

# KBody: Towards general, robust, and aligned monocular whole-body estimation

Nikolaos Ziouli<sup>1</sup>, James F. O’Brien<sup>\*1,2</sup>

<sup>1</sup> Klothed Technologies Inc., <sup>2</sup> UC Berkeley

<https://klothed.github.io/KBody>

## Abstract

*KBody* is a method for fitting a low-dimensional body model to an image. It follows a predict-and-optimize approach, relying on data-driven model estimates for the constraints that will be used to solve for the body’s parameters. Acknowledging the importance of high quality correspondences, it leverages “virtual joints” to improve fitting performance, disentangles the optimization between the pose and shape parameters, and integrates asymmetric distance fields to strike a balance in terms of pose and shape capturing capacity, as well as pixel alignment. We also show that generative model inversion offers a strong appearance prior that can be used to complete partial human images and used as a building block for generalized and robust monocular body fitting. Project page: <https://klothed.github.io/KBody>.

## 1. Introduction









Machine perception of humans in images has seen remarkable progress recent years. This rapid advance has been the combined result of datasets like MS-COCO [44], the evolution of data-driven methods [24, 62], and modern parametric human body representations that are compact and continuously differentiable [46, 71]. Estimating the parameters of a dynamic human body is a cornerstone for human-centric applications such as virtual try-on based e-commerce [39], avatar creation for virtual presence [43], and performance analysis for virtual coaching [19]. Multi-view configurations offer robust estimations in challenging conditions [13]. This is a result of strongly-constraining the problem, and the same cannot be said for the ill-posed monocular case, which is nonetheless, the foundation of many consumer-facing products.

Despite the significant progress, high-quality monocular human body estimation in-the-wild remains elusive due to the challenges arising from the problem formulation itself and the limitations of available constraints. From an abstracted point of view, estimating the human body from a

\*Corresponding author: [james@getklothed.com](mailto:james@getklothed.com)

	SMPLify-X [51]	PyMAF-X [73]	SHAPY [14]	KBody (Ours)
Pose	✓	✓✓	✓	✓✓
Shape	✓	✗	✓✓	✓✓
Pixel	✓	✓	✗	✓✓

Full body				
				

**Figure 1.** Flexible, pixel aligned, accurate body pose and shape capture is the challenging, yet ultimate goal of monocular expressive body fitting. KBody is a general approach that improves the balance between all 3 traits using a *predict-and-optimize* approach while also gracefully handling partial images.

single image corresponds to estimating the articulation parameters  $\theta \in \mathbb{SO}(3)^P$ , the shape parameters  $\beta \in \mathbb{R}^B$  and the global transformation  $\mathbf{T} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$ . These parameters reconstruct the human mesh  $(\mathbf{V}, \mathbf{F}) = \mathcal{H}(\theta, \beta, \mathbf{T})$  via the body function  $\mathcal{H}$ . Two dominant classes of approach exist. The first fits the body by minimizing an objective [10, 51]:

$$\operatorname{argmin}_{\theta, \beta, \mathbf{T}} \mathcal{E}_{data} + \mathcal{E}_{prior}, \quad (1)$$

that includes a data fitting term,  $\mathcal{E}_{data}$ , and  $\mathcal{E}_{prior}$ , an important prior regularization term to prevent degenerate solutions and provide additional constraints to alleviate the ill-posedness of the problem. The constraints involved in the data term most typically include 2D keypoints [51], that are typically inferred by a data-driven method [12]. While the prior term helps, 2D keypoints usually lead to solutions that suffer from monocular ambiguity, producing poor re-

sults from a 3D accuracy perspective. The second class of approach consists of data-driven methods that encode a learned prior in the parameters,  $\chi$ , of a neural network,  $f$ , and perform monocular inference:

$$(\theta, \beta, \mathbf{T}, \pi) = f_{\chi}(\mathbf{I}), \quad (2)$$

with  $\pi$  being the – typically weak perspective / orthographic – projection parameters that best explain the image content using the estimated parameters. As the neural network function  $f_{\chi}$  is supervised, it preserves 3D awareness but usually suffers from predictions with poor pixel alignment, and bias due to the long tailed distribution of data [53].

Another challenge that also hinders high-quality pixel alignment is the conflict between the pose  $\theta$  and shape  $\beta$ , that are entangled though  $\mathcal{H}$ . Early works [15, 51] focused on the difficult problem of pose capturing foremost, with proper shape being an unaccomplished side-objective. Yet as progress was made, it became evident that inaccurate shape was hindering further advances, with more recent works [14, 16] focusing on higher quality shape capture, but seemingly, at the cost of poorer pose estimation.

In-the-wild images introduce additional challenges, some of which are only partly addressed by complex augmentation schemes [67, 68], and others, like missing information in partial human images, which is prevalent in some domains, are challenging to overcome. For fitting approaches the prior terms are not sufficient to regularize the optimization process when keypoints are missing, while data-driven methods can only extrapolate up to the training data distribution’s capacity. Overall, achieving pixel-aligned estimates that are metrically correct (in world scale, not up to an unknown scale factor), and doing so robustly for a wide range of inputs remains a significant challenge.

In this work, we present a general framework for estimating whole-body human parameters from a single image. Our goal is to deliver robust estimates, for a variety of inputs, while preserving pixel alignment and proper 3D estimations as much as possible, as well as to capture shape cues and pose information simultaneously, as seen in Fig. 1. More specifically, our contributions are the following:

- We improve fitting quality by introducing virtual joints, adapted to fit the estimated data, and allowing for smooth interplay with silhouette constraints, expressed as an asymmetric distance field. We additionally show how disentangling the optimization process allows for improved joint shape and pose estimates.
- We present an appearance prior based approach to handle images with missing information by completing them in a structurally plausible manner. Plausibility is enough for inferring constraints on the hallucinated parts which enable higher quality fits on partial images.

## 2. Related Work

Estimating parametric human models from images is a rapidly evolving area forming a complex landscape of data, models, and training strategies, as discussed in a recent survey [50] and benchmark [64] papers. Several parametric human body models, including STAR [48], GHUM [4, 71] and most recently SUPR [49] have been released, but we will focus on the expressive variant of SMPL, SMPL-X [51].

