

PLIKS: A Pseudo-Linear Inverse Kinematic Solver for 3D Human Body Estimation

Karthik Shetty^{1,2} Annette Birkhold² Srikrishna Jaganathan^{1,2} Norbert Strobel^{2,3}

Markus Kowarschik² Andreas Maier¹ Bernhard Egger¹

¹FAU Erlangen-Nürnberg, Erlangen, Germany ²Siemens Healthineers AG, Forchheim, Germany

³University of Applied Sciences Würzburg-Schweinfurt, Germany

karthik.shetty@fau.de

Abstract

We introduce *PLIKS* (Pseudo-Linear Inverse Kinematic Solver) for reconstruction of a 3D mesh of the human body from a single 2D image. Current techniques directly regress the shape, pose, and translation of a parametric model from an input image through a non-linear mapping with minimal flexibility to any external influences. We approach the task as a model-in-the-loop optimization problem. *PLIKS* is built on a linearized formulation of the parametric SMPL model. Using *PLIKS*, we can analytically reconstruct the human model via 2D pixel-aligned vertices. This enables us with the flexibility to use accurate camera calibration information when available. *PLIKS* offers an easy way to introduce additional constraints such as shape and translation. We present quantitative evaluations which confirm that *PLIKS* achieves more accurate reconstruction with greater than 10% improvement compared to other state-of-the-art methods with respect to the standard 3D human pose and shape benchmarks while also obtaining a reconstruction error improvement of 12.9 mm on the newer AGORA dataset.

1. Introduction

Estimating human surface meshes and poses from single images is one of the core research directions in computer vision, allowing for multiple applications in computer graphics, robotics and augmented reality [16, 50]. Since humans have complex body articulations and the scene parameters are typically unknown, we are essentially dealing with an ill-posed problem that is difficult to solve in general.

Thanks to models such as SMPL [38] and SMPL-X [48] additional constraints on body shape and pose became available. They made the problem somewhat more tractable. Most state-of-the-art methods [21, 23, 27, 30, 56] directly regress the shape and pose parameters from a given input image. These approaches rely completely on neural net-

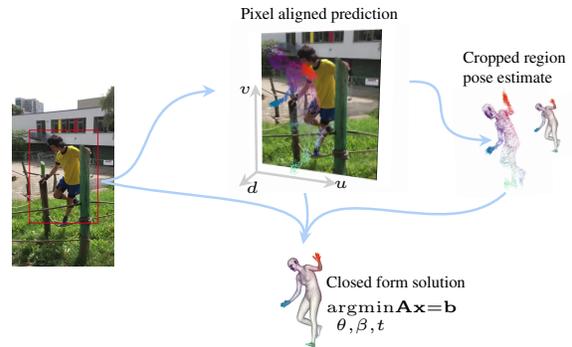


Figure 1. Network predicts a pixel-aligned vertex map (u, v, d) which is used to obtain an initial pose estimate. Then a closed-form solution is made use of to solve the Inverse kinematics between the 2D pixel-aligned vertex map (u, v) and a pseudo-parametric model given the detected bounding-box camera intrinsic and initial pose estimate.

works, while making several assumptions about the image generation process. One typical assumption is the use of a simplified camera model such as the weak perspective camera. In this scenario, the camera is assumed to be far away from the subject, which is generally realized by setting a large focal length constant for all images. A weak perspective camera can be described based on three parameters, two with respect to translation in the horizontal and vertical directions, and the third being scale. While these methods can estimate plausible shape and pose parameters, it can happen that the resulting meshes are either misaligned in the 2D image space or in the 3D object space. This is because the underlying optimization problem is often not constrained enough such that it is difficult for the underlying networks to optimize between the 2D re-projection loss and the 3D loss.

Some existing methods [26, 28, 34] propose a workaround by tackling the problem using a hybrid

approach involving learning-based and optimization-based techniques while incorporating a full perspective camera [26]. Optimization-based approaches are, however, prone to local minima, and they are computationally expensive. In [28], the authors propose to regress the SMPL parameters by conditioning on features from a CamCalib network meant to predict the camera parameters. Unfortunately, this camera prediction network needs a specialized dataset to train on, which is very hard to acquire in practice. It also prevents end-to-end learning.

On the other hand, recent non-parametric or model-free approaches [36, 43] directly regress the mesh vertex coordinates based on their 2D projections, aligning well to the input image. However, by ignoring the effects of a perspective camera, even these methods suffer from the same limitations as the parametric models.

Motivated by the above observations, we present a novel approach, named PLIKS, for 3D human shape and pose estimation that incorporates the perspective camera while analytically solving for all the parameters of the parametric model. The pipeline of our approach comprises of two modules, namely the mesh regressor and PLIKS. The mesh regressor provides a mapping between an image and the 3D vertices of the SMPL model. Given a single image, any off-the-shelf Convolution Neural Network (CNN) can be used for feature extraction. The extracted features can then be used to obtain a mesh representation either by using 1D CNNs [43], GraphCNNs [31], or even transformers [36]. This way, correspondences to the image space can be found and a relative depth estimate can be computed. From the image-aligned mesh prediction, we can roughly estimate the rotations with respect to a template mesh in canonical space with the application of Inverse Kinematics (IK), denoted in this work as the Approximate Rotation Estimator (ARE). Finally, we reformulate the SMPL model as a linear system of equations, with which we can use the 2D pixel-aligned vertex maps and any known camera intrinsic parameters to fully estimate the model without the need for any additional optimization. As our approach is end-to-end differentiable and fits the model within the training loop, it is self-improving in nature. The proposed approach is benchmarked against various 3D human pose and shape datasets, and significantly outperforms other state-of-the-art approaches.

To summarize, the contribution of our paper is the following: (1) We bridge the gap between the 2D pixel-aligned vertex maps and the parametric model by reformulating the SMPL model as a linear system of equations. Since the proposed approach is fully differentiable, we can perform end-to-end training. (2) We propose a 3D human body estimation framework that reconstructs the 3D body without relying on weak-perspective assumptions. (3) We show that our approach can improve upon other state-of-the-art methods

when evaluated across various 3D human pose and shape benchmarks.

2. Related Work

Recovering human pose and shape from monocular images has been extensively studied using both model-based [23, 27, 56] and model-free approaches [12, 36, 43]. Model-based methods estimate the parameters of a parametric body model such as SMPL [38] based on a single input RGB image. The use of parametric body models makes it possible to enforce strong statistical priors of the human mesh. Model-based methods can be further split into optimization and regression techniques. Optimization-based approaches [4, 7] make use of 2D keypoints estimated by a Deep Neural Network (DNN) which are iteratively fit with the SMPL model. These methods however are sensitive to initialization and are susceptible to local minima. Regression-based techniques based on DNNs directly estimate the pose and shape parameters [11, 23, 45, 49, 56, 58]. However, these approaches typically require a large amount of training data. SPIN [30] provided a revolutionary architecture combining optimization-based techniques with regression-based methods. This allowed for much stronger supervision and improved performance on mesh accuracy. However, due to the difficulty in directly estimating the mapping from a single image to the shape and pose space, the mesh alignment with respect to the input image is often imperfect [19, 32]. Model-free approaches directly regress the vertices based on intermediate representations, such as IUV maps [12, 60, 63], 2D/3D heatmaps [19, 32], silhouette [53, 61], and direct vertex regression [31, 35, 36, 43] where correspondence is established between the model and the input image.

To overcome the issues of the learning-based and optimization-based approaches, hybrid techniques combining both approaches have been proposed [8, 19, 32, 33]. In contrast to SPIN [30], here we have a closed-form solution. First, the learning-based approach localizes 3D human joint coordinates [5, 6, 18, 41], utilizing volumetric heatmaps or GraphCNNs for the target representation. From the localized 3D joints, the swing rotations are then determined analytically, whereas the twist rotations, shape parameters, and root translation are either predicted by a network or iteratively optimized [19, 32].

Typical human reconstruction methods [23, 30] take a cropped input image while using 3D predictions projected onto the cropped image for 2D supervision. This, however, ignores the effects of perspective warping. This happens when the cropped image is off-center resulting in inaccurate rotations [62]. To overcome this, some approaches [19, 26] perform iterative optimization [4] after the initial DNN-based predictions. SPEC [28] proposes to condition the image features from a camera calibration network. How-

ever, a few methods include perspective warping during the cropping process and add an implicit camera rotation as post-processing [39, 62]. CLIFF [34] addresses the problem by incorporating the bounding box information into the cropped image.

3. Methodology

In this section, we present our network architecture which includes an analytical solver for inverse kinematics and is end-to-end trainable. As illustrated in Fig. 2, our network consists of two parts, first, a mesh regressor and second, a Pseudo-Linear Inverse Kinematic Solver (PLIKS). In Sec. 3.1, we go over the SMPL model along with its forward kinematics process, and in Sec. 3.2, we explain the full pipeline for mesh reconstruction.

3.1. Parametric Mesh Representation

We use the Skinned Multi-Person Linear (SMPL) model to parameterize the human body [38]. The SMPL model is a statistical parametric function $\mathcal{M}(\beta, \theta; \Phi)$. The output of this function is a triangulated surface mesh with $N = 6890$ vertices. The shape parameters β are represented by a low dimensional principal component which maps the linear basis \mathbf{B} from $\mathbb{R}^{|\beta|} \mapsto \mathbb{R}^{3N}$, representing offsets to the average mesh \bar{x}_m as $x = \bar{x}_m + \beta\mathbf{B}$. The pose of the model is defined with the help of a kinematics chain involving a set of relative rotation vectors $\theta = [\theta_1, \dots, \theta_K] \in \mathbb{R}^{K \times 3}$ made up of $K = 24$ joints represented using axis-angle rotations. Additional model parameters summarized as Φ are involved in the deformation process of the SMPL model. They are used as joint regressor $\mathcal{J} \in \mathbb{R}^{K \times 3N}$, blend weights $\mathcal{W} \in \mathbb{R}^{K \times N}$ and shape deformations conditioned on the body pose $B_P(\theta)$. Starting from a mean template mesh, the desired body mesh is obtained by applying forward kinematics based on the relative rotations θ and shape deformations β . The 3D body joints can be obtained by a linear combination of the mesh vertices using any desired linear regressor $\mathcal{J}'(\mathcal{M}(\beta, \theta; \Phi))$.

