

Abstract

Dynamic multi-person mesh recovery has been a hot topic in 3D vision recently. However, few works focus on the multi-person motion capture from uncalibrated cameras, which mainly faces two challenges: the one is that inter-person interactions and occlusions introduce inherent ambiguities for both camera calibration and motion capture; The other is that a lack of dense correspondences can be used to constrain sparse camera geometries in a dynamic multi-person scene. Our key idea is incorporating motion prior knowledge into simultaneous optimization of extrinsic camera parameters and human meshes from noisy human semantics. First, we introduce a physics-geometry consistency to reduce the low and high frequency noises of the detected human semantics. Then a novel latent motion prior is proposed to simultaneously optimize extrinsic camera parameters and coherent human motions from slightly noisy inputs. Experimental results show that accurate camera parameters and human motions can be obtained through one-stage optimization. The codes will be publicly available at <https://www.yangangwang.com>.

*Corresponding author. E-mail: yangangwang@seu.edu.cn. All the authors from Southeast University are affiliated with the Key Laboratory of Measurement and Control of Complex Systems of Engineering, Ministry of Education, Nanjing, China.

1. Introduction

Recovering multiple human motions from video is essential for many applications, such as social behavior understanding, sports broadcasting, virtual reality applications, etc. Numerous previous works have been aimed at capturing multi-person motions from multi-view input via geometry constraints [2, 16, 9, 38, 62, 29] or optimization-based model fitting [61, 35, 40, 34, 59]. While these works have made remarkable advances in multi-person motion capture, they all rely on accurate calibrated cameras to build view-view and model-view consistency. Few works focus on multi-person motion capture from uncalibrated cameras. [47] constructs a two-stage framework that first calibrates the camera using the static geometry from the background and then generates 3D human models from dynamic object reconstruction and segmentations. [17] utilizes the similarity of the estimated 3D poses in each view to find pose pairs and refines them in the global coordinate system. However, these methods require a large space distance among the target people and can not capture interactive human bodies.

In this paper, we address the problem of directly recovering multiple human bodies with unknown extrinsic camera parameters. There are two main challenges. The first one is that inter-person interactions and occlusions introduce inherent ambiguities for both camera calibration and motion reconstruction. The ambiguous low-level vi-

sual features lead to severe low and high frequency noises in detected human semantics (*e.g.*, 2D pose [3], appearance [35]), which causes extreme difficulty in establishing view-view and model-view consistency. The other is that a lack of sufficient local image features (*e.g.*, SIFT [43]) can be used to constrain sparse camera geometries in a dynamic multi-person scene.

To tackle the obstacles, our key-idea is to **use motion prior knowledge to assist the simultaneous recovery of camera parameters and dynamic human meshes from noisy human semantics**. We introduce a physics-geometry consistency to reduce the low and high-frequency noises of the detected multi-person semantics. Then a latent motion prior is proposed to recover multiple human motions with extrinsic camera parameters from partial and slightly noisy multi-person 2D poses. As shown in Fig.2, the multi-view 2D poses from off-the-shelf 2D pose detection [18, 7] and tracking [66] contain high-frequency 2D joint jitter and low-frequency identity error. Without proper camera parameters, we can not filter out the noises by epipolar constraint [2, 9]. However, we found that the triangulated skeleton joint trajectories are continuous, even though the camera parameters are inaccurate. Based on this observation, we propose a physics-geometry consistency and construct a convex optimization to combine kinetic energy prior and epipolar constraint to reduce the high and low frequency noises.

Simultaneously optimizing extrinsic camera parameters and multi-person motions from the filtered and slightly noisy 2D poses is a highly non-convex problem. We then introduce a compact latent motion prior to jointly recover temporal coherent human motions and accurate camera parameters. We adopt a variational autoencoder [30] (VAE) architecture for our motion prior. Different from existing VAE-based motion models [41, 44, 39], we use bidirectional GRU [10] as backbone and design a latent space both considering local kinematics and global dynamics. Therefore, our latent prior can be trained on a limited amount of short motion clips [45] and be used to optimize long sequences. While the motion prior can generate diverse and temporal coherent motions, it is not robust to noises in motion optimization. We found that linearly interpolating the latent code of VPoser [48] will produce consecutive poses. Inspired by this, we propose a local linear constraint on motion latent code in model training and optimization. This constraint ensures motion prior to produce coherent motions from noisy input. In addition, to keep local kinematics, a skip-connection between explicit human motion and latent motion code is incorporated in the model. Using the noisy 2D poses as constraints, we can recover human motions and camera parameters by simultaneously optimizing the latent code and cameras.

The main contributions of this work are summarized as

follows.

- We propose a framework that directly recovers multi-person human motions with accurate extrinsic camera parameters from sparse multi-view cameras.
- We propose a physics-geometry consistency to reduce the notorious low and high frequency noises in detected human semantics.
- We propose a human motion prior that contains both local kinematics and global dynamics, which can be trained on limited short motion clips and be used to optimize temporal coherent long sequences.

2. Related Work

Multi-view Human pose and shape estimation. Reconstructing human pose and shape from multi-view inputs has been a long-standing problem in 3D vision. [40] reconstructs interactive multi-person with manually specified masks. To avoid manual operations, the color [46, 59], appearance [35], location [34] and other cues of human are utilized to build the spatio-temporal correspondences, thus realizing optimization-based model fitting. In contrast, [2, 3, 38, 62, 6, 29] firstly establish view-view correspondences via detected 2D poses and geometric constraints and then reconstruct through triangulation or optimization. [16] considers geometric and appearance constraints simultaneously. However, these methods all rely on accurate camera parameters. Besides, 2D poses and appearance can be easily affected by partial occlusion, which is very common in multi-person interaction sceneries.

To recover multiple human meshes from uncalibrated cameras, [47] first calibrates the camera using the static geometry from the background and then generates 3D human models from dynamic object reconstruction. [17] realizes reconstruction via the similarity of the detected 3D poses from different views. However, these methods require a large space distance among the target people and can not capture interactive human bodies.