### 2.1. Single-shot estimation methods

Pioneering the transition from keypoint estimation to full-body estimation involved the direct regression of low-dimensional body parameters from a single image [28]. The method was supervised using keypoint annotations and thus, end-to-end training was achieved after also regressing the camera parameters that would project the articulated body joints to correct positions. Regularization was applied in the form of a discriminator for the estimated pose and shape, so as to match a realistic distribution made available as a corpus of fit human scans. Various extensions were later proposed, integrating inverse kinematics [40], topological priors [47], and external camera estimation [35] to improve pose estimation performance. While the latter two approaches use silhouettes in their training schemes, they remain an intermediate representation for skeletonization [47] or they include clothing layers [35].

Initial efforts only regressed pure body parameters (*i.e.* SMPL), which unfortunately disregards details like hands and faces. ExPose [15] included regressing parameters for the hands and face. FrankMoCap [55] built an efficient system, achieving real-time rates. ExPose was extended to PIXIE [18], which had separate experts for the body, hands and face that were optimally combined to improve results. More recently, PyMAF-X [73] builds on the iterative nature of these models (*e.g.* [15, 28]) but instead of using global features at a single scale, PyMAF-X uses a pyramid of features, including finer-grained ones, achieving higher quality pixel alignment than other approaches.

Taking another direction, SHAPY [14] focuses on shape estimation using model agency annotations for shape measurements. Having been trained with this supervision, it is capable of regressing metric-scale shapes. SHAPY’s pose estimation performance is not at the same level of PyMAF-X, but its capacity to output metric-scale shapes heavily compensates.

### 2.2. Iterative optimization methods

SMPLify [10] was the seminal work that fit the SMPL body to a single image, showing the effectiveness of having priors for both the pose and shape alike. SMPLify was later extended to use annotated silhouettes in its iterative optimization scheme, with the goal of improving dataset annotations [38]. Using an L1 silhouette objective allowed

for capturing human performances in video [23] using differentiable soft-rasterization [45, 52], and improved results when combined with a differentiable ray-tracer [42] and part-based masks [6]. In a follow up work, it was extended to SMPLify-X [51], adding details like hands and face, as well as a learned prior, VPoser [51]. Similarly, to improve shape capturing for use within forensic contexts [63], an L2 mask loss was added into the optimization scheme through a differentiable renderer [32].

While orthogonal improvements like better priors (e.g. Pose-NDF [65]) can improve fitting performance, results ultimately rely on the constraints  $\mathbf{k}$  and (optionally)  $\mathbf{S}$  [22, 33]. Another important component is the initialization of the optimization which can significantly affect convergence due to the ill-posedness of monocular fit. One solution [27] to this uses 3D keypoint estimates as constraints, and iteratively refines the estimate via forward kinematics.

Finally, a relatively recent and novel direction combines data-driven models and optimization techniques. HUND [72] learns a recurrent model that is learned to optimize a recurrent (initial) state and alignment errors iteratively, which proves to be faster than traditional optimization approaches. The same applies to LVD [16] that learns descent updates for each body vertex so as to predict the depicted human mesh. Both approaches are limited by their training data compared to other optimization techniques. Particularly so for LVD, which is trained on 3D human scans, a data category that is hard to acquire at scale.

Last, test-time optimization is a relatively new field that finetunes an entire model on a specific target sample using predicted constraints like keypoints and silhouettes. Through this technique and the use of separate model and parameter steps, as well as silhouette constraints, a recent work [41] has demonstrated improved shape estimation.

### 3. Approach

Our overall framework is presented in Fig. 2, and thematically split in two distinct stages, both following a *predict-then-optimize* approach. First, there is an optional completion preprocessing step on the left for use with partial images of humans. Second, is a process for high quality monocular fitting using 2D constraints on the right. While the illustration follows the processing order from left-to-right, we will first present the fitting on Sec. 3.1, and then the optional completion preprocessing step on Sec. 3.2.

#### 3.1. Pixel-aligned Shape-aware 3D Body Fitting

Similar to prior approaches, we fit a parametric body model to image-domain constraints by minimizing Eq. (1), using the same prior terms as SMPLify-X [51], but with a disentangled optimization process (Sec. 3.1.1), while also using virtual joints in the projected keypoints objective (Sec. 3.1.2), and adding a silhouette-based objective

(Sec. 3.1.3):

$$\mathcal{E}_{data} = \underbrace{\lambda_k(\mathcal{E}_{rj} + \mathcal{E}_{vj})}_{\text{keypoints}} + \underbrace{\lambda_m \mathcal{E}_{mask} + \lambda_d \mathcal{E}_{adf}}_{\text{silhouette}}, \quad (3)$$