Simplified SMPL Although the SMPL model described above is intrinsically linear, it cannot be defined as a linear system of equations to solve the IK due to the forward kinematics and pose-related shape deformations $B_P(\theta)$. To deal with this problem, we make use of a simplified model, where we ignore the pose-related shape deformations $B_P(\theta)$ in the formulation. We further split the SMPL mesh into $K=24$ segments $s_{k \in (1, K)}$ corresponding to the influence of the 24 rotation vectors θ of the SMPL model as shown in Fig. 3. A vertex on the mesh with index i belongs to a particular segment that has the maximum blend weight of all the joints, i.e., $\arg \max_{k \in (1, K)} \mathcal{W}(i)$. For ease of notation, we represent a set of vertices/indices corresponding

to a segment s_k with the superscript k . Here, we also consider the individual segments to be a set of rigid 3D points which enables us to orient these segments from the template mesh \bar{x}_m^k to a network predicted mesh \mathbf{X}^k to obtain an initial rotation estimate using ARE explained in the following section.

3.2. Pseudo-Linear Inverse Kinematic Solver

Solving Inverse Kinematics (IK) from 2D correspondences is a challenging task due to the inherent non-linearity that exists in the parameterized mesh generation process. We propose to solve IK from 2D pixel-aligned vertex inputs by building a linear system of equations of the form $\mathbf{A}\mathbf{x}=\mathbf{b}$. We add shape constraints to obtain the optimal world pose $\hat{\theta}$, shape β , and world translation \mathbf{t} . Linear-Least-Squares is used to estimate the optimal parametric solution making the entire pipeline end-to-end differentiable. From the world pose $\hat{\theta}$, the relative rotations θ can be inferred by recursively solving the kinematic tree. For the linear system, we assume the rotations as a first-order Taylor approximation [3] with the individual segments of the predicted mesh model approximately oriented $\tilde{\theta}$ along the optimal solution. Based on this assumption, using the rotations estimated from ARE should provide the exact solution in a single optimization step.

Regression Network Assuming that some form of dense vertex correspondences exists that facilitates mapping between 3D vertices to pixels in the 2D image plane, we can then incorporate the IK solver into the network. As shown in Fig. 2, our architecture is comprised of an encoder, a mesh regressor, and a parameter regressor. The encoder is based on HRNet-W32 [55], the mesh regressor is based on MeshNet [43], whereas the parameter regressor is a set of fully connected layers.

The encoder acts as a feature extractor $\mathbf{F} \in \mathbb{R}^{C \times c' \times c'}$, with a channel dimension $C=480$, height and width $c'=58$. It takes a cropped image $\mathbf{I} \in \mathbb{R}^{3 \times 224 \times 224}$ of a person as input. Similar to MeshNet [43], we make use of 1D convolutions to generate four feature vectors

$$\mathbf{P} = \{\{\mathbf{P}^u, \mathbf{P}^v, \mathbf{P}^d\} \in \mathbb{R}^{N \times 58}, \mathbf{w} \in \mathbb{R}^N\},$$

for each mesh vertex. Here, \mathbf{P}^u and \mathbf{P}^v represent the features along the u and v axis, \mathbf{P}^d represents the features along the root-normalized depth, and \mathbf{w} represents the weighting factor. They are obtained as follows

$$\begin{aligned} \mathbf{P}^u &= f_u^{1D}(\text{avg}^u(\mathbf{F})), \\ \mathbf{P}^v &= f_v^{1D}(\text{avg}^v(\mathbf{F})), \\ \mathbf{P}^d &= f_d^{1D}(\psi_d(\text{avg}^{u,v}(\mathbf{F}))), \\ \mathbf{w} &= \sigma(\text{avg}^d(f_d^{1D}(\psi_d(\text{avg}^{u,v}(\mathbf{F}))))). \end{aligned} \tag{1}$$

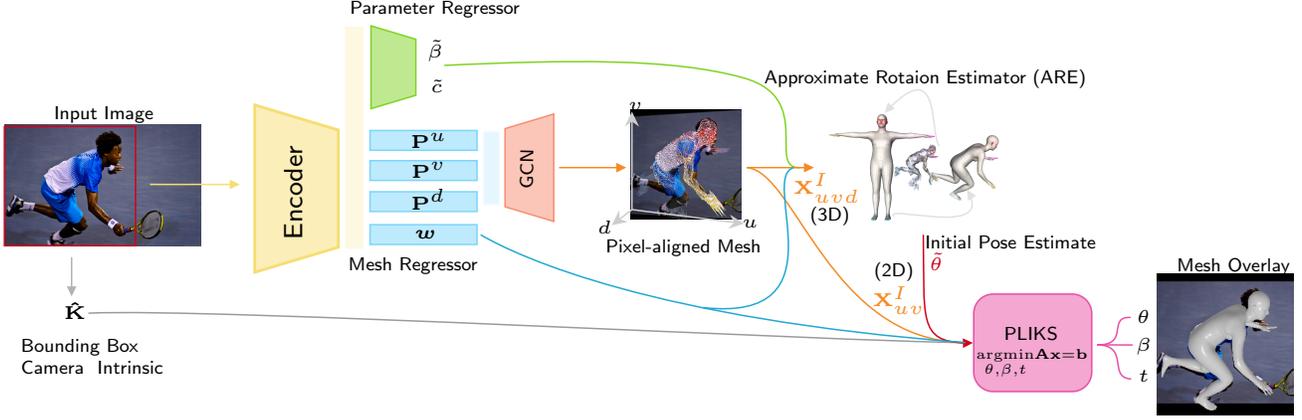


Figure 2. **Overview of the proposed framework:** For a given RGB image input, the mesh regressor generates a mapping to the 3D mesh vertices aligned to the image plane. The ARE module provides rotation estimates in world space for each joint segment of the SMPL. The reconstructed body mesh is obtained via an analytical solution through the PLIKS process, fed with $\tilde{\theta}$ comprising the segments’ initial world rotations and the 2D mesh projections \mathbf{X}_{uv}^I as inputs.

Here $f_i^{1D}(\cdot)$ and $\psi_i(\cdot)$ represent a 1D convolution along the i -th dimension, which converts the features \mathbf{F} from $C \times c' \rightarrow N \times c'$ and $C \times 1 \rightarrow (c' \times C)^T$, respectively. The average function $\text{avg}^i(\cdot)$ averages the features along the i -th dimension. For the weighting factor w , we average across the channel dimension, followed by applying the sigmoid activation function $\sigma(\cdot)$ [44]. More details about the weighting function w are explained in the following sections.

We then concatenate and process $\tilde{\mathbf{P}} = \{\mathbf{P}^u, \mathbf{P}^v, \mathbf{P}^d\}$ using a graph convolution network (GCN) to predict $uvd \in \mathbb{R}^{N \times 3}$. We use GCN rather than the heatmap-based approach from MeshNet [43] to avoid any truncation-based artifacts, which are prominent when partial images of humans are supplied as inputs. For the GCN we use the formulation from Kipf *et al.* [25], defined as $\mathbf{G} = \sigma(\tilde{\mathbf{A}}\tilde{\mathbf{P}}\mathbf{W})$, where $\tilde{\mathbf{A}} \in \mathbb{R}^{N \times N}$ denotes the graph adjacency matrix, $\mathbf{W} \in \mathbb{R}^{c' \times l}$ denotes the trainable weights with l as the out-



Figure 3. We split the SMPL model into 24 segments based on the blend weights. These segments help in determining the initial rotation estimate in the ARE pipeline. Here the segments are assumed to be rigid.

put channel dimension, and σ denotes the ReLU [1] activation function. We make use of 3 GCNs in series, with channel sizes 64, 32, and 3 respectively. The final output $\mathbf{X}_{uvd}^I \in \mathbb{R}^{N \times 3}$ acts as the vertex correspondence in the image coordinate system.

The parameter regressor is a set of fully connected layers to obtain an approximate shape $\tilde{\beta}$ and a weak perspective camera $\tilde{c} \in \mathbb{R}^3$, which are used to determine the approximate world rotations later. We make use of the estimated depth \tilde{c}_d from the camera prediction to obtain a mesh in the world coordinate system \mathbf{X} as follows,

$$\mathbf{X} = (\mathbf{K}^{-1}\mathbf{X}_{uv1}^I) \cdot (\mathbf{X}_d^I + \tilde{c}_d), \quad (2)$$

where $\mathbf{K} \in \mathbb{R}^{3 \times 3}$ is an intrinsic matrix with fixed focal length of 1000.