Extrinsic camera calibration. Conventional camera calibration methods rely on specific tools (*e.g.*, checkerboard[63] and one-dimensional objects[64]). Except for the complex calibration process, it leads to two separate stages for calibration and reconstruction. [26, 47, 69] propose more convenient methods that directly use image features from static background (*e.g.*, SIFT [43]) to calibrate the camera. However, the dynamic human bodies occupy the most proportion of the image pixels in multi-person scenarios. To handle this obstacle, [50, 12, 8, 50, 13] obtain structure cues and estimate camera parameters from the semantics of the scene (*e.g.*, lines of the basketball court). [24, 55] estimate the extrinsic camera parameters from the tracked human trajectories in more general multi-person scenes. [52, 4, 5] extract frontier points of the silhouette and recover epipolar geometry by using points between

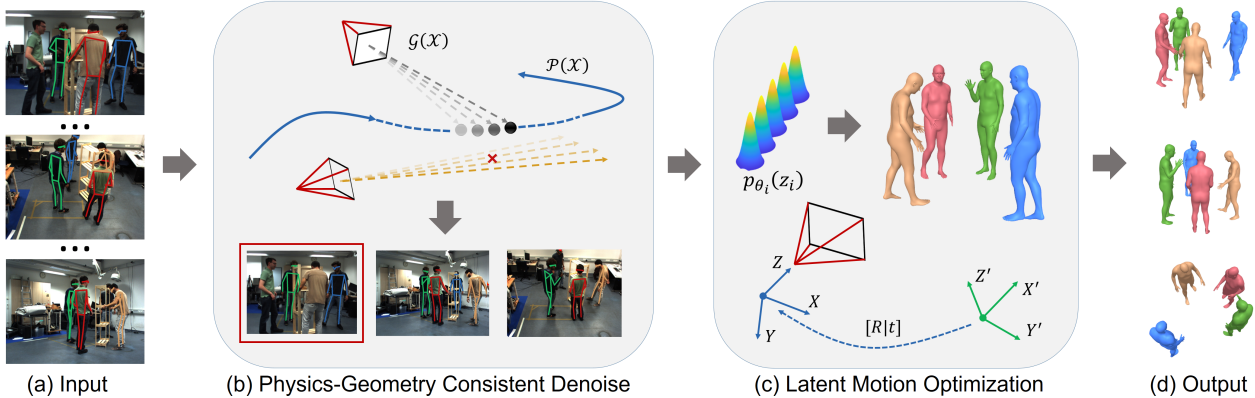


Figure 2: Overview of our method. Since directly optimizing cameras and human motions from noisy detections (a) always lead to suboptimal solutions, we first introduce a physics-geometry consistency (b) to reduce high and low frequency noises in the detected human semantics. Then, to recover from the filtered partial and slightly noisy inputs (b), we incorporate a novel latent motion prior to the optimization framework (c) to obtain accurate camera parameters and coherent human motions (d).

different perspectives. Nevertheless, getting accurate human segmentations from in-the-wild images itself is a challenging problem. [15] realizes camera calibration by using the depth camera in an indoor scene to extract the skeleton. [49, 20, 54] and [21] use detected human 2D joints and mesh respectively to calibrate the camera, further simplifying the calibration device. State-of-the-art 2D/3D pose estimation frameworks [18, 7, 32] can hardly get accurate 2D/3D keypoints in multi-person scenes, and such methods cannot be directly applied to multi-person cases. To reduce the ambiguities generated by human interactions and occlusions, we propose a physics-geometry consistent denoising framework and a robust latent motion prior to remove the noises, realizing multi-person reconstruction and extrinsic camera calibration in an end-to-end way.

Motion prior. Traditional marker-less motion capture relies on massive views to provide sufficient visual cues [29, 57, 14]. To reconstruct from sparse cameras, [67, 35] employ the euclidean distance of poses in adjacent frames as the regularization term, which may limit the dynamics of the reconstructed motions. Thus, applying strong and compact motion prior in motion capture has attracted wide attention. The simple and feasible motion priors (*e.g.*, Principal Component Analysis [51], Low-dimensional Non-linear Manifolds [27, 19]) lack expressiveness and are not robust to noises. Historically, Gaussian Process Latent Variable Model (GPLVM) [33, 60, 37, 36] succeed in modeling human motions [58, 56] since it takes uncertainties into account, but is difficult to make a smooth transition among mixture models. [25] uses low-dimensional Discrete Cosine Transform (DCT) basis [1] as the temporal prior to capture human motions. With the development of deep learning, VIBE [31] trains a discriminator to determine the quality of motion, but one-dimensional variables can hardly describe dynamics. [41] and [44, 65] train VAEs based on Temporal Convolutional Networks(TCN) and Recurrent Neural

Network(RNN) respectively and represent motion with latent code. However, both of these two methods use latent code in a fixed dimension, which is not suitable for dealing with sequences of varying lengths. [39] constructs a conditional variational autoencoder (cVAE) to represent motions of the two adjacent frames. Although this structure solves the problem of sequence length variation, it can only model sequence information of the past, which is not suitable for optimizing the whole sequence.

In this paper, we propose a motion prior that contains local kinematics and global dynamics of the motion. The structure of the model makes it is suitable for large-scale variable-length sequence optimization.

3. Method

Our goal is to recover both multi-person motions and extrinsic camera parameters simultaneously from multi-view videos. Firstly, we propose a physics-geometry consistency to reduce the high and low frequency noises in the detected human semantics (Sec.3.2). Then, we introduce a robust latent motion prior (Sec.3.3), which contains human dynamics and kinematics, to assist estimation from noisy inputs. Finally, with the trained motion prior, we design an optimization framework to recover accurate extrinsic camera parameters and human motions from multi-view uncalibrated videos (Sec.3.4).

3.1. Preliminaries

Human motion representation. We adopt SMPL [42] to represent human motion, which consists of the shape $\beta \in \mathbb{R}^{10}$, pose $\theta \in \mathbb{R}^{72}$ and translation $\mathcal{T} \in \mathbb{R}^3$. To generally learn human dynamics and kinematics from training data, we separate global rotation $\mathcal{R} \in \mathbb{R}^{T \times 3}$, translation \mathcal{T} and human shape β when constructing the motion prior. Moreover, we use the more appropriate continuous 6D rotation representation [68] for the prior. Finally, a motion that

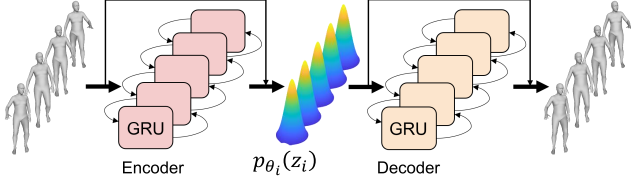


Figure 3: The motion prior is a symmetrical encoder-decoder network, which compactly models human dynamics and kinematics. The prior can be trained on short clips and be used to fit long sequences.

contains T frames is represented as $\mathcal{X} \in \mathbb{R}^{T \times 138}$.

2D pose detection and camera initialization. We first use off-the-shelf 2D pose estimation [18] and tracking framework [66] to get tracked 2D poses for each person. Then, we estimate initial camera extrinsic parameters for the denoising framework Sec.3.2. We obtain the fundamental matrix from multi-view 2D poses in the first frame using epipolar geometry with known intrinsic parameters. Then the initial extrinsic parameters can be decomposed from it. Since the 2D poses are noisy, a result selection is used to ensure robustness. The details can be found in the Sup. Mat.