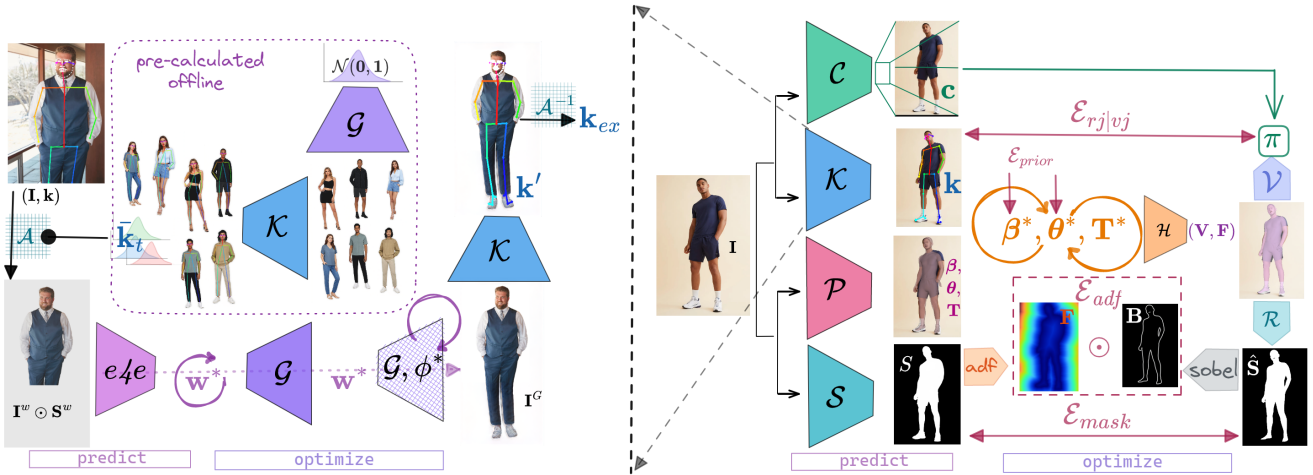
where  $\mathcal{E}_{rj|vj} = \varrho(\mathbf{k}, \pi(\mathbf{j}_{rj|vj}))$  is the Geman-McClure penalty function [21] for the regular joints,  $\mathbf{j}_{rj}$ , and virtual joints,  $\mathbf{j}_{vj}$ , matching them to the corresponding keypoints  $\mathbf{k}$  via the projection function  $\pi$  of given camera model. The parameters of the camera can be known (when available in image metadata), estimated (e.g. using a camera parameters estimation model), or fixed (when no information is available).  $\mathcal{E}_{mask} = \sum_{\Omega} \|\mathbf{S} - \hat{\mathbf{S}}\|_1$  is an L1 silhouette overlay term defined on the image domain  $\Omega$ , between an inferred silhouette  $\mathbf{S}$  and the body model’s silhouette  $\hat{\mathbf{S}} = \mathcal{R}(\mathbf{V}, \mathbf{F})$  rendered via a differentiable rendering function  $\mathcal{R}$ . The virtual joints calculation and the asymmetric distance field term,  $\mathcal{E}_{adf}$ , are described in the following subsections.

#### 3.1.1 Disentangled Optimization

Prior monocular human body fitting works perform a staged optimization of Eq.(1), where each stage adds a layer of complexity in the optimization (e.g. details like fingers), and also anneals the constraints’ weights [10, 51] across stages. Initial estimates of global parameters  $\mathbf{T}$  have also been included as a first stage [10, 51], but sensitivity to localisation of the torso joints has led to alternatives [38]. To overcome sensitivity to initialization we use a data-driven initial estimate which serves as a good initial starting point.

However, all prior work up to now optimize both  $\beta$  and  $\theta$  simultaneously at each iteration  $i$  of each stage  $s$ :  $(\beta_{i+1}^s, \theta_{i+1}^s) = (\beta_i^s + \Delta\beta_{s_i}, \theta_i^s + \Delta\theta_i^s)$ . These two sets of parameters are entangled by the human body function  $\mathcal{H}$  that allows for their joint optimization. While this is effective with a 3D objective that is conditioned on the same domain where the function  $\mathcal{H}$  exists, it is much less effective in the monocular 2D case that comes with inherent 3D ambiguities. As a result, optimization is dominated by the pose updates  $\Delta\theta$ . This imbalance is evident in both keypoint-only optimization approaches [51] as well as data-driven models trained with only keypoint losses [15, 36, 73]. Both tend to produce shape coefficients biased towards the zero mean vector. More recent shape-aware approaches either optimize in 3D [16] or use 3D losses during training [14].

Seeking to improve our optimization loop, we separate the parameter updates of the shape  $\beta$  and pose  $\theta$  components in an alternating fashion for stage  $s$ :  $(\beta_i^s, \theta_{i+1}^s) = (\beta_{i-1}^s + \Delta\beta_{i-1}^s, \theta_i^s + \Delta\theta_i^s)$ . Similar to block coordinate optimization, the shape  $\beta$  parameters are only updated in even iterations  $i$ , while the pose parameters  $\theta$  are only updated in odd iterations  $i + 1$ . This method exhibits significantly better joint optimization of these parameters even in the highly



**Figure 2.** The KBody framework considers 2 stages, an optional image-based body completion on the left, and a general body fitting on the right. Keypoints  $\mathbf{k}$ , silhouette  $\mathbf{S}$  and (optionally) camera  $\mathbf{c}$  constraints are predicted from the respective models  $\mathcal{K}$ ,  $\mathcal{S}$  and  $\mathcal{C}$ . Then, an initial state  $\beta, \theta, \mathbf{T}$  predicted by  $\mathcal{P}$  is iteratively optimized to fit these constraints using the rendering  $\mathcal{R}$ , virtual joint  $\mathcal{V}$ , and camera-conditioned projection  $\pi$  functions. When identifying partial keypoints  $\mathbf{k}$ , the optional step on the left produces extrapolated keypoints  $\mathbf{k}_{ex}$  to improve fits on partial images. After properly aligning the masked image  $\mathbf{I}^w \odot \mathbf{S}^w$  using  $\mathbf{k}$  and the distribution  $\bar{\mathbf{k}}_t$  expected by the generative model, an initial inversion vector  $\mathbf{w}$ , estimated by a single-shot inversion model  $e4e$ , is iteratively refined twice, first on the  $\mathcal{W}$  latent space and then on the manifold  $\mathcal{G}_\phi$  using the warped masked partial image as constraint.

ill-posed monocular case. However, it is critical that the global pose is close to the minima, meaning that such a disentangled optimization stage can only be introduced later in the optimization process. Alternatively, one could introduce scaling factors on the parameters and loss function so that  $\theta$  and  $\beta$  would be well balanced, but it is unclear how the scaling factors would be computed and then updated as the optimization progresses.

### 3.1.2 Virtual Joints

An iterative fitting approach crucially relies on high quality correspondences. Defining proper joint locations on the body to match the keypoint estimates has troubled past approaches, with the hip joints ignored from the optimization [51], or regressed via empirically defined and manually created joint regressor functions [36]. However, the location of the keypoints  $\mathbf{k}$  are typically inferred from a data-driven model which aggregates numerous annotations and thus, includes their biases as well. Recent works that acknowledge this have resorted to learning a joint regressor for a specific dataset [25] which comes with new challenges like properly constraining the joints' locations inside the human body.