Approximate Rotation Estimator (ARE) Given two sets of corresponding points in 3D space, it is possible to obtain an optimal rotation as a closed-form solution using the Kabsch solver [22]. For a given segment from network mesh prediction \mathbf{X}^k we make use of the Kabsch solver to determine the rotation that a same segment from the template mesh needs to go through from its rest pose. As the mesh prediction \mathbf{X}^k can represent a wide range of human shapes, we make use of the shape predictions $\tilde{\beta}$ on the mean shape as $\mathbf{x}^k = \bar{\mathbf{x}}_m^k + (\tilde{\beta}\mathbf{B})^k$. The pose solver minimizes the squared distances between a set of 3D correspondences to obtain an optimal pose as follows,

$$(\mathbf{U}\Sigma\mathbf{V}^T)_k = \mathbf{C}_k = \sum_{(\mathbf{X}^k, \bar{\mathbf{X}}^k)} (\mathbf{X}^k - \bar{\mathbf{X}}^k) (\mathbf{x}^k - \bar{\mathbf{x}}^k)^T \mathbf{w}^k \mathcal{W}^k. \quad (3)$$

Here C_k is the covariance matrix between the correspondences. In this context, $\bar{\mathbf{x}}^k$ and $\bar{\mathbf{X}}^k$ represent the means over the associated point sets. The world rotation for the k -th segment $\hat{\boldsymbol{\theta}}_k = (\mathbf{V}\mathbf{U}^T)_k$ is obtained by applying singular value decomposition (SVD) over the covariance matrix. Since the SVD is differentiable, gradients can be back-propagated during the training process. Due to the blend skinning \mathcal{W} process, a vertex may have rotation influence from one or more joints. To tackle this, we multiply the squared distances in Eq. (3) with a weight term w^k from Eq. (1) and the maximum blend influence \mathcal{W}^k . The weight is learned in an unsupervised manner to offset any influence the blending process contributes.

The obtained rotations are in the world space, whereas the SMPL model expects relative rotations in its axis space. The rotation around the pelvis $\hat{\boldsymbol{\theta}}_{k=1}$ corresponds to the global root rotation. The relative rotations for the other joints can be recursively based on the parent rotation following a pre-defined kinematic tree for the SMPL model. To this end, let $\hat{\mathbf{R}}_k = \mathcal{R}(\hat{\boldsymbol{\theta}}_k) \in \mathbb{SO}(3)$ represent the estimated world rotation determined for segment k . Then the relative rotation for segment k is

$$\hat{\boldsymbol{\theta}}_k = \hat{\mathbf{R}}_k = \tilde{\mathbf{R}}_{p(k)}^{-1} \tilde{\mathbf{R}}_k, \quad (4)$$

where $p(k)$ represents the parent joint of k and $\hat{\boldsymbol{\theta}}_k$ the estimated relative pose.

Pseudo-Linear Inverse Kinematic Solver (PLIKS) An approximate SMPL model projected onto the image plane, with no pose-related blend shapes can be represented as

$$\mathbf{i}^k = \hat{\mathbf{K}} \sum_{j=1}^K \mathcal{W}_j^k \left(\Delta \mathbf{R}_k \tilde{\mathbf{R}}_k (\bar{\mathbf{x}}_m^k + \beta \mathbf{B}^k) + \mathbf{t}_k \right), \quad (5)$$

Here, the superscript k refers to the k -th's subset of vertices defined in Sec. 3.1. The matrix $\hat{\mathbf{K}} \in \mathbb{R}^{3 \times 4}$ represents a perspective projection matrix taking into account the affine transformation (only crop and resize) for the image fed into the network as

$$\hat{\mathbf{K}} = \begin{bmatrix} f_x & 0 & p_x & 0 \\ 0 & f_y & p_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}.$$

Further, $\tilde{\mathbf{R}}_k$ represents the approximate world rotation obtained from the previous step, $\Delta \mathbf{R}_k$ is defined as the additional rotation required to get an optimal solution, \mathbf{t}^k represents the joint translation, and \mathcal{W} are the blend weights. By combining all segments, we get $\mathbf{i} = \mathbf{X}_{uv1}^I \in \mathbb{R}^3$. This represents the correspondence between mesh vertices, defined in homogeneous coordinates, and the pixels in the image plane.

Since the additional rotation required is considered to be small we linearize the rotation matrix which needs to be determined based on Taylor expansion with angles α , β , γ along the x , y and z axes, respectively, as follows

$$\Delta \mathbf{R}_k = \begin{bmatrix} 1 & -\gamma_k & \beta_k & 0 \\ \gamma_k & 1 & -\alpha_k & 0 \\ -\beta_k & \alpha_k & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

Further, we simplify the projected SMPL model in Eq. (5) by making some additional assumptions using the definitions introduced in Eq. (6). For the \mathbf{x}_r^k term we assume that the majority of the rotation is significantly effected by the rotation with respect to the primary segment, i.e., we ignore the impact of neighbouring rotations as they are usually minuscule. This assumption is also made for the term \mathbf{B}_r^k by assuming that for small rotations $\Delta \mathbf{R}_k \tilde{\mathbf{R}}_k \beta \mathbf{B}^k \approx \tilde{\mathbf{R}}_k \beta \mathbf{B}^k$.

$$\mathbf{i}^k = \hat{\mathbf{K}} \left(\underbrace{\Delta \mathbf{R}_k \sum_{j=1}^K \mathcal{W}_j^k \tilde{\mathbf{R}}_k \bar{\mathbf{x}}_m^k}_{\mathbf{x}_r^k} + \beta \underbrace{\sum_{j=1}^K \mathcal{W}_j^k \tilde{\mathbf{R}}_k \mathbf{B}^k}_{\mathbf{B}_r^k} + \mathbf{t}_k \underbrace{\sum_{j=1}^K \mathcal{W}_j^k}_{\mathbf{W}_r^k} \right). \quad (6)$$

The equation is now linear, with 154 unknown parameters corresponding to $\Delta \mathbf{R}_k$, β and \mathbf{t}_k . If the focal lengths f_x , f_y and principal points p_x , p_y are not known, we assume a fixed focal length of 1000 mm and the image center as the principal point. Note that the fixed values here apply to the entire image, i.e., not to the cropped and resized input fed to the network.

Using Direct Linear Transform (DLT) [13] we can rewrite Eq. (6) in the form $\mathbf{A}\mathbf{x} = \mathbf{b}$. The optimal parameters can be obtained by minimizing the analytical error using linear least square defined by $\mathbf{x} = \mathbf{A}^+ \mathbf{b}$, where \mathbf{A}^+ represents the pseudo-inverse of \mathbf{A} . As the pseudo-inverse of a tall matrix is differentiable [9], gradients can be back-propagated during the training process. One of the major drawbacks when using DLT is that it minimizes the analytical loss rather than the geometric loss. One option to overcome this is to use Iterative Reweighted Least Squares (IRLS) [10], which robustly minimizes the objective function in an iterative manner by reweighing the geometric loss. However, we make use of the network predicted weighting w^k to weight the correspondences [29, 51]. To further enforce the predicted shape to be close to the mean shape, we add an additional constraint such that $\|\beta\|_2 \approx 0$ with a regularizing weight ω_β as

$$\operatorname{argmin}_{\Delta \mathbf{R}_k, \beta, \mathbf{t}_k} \left\| \mathbf{w}^k \left(\mathbf{i}^k - \hat{\mathbf{K}} (\Delta \mathbf{R}_k \mathbf{x}_r^k + \beta \mathbf{B}_r^k + \mathbf{t}_k \mathbf{W}_r^k) \right) \right\|_2 + \omega_\beta \|\beta\|_2. \quad (7)$$

To get the final relative pose $\boldsymbol{\theta}_k$ for each joint k , we apply Eq. (4) on the obtained world rotations $\Delta \mathbf{R}_k \tilde{\mathbf{R}}_k$. The

Method	3DPW (14)			Human3.6M (14)		MPI-INF-3DHP (17)		
	PA-MPJPE↓	MPJPE↓	PVE↓	PA-MPJPE↓	MPJPE↓	PCK↑	AUC↑	MPJPE↓
HMR [23]	81.3	130.0	-	56.8	88.0	-	-	-
SPIN [30]	59.2	96.9	116.4	41.1	-	76.4	37.1	105.2
I2L+ [43]	58.6	93.2	-	41.7	55.7	-	-	-
EFT† [21]	51.6	-	-	44.0	-	-	-	-
ROMP† [56]	47.3	76.7	93.4	-	-	-	-	95.1
PARE† [27]	46.5	74.5	88.6	-	-	-	-	-
Mesh Graphormer†+ [36]	45.6	74.7	87.7	34.5	51.2	-	-	-
HybrIK† [32]	45.3	74.1	86.5	33.6	55.4	87.5	46.9	93.9
CLIFF† [34]	43.0	69.0	81.2	32.7	47.1	-	-	-
PLIKS†	42.8	66.9	82.6	34.7	49.3	91.8	52.3	72.3
PLIKS‡	40.1	63.5	76.7	34.9	46.8	93.1	53.0	69.8
PLIKS‡(HR48)	38.5	60.5	73.3	34.5	47.0	93.9	54.1	67.6

Table 1. Benchmark of state-of-the-art models on 3DPW, Human3.6M, and MPI-INF-3DHP datasets. All units are in mm. Here + represents non-parametric methods. † means that the network was additionally trained with 3DPW. ‡ means that the network was additionally trained with 3DPW and AGORA.

	3DPW		AGORA	
	MPJPE↓	PVE↓	MPJPE↓	PVE↓
ARE	23.0	24.9	52.0	55.3
PLIKS	1.3	1.5	8.4(1.1)	10.3(1.7)

Table 2. Ground truth errors to validate the closed-form solution of PLIKS. Results in brackets represent running the PLIKS module twice. All units are in mm.

global translation for the camera system $\hat{\mathbf{K}}$ is obtained from the root joint as $\mathbf{t}_1 - (\mathbf{j}_1 - \Delta \mathbf{R}_1 \hat{\mathbf{R}}_1 \mathbf{j}_1)$, where \mathbf{j}_1 is the root joint for a rest pose mesh with shape coefficient β .

4. Experiments

Following previous works, the base PLIKS model is trained on a combination of 3D datasets (Human3.6M [17] and MPI-INF-3DHP [39]), and 2D dataset (COCO [37]) with pseudo-GT labels obtained from EFT [21]. We evaluate on Human3.6M [17], 3DPW [59], MPI-INF-3DHP [39], AGORA [47], MuPoTs-3D [40], and 3DOH50K [65]. During the evaluation, we highlight our results for networks trained with any additional datasets.