3.2. Physics-geometry Consistent Denoising

Due to the inherent ambiguities in inter-person interactions and occlusions, state-of-the-art pose detection and tracking methods [18, 7, 53, 66] can hardly get the precise 2D poses with accurate identity from in-the-wild videos. The drift and jitter generated by pose detection are often high-frequency, while identity error generated by pose tracking is low-frequency. The mixture of the two types of noises is notorious in multi-person mesh recovery. To solve this obstacle, we propose a physics-geometry consistency to reduce both high and low frequency noises in 2D poses from each view.

Supposing the target person is detected in V views, our goal is to remove the noisy detections that do not satisfy the physics-geometry consistency. Theoretically, despite that the camera parameters are not accurate, the triangulated skeleton joint trajectories from 2D poses with accurate identity are continuous. So we first utilize a set of optical rays, which come from the optical center of the camera and pass through corresponding 2D joint coordinates, to construct a physical constraint. For view i , the ray in the plücker coordinates is represented as (n_i, l_i) . Given the skeleton joint positions of the previous frame x_{t-1} , the optical rays should be close to x_{t-1} . We represent the distance between x_{t-1} and the rays as:

$$\mathcal{L}_p^i = \|x_{t-1} \times n_i - l_i\|. \quad (1)$$

The rays generated by the wrong detection will produce an out-of-range physical cost \mathcal{L}_p . However, with only the above physical constraint, the system may get the wrong results in inter-person occlusion cases. Consequently, we fur-

ther propose an additional geometric constraint. We enforce the rays from view i and view j to be coplanar precisely:

$$\mathcal{L}_g^{i,j} = n_i^T l_j + n_j^T l_i. \quad (2)$$

We combine these two constraints as the physics-geometry consistency. We then follow [23] to filter out incorrect detections with the physics-geometry consistency. The physical cost and geometric cost of different views are represented in matrices \mathcal{P} and \mathcal{G} .

$$\begin{cases} \mathcal{P}_{i,j} = \mathcal{L}_p^i + \mathcal{L}_p^j \\ \mathcal{G}_{i,j} = \mathcal{L}_g^{i,j} \end{cases}, \quad (3)$$

where $\mathcal{P}_{i,j}$ and $\mathcal{G}_{i,j}$ are physical cost and geometric cost of view i and view j . We use a positive semidefinite matrix $\mathcal{M} \in \{0, 1\}^{v \times v}$ to represent the correctness of correspondences among different views. Our goal is to solve \mathcal{M} , which minimizes the physics-geometry consistency cost:

$$\arg \min_{\mathcal{M}} f(\mathcal{M}) = -c_g \langle \mathcal{G}, \mathcal{M} \rangle - c_p \langle \mathcal{P}, \mathcal{M} \rangle, \quad (4)$$

where c_g, c_p are 0.7 and 0.3 in our experiment. $\langle \cdot \rangle$ denotes the hadamard product. Finally, we use the estimated \mathcal{M} to extract accurate detections.

The skeleton joint position of the start frame x_0 is triangulated with the queries of pose tracking [66]. We triangulate x_t with filtered results and use it to calculate the physical consistency cost in the next frame. The filtered 2D poses will be used in Eqn.(13) to find optimal motions. More details can be found in Sup. Mat.

3.3. Latent Motion Prior

Simultaneous optimization of multi-person motions and camera parameters from slightly noisy 2D poses is a highly non-convex problem and is likely to fall into the local minima. To address this challenge, we design a compact VAE-based latent motion prior to obtain accurate and temporal coherent motions. The prior has three strengths. 1) It contains compact dynamics and kinematics to reduce computational complexity. 2) It can be trained on short motion clips and applied to long sequence fitting. 3) The latent local linear constraint ensures robustness to noisy input. The details are described as following.

Model architecture. Our network is based on VAE [30], which shows great power in modeling motions [39, 44]. As shown in Fig.3, the encoder consists of a bidirectional GRU, a mean and variance encoding network with a skip-connection. The decoder has a symmetric network structure. Different from previous work [39], the bidirectional GRU ensures that the prior is able to see all the information from the entire sequence and that the latent code can represent global dynamics. However, the latent prior encoded only by features extracted from GRU is difficult to

reconstruct accurate local fine-grained poses when used for large-scale sequence optimization. Thus, we construct a skip-connection for the encoder and decoder, respectively, allowing the latent prior to accurately capture the refined kinematic poses and the global correlation between them. Besides, we design the latent code $\mathbf{z} \in \mathbb{R}^{T \times 32}$ whose frame length T is corresponding to the input sequence. Thus, our prior can be trained on a limited amount of short motion clips [45] and be applied to long sequence fitting.

Training. In the training phase, a motion \mathcal{X} is fed into the encoder to generate mean $\mu(\mathcal{X})$ and variance $\sigma(\mathcal{X})$. The sampled latent code $\mathbf{z} \sim q_\phi(\mathbf{z} | \mu(\mathcal{X}), \sigma(\mathcal{X}))$ is then decoded to get the reconstructed motion $\hat{\mathcal{X}}$. The reparameterization trick [30] is adopted to achieve gradient backpropagation. We train the network through maximizing the Evidence Lower Bound (ELBO):

$$\log p_\theta(\mathcal{X}) \geq \mathbb{E}_{q_\phi} [\log p_\theta(\mathcal{X} | \mathbf{z})] - D_{\text{KL}}(q_\phi(\mathbf{z} | \mathcal{X}) \| p_\theta(\mathbf{z})). \quad (5)$$

The specific loss function is:

$$\mathcal{L}_{vae} = \mathcal{L}_{6d} + \mathcal{L}_v + \mathcal{L}_{kl} + \mathcal{L}_{\text{linear}} + \mathcal{L}_{\text{reg}}, \quad (6)$$

where \mathcal{L}_{6d} and \mathcal{L}_v are:

$$\mathcal{L}_{6d} = \sum_{t=1}^T \left\| \mathcal{X}_t - \hat{\mathcal{X}}_t \right\|^2, \quad (7)$$

$$\mathcal{L}_v = \sum_{t=1}^T \left\| \mathcal{V}_t - \hat{\mathcal{V}}_t \right\|^2, \quad (8)$$

where \mathcal{V}_t is the deformed SMPL vertices of frame t . This term guarantees that the prior learns high fidelity local details.

$$\mathcal{L}_{kl} = KL(q(\mathbf{z} | \mathcal{X}) \| \mathcal{N}(0, I)), \quad (9)$$

which enforces its output to be near the Gaussian distribution. The regularization term, which ensures the network will not be easily overfitted:

$$\mathcal{L}_{\text{reg}} = \|\phi\|_2^2. \quad (10)$$

Although applying the above constraints can produce diverse and temporal coherent motions, it is not robust to noisy 2D poses. The jitter and drift of 2D poses and identity error will result in an unsmooth motion. Inspired by the interpolation of VPoser [48], we add a local linear constraint to enforce a smooth transition on latent code:

$$\mathcal{L}_{\text{linear}} = z_{t+1} - 2z_t + z_{t-1}. \quad (11)$$

When the motion prior is applied in long sequence fitting, the parameters of the decoder are fixed. The latent code is decoded to get the motion $\hat{\mathcal{X}} \in \mathbb{R}^{T \times 138}$.