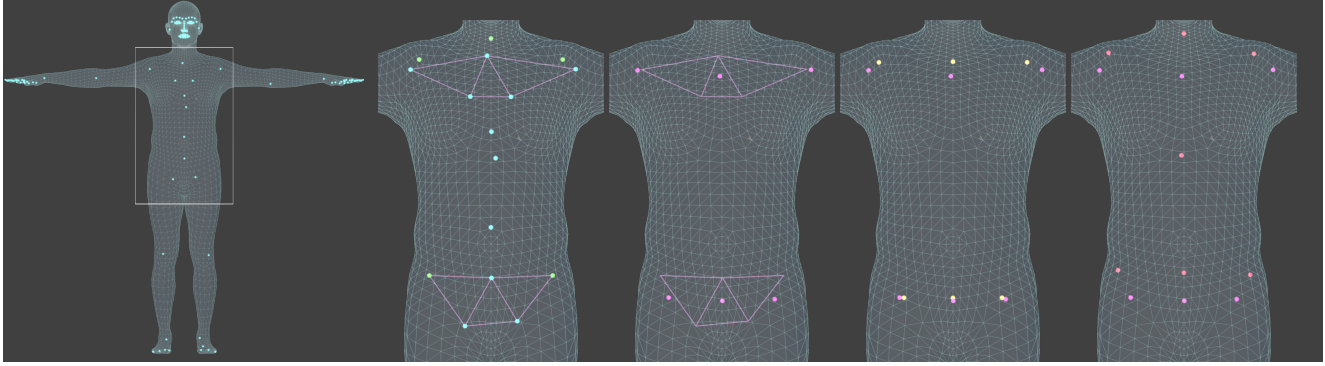
Our approach also seeks to identify better matching locations, but not for a specific dataset, instead matching the inference distribution of a pre-trained 2D keypoint estimator. We introduce the concept of *virtual joints*  $\mathbf{j}_{vj} = \mathcal{V}(\mathbf{b}, \mathbf{j}_s)$ , by parameterizing joint locations as a linear combination of weights  $\mathbf{b}$  and pre-defined (empirically or anthropomor-

phically) joint subsets  $\mathbf{j}_s, s \in [1, \dots, S]$ . More specifically, we focus on the more ambiguous torso joints, which carry a two-fold importance, **i)** they are high in the kinematic chain, and thus, highly influential of the articulated body fit, and **ii)** they are highly dependent on human shape, and thus, are necessary to avoid cross data-term conflicts between the keypoint and the silhouette terms.

Virtual joint localisation is restricted to planes formed by joint triangles (*i.e.*  $S = 3$ ), illustrated in Fig. 3, using a barycentric formulation for the virtual joints. This allows for the reduction of the number of weights  $\mathbf{b}$  to 2 (or 1 for joints that need to lie on one of the triangle's altitudes), by exploiting  $\sum_{b \in \mathbf{b}} b = 1$ . While this relies on a non-holding rigidity assumption for the joints subset, albeit relaxed in the torso area, the goal is to better localize joints matching those inferred by a 2D estimation model, which itself exhibits limited expressivity at the torso. Finding the best matching locations is an one-off process that involves fitting a variety of pre-defined poses to inferred keypoints and identifying the best performing weights using a performance indicator.

### 3.1.3 Asymmetric Distance Fields

Human body fitting requires both pose and shape parameters that eventually get mapped to 3D or 2D joints, the latter being a function of the reconstructed mesh vertices. Dense representations like silhouettes or distance maps have been used even in earlier parametric model fitting approaches



**Figure 3.** From left to right: **i)** the SMPL-X body surface and joints, **ii)** the inset torso with the barycentric parameterization comprising the triangles formed by raw and manually picked [36] joints, **iii)** our best-estimated virtual joints, and their comparison with **iv)** manually picked openpose joints [7, 8] and **v)** the learned regressor joints fit to Human3.6M [25]. As illustrated, the virtual joints can extrapolate to exterior triangle locations by using negative barycentric weights.

[5, 60], and lately in approaches involving both Eq.(1) or Eq.(2) to offer a less domain sensitive (proxy) representation for synthetic data training [9, 57, 58], topological objectives [47], pose refinement [23] and better shape (and pose) estimates [26, 38, 63].

Optimization approaches typically rely on differentiable rendering [32, 52] and a per-pixel L1/2 loss between a constraint input silhouette and the body rendered one. This loss is inefficient, suffering from an irregular loss landscape and the lack of directional information for parameter updates [47]. The approach is also highly susceptible to body and estimated silhouette inconsistencies, usually derived from hair, clothing, background mixing or inference uncertainty.

Instead our silhouette term supplements the per-pixel mask alignment objective with a boundary-based distance-field objective by first extracting the boundary  $\mathbf{B} = G \circledast \hat{\mathbf{S}}$  of the rendered fit body mesh silhouette  $\hat{\mathbf{S}}$  using convolutional edge extraction kernel  $G$  [61]. The result is used to derive a Chamfer distance objective using a distance field  $\mathbf{F}$  via the Hadamard product  $\mathcal{E}_{adf} = \sum_{\Omega} \mathbf{B} \odot \mathbf{F}$  summed over the pixel domain  $\Omega$  which is minimized when the two silhouette boundaries align. This is a more efficient alternative compared to nearest-neighbor queries [26, 38] but lacks the symmetric component of Chamfer distance, *i.e.* the distance from the inferred silhouette  $\mathbf{S}$  to the fit rendered one  $\hat{\mathbf{S}}$ , which is nonetheless noisier [3, 26, 38].

Indeed, for generalized fitting where hair, clothing and silhouette estimation artifacts come into play, the silhouette-based term tends to become noisy and hinder optimization or produce unrealistic shape estimates. To overcome this, we calculate an asymmetric distance field (ADF) defined over the entire image domain  $\Omega$ :

$$\mathbf{F} = \lambda_o D(\mathbf{S}) \odot \bar{\mathbf{S}} + \lambda_i D(\bar{\mathbf{S}}) \odot \mathbf{S}, \quad (4)$$

with  $D(\cdot)$  being the distance field function and  $\bar{\mathbf{S}}$  denotes pixel-wise binary inversion. While [3] downscales the nois-

ier symmetric Chamfer objective, we completely disregard it and provide explicit control over pushing the body inwards or outwards with respect to the silhouette by respectively controlling the outer and inner distance field scaling factors  $\lambda_o$  and  $\lambda_i$ . For blendshape models such as SMPL(-X), downweighting the inner field and/or upweighting the outer one, heavily restricts the body shape inside the silhouette while still allowing for greater freedom in not exactly matching the boundary in its entirety.