4.1. PLIKS Error Analysis

We measure the effective reconstruction capability of ARE and PLIKS with ground truth conditions for the 3DPW [59] test set and the AGORA [47] validation set. For the ARE module experiment, we provide the GT SMPL mesh vertex projections and the GT shape parameters. Additionally, we set the root depth as 7 m for all the images as the SMPL model in its template pose can be well represented inside a 224×224 image when assuming weak perspective settings with a focal length of 1 m. For the PLIKS module experiment, we only provide the GT SMPL mesh vertex projections and a root depth of 7 m to the ARE module along with the GT bounding-box camera intrinsics as inputs. We also set the weights for the shape regular-

	MRPE	MRPE _x	MRPE _y	MRPE _z
Baseline [26, 52]	267.8	27.5	28.3	261.9
RootNet [41]	120.0	23.3	23.0	108.1
PLIKS	135.5	18.0	14.1	128.9
PLIKS‡	96.1	16.1	14.7	88.5

Table 3. Evaluation on Human3.6M dataset, with respect to the MRPE and MRPE on x , y and z axis. ‡ means the network was additionally trained with 3DPW [59] and AGORA [47].

	Matched People		All People
	PCK _{abs} ↑	PCK _{root} ↑	PCK _{abs} ↑
Rootnet [41]	31.8	31.0	31.5
VirtualPose [54]	47.0	53.5	44.0
PLIKS‡	44.8	55.7	44.2

Table 4. Absolute PCK evaluation on the MuPoTs-3D [40] dataset.

izer $\omega_\beta=0$. We measure the Mean Per Joint Position Error (MPJPE), and the Per Vertex Error (PVE). Following previous works [23, 27, 30], we use the LSP joint regressor [30] to determine the 14 joints which are regressed from the body mesh. We report the results in Table 2 showing that the assumptions made in Eq. (6) are reasonable providing a stable and accurate fit. We observe a larger error on the AGORA dataset with the PLIKS module, due to larger perspective warping effects, which results in incorrect global rotation estimation from the ARE module. As PLIKS uses a linearized rotation matrix, running the PLIKS module twice, that is iteratively optimizing the rotation estimates provides a more accurate fit. Though ARE is not an accurate IK solver, we can safely conclude about the drawbacks of using a weak perspective camera model.

4.2. Comparison with the State-of-the-art

We compare our method with previous human mesh reconstruction approaches based on Human3.6M, MPI-INF-3DHP, and 3DPW datasets. To remain consistent with previous approaches, we use the LSP regressor [30] to obtain the 14 joints for Human3.6M and 3DPW datasets, while we

use 17 joints with the Human3.6M regressor [30] for the MPI-INF-3DHP dataset. We measure Procrustes Aligned Mean Per Joint Position Error (PA-MPJPE), Percentage of Correct Keypoints (PCK), and Area Under Curve (AUC) on the 3D pose results.

In Table 1, we provide quantitative results for previous 3D human shape and pose estimation results. We use the best results reported in all the other works for our comparison. Our network outperforms all previous state-of-the-art techniques on MPI-INF-3DHP and 3DPW datasets along with the lowest MPJPE on Human3.6M. Further fine-tuning with the AGORA dataset results in a significant performance improvement. We also provide results of our network when trained on a larger HRNet-W48 [55] backbone following previous works [34, 36, 64]. We provide qualitative results of our approach comparing prior methods [32, 34] in Fig. 4. Though all methods shown in Fig. 4 align well in the 2D projected space, it is apparent that the prior methods fail to align well in 3D space due to perspective warping. The effect of perspective warping is observed significantly when the human is off-centered. On the MPI-INF-3DHP dataset, the advantage of incorporating PLIKS is more pronounced since the images are captured with a wide Field-of-View (FOV) ($66^\circ \sim 90^\circ$) showing a 26.3 mm MPJPE improvement.

In Table 5, we list quantitative results on the AGORA benchmark. It can be seen that the proposed approach significantly outperforms all prior methods. AGORA additionally measures the F1 score to get Normalized Mean Joint Error (NMJE) and Normalized Mean Vertex Error (NMVE), which penalizes missed or faulty human detection. Here the human bounding box and camera parameters are not provided for the test set. We make use of YOLO-v5 [20, 42] for person detection. From Table 5, PLIKS was fine-tuned only on AGORA, whereas, PLIKS[†] was trained on all the 2D and 3D datasets which helps in generalizing to real-world images. As the focal length is not available during the test, we set $f=1$ m for all images. Making use of the estimated focal length from CamCalib [28] on PLIKS[†] shows a 3.4 mm NMJE improvement (PLIKS[‡]).

Most methods reported in Table 1 generate root relative meshes using a weak perspective camera. As our results are in absolute camera coordinates, we present root localization results in Table 3. We report previous results on Human3.6M provided from [41] on the mean root position error (MRPE). Here baseline [39, 52] depth is obtained using the least squares fit from the predicted 3D joints and its corresponding 2D projections. Though [41] has explicitly been designed to learn the absolute depth, we observe comparable results along the depth while performing better in the horizontal and vertical directions. Using 3DPW and AGORA for training results in significant improvement of the root localization error. We also report the Absolute PCK

Method	F1 \uparrow	NMVE \downarrow	NMJE \downarrow	MVE \downarrow	MPJPE \downarrow
SPIN [30]	0.77	193.4	199.2	148.9	153.4
SPEC [28]	0.84	113.6	118.8	103.4	108.1
BEV [57]	0.93	108.3	113.2	100.7	105.3
H4W [42]	0.94	90.2	95.5	84.8	89.8
CLIFF [34]	0.91	83.5	89.0	76.0	81.0
PLIKS	0.94	76.8	81.5	72.2	76.6
PLIKS [†]	0.94	78.3	83.0	73.6	78.0
PLIKS [‡]	0.94	74.9	79.6	70.4	74.8
PLIKS [‡] (HR48)	0.94	71.6	76.1	67.3	71.5

Table 5. Reconstruction errors on the AGORA test set. All results are taken from the official evaluation platform. PLIKS[†] was trained with all 2D and 3D datasets. PLIKS[‡] represents the usage of CamCalib [28] for focal length estimation during evaluation.

result on the MuPoTs-3D [40] dataset in Table 4. Note that the results are not strictly comparable due to differences in network backbones, the training datasets, and the joint prediction techniques.

4.3. Occlusion Analysis

To validate the stability under occlusion, we evaluate PLIKS on the object-occluded benchmark dataset 3DOH50K [65] and the person-occluded dataset 3DPW-OCC [59]. Following [27], we use only COCO, Human3.6M, and 3DOH datasets for training. To arrive at a fair comparison, we also train a network with ResNet50 [14] as the backbone. We observe better occlusion performance on both backbones from Table 6 for our approach. This can be attributed to the fact, that our approach is intrinsically trained on dense correspondences.

4.4. Ablation Studies

In this context, we evaluate the individual components of our approach. For all experiments, we use only the COCO, Human3.6M, and MPI-INF-3DHP datasets for training and evaluation on all 3D datasets. We use the 3DPW validation set to determine the best model. In Table 7, we provide all the results for our ablation studies. The ARE module on its own can reconstruct 3D models provided the shape and camera parameters are regressed by the network. Hence, we evaluate its performance without the PLIKS module. Here, we observe that it performs poorly on all metrics for all datasets except for the PA-MPJPE of 3DPW. Next, on the same ARE module, we replace the HRNet-32 backbone with a ResNet-50 backbone. In this setting, similar to the

	3DPW-OCC		3DOH	
	MPJPE \downarrow	PA-MPJPE \downarrow	MPJPE \downarrow	PA-MPJPE \downarrow
DOH [65]	-	72.2	-	58.5
EFT [21]	94.4	60.9	75.2	53.1
PARE [27]	90.5	56.6	63.3	44.3
PLIKS(R50)	86.1	53.2	51.5	39.3

Table 6. Evaluation on occlusion datasets 3DPW-OCC and 3DOH.

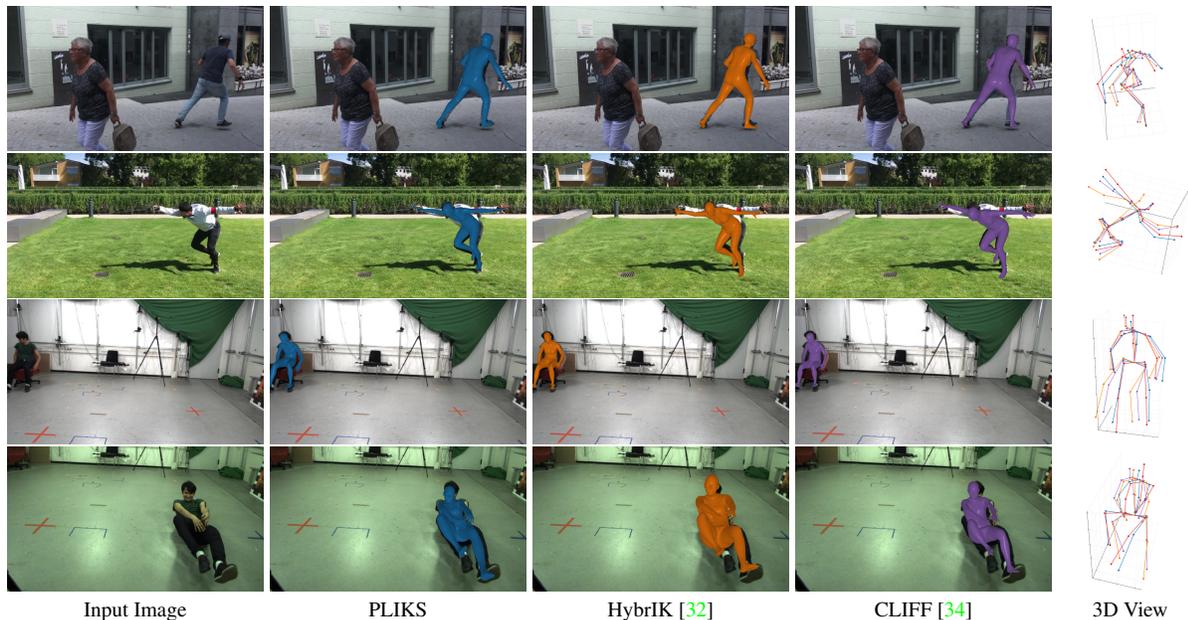


Figure 4. Qualitative results from 3DPW (top 2), MPI-INF-3DHP (bottom 2) datasets. The last column shows our shows the pelvis centred 3D joint locations. Here, the ground truth joint locations are represented in red with PLIKS in blue, HybrIK in orange and CLIFF in purple.