3.4. Joint Optimization of Motions and Cameras

Optimization variables. Different from traditional structure-from-motion (SFM), which lacks structural constraints between 3D points and is not robust to noisy input. We directly optimize the motion prior, so that the entire motions are under inherent kinematic and dynamic constraints. The optimization variables of V views videos that contain N people are $\{(\beta, \mathbf{z}, \mathcal{R}, \mathcal{T})_{1:N}, \mathcal{E}_{1:V}\}$. The $\mathcal{E} \in \mathbb{R}^6$ is camera extrinsic parameter that contains rotation and translation.

Objective. We formulate the objective function as following:

$$\arg \min_{(\beta, \mathbf{z}, \mathcal{R}, \mathcal{T})_{1:N}, \mathcal{E}_{1:V}} \mathcal{L} = \mathcal{L}_{\text{data}} + \mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{pen}}, \quad (12)$$

where the data term is:

$$\mathcal{L}_{\text{data}} = \sum_{v=1}^V \sum_{n=1}^N \sigma_v^n \rho(\Pi_{\mathcal{E}_v}(\mathbf{J}^n) - \mathbf{p}_v^n) \quad (13)$$

where ρ is the robust Geman-McClure function [22]. \mathbf{p} , σ are the filtered 2D poses and its corresponding confidence. \mathbf{J} is the skeleton joint position generated by model parameters.

Besides, the regularization term is:

$$\mathcal{L}_{\text{prior}} = \sum_{n=1}^N \|\mathbf{z}_n\|^2 + \sum_{n=1}^N \|\beta_n\|^2 + \sum_{n=1}^N \mathcal{L}_{\text{linear}}. \quad (14)$$

$\mathcal{L}_{\text{linear}}$ is the same as Eqn.(11). We further apply a collision term based on differentiable Signed Distance Field (SDF) [28] to prevent artifacts generated from multi-person interactions.

$$\mathcal{L}_{\text{pen}} = \sum_{j=1}^N \sum_{i=1, i \neq j}^N \sum_{vt \in \mathcal{V}_j} -\min(\text{SDF}_i(vt), 0), \quad (15)$$

where $\text{SDF}(vt)$ is the distance from sampled vertex vt to the human mesh surface.

4. Experiments

In this section, we conduct several evaluations to demonstrate the effectiveness of our method. The comparisons in Sec.4.1 show that our method can recover multiple human bodies from uncalibrated cameras and achieves state-of-the-art. Then, we prove that the accurate extrinsic camera parameters can be obtained from joint optimization. Finally, several ablations in Sec.4.3 are conducted to evaluate key components. The details of the datasets that are used for training and testing can be found in the Sup. Mat.

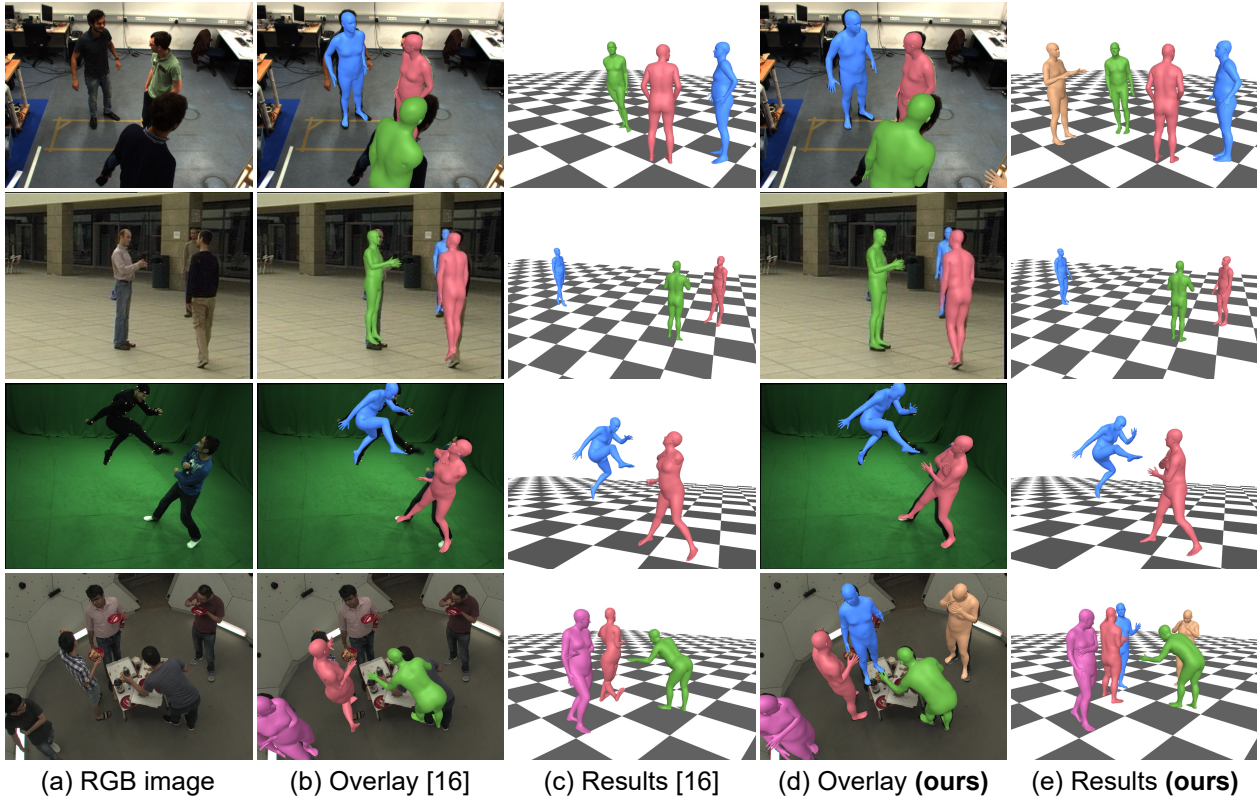


Figure 4: Qualitative comparison with [16]. Due to the mismatched 2D pose and a lack of prior knowledge, [16] fails on these cases while our method obtains accurate results with the proposed motion prior and physics-geometry consistency.