### 3.2. Structurally Plausible Human Completion

Generalized human body estimation is greatly challenged by partial human images, be it either data-driven estimates or optimization-based approaches, as the partial context and lack of annotated data reduce the prediction accuracy of single-shot estimates and the quality of the constraints for iterative fitting. We integrate partial human image completion in our fitting framework as an optional step which can be easily identified post keypoint estimation. The goal is to complement the partial image inferred keypoints  $\mathbf{k}_{in}$  with high confidence estimates for the invisible keypoints  $\mathbf{k}_{ex} = \mathbf{k}_c \setminus \mathbf{k}_{in}$ ,  $\mathbf{k}_c$  being the inferred keypoints on the completed image. This gracefully benefits the fitting process as the projection function  $\pi$  is not confined in the original image domain  $\Omega$ , and can optimize for the combined set  $\mathbf{k} = \mathbf{k}_{in} \uplus \mathbf{k}_{ex}$ .

While image inpainting could be a proper technique, we find that its generalized nature hurts the structural plausibility and quality of the results, mainly due to the extended nature of partial human image completions. Human specific solutions either only focus on visible body part completion [75] or rely on intermediate models producing high level extrapolated completions [70] as secondary inputs. Instead, we invert the partial images to the latent space of a generative appearance prior [31] learned using clothed humans [20], carrying more appropriate human structural

and shape-aware semantics, and producing high quality extended completions.

### 3.2.1 Partial image alignment

The StyleGAN variants [29–31] are a pure appearance-based family of generative models. They have been shown to be highly sensitive to affine transformations [2] even when considering faces whose articulation, and thus spatial variance, is limited compared to full body images. Therefore, for all StyleGAN face inversion techniques, centralization is a necessary preprocessing step. To overcome this challenge, which is pronounced when considering full bodies, StyleGAN-Human [20] investigated various alignment techniques and showed that mid-body alignment behaved better. In our case we seek to align partial images, making mid-body alignment not an option, and necessitating the design of a different alignment strategy.

Given that most partial images are either missing heads or the bottom half, and typically include some torso joints, we seek to identify the likely positions of these joints on generated samples. We generate  $M = 20000$  samples from a normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{1})$  using a truncation factor [11, 34]  $\psi = 0.5$  to enforce a proper distribution, and run the keypoint detector on each sample. This derives the expected (*i.e.* mean) location of these joints  $\bar{\mathbf{k}}_t$ . We can then align the corresponding torso joints of the partial image with these expected locations by calculating an affine transformation  $\mathbf{A} = \alpha(\bar{\mathbf{k}}_t, \mathbf{k}_t)$ , comprising a translation and scale. This in turn drives an image warping operation  $\mathbf{I}^w = \mathcal{A}(\mathbf{I} \odot \mathbf{S}, \alpha)$  that aligns the masked image to the StyleGAN-Human resolution  $(w, h) = (512, 1024)$  and the partial image content to the model’s expectation from an appearance reconstruction perspective.

### 3.2.2 Generative appearance-based completion

Using a pre-trained StyleGAN-Human model [20] our goal is to invert an aligned partial image  $\mathbf{I}^w$  into its intermediate latent space  $\mathbf{w} \in \mathcal{W}$  that reconstructs a full, completed image matching the partial image’s appearance. We initialize  $\mathbf{w}$  with a single-shot estimate from a pre-trained inversion model [66] (*e4e*). As the input image is partial, this only serves as a rough initialization into  $\mathcal{W}$ , which we further refine following [31] by minimizing:

$$\operatorname{argmin}_{\mathbf{w}} \lambda_{L1} \mathcal{E}_{L1} + \lambda_{VGG} \mathcal{E}_{VGG} + \lambda_{reg} \mathcal{E}_{reg}, \quad (5)$$

where  $\mathcal{E}_{L1}$  is a  $L1$  reconstruction loss,  $\mathcal{E}_{VGG}$  is a VGG [59] based perceptual feature loss, and  $\mathcal{E}_{reg}$  is a noise regularization term. These terms are described in [31] but in our implementation we adapt the pixel-based losses to focus only the original content of the aligned partial image. This is achieved by using the warped silhouette image  $\mathbf{S}^w$  to mask

the generated  $(\mathbf{I}^G \odot \mathbf{S}^w)$  and target images  $(\mathbf{I}^w \odot \mathbf{S}^w)$  prior to error term calculation.

Nonetheless this process is not sufficient to plausibly complete the image but is rather a way to quickly initialize the latent space. Plausible completion is achieved by relying on generator fine-tuning with latent space regularization [54]. Essentially, this process is a test-time optimization of the StyleGAN generator manifold around the latent space point  $\mathbf{w}^*$ . We optimize the generator  $\mathcal{G}$  parameters  $\phi$ :

$$\operatorname{argmin}_{\phi} \lambda_{L2} \mathcal{E}_{L2} + \lambda_{LPIPS} \mathcal{E}_{LPIPS} + \lambda_R \mathcal{E}_R + \lambda_D \mathcal{E}_D, \quad (6)$$

where we complement the LPIPS [74] and space regularization ( $\mathcal{E}_R$ ) terms used in [54] with an  $L2$  reconstruction term to accelerate convergence and a discriminator loss ( $\mathcal{E}_D$ ) to improve the global coherence of our results. Notably, the reconstruction and perceptual image domain losses are calculated only at the valid warped image  $\mathbf{I}^w$  domain, denoted by the mask  $\mathbf{M}^c$ , including the partial image’s background which we found to significantly improve convergence. Thus, only the areas to be completed are masked to not participate in error and gradient calculations. The high quality inversion capacity of the pivotal tuning technique ensures photometric convergence on the valid image regions. While image editing is not our goal, we find that the space regularization and discriminator terms help in producing structurally and photometrically plausible extended completions, acting as global regularizers.