Module	Camera	3DPW		Human3.6M		MPI-INF-3DHP	
		MPJPE↓	PA-MPJPE↓	MPJPE↓	PA-MPJPE↓	MPJPE↓	PA-MPJPE↓
ARE	(wP)	83.6	50.5	53.6	39.4	94.6	61.4
ARE(R50)	(wP)	87.1	53.2	54.5	40.5	95.9	61.7
PLIKS	(wP)	85.5	50.6	51.3	36.2	93.4	63.2
PLIKS	$f=\sqrt{W^2+H^2}$ (P)	81.8	51.7	47.8	34.8	83.2	61.4
PLIKS	Known (P)	81.8	51.9	48.9	34.8	76.9	60.6

Table 7. Ablation studies by varying network setting. Here, (wP) and (P) refers to the weak-perspective and perspective camera model respectively.

occlusion test, HRNet performs better than ResNet. Next, we evaluate the importance of accounting for camera intrinsics by training a network with the PLIKS module but with all cropped images given as input to the network having a fixed focal length of $f=0.3$ m. The low focal length value violates the assumption of a weak-perspective setting. We observe slight improvements compared to that of ARE on the Human3.6M, whereas we see on-par results with the other datasets. Finally, we determine the effect of using a focal length of $f=\sqrt{W^2+H^2}$ proposed in [26] only during inference. While there is no significant difference in 3DPW and Human3.6M, we observe a drop in performance for the MPI-INF-3DHP dataset compared to PLIKS with known camera intrinsics. This is due to the actual FOV being larger than the estimated FOV. Our ablations clearly demonstrate the importance of incorporating camera intrinsics into a network.

5. Conclusion

In this paper, we bridge the gap between 2D correspondence and body mesh estimation by creating a pipeline from

which the inverse kinematics can be solved in closed form. This can be effectively leveraged to use the full perspective projection rather than having to rely on weak-perspective counterparts. Our approach yields considerable improvements in both root-relative and absolute-3D estimation for human pose estimation. PLIKS further enables our method to be fully differentiable facilitating end-to-end training. We validated the effectiveness of our method on various 3D pose and shape datasets and achieved state-of-the-art on multiple benchmarks. Due to the inherent nature of a built-in solver, we can extend our work with additional constraints like using multi-view systems for even higher reconstruction accuracy or temporal constraints for smooth video reconstruction should be feasible and foster future research in this direction.

Disclaimer The concepts and information presented in this article are based on research and are not commercially available.

References

- [1] Abien Fred Agarap. Deep learning using rectified linear units (relu), 2018. cite arxiv:1803.08375Comment: 7 pages, 11 figures, 9 tables. [4](#)
- [2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014. [12](#), [17](#)
- [3] Timothy Barfoot, James R. Forbes, and Paul T. Furgale. Pose estimation using linearized rotations and quaternion algebra. *Acta Astronautica*, 68(1):101–112, 2011. [3](#)
- [4] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In *European conference on computer vision*, pages 561–578. Springer, 2016. [2](#)
- [5] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Pose2mesh: Graph convolutional network for 3d human pose and mesh recovery from a 2d human pose. In *European Conference on Computer Vision (ECCV)*, 2020. [2](#)
- [6] Hongsuk Choi, Gyeongsik Moon, JoonKyu Park, and Kyoung Mu Lee. Learning to estimate robust 3d human mesh from in-the-wild crowded scenes. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [2](#)
- [7] Qi Fang, Qing Shuai, Junting Dong, Hujun Bao, and Xiaowei Zhou. Reconstructing 3d human pose by watching humans in the mirror. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12814–12823, 2021. [2](#)
- [8] Georgios Georgakis, Ren Li, Srikrishna Karanam, Terrence Chen, Jana Košecká, and Ziyang Wu. Hierarchical kinematic human mesh recovery. In *European Conference on Computer Vision*, pages 768–784. Springer, 2020. [2](#)
- [9] G. H. Golub and V. Pereyra. The differentiation of pseudo-inverses and nonlinear least squares problems whose variables separate. *SIAM Journal on Numerical Analysis*, 10(2):413–432, 1973. [5](#)
- [10] P. J. Green. Iteratively reweighted least squares for maximum likelihood estimation, and some robust and resistant alternatives. *Journal of the Royal Statistical Society. Series B (Methodological)*, 46(2):149–192, 1984. [5](#)
- [11] Riza Alp Guler and Iasonas Kokkinos. Holopose: Holistic 3d human reconstruction in-the-wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10884–10894, 2019. [2](#)
- [12] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7297–7306, 2018. [2](#)
- [13] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, New York, NY, USA, 2 edition, 2003. [5](#)
- [14] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. [7](#)
- [15] Nikolas Hesse, Sergi Pujades, Javier Romero, Michael J. Black, Christoph Bodensteiner, Michael Arens, Ulrich G. Hofmann, Uta Tacke, Mijna Hadders-Algra, Raphael Weinberger, Wolfgang Müller-Felber, and A. Sebastian Schroeder. Learning an infant body model from RGB-D data for accurate full body motion analysis. In *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Springer, 2018. [14](#)
- [16] David C. Hogg. Model-based vision: a program to see a walking person. *Image Vis. Comput.*, 1:5–20, 1983. [1](#)
- [17] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Trans. Pattern Anal. Mach. Intell.*, 36(7):1325–1339, 2014. [6](#), [12](#)
- [18] Umar Iqbal, Pavlo Molchanov, and Jan Kautz. Weakly-supervised 3d human pose learning via multi-view images in the wild. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5242–5251, 2020. [2](#)
- [19] Umar Iqbal, Kevin Xie, Yunrong Guo, Jan Kautz, and Pavlo Molchanov. Kama: 3d keypoint aware body mesh articulation. *2021 International Conference on 3D Vision (3DV)*, pages 689–699, 2021. [2](#)
- [20] Glenn Jocher, Ayush Chaurasia, Alex Stoken, Jirka Borovec, NanoCode012, Yonghye Kwon, TaoXie, Kalen Michael, Jiacong Fang, Imyhxy, , Lorna, Colin Wong, (Zeng Yifu), Abhiram V, Diego Montes, Zhiqiang Wang, Cristi Fati, Jebastin Nadar, Laughing, UnglvKitDe, Tkianai, YxNONG, Piotr Skalski, Adam Hogan, Max Strobel, Mrinal Jain, Lorenzo Mammana, and Xylieong. ultralytics/yolov5: v6.2 - yolov5 classification models, apple m1, reproducibility, clearml and deci.ai integrations, 2022. [7](#), [13](#)
- [21] Hanbyul Joo, Natalia Neverova, and Andrea Vedaldi. Exemplar fine-tuning for 3d human pose fitting towards in-the-wild 3d human pose estimation. In *3DV*, 2020. [1](#), [6](#), [7](#), [12](#), [13](#)
- [22] W. Kabsch. A solution for the best rotation to relate two sets of vectors. *Acta Crystallographica Section A*, 32(5):922–923, 1976. [4](#)
- [23] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7122–7131. IEEE Computer Society, 2018. [1](#), [2](#), [6](#), [12](#)
- [24] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. [13](#)
- [25] Thomas N. Kipf and Max Welling. Semi-Supervised Classification with Graph Convolutional Networks. In *Proceedings of the 5th International Conference on Learning Representations, ICLR '17, 2017*. [4](#)
- [26] Imry Kissos, Lior Fritz, Matan Goldman, Omer Meir, Eduard Oks, and Mark Kliger. Beyond weak perspective for monocular 3d human pose estimation. In *Computer Vision – ECCV 2020 Workshops: Glasgow, UK, August 23–28*,