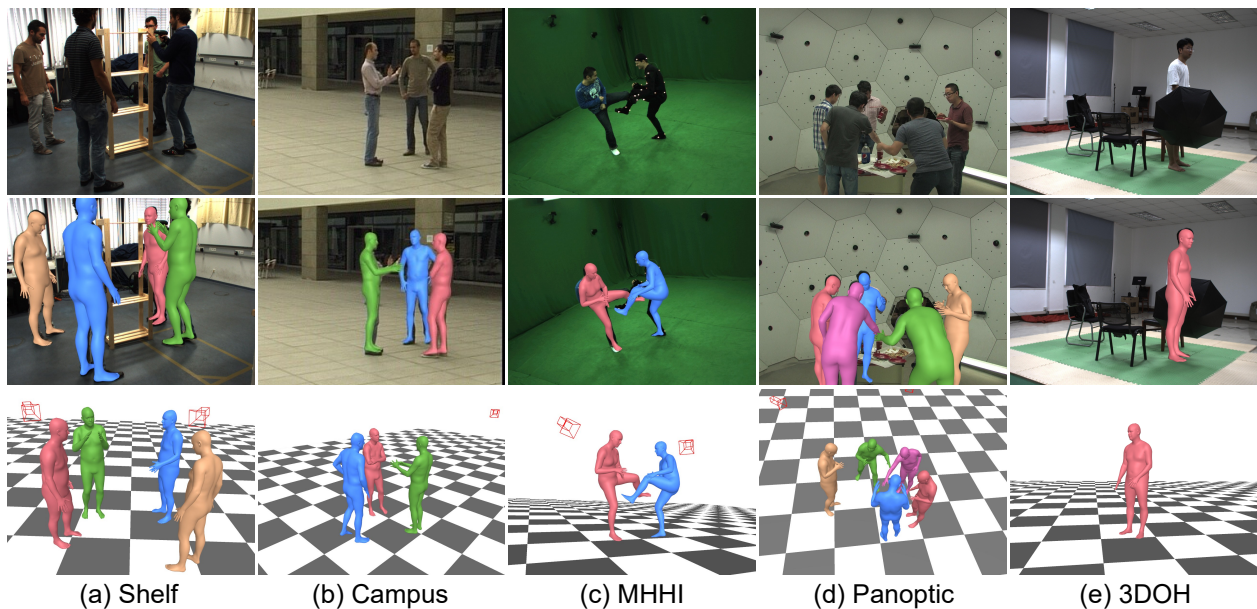


Figure 5: The estimated results on different datasets. Our methods can obtain temporal coherent human motions and accurate extrinsic cameras parameters simultaneously from multi-view uncalibrated videos.

4.1. Multi-person Motion Capture

We first conducted qualitative and quantitative comparisons on Campus and Shelf datasets. To the best of our knowledge, no method has ever recovered human meshes

on these datasets. We compared several baseline methods that regress 3D poses. [2] and [3] introduce 3D pictorial structure for multi-person 3D pose estimation from multi-view images and videos respectively. [6, 16, 9, 62, 11] are recent works based on calibrated cameras. The quantita-

Method	Campus			Shelf		
	A1	A2	A3	A1	A2	A3
Belagiannis <i>et al.</i> [2]	82.0	72.4	73.7	66.1	65.0	83.2
Belagiannis <i>et al.</i> [3]	93.5	75.7	85.4	75.3	69.7	87.6
Bridgeman <i>et al.</i> [6]	91.8	92.7	93.2	99.7	92.8	97.7
Dong <i>et al.</i> [16]	97.6	93.3	98.0	98.9	94.1	97.8
Chen <i>et al.</i> [9]	97.1	94.1	98.6	99.6	93.2	97.5
Zhang <i>et al.</i> [62]	-	-	-	99.0	96.2	97.6
Chu <i>et al.</i> [11]	98.4	93.8	98.3	99.1	95.4	97.6
VPoser-t [48]	97.3	93.5	98.4	99.8	94.1	97.5
Ours	97.6	93.7	98.7	99.8	96.5	97.6

Table 1: Comparison with baseline methods that estimate multi-person 3D poses. The numbers are the percentage of correctly estimated parts (PCP). The proposed method achieves state-of-the-art on some metrics. VPoser-t is a combination of VPoser [48].

tive results shown in Tab.1 demonstrate that our method achieves state-of-the-art on Campus and Shelf datasets in terms of PCP. Since only a few works target to multi-person mesh recovery task from multi-view input, we compared with EasyMocap* which fits SMPL model to the 3D pose estimated by [16]. Row 2 and row 4 of Fig.4 show that [16] produces the wrong result due to partial occlusion, while our method generates accurate poses with physics-geometry consistency. Besides, our method obtains more natural and temporal coherent results even for challenging poses since the proposed motion prior provides local kinematics and global dynamics.

We then evaluated our method on MHHI dataset. [40, 34, 35] can reconstruct closely interacting multi-person meshes from multi-view input, but all these works rely on accurate calibrated camera parameters. We conducted quantitative comparisons with these methods in Tab.2. The numbers are the mean distance with standard deviation between the tracked 38 markers and its paired 3D vertices in *mm*. In the single-view case, since the motion prior provides additional prior knowledge, our method generates far more accurate results than [34]. In addition, the proposed approach achieves competitive results with the least views.

To further demonstrate the effectiveness of the proposed method in single-view occluded situations, we show the qualitative results on 3DOH in Fig.5. Our method can recover complete and reasonable human bodies from partial observation with the local kinematics and global dynamics in the motion prior. More qualitative and quantitative results on single-person datasets can be found in Sup. Mat.

4.2. Camera Calibration Evaluation

We then qualitatively and quantitatively evaluate the estimated camera parameters. Since there exists a rigid transformation between the predicted camera parameters and the ground-truth provided in the datasets, we follow [12] to apply rigid alignment to the estimated cameras. We first com-

*<https://github.com/zju3dv/EasyMocap>

Method	1 view	2 views	4 views	8 views	12 views
Liu <i>et al.</i> [40]	-	-	-	-	51.67
Li <i>et al.</i> [34]	1549.88	242.27	58.42	48.57	43.30
Li <i>et al.</i> [35]	-	63.93	37.88	32.73	30.35
VPoser-t [48]	158.33	60.02	38.46	32.11	31.48
Ours	140.96	58.04	37.86	30.92	29.83

Table 2: Quantitative comparison with multi-person mesh recovery methods on MHHI dataset. The numbers are the mean distance with standard deviation between markers and its paired 3D vertices in *mm*.