## 4. Results

We refer to the approach depicted in Fig. 2 as KBody and implement it using SMPL-X [51] as the body model  $\mathcal{H}$ , OpenPose [12] as the 2D keypoint  $\mathbf{k}$  estimator model  $\mathcal{K}$ , MODNet [33] as the silhouette  $\mathbf{S}$  estimator  $\mathcal{S}$  after thresholding the estimated matte at 0.85, ExPose [15] as the initial body parameters  $(\beta, \theta, \mathbf{T})$  predictor  $\mathcal{P}$ , and the CamCalib model  $\mathcal{C}$ , presented in SPEC [35], as the camera parameter  $\mathbf{c}$  estimator. For the subsequent optimization we rely on the limited memory BFGS optimizer [69] with strong Wolfe line search and a budget of 30 iterations. Similar to prior work we perform annealed optimization with the early stages using stronger regularization to make the objective function more convex, and then progressively reduce the regularization term weights and increase the data terms of the details (hands, face). For the differentiable mesh rendering we use a high-performance rasterization based implementation [37]. For the pose prior and regularization terms we use the same as SMPLify-X [51], but relax the latter’s weights as the initialization and silhouette constraints provide extra prior knowledge about the pose and shape.

First we validate the effectiveness of the virtual joint localization by running only two stages of fitting to  $\mathcal{K}$  after

initializing with  $\mathcal{P}$ , with only the second stage optimizing the details, and without involving  $\mathcal{S}$  or  $\mathcal{C}$  for a fair comparison with other works. We use the EHF [51] dataset as well as a manually collected set of plain background human photos whose foreground masks are estimated in high-quality using a background removal service [1]. We run a hierarchical and empirically defined search to identify the parameters  $\mathbf{b}$  by fitting to the keypoints estimated by  $\mathcal{K}$ . Performance is measured using an indicator combining the keypoints’ RMSE and the IoU using the service generated masks, defined as  $(1 - IoU) \times RMSE$ . After estimating the barycentric coordinates  $\mathbf{b}$  resulting in the best fits, we conduct an experiment on EHF that is presented in Tab. 1. Performance is assessed via procrustes-aligned vertex-to-vertex error on the SMPL-X body’s vertices (PA-V2V-X) [51]. As also shown in Pose-NDF [65] and the first 3 rows of Tab. 1, simply optimizing the initial estimates of a data-driven model does not necessarily lead to improved fits. Using better priors like GAN-S [17] and Pose-NDF [65] slightly improves results over the baseline single-shot model ExPose [15], while a manually selected joint regressor [7, 8] (Fig. 3 iv) does not result in improved fits. The virtual joints produce the most significant gain, showcasing the importance of higher quality correspondences between the estimated keypoints used as constraints and body’s joints. It should be noted that, apart from the last 2 rows, bad joint-to-keypoint correspondences (e.g. hips) are ignored during optimization.

Initialization	Optimization	Joints	Prior	PA-V2V-X↓
✗	SMPLify-X [51]	$\mathbf{j}_{rj}$	VPoser [51]	60.3 mm
ExPose [15]	✗	$\mathbf{j}_{rj}$	✗	54.8 mm
ExPose [15]	SMPLify-X [51]	$\mathbf{j}_{rj}$	VPoser [51]	67.2 mm
✗	SMPLify-X [51]	$\mathbf{j}_{rj}$	PoseNDF [65]	57.4 mm
ExPose [15]	SMPLify-X [51]	$\mathbf{j}_{rj}$	PoseNDF [65]	53.8 mm
ExPose [15]	SMPLify-X [51]	$\mathbf{j}_{rj}$	GAN-S [17]	54.1 mm
ExPose [15]	SMPLify-X [51]	$\mathbf{j}_{op}$ [7, 8]	VPoser [51]	57.5 mm
ExPose [15]	SMPLify-X [51]	$\mathbf{j}_{rj vj}$	VPoser [51]	49.3 mm

**Table 1.** Virtual joints improvement analysis on EHF [51]. The columns indicate parameter initialization and optimization, which joints are optimized, and with which pose prior.

Next we evaluate our approach when adding the silhouette constraints using  $\mathcal{S}$  and the disentangled optimization for improved shape estimation, adding one such stage and then a final stage for detail (hands, face) capture. We perform two experiments, first using EHF to assess performance for pose capture as a single subject is used, and second, using SSP3D [56] that includes higher shape variance. For both experiments we employ pixel-based IoU and use the PA-V2V metric for general body pose estimation and the PVE-T-SC metric [56] for shape estimation. Likewise we use the SMPL meshes instead of SMPL-X to reduce the effect of the densely sampled head, after converting SMPL-X fits to SMPL meshes using pre-calculated mesh-

to-mesh vertex transfer maps [48]. Comparisons against optimization [16, 51] and single-shot [14, 15, 73] based approaches are given, with some of the latter focusing on shape [14], and others on expressive pose [15, 73]. Tab. 2 (left) shows that our *predict-then-optimize* approach outperforms the other methods with respect to pose estimates. PyMAF-X is a robust pose estimator when considering an average shaped subject, while LVD suffers due to its limited training data. (accordingly, LVD is omitted from the remainder of the experiments). On the contrary, on SSP3D the shape-aware SHAPY method offers better performance than PyMAF-X as presented in Tab. 2 (right). Still, our approach produces the best results in terms of pixel alignment and shape estimation, while also showing the benefit of disentangled optimization on shape capturing performance.

Method	EHF [51]		SSP3D [56]	
	PA-V2V↓	IoU↑	PVE-T-SC↓	IoU↑
ExPose [15]	71.7 mm	84.72%	33.0 mm	71.00%
LVD [16]	131.7 mm	-	-	-
SMPLify-X [51]	95.9 mm	81.46%	33.9 mm	76.60%
PyMAF-X [73]	66.6 mm	85.57%	30.6 mm	75.87%
SHAPY [14]	71.1 mm	81.29%	29.3 mm	72.65%
KBody (w/o $\mathcal{C}$ & §3.1.1)	-	-	28.1 mm	77.87%
KBody (w/o $\mathcal{C}$ )	64.2 mm	87.72%	25.6 mm	80.35%

**Table 2.** Results on the the EHF [51] & SSP3D [56] datasets.

We also present results on the validation set of the HBW dataset [14] that uses pre-scanned shapes and an assortment of in-the-wild images of the same persons to assess body shape estimation. However, up to now all results were presented using an arbitrary camera, in line with prior work for a fair comparison. As shown in Tab. 3 SHAPY outperforms all methods but this is reasonable as it was trained with metric scale supervision, whereas all other approaches were not. Still, our approach compares favorably to the remaining methods while offering higher quality pixel alignment than all alternatives. We also ablate the effect of estimating the camera’s parameters through  $\mathcal{C}$ , which naturally improves metric-scale performance.