- 2020, *Proceedings, Part II*, page 541–554, Berlin, Heidelberg, 2020. Springer-Verlag. [1](#), [2](#), [6](#), [8](#)
- [27] Muhammed Kocabas, Chun-Hao P. Huang, Otmar Hilliges, and Michael J. Black. PARE: Part attention regressor for 3D human body estimation. In *Proceedings International Conference on Computer Vision (ICCV)*, pages 11127–11137. IEEE, Oct. 2021. [1](#), [2](#), [6](#), [7](#), [12](#)
- [28] Muhammed Kocabas, Chun-Hao P. Huang, Joachim Tesch, Lea Müller, Otmar Hilliges, and Michael J. Black. SPEC: Seeing people in the wild with an estimated camera. In *Proc. International Conference on Computer Vision (ICCV)*, pages 11035–11045, Oct. 2021. [1](#), [2](#), [7](#), [13](#)
- [29] Tatsuro Koizumi and William A. P. Smith. “look ma, no landmarks!” – unsupervised, model-based dense face alignment. Berlin, Heidelberg, 2020. Springer-Verlag. [5](#)
- [30] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. [1](#), [2](#), [6](#), [7](#), [12](#)
- [31] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *CVPR*, 2019. [2](#)
- [32] Jiefeng Li, Chao Xu, Zhicun Chen, Siyuan Bian, Lixin Yang, and Cewu Lu. Hybrik: A hybrid analytical-neural inverse kinematics solution for 3d human pose and shape estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3383–3393, 2021. [2](#), [6](#), [7](#), [8](#), [12](#)
- [33] Runze Li, Srikrishna Karanam, Ren Li, Terrence Chen, Bir Bhanu, and Ziyang Wu. Learning local recurrent models for human mesh recovery. In *2021 International Conference on 3D Vision (3DV)*, pages 555–564. IEEE, 2021. [2](#)
- [34] Zhihao Li, Jianzhuang Liu, Zhensong Zhang, Songcen Xu, and Youliang Yan. Cliff: Carrying location information in full frames into human pose and shape estimation. *arXiv preprint arXiv:2208.00571*, 2022. [1](#), [3](#), [6](#), [7](#), [8](#)
- [35] Kevin Lin, Lijuan Wang, and Zicheng Liu. End-to-end human pose and mesh reconstruction with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1954–1963, 2021. [2](#)
- [36] Kevin Lin, Lijuan Wang, and Zicheng Liu. Mesh graphormer. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12919–12928, 2021. [2](#), [6](#), [7](#)
- [37] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2014. cite arxiv:1405.0312Comment: 1) updated annotation pipeline description and figures; 2) added new section describing datasets splits; 3) updated author list. [6](#), [12](#), [17](#)
- [38] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM Trans. Graphics (Proc. SIGGRAPH Asia)*, 34(6):248:1–248:16, Oct. 2015. [1](#), [2](#), [3](#)
- [39] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3DV*, pages 506–516. IEEE Computer Society, 2017. [3](#), [6](#), [7](#), [12](#)
- [40] Dushyant Mehta, Oleksandr Sotnychenko, Franziska Mueller, Weipeng Xu, Srinath Sridhar, Gerard Pons-Moll, and Christian Theobalt. Single-shot multi-person 3d pose estimation from monocular rgb. In *2018 International Conference on 3D Vision (3DV)*, pages 120–130. IEEE, 2018. [6](#), [7](#), [12](#), [17](#)
- [41] Gyeongsik Moon, Juyong Chang, and Kyoung Mu Lee. Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In *The IEEE Conference on International Conference on Computer Vision (ICCV)*, 2019. [2](#), [6](#), [7](#)
- [42] Gyeongsik Moon, Hongsuk Choi, and Kyoung Mu Lee. Accurate 3d hand pose estimation for whole-body 3d human mesh estimation. In *Computer Vision and Pattern Recognition Workshop (CVPRW)*, 2022. [7](#)
- [43] Gyeongsik Moon and Kyoung Mu Lee. I2l-meshnet: Image-to-lixel prediction network for accurate 3d human pose and mesh estimation from a single rgb image. *ArXiv*, abs/2008.03713, 2020. [2](#), [3](#), [4](#), [6](#), [12](#)
- [44] Chigozie Nwankpa, W. Ijomah, Anthony Gachagan, and Stephen Marshall. Activation functions: Comparison of trends in practice and research for deep learning. 12 2020. [4](#)
- [45] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *2018 international conference on 3D vision (3DV)*, pages 484–494. IEEE, 2018. [2](#)
- [46] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. In *NIPS 2017 Workshop on Autodiff*, 2017. [12](#)
- [47] Priyanka Patel, Chun-Hao P. Huang, Joachim Tesch, David T. Hoffmann, Shashank Tripathi, and Michael J. Black. Agora: Avatars in geography optimized for regression analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13468–13478, June 2021. [6](#), [12](#), [13](#), [14](#)
- [48] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed A. A. Osman, Dimitrios Tzionas, and Michael J. Black. Expressive body capture: 3D hands, face, and body from a single image. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 10975–10985, 2019. [1](#)
- [49] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 459–468, 2018. [2](#)
- [50] Liliana Lo Presti and Marco La Cascia. 3d skeleton-based human action classification: A survey. *Pattern Recognition*, 53:130–147, May 2016. [1](#)
- [51] Rene Ranftl and Vladlen Koltun. Deep fundamental matrix estimation. In *The European Conference on Computer Vision (ECCV)*, 2018. [5](#)

- [52] Gregory Rogez, Philippe Weinzaepfel, and Cordelia Schmid. LCR-Net: Localization-Classification-Regression for Human Pose. In *CVPR*, 2017. 6, 7
- [53] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. Synthetic training for accurate 3d human pose and shape estimation in the wild. In *British Machine Vision Conference (BMVC)*, September 2020. 2
- [54] Jiajun Su, Chunyu Wang, Xiaoxuan Ma, Wenjun Zeng, and Yizhou Wang. Virtualpose: Learning generalizable 3d human pose models from virtual data. *ArXiv*, abs/2207.09949, 2022. 6
- [55] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. 3, 7, 12
- [56] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Black Michael J., and Tao Mei. Monocular, one-stage, regression of multiple 3d people. In *ICCV*, 2021. 1, 2, 6
- [57] Yu Sun, Wu Liu, Qian Bao, Yili Fu, Tao Mei, and Michael J Black. Putting people in their place: Monocular regression of 3d people in depth. In *CVPR*, 2022. 7, 14, 15, 16
- [58] Hsiao-Yu Tung, Hsiao-Wei Tung, Ersin Yumer, and Katerina Fragkiadaki. Self-supervised learning of motion capture. *Advances in Neural Information Processing Systems*, 30, 2017. 2
- [59] Timo von Marcard, Roberto Henschel, Michael J. Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3D human pose in the wild using IMUs and a moving camera. In *European Conference on Computer Vision (ECCV)*, volume Lecture Notes in Computer Science, vol 11214, pages 614–631. Springer, Cham, Sept. 2018. 6, 7, 12, 13, 17
- [60] Yuanlu Xu, Song-Chun Zhu, and Tony Tung. Denserac: Joint 3d pose and shape estimation by dense render-and-compare. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7759–7769, 2019. 2
- [61] Kaibing Yang, Renshu Gu, Maoyu Wang, Masahiro Toyoura, and Gang Xu. Lasor: Learning accurate 3d human pose and shape via synthetic occlusion-aware data and neural mesh rendering. *IEEE Transactions on Image Processing*, 31:1938–1948, 2022. 2
- [62] Frank Yu, Mathieu Salzmann, Pascal Fua, and Helge Rhodin. Pcls: Geometry-aware neural reconstruction of 3d pose with perspective crop layers. *CVPR*, 2021. 2, 3
- [63] Wang Zeng, Wanli Ouyang, Ping Luo, Wentao Liu, and Xiaogang Wang. 3d human mesh regression with dense correspondence. In *CVPR*, 2020. 2
- [64] Hongwen Zhang, Yating Tian, Yuxiang Zhang, Mengcheng Li, Liang An, Zhenan Sun, and Yebin Liu. Pymaf-x: Towards well-aligned full-body model regression from monocular images. *arXiv preprint arXiv:2207.06400*, 2022. 7
- [65] Tianshu Zhang, Buzhen Huang, and Yangang Wang. Object-occluded human shape and pose estimation from a single color image. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 6, 7, 12, 17

Appendix

6. Introduction

In this material, we provide implementation details and analysis of focal lengths and regularizers for our method. We further discuss the benefits of using a solver for human pose estimation utilizing constraints. Additionally, we present more qualitative results, to show the performance of PLIKS and to explore its failure scenarios.

6.1. Datasets

COCO: COCO [37] is a large-scale in-the-wild 2D key-point dataset. We use this for training. We make use of pseudo-ground truth SMPL annotations provided by EFT [21].

MPI-INF-3DHP: MPI-INF-3DHP is an indoor multi-view and outdoor scene dataset for 3D human pose estimation. We make use of SMPL multi-view fits by SPIN [30]. We use this for training and evaluation.

Human3.6M: Human3.6M [17] is an indoor, multi-view 3D human pose estimation dataset. We follow the standard practice [23, 30] where subjects S1, S5, S6, S7, and S8 are used for training while S9 and S11 are the test subjects. We follow Protocol 2 using only the front-facing cameras.

3DPW: 3DPW [59] is a challenging outdoor benchmark for 3D pose and shape estimation. To get a fair comparison with previous state-of-the-art [27, 32], we use 3DPW training data for 3DPW experiments. We make use of a subset of this dataset 3DPW-OCC following [27] for the occlusion benchmark.

AGORA: AGORA [47] is a synthetic dataset with accurate SMPL models fitted to 3D scans. The test set is not publicly available, here the evaluation is performed on the official platform. For both training and testing, we use the images of resolution 1280×720 .

3DOH: 3DOH [65] is an object-occluded dataset. We use this to train and evaluate only for occlusion benchmark.

MuPoTs-3D: MuPoTs-3D [40] is a mixed indoor and outdoor multi-person dataset consisting of 20 sequences showing people performing various actions and interactions. We use this for evaluating the absolute root error.

6.2. Network Training

As the entire pipeline is differentiable, the network is trained end-to-end. We split the training into two steps, pre-training (ARE) and training (PLIKS) to accelerate the network training speed. In pre-training, we train exclusively with the ARE module, and optimize only with respect to the mesh and network predicted parameters $(\tilde{\beta}, \tilde{\theta}_k)$ by minimizing,

$$L = \omega_1 L_\theta + \omega_2 L_\beta + \omega_3 L_{2d} + \omega_4 L_{3d} + \omega_5 L_M. \quad (8)$$

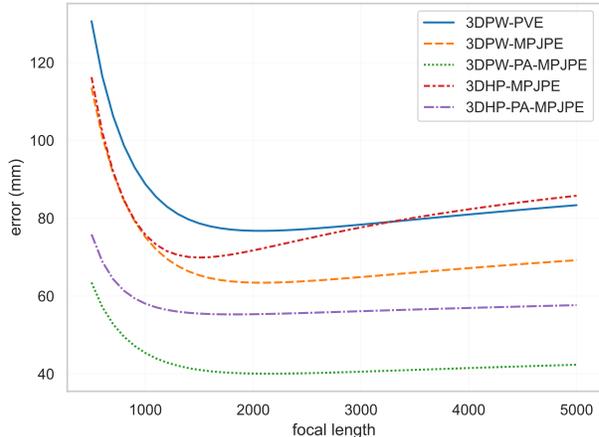


Figure 5. Impact on focal length on estimation errors when using the 3DPW [59] and MPI-INF-3DHP [39] dataset.