Method	Panoptic Dataset			Shelf Dataset		
	Pos.	Ang.	Reproj.	Pos.	Ang.	Reproj.
PhotoScan	505.02	35.29	188.18	-	-	-
initial	3358.51	44.30	637.21	1532.42	26.86	79.34
w/o P-G consis. + <i>opt cam.</i>	178.78	1.10	23.00	29.09	0.68	18.88
VPoser-t [48] + <i>opt cam.</i>	118.88	0.64	22.76	34.30	0.59	18.83
MotionPrior + <i>opt cam.</i>	101.25	0.59	22.69	23.18	0.52	18.70

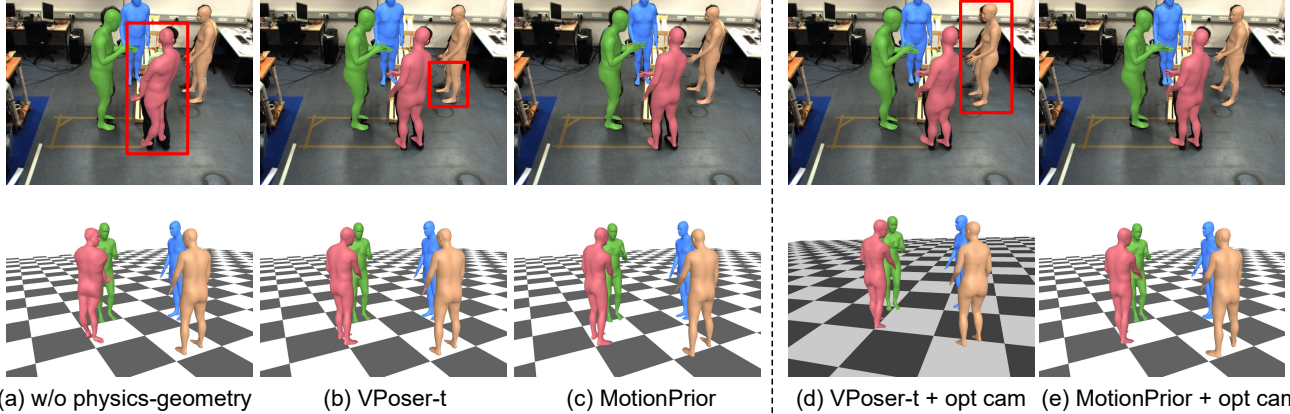
Table 3: Evaluation of the estimated camera. The *Pos.* and *Ang.* are position error and angle error between predicted cameras and ground-truth camera parameters. The units are *mm* and *deg*, respectively. The *Reproj.* is re-projection error in pixel. The initial is the coarse camera parameters estimated from Sec.3.1. + *opt cam.* denotes simultaneously optimize cameras and human motions.

Method	MHHI		Shelf		
	Mean	Std	A1	A2	A3
VPoser-t [48]	31.48	11.54	99.8	94.1	97.5
w/o P-G consis.	32.31	12.17	92.4	89.8	91.6
w/o local linear	30.25	11.07	99.8	95.4	97.3
MotionPrior	29.83	9.87	99.8	96.5	97.6
VPoser-t [48] + <i>opt cam.</i>	43.72	19.57	97.4	89.7	89.7
w/o P-G consis.+ <i>opt cam.</i>	49.34	24.37	91.5	86.7	88.6
w/o local linear + <i>opt cam.</i>	35.25	17.07	97.5	90.4	93.3
MotionPrior + <i>opt cam.</i>	34.44	10.57	98.4	91.5	94.4

Table 4: Ablation on physics-geometry consistency and our motion prior. *opt cam.* denotes simultaneously optimize cameras and human motions.

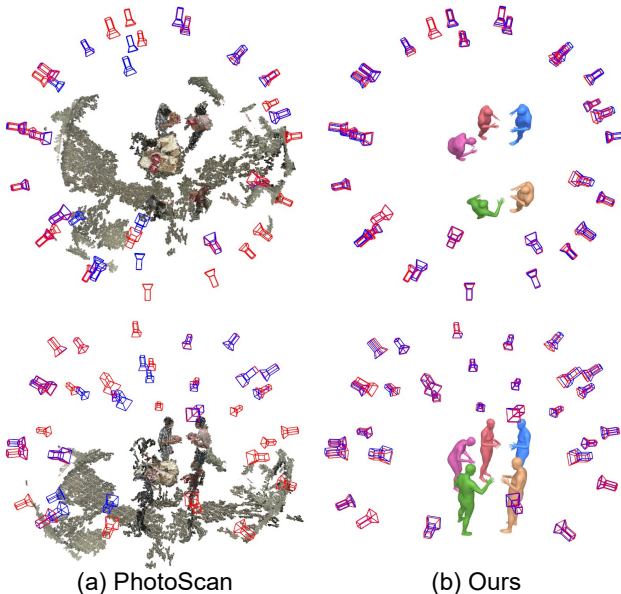
pared with PhotoScan[†], which is a commercial software that reconstructs 3D point clouds and cameras. As shown in Tab.3, PhotoScan fails to work for sparse inputs (Shelf dataset) since it relies on the dense correspondences between each view. We evaluate the results with position error, angle error, and re-projection error. Under relatively massive views, our method outperforms PhotoScan in all metrics. Fig.7 shows the results on Panoptic dataset with 31 views. The cameras in red and blue colors are the ground-truth and the predictions, respectively. PhotoScan only captures part of the cameras with low accuracy. On the contrary, our method successfully estimates all the cameras with complete human meshes. We then compared with the initial extrinsic parameters estimated in Sec.3.1. After joint optimization, the final results gain significant improvement. Our method achieves better performance both from massive and sparse inputs with the physics-geometry consistency and the motion prior.

[†]<https://www.agisoft.com/>



(a) w/o physics-geometry (b) VPoser-t (c) MotionPrior (d) VPoser-t + opt cam (e) MotionPrior + opt cam

Figure 6: Ablation on physics-geometry consistency and our motion prior. Without physics-geometry consistency, it can not obtain accurate motion due to the influence of noises. Since the lack of motion dynamics, the VPoser-t is hard to estimate plausible cameras and motions when the cameras are not provided.



(a) PhotoScan (b) Ours

Figure 7: PhotoScan can not work on sparse inputs. We conducted a comparison with PhotoScan on Panoptic with 31-views input. Our method accurately estimates all camera extrinsic parameters from noisy human semantics, while PhotoScan gets only a part of cameras.

4.3. Ablation Study

Physics-geometry consistency. We conducted ablation on the physics-geometry consistency to reveal its importance of removing the noises in the human semantics. Fig.6 illustrates that without the consistency, the reconstruction is unnatural due to the noisy detections. As shown in Tab.4, without the proposed consistency, the mean distance error of joint optimization increases 12.42, demonstrating its significance.

Motion prior. VPoser-t is a combination of [48] which lacks global dynamics. We first compared it to illustrate the superiority of the proposed motion prior. Tab.4 shows that

the standard variance of our method on MHHI is smaller since the motion prior models the temporal information. Tab.3, Tab.4 and Fig.6 demonstrate that due to the lack of temporal constraints, VPoser-t is more sensitive to the noisy detections. The local linear constraint ensures a smooth transition between each frame of the latent code. We then removed the local linear constraint when training the motion prior. In Tab.4, without local linear constraint, although the mean distance error of joint optimization on MHHI dataset is small, the standard variance of which is large. Thus, the results prove that the constraint is effective in modeling temporal coherent motions.