KBody’s efficacy is qualitatively illustrated in Fig. 1 using images collected online. For these representative examples, KBody provides more balanced solutions, capturing pose and shape in high-quality for both heavy and lighter subjects, while also achieving good pixel alignment. With

Method	Height↓	Chest↓	Waist↓	Hips↓	P2P <sub>20k</sub> ↓	IoU (%)↑
ExPose [15]	75	91	93	91	36	80.50
SMPLify-X [51]	121	133	150	62	41	83.68
PyMAF-X [73]	100	74	90	64	34	80.36
SHAPY [14]	62	58	83	63	24	77.40
KBody (w/o $\mathcal{C}$ )	79	81	96	70	32	84.40
KBody	78	70	88	61	30	85.19

**Table 3.** Quantitative results on the HBW (val) [14] dataset.

respect to partial images, we provide a qualitative evaluation in Fig. 5 that shows how our inversion-based completion can handle missing head and/or lower body information. Moreover, Fig. 6 presents an ablation of the ADF objective and its benefits to clothed estimation.

Finally, an extended set of 112 full, 78 partial, and 32 ADF ablations of randomly selected in-the-wild examples can be found in our supplemental material and [project page](#).



**Figure 4.** Left-to-right: SMPLify-X [51] (light green), PyMAF-X [73] (purple), SHAPY [14] (green) and KBody (pink).

## 5. Conclusion

In this work we have presented KBody a general method for monocular body fitting. KBody can handle partial images gracefully via a generative completion stage and employs a multi-constraint fitting approach that delivers high-quality fits, and a balanced performance across pose, shape and image alignment performance. While the conflicts between pose and shape performance as well as world scale outputs and image alignment remain to be solved, we believe KBody is a step towards that direction as it shows that it is necessary for single-shot and iterative approaches to co-exist. However, relying on externally estimated constraints limits applicability to situations where these models under-



**Figure 5.** Partial image qualitative results. Same scheme as Fig. 4.



**Figure 6.** KBody fitting results **with** (top) and **w/o** (bottom) ADF.

perform. Novel solutions for (black-box) uncertainty estimation and multi-modal solving are required. Still, improving 2D estimation models is more practical than acquiring 3D data for supervision [16] or a wide-range of images and corresponding measurements [14]. Finally, relying on an image-based appearance prior for completion comes with limitations for non frontal facing images which can be addressed in the future with 3D aware generative models.



## References

- [1] Remove Image Background 100% Automatically and Free. <https://www.remove.bg/>. 7
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2StyleGAN: How to embed images into the StyleGAN latent space? In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4432–4441, 2019. 6
- [3] Thiemo Alldieck, Marc Kassubeck, Bastian Wandt, Bodo Rosenhahn, and Marcus Magnor. Optical flow-based 3d human motion estimation from monocular video. In *German Conference on Pattern Recognition*, pages 347–360. Springer, 2017. 5
- [4] Thiemo Alldieck, Hongyi Xu, and Cristian Sminchisescu. imGHUM: Implicit generative models of 3d human shape and articulated pose. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5461–5470, 2021. 2
- [5] Alexandru O. Balan, Leonid Sigal, Michael J Black, James E Davis, and Horst W Haussecker. Detailed human shape and pose from images. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2007. 5
- [6] Munkhtulga Battogtokh and Rita Borgo. Simple Techniques for a Novel Human Body Pose Optimisation Using Differentiable Inverse Rendering. *Eurographics 2022-Short Papers*, pages 65–684, 2022. 3
- [7] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Combining implicit function learning and parametric models for 3d human reconstruction. In *European Conference on Computer Vision*, pages 311–329. Springer, 2020. 5, 7
- [8] Bharat Lal Bhatnagar, Cristian Sminchisescu, Christian Theobalt, and Gerard Pons-Moll. Loopreg: Self-supervised learning of implicit surface correspondences, pose and shape for 3d human mesh registration. *Advances in Neural Information Processing Systems*, 33:12909–12922, 2020. 5, 7
- [9] Benjamin Biggs, Thomas Roddick, Andrew Fitzgibbon, and Roberto Cipolla. Creatures great and small: Recovering the shape and motion of animals from video. In *Asian Conference on Computer Vision*, pages 3–19. Springer, 2018. 5
- [10] 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. 1, 2, 3
- [11] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large Scale GAN Training for High Fidelity Natural Image Synthesis. In *International Conference on Learning Representations*, 2018. 6
- [12] Z. Cao, G. Hidalgo Martinez, T. Simon, S. Wei, and Y. A. Sheikh. OpenPose: Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 1, 6
- [13] Wei Cheng, Su Xu, Jingtian Piao, Chen Qian, Wayne Wu, Kwan-Yee Lin, and Hongsheng Li. Generalizable Neural Performer: Learning Robust Radiance Fields for Human Novel View Synthesis. *arXiv preprint arXiv:2204.11798*, 2022. 1
- [14] Vasileios Choutas, Lea Müller, Chun-Hao P Huang, Siyu Tang, Dimitrios Tzionas, and Michael J Black. Accurate 3D Body Shape Regression Using Metric and Semantic Attributes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2718–2728, 2022. 1, 2, 3, 7, 8
- [15] Vasileios Choutas, Georgios Pavlakos, Timo Bolkart, Dimitrios Tzionas, and Michael J Black. Monocular expressive body regression through body-driven attention. In *European Conference on Computer Vision*, pages 20–40. Springer, 2020. 2, 3, 6, 7
- [16] Enric Corona, Gerard Pons-Moll, Guillem Alenyà, and Francesc Moreno-Noguer. Learned Vertex Descent: A New Direction for 3D Human Model Fitting. *arXiv preprint arXiv:2205.06254*, 2022. 2, 3, 7, 8
- [17] Andrey Davydov, Anastasia Remizova, Victor Constantin, Sina Honari, Mathieu Salzmann, and Pascal Fua. Adversarial parametric pose prior. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10997–11005, 2022. 7
- [18] Yao Feng, Vasileios Choutas, Timo Bolkart, Dimitrios Tzionas, and Michael J Black. Collaborative regression of expressive bodies using moderation. In *2021 International Conference on 3D Vision (3DV)*, pages 792–804. IEEE, 2021. 2
- [19] Mihai Fieraru, Mihai Zanfir, Silviu Cristian Pirlea, Vlad Olaru, and Cristian Sminchisescu. Aifit: Automatic 3d human-interpretable feedback models for fitness training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9919–9928, 2021. 1
- [20] Jianglin Fu, Shikai Li, Yuming Jiang, Kwan-Yee Lin, Chen Qian, Chen Change Loy, Wayne Wu, and Ziwei Liu. StyleGAN-HUMAN: A data-centric odyssey of human generation. In *European Conference on Computer Vision*, pages 1–19. Springer, 2022. 5, 6
- [21] Stuart Geman. Statistical methods for tomographic image reconstruction. *Bull. Int. Stat. Inst.*, 4:5–21, 1987. 3
- [22] Rıza 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. 3
- [23] Chen Guo, Xu Chen, Jie Song, and Otmar Hilliges. Human performance capture from monocular video in the wild. In *2021 International Conference on 3D Vision (3DV)*, pages 889–898. IEEE, 2021. 3, 5
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1
- [25] Eric Hedlin, Helge Rhodin, and Kwang Moo Yi. A Simple Method to Boost Human Pose Estimation Accuracy by Correcting the Joint Regressor for the Human3.6m Dataset. *arXiv preprint arXiv:2205.00076*, 2022. 4, 5
- [26] Yinghao Huang, Federica Bogo, Christoph Lassner, Angjoo Kanazawa, Peter V Gehler, Javier Romero, Ijaz Akhter, and Michael J Black. Towards accurate marker-less human shape and pose estimation over time. In *2017 international conference on 3D vision (3DV)*, pages 421–430. IEEE, 2017. 5