Following previous work [23, 30, 43], we employ standard mesh losses to supervise the training process. Here, L_θ is the L2 loss between the predicted pose and ground truth (GT) pose. Similarly, L_β is the L2 loss between the predicted shape and GT shape. L_{2d} , L_{3d} and L_M are the L1 loss between predictions and GT 2D joint re-projection, 3D joints and, the mesh vertex in image space respectively. To supervise the 2D annotations, the predicted 3D joints are projected by the weak-perspective camera \tilde{c} as predicted by the network.

During training we make use of the PLIKS module. Due to the presence of the linear solver in PLIKS, we observe numerical instability in the early stages of training, i.e. the pixel-aligned vertex predictions are not adequately consistent for the solver, making the reconstruction ill-posed. To keep the error within bounds, we add strong shape and pose regularizers for two epochs. In this stabilization period, the shape regularizer ω_β exponentially decays from 1 to 0.1. We further add a pose-constraint to the objective function of PLIKS (Eq. (9)), such that $\omega_\theta \sum |\Delta \mathbf{R}_k| \approx \mathbf{I}$. As a consequence, the additional rotation $\Delta \mathbf{R}_k$ obtained during the stabilization period is constrained to be close to zero. Similar to ω_β , we decay ω_θ from 1 to 0. For training, we use the same objective function from Eq. (8) to minimize the mesh and the analytically predicted parameters (β, θ_k) .

$$\operatorname{argmin}_{\Delta \mathbf{R}_k, \beta, \theta_k} \left\| \mathbf{w}^k \left(\mathbf{i}^k - \hat{\mathbf{K}}(\Delta \mathbf{R}_k \mathbf{x}_r^k + \beta \mathbf{B}_r^k + \mathbf{t}_k \mathbf{W}_r^k) \right) \right\|_2 + \omega_\beta \|\beta\|_2. \quad (9)$$

6.3. Implementation Details

PyTorch [46] is used for implementation. For all our experiments we initialize the HRNet [55] backbone with weights pre-trained on the MPII [2] dataset, which exhibits

faster convergence during training. We use the Adam optimizer [24] with a mini-batch size of 32. The learning rate at pre-training is set to $1e^{-4}$, whereas, while training the entire pipeline it is initialized to $5e^{-5}$. The network is pre-trained for 20 epochs, stabilized for 2 epochs, and then finally trained for further 30 epochs. We set the learning rate to $1e^{-5}$ while fine-tuning with the 3DPW [59] or AGORA [47] dataset. For fine-tuning, we use the previous pre-trained network as the starting point. This is to accelerate convergence and correct the 3D inaccuracies from the pseudo-GT labels. It takes around 3-5 days to train on a single NVIDIA Tesla V-100-16GB GPU. We set ω_1 , ω_2 , ω_3 , ω_4 , and ω_5 to 1, 0.05, 4, 8, and 4, respectively. As the pseudo-GT labels from EFT [21], are defined with respect to weak-perspective projection, we reduce ω_1 , ω_2 , ω_4 , and ω_5 by a factor of 0.1 for the 2D dataset.

7. Ablations

Here we discuss the effects of shape regularizer and effects of focal length estimation.

7.1. Regularizer

To demonstrate the importance of a strong regularizer, we perform a similar experiment (from Sec 4.1) where we add random noise to the GT of the mesh vertices from the 3DPW [59] test set. Here we vary the shape regularizer weights ω_β and observe the final MPJPE obtained. From Table 8, it is evident that larger weights for ω_β is more robust to noise. However, training the network using larger weights has its own drawbacks as shown in Figure 6. The network forces the shape components β to always be close to zero. As the shape β is determined by a solver, it enables us to switch to a male, female, or neutral model seamlessly by replacing the shape coefficients \mathbf{B} during inference. For our training, we set $\omega_\beta = 0.1$, as this is a good mixture between stability and shape variations.

7.2. Focal Length

We conduct experiments on the 3DPW and MPI-INF-3DHP test sets by varying the focal lengths. As shown in Fig. 5, PLIKS is robust to a wide range of focal lengths when the FOV is small (e.g., 3DPW), but it suffers from the effects of perspective warping on large focal lengths for wide FOV images (e.g., MPI-INF-3DHP). Using CamCalib [28] on the MPI-INF-3DHP to determine the FOV and consequently the focal length of the image, we could only obtain a reduction in MPJPE of 72.01 mm, i.e., a drop of just 3%. In particular, when the ground truth camera matrix is known, our approach can be expected to yield optimal performance.

	± 10 mm	± 20 mm	± 30 mm
$\omega_\beta = 2.0$	17.3	22.6	34.9
$\omega_\beta = 1.0$	18.4	37.3	64.0
$\omega_\beta = 0.1$	122	254	322.1

Table 8. Ground truth errors in the presence of per vertex noise ranging from $\pm 10mm$ to $\pm 30mm$ and the effect of using a shape regularizer, ω_β .

8. Qualitative Results

In this section, we show comparisons to SOTA methods on AGORA and provide more qualitative results on various other datasets.

8.1. Qualitative Comparison

We display several examples of PLIKS on the AGORA test set in Fig. 7. We use YOLO [20] for the bounding box estimation and CamCalib [28] for the focal length estimation. The images demonstrate that PLIKS performs better than previous approaches, by aligning the bodies well in 3D as well as 2D.

8.2. Inference Modification

As mentioned in the main paper, one of the strengths of our method is the application of constraints during inference. Here, we discuss a proof-of-concept for two use cases, where we show the benefits of using a solver without any retraining of the network. We discuss dynamic shape and translation constraints.

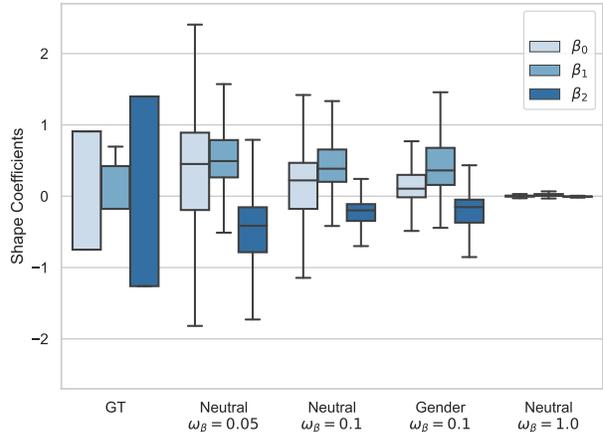


Figure 6. Effect of the regularizer weight used during training on the shape coefficients. Picking a higher ω_β reduces the error, but causes the network to output meshes getting progressively closer to the identity representation. Here neutral represents the neutral SMPL model, and gender refers to the gender-specific model on the 3DPW [59] dataset.