5. Conclusion

This paper proposes a framework that directly recovers human motions and extrinsic camera parameters from sparse multi-view video cameras. Unlike previous work, which fails to establish view-view and model-view corresponds, we introduce a physics-geometry consistency to reduce the low and high frequency noises of the detected human semantics. In addition, we also propose a novel latent motion prior to jointly optimize camera parameters and coherent human motions from slightly noisy inputs. The proposed method simplifies the conventional multi-person mesh recovery by incorporating the calibration and reconstruction into a one-stage optimization framework.

Acknowledgments. The authors would like to thank Professor Yebin Liu and Professor Kun Li for sharing the data. This work was supported in part by National Key R&D Program of China (No. 2018YFB1403900), in part by National Natural Science Foundation of China (No. 61806054), in part by Natural Science Foundation of Jiangsu Province (No. BK20180355), Young Elite Scientist Sponsorship Program by the China Association for Science and Technology and "Zhishan Young Scholar" Program of Southeast University.

References

- [1] I. Akhter, T. Simon, S. Khan, I. Matthews, and Y. Sheikh. Bilinear spatiotemporal basis models. *ACM Transactions on Graphics (TOG)*, 31(2):1–12, 2012. 3
- [2] V. Belagiannis, S. Amin, M. Andriluka, B. Schiele, N. Navab, and S. Ilic. 3d pictorial structures for multiple human pose estimation. In *CVPR*, 2014. 1, 2, 6, 7
- [3] V. Belagiannis, S. Amin, M. Andriluka, B. Schiele, N. Navab, and S. Ilic. 3d pictorial structures revisited: Multiple human pose estimation. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):1929–1942, 2015. 2, 6, 7
- [4] G. Ben-Artzi, Y. Kasten, S. Peleg, and M. Werman. Camera calibration from dynamic silhouettes using motion barcodes. In *CVPR*, 2016. 2
- [5] E. Boyer. On using silhouettes for camera calibration. In *ACCV*, 2006. 2
- [6] L. Bridgeman, M. Volino, J.-Y. Guillemaut, and A. Hilton. Multi-person 3d pose estimation and tracking in sports. In *CVPR Workshops*, 2019. 2, 6, 7
- [7] Z. Cao, G. Hidalgo, T. Simon, S.-E. Wei, and Y. Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *IEEE transactions on pattern analysis and machine intelligence*, 43(1):172–186, 2019. 2, 3, 4
- [8] J. Chen and J. J. Little. Sports camera calibration via synthetic data. In *CVPR Workshops*, 2019. 2
- [9] L. Chen, H. Ai, R. Chen, Z. Zhuang, and S. Liu. Cross-view tracking for multi-human 3d pose estimation at over 100 fps. In *CVPR*, 2020. 1, 2, 6, 7
- [10] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *EMNLP*, 2014. 2
- [11] H. Chu, J.-H. Lee, Y.-C. Lee, C.-H. Hsu, J.-D. Li, and C.-S. Chen. Part-aware measurement for robust multi-view multi-human 3d pose estimation and tracking. In *CVPR*, 2021. 6, 7
- [12] A. Cioppa, A. Delière, F. Magera, S. Giancola, O. Barnich, B. Ghanem, and M. Van Droogenbroeck. Camera calibration and player localization in soccer-net-v2 and investigation of their representations for action spotting. *arXiv preprint arXiv:2104.09333*, 2021. 2, 7
- [13] L. Citraro, P. Márquez-Neila, S. Savarè, V. Jayaram, C. Dubout, F. Renaut, A. Hasfura, H. B. Shitrit, and P. Fua. Real-time camera pose estimation for sports fields. *Machine Vision and Applications*, 31(3):1–13, 2020. 2
- [14] A. Collet, M. Chuang, P. Sweeney, D. Gillett, D. Evseev, D. Calabrese, H. Hoppe, A. Kirk, and S. Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (ToG)*, 34(4):1–13, 2015. 3
- [15] K. Desai, B. Prabhakaran, and S. Raghuraman. Skeleton-based continuous extrinsic calibration of multiple rgb-d kinect cameras. In *Proceedings of the 9th ACM Multimedia Systems Conference*, 2018. 3
- [16] J. Dong, W. Jiang, Q. Huang, H. Bao, and X. Zhou. Fast and robust multi-person 3d pose estimation from multiple views. In *CVPR*, 2019. 1, 2, 6, 7
- [17] S. Ershadi-Nasab, S. Kasaei, and E. Sanaei. Uncalibrated multi-view multiple humans association and 3d pose estimation by adversarial learning. *Multimedia Tools and Applications*, 80(2):2461–2488, 2021. 1, 2
- [18] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu. Rmpe: Regional multi-person pose estimation. In *ICCV*, 2017. 2, 3, 4
- [19] J. Gall, A. Yao, and L. Van Gool. 2d action recognition serves 3d human pose estimation. In *ECCV*, 2010. 3
- [20] N. Garau and N. Conci. Unsupervised continuous camera network pose estimation through human mesh recovery. In *Proceedings of the 13th International Conference on Distributed Smart Cameras*, 2019. 3
- [21] N. Garau, F. G. De Natale, and N. Conci. Fast automatic camera network calibration through human mesh recovery. *Journal of Real-Time Image Processing*, 17(6):1757–1768, 2020. 3
- [22] S. Geman. Statistical methods for tomographic image reconstruction. *Bull. Int. Stat. Inst.*, 4:5–21, 1987. 5
- [23] Q.-X. Huang and L. Guibas. Consistent shape maps via semidefinite programming. In *Computer Graphics Forum*, volume 32, pages 177–186. Wiley Online Library, 2013. 4
- [24] S. Huang, X. Ying, J. Rong, Z. Shang, and H. Zha. Camera calibration from periodic motion of a pedestrian. In *CVPR*, 2016. 2
- [25] 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 *3DV*, 2017. 3
- [26] R. Inomata, K. Terabayashi, K. Umeda, and G. Godin. Registration of 3d geometric model and color images using sift and range intensity images. In *International Symposium on Visual Computing*, 2011. 2
- [27] T. Jaeggli, E. Koller-Meier, and L. Van Gool. Learning generative models for multi-activity body pose estimation. *International Journal of Computer Vision*, 83(2):121–134, 2009. 3
- [28] W. Jiang, N. Kolotouros, G. Pavlakos, X. Zhou, and K. Daniilidis. Coherent reconstruction of multiple humans from a single image. In *CVPR*, 2020. 5
- [29] H. Joo, T. Simon, X. Li, H. Liu, L. Tan, L. Gui, S. Banerjee, T. Godisart, B. Nabbe, I. Matthews, et al. Panoptic studio: A massively multiview system for social interaction capture. *IEEE transactions on pattern analysis and machine intelligence*, 41(1):190–204, 2017. 1, 2, 3
- [30] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 2, 4, 5
- [31] M. Kocabas, N. Athanasiou, and M. J. Black. Vibe: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 3
- [32] N. Kolotouros, G. Pavlakos, M. J. Black, and K. Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 3
- [33] N. Lawrence and A. Hyvärinen. Probabilistic non-linear principal component analysis with gaussian process latent variable models. *Journal of machine learning research*, 6(11), 2005. 3

- [34] K. Li, N. Jiao, Y. Liu, Y. Wang, and J. Yang. Shape and pose estimation for closely interacting persons using multi-view images. In *Computer Graphics Forum*, 2018. 1, 2, 7
- [35] K. Li, Y. Mao, Y. Liu, R. Shao, and Y. Liu. Full-body motion capture for multiple closely interacting persons. *Graphical Models*, 110:101072, 2020. 1, 2, 3, 7
- [36] R. Li, T.-P. Tian, and S. Sclaroff. Simultaneous learning of nonlinear manifold and dynamical models for high-dimensional time series. In *ICCV*, 2007. 3
- [37] R. Li, T.-P. Tian, S. Sclaroff, and M.-H. Yang. 3d human motion tracking with a coordinated mixture of factor analyzers. *International Journal of Computer Vision*, 87(1-2):170, 2010. 3
- [38] J. Lin and G. H. Lee. Multi-view multi-person 3d pose estimation with plane sweep stereo. In *CVPR*, 2021. 1, 2
- [39] H. Y. Ling, F. Zinno, G. Cheng, and M. Van De Panne. Character controllers using motion vaes. *ACM Transactions on Graphics (TOG)*, 39(4):40–1, 2020. 2, 3, 4
- [40] Y. Liu, J. Gall, C. Stoll, Q. Dai, H.-P. Seidel, and C. Theobalt. Markerless motion capture of multiple characters using multi-view image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2720–2735, 2013. 1, 2, 7
- [41] S. Lohit, R. Anirudh, and P. Turaga. Recovering trajectories of unmarked joints in 3d human actions using latent space optimization. In *WACV*, 2021. 2, 3
- [42] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. 3
- [43] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60(2):91–110, 2004. 2
- [44] Z. Luo, S. A. Golestaneh, and K. M. Kitani. 3d human motion estimation via motion compression and refinement. In *ACCV*, 2020. 2, 3, 4
- [45] N. Mahmood, N. Ghorbani, N. F. Troje, G. Pons-Moll, and M. J. Black. Amass: Archive of motion capture as surface shapes. In *ICCV*, 2019. 2, 5
- [46] J. R. Mitchelson and A. Hilton. Simultaneous pose estimation of multiple people using multiple-view cues with hierarchical sampling. In *BMVC*, 2003. 2
- [47] A. Mustafa, H. Kim, J.-Y. Guillemaut, and A. Hilton. General dynamic scene reconstruction from multiple view video. In *ICCV*, 2015. 1, 2
- [48] 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 *CVPR*, 2019. 2, 5, 7, 8
- [49] J. Puwein, L. Ballan, R. Ziegler, and M. Pollefeys. Joint camera pose estimation and 3d human pose estimation in a multi-camera setup. In *ACCV*, 2014. 3
- [50] L. Sha, J. Hobbs, P. Felsen, X. Wei, P. Lucey, and S. Ganguly. End-to-end camera calibration for broadcast videos. In *CVPR*, 2020. 2
- [51] H. Sidenbladh, M. J. Black, and D. J. Fleet. Stochastic tracking of 3d human figures using 2d image motion. In *ECCV*, 2000. 3
- [52] S. N. Sinha and M. Pollefeys. Camera network calibration and synchronization from silhouettes in archived video. *International journal of computer vision*, 87(3):266–283, 2010. 2
- [53] K. Sun, B. Xiao, D. Liu, and J. Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. 4
- [54] 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 *CVPR Workshops*, 2018. 3
- [55] Z. Tang, Y.-S. Lin, K.-H. Lee, J.-N. Hwang, and J.-H. Chuang. Esther: Joint camera self-calibration and automatic radial distortion correction from tracking of walking humans. *IEEE Access*, 7:10754–10766, 2019. 2
- [56] R. Urtasun, D. J. Fleet, and P. Fua. 3d people tracking with gaussian process dynamical models. In *CVPR*, 2006. 3
- [57] D. Vlasic, P. Peers, I. Baran, P. Debevec, J. Popović, S. Rusinkiewicz, and W. Matusik. Dynamic shape capture using multi-view photometric stereo. In *SIGGRAPH Asia*, 2009. 3
- [58] J. M. Wang, D. J. Fleet, and A. Hertzmann. Gaussian process dynamical models for human motion. *IEEE transactions on pattern analysis and machine intelligence*, 30(2):283–298, 2007. 3
- [59] C. Wu, C. Stoll, L. Valgaerts, and C. Theobalt. On-set performance capture of multiple actors with a stereo camera. *ACM Transactions on Graphics (TOG)*, 32(6):1–11, 2013. 1, 2
- [60] A. Yao, J. Gall, L. V. Gool, and R. Urtasun. Learning probabilistic non-linear latent variable models for tracking complex activities. *Advances in Neural Information Processing Systems*, 24:1359–1367, 2011. 3
- [61] G. Ye, Y. Liu, N. Hasler, X. Ji, Q. Dai, and C. Theobalt. Performance capture of interacting characters with handheld kinects. In *ECCV*, 2012. 1
- [62] Y. Zhang, L. An, T. Yu, X. Li, K. Li, and Y. Liu. 4d association graph for realtime multi-person motion capture using multiple video cameras. In *CVPR*, 2020. 1, 2, 6, 7
- [63] Z. Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000. 2
- [64] Z. Zhang. Camera calibration with one-dimensional objects. *IEEE transactions on pattern analysis and machine intelligence*, 26(7):892–899, 2004. 2
- [65] Z. Zhao, X. Zhao, and Y. Wang. Travelnet: Self-supervised physically plausible hand motion learning from monocular color images. In *ICCV*, 2021. 3
- [66] K. Zhou, Y. Yang, A. Cavallaro, and T. Xiang. Omni-scale feature learning for person re-identification. In *ICCV*, 2019. 2, 4
- [67] X. Zhou, M. Zhu, G. Pavlakos, S. Leonardos, K. G. Derpanis, and K. Daniilidis. Monocap: Monocular human motion capture using a cnn coupled with a geometric prior. *IEEE transactions on pattern analysis and machine intelligence*, 41(4):901–914, 2018. 3
- [68] Y. Zhou, C. Barnes, J. Lu, J. Yang, and H. Li. On the continuity of rotation representations in neural networks. In *CVPR*, 2019. 3

- [69] D. Zou and P. Tan. Coslam: Collaborative visual slam in dynamic environments. *IEEE transactions on pattern analysis and machine intelligence*, 35(2):354–366, 2012. 2