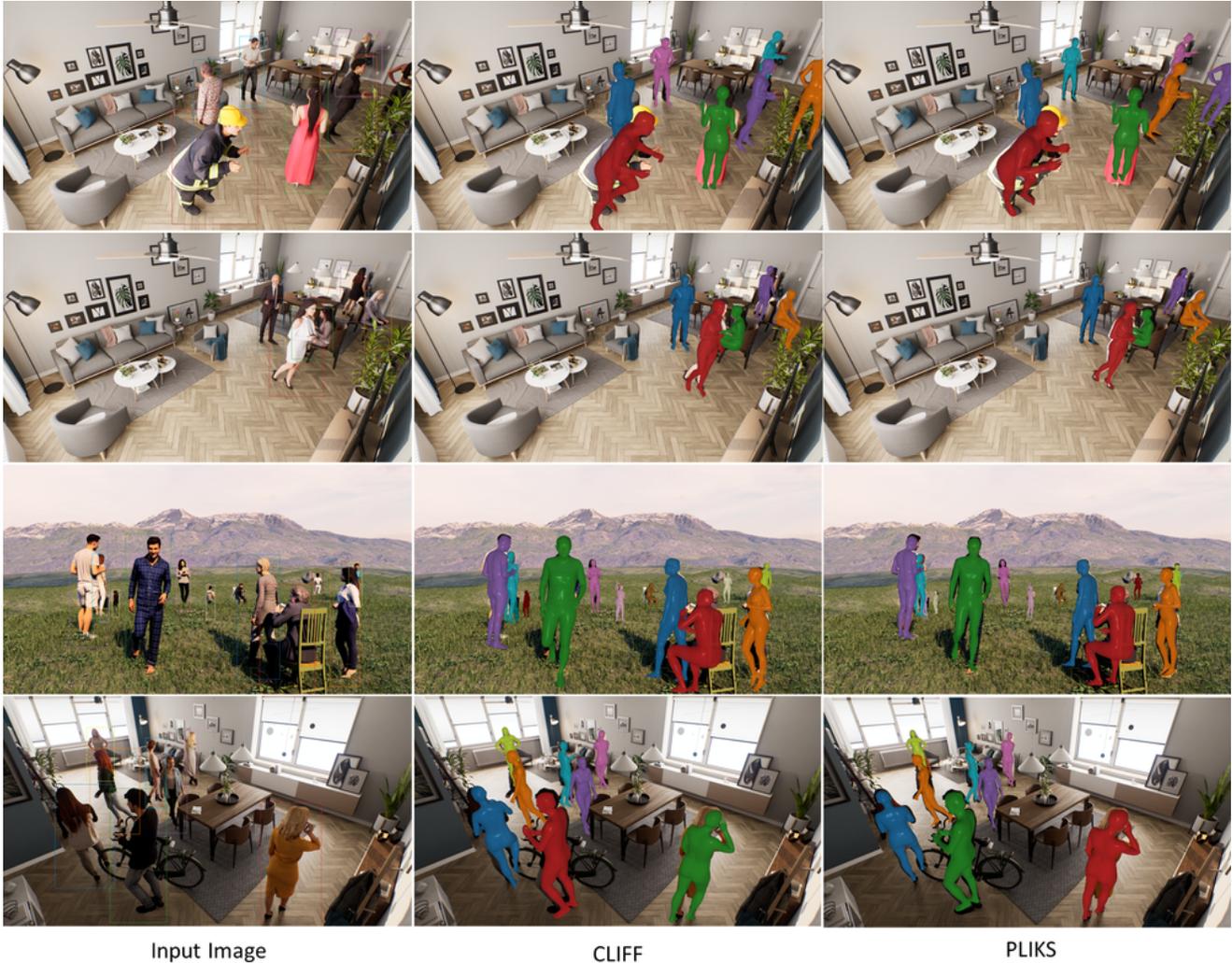


Figure 7. Qualitative results from AGORA test set.

Dynamic Shape Although our network was trained only on a neutral SMPL model with 10 shape components, it can make use of other shape models during inference if they follow the same design principle as SMPL. As an extreme scenario, we show the application using the kid-SMPL model [47, 57]. The kid-SMPL is an extended version of the SMPL model supporting children by linearly blending the SMPL and Skinned Multi-Infant Linear Model (SMIL) [15] by a weighting factor $\alpha \in (0, 1)$ [47]. Here, larger weights represents infants, while smaller weights are associated with adults. For simplicity, we denote the kid-SMPL model as having 11 shape components.

Qualitative results of using the kid-SMPL model on the Relative Human (RH) dataset [57] are shown in Fig. 8. The only modification performed was adapting the shape coefficients \mathbf{B}_r^k in Eq. 9 from the SMPL to their kid-SMPL counterparts. In that context, we further empirically set ω_β to 0.5. From the RH dataset we employ the GT age classifier,

i.e., we use SMPL for adults, and kid-SMPL for child or infant. We observe visually satisfactory results, with sufficiently reliable depth reasoning. A top-down approach [57] or a simple age classifier could be designed to determine the age as a future work.

Translation Constraints Previous examples of just using dynamic shapes is not a complete solution, due to the ill-posed nature of the problem. This is quite evident from the fifth column of Fig. 9. As a proof-of-concept, we show the application of translation constraints during inference. We add a simple depth constraint to Eq. 9 as $\omega_t t_{0,z}^k = \omega_t t_{0,z}^a$. Here, $t_{0,z}^a$ is the root depth of the adult in the image, and $t_{0,z}^k$ is the constrained setting for the root-depth of the kids in the image, with ω_t being a weighting factor. We make the assumption that the children in the images are standing close to the adults. The solver optimizes the shape such



Figure 8. Example images with dynamic shape during inference. Set of input images, overlay and, 3D view.

that the translation constraint is satisfied. We empirically set ω_t to 0.2. Though, strictly not comparable, we visualize the results of BEV [57] in Fig. 9. There, all images are in fact from the RH training set on which BEV was trained. We quickly add that this is not a real-world solution to the problem, but it emphasizes the importance of using constraints during inference or training. As future work, one could make use of the RH dataset with the depth-level information by adding a top-down approach [57] for better constraints.

8.3. Failure Mode

In Fig. 10, we show a few examples where PLIKS fails to reconstruct reasonable human body poses. The failure cases range from (a) too many people in the crop, (b) extreme poses not seen in training, and (c) extreme occlusion.

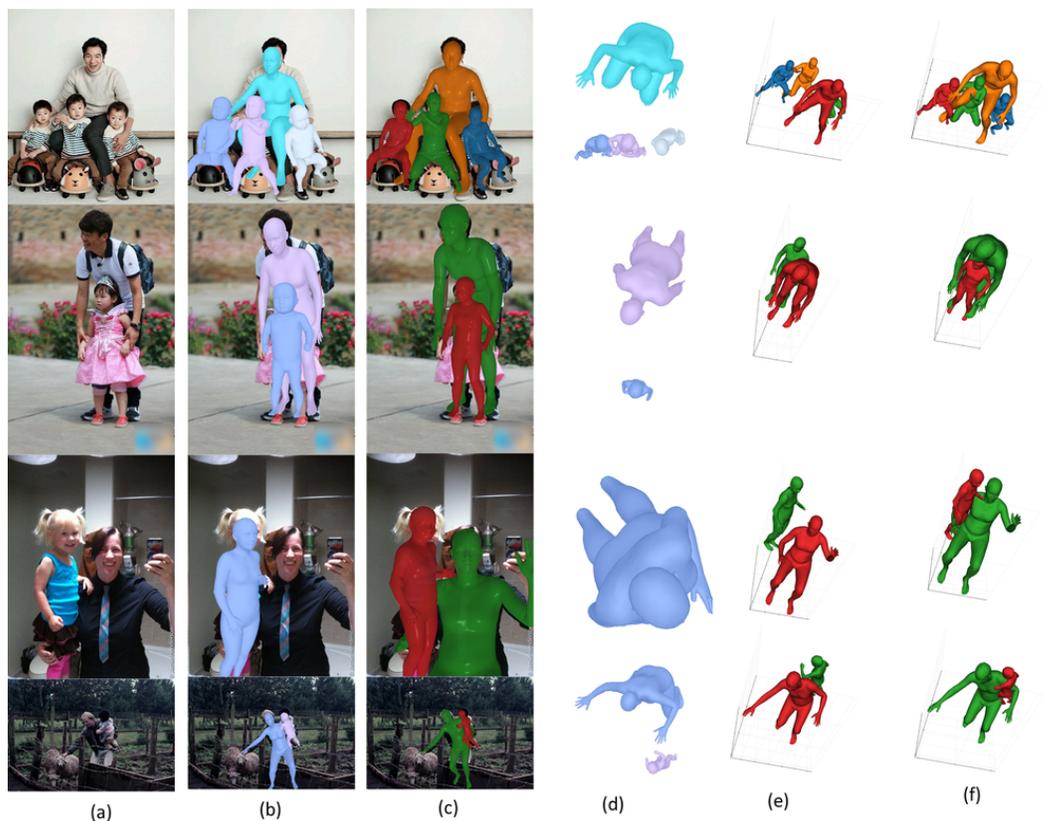


Figure 9. Example images with translation constraints during inference. (a) Input Image, (b,c) 3D overlay from BEV [57] and PLIKS respectively, (d) 3D view of the model from BEV [57], (e,f) 3D view of the model from PLIKS without and with using the translation constraint.

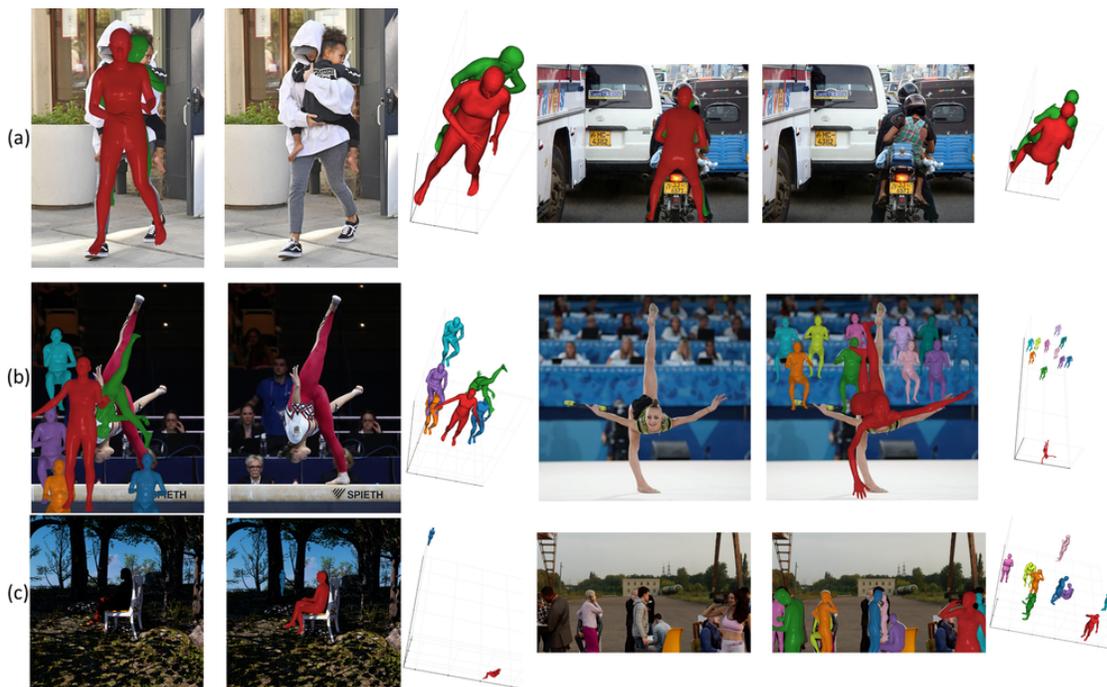


Figure 10. Example of failure cases.



Figure 11. Additional qualitative results of PLIKS from COCO [37], MPII [2], 3DPW [59], 3DOH [65] and MuPoTs-3D [40]. Set of challenging input images, overlay and, 3D view.