- [27] Umar Iqbal, Kevin Xie, Yunrong Guo, Jan Kautz, and Pavlo Molchanov. KAMA: 3d keypoint aware body mesh articulation. In *2021 International Conference on 3D Vision (3DV)*, pages 689–699. IEEE, 2021. **3**
- [28] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7122–7131, 2018. **2**
- [29] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34:852–863, 2021. **6**
- [30] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. **6**
- [31] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and Improving the Image Quality of StyleGAN. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. **5, 6**
- [32] Hiroharu Kato, Yoshitaka Ushiku, and Tatsuya Harada. Neural 3d mesh renderer. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3907–3916, 2018. **3, 5**
- [33] Zhanghan Ke, Jiayu Sun, Kaican Li, Qiong Yan, and Rynson WH Lau. Modnet: Real-time trimap-free portrait matting via objective decomposition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 1140–1147, 2022. **3, 6**
- [34] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018. **6**
- [35] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11035–11045, 2021. **2, 6**
- [36] 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 *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2252–2261, 2019. **3, 4, 5**
- [37] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering. *ACM Transactions on Graphics (TOG)*, 39(6):1–14, 2020. **6**
- [38] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler. Unite the people: Closing the loop between 3d and 2d human representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6050–6059, 2017. **2, 3, 5**
- [39] Sangyun Lee, Gyojung Gu, Sunghyun Park, Seunghwan Choi, and Jaegul Choo. High-Resolution Virtual Try-On with Misalignment and Occlusion-Handled Conditions. In *European Conference on Computer Vision*, pages 204–219. Springer, 2022. **1**
- [40] 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**
- [41] Ren Li, Meng Zheng, Srikrishna Karanam, Terrence Chen, and Ziyang Wu. Everybody Is Unique: Towards Unbiased Human Mesh Recovery. *arXiv preprint arXiv:2107.06239*, 2021. **3**
- [42] Tzu-Mao Li, Miika Aittala, Frédo Durand, and Jaakko Lehtinen. Differentiable monte carlo ray tracing through edge sampling. *ACM Transactions on Graphics (TOG)*, 37(6):1–11, 2018. **3**
- [43] Zhe Li, Zerong Zheng, Hongwen Zhang, Chaonan Ji, and Yebin Liu. AvatarCap: Animatable Avatar Conditioned Monocular Human Volumetric Capture. In *European Conference on Computer Vision*, pages 322–341. Springer, 2022. **1**
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. **1**
- [45] Shichen Liu, Tianye Li, Weikai Chen, and Hao Li. Soft rasterizer: A differentiable renderer for image-based 3d reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7708–7717, 2019. **3**
- [46] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. **1**
- [47] Ramesha Rakesh Mugaludi, Jogendra Nath Kundu, Varun Jampani, et al. Aligning silhouette topology for self-adaptive 3D human pose recovery. *Advances in Neural Information Processing Systems*, 34:4582–4593, 2021. **2, 5**
- [48] Ahmed AA Osman, Timo Bolkart, and Michael J Black. STAR: Sparse trained articulated human body regressor. In *European Conference on Computer Vision*, pages 598–613. Springer, 2020. **2, 7**
- [49] Ahmed AA Osman, Timo Bolkart, Dimitrios Tzionas, and Michael J Black. SUPR: A Sparse Unified Part-Based Human Representation. *arXiv preprint arXiv:2210.13861*, 2022. **2**
- [50] Hui En Pang, Zhongang Cai, Lei Yang, Tianwei Zhang, and Ziwei Liu. Benchmarking and Analyzing 3D Human Pose and Shape Estimation Beyond Algorithms. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022. **2**
- [51] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10975–10985, 2019. **1, 2, 3, 4, 6, 7, 8**
- [52] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv preprint arXiv:2007.08501*, 2020. **3, 5**

- [53] Jiawei Ren, Mingyuan Zhang, Cunjun Yu, and Ziwei Liu. Balanced MSE for Imbalanced Visual Regression. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7926–7935, 2022. [2](#)
- [54] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022. [6](#)
- [55] Yu Rong, Takaaki Shiratori, and Hanbyul Joo. FrankMoCap: A monocular 3d whole-body pose estimation system via regression and integration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1749–1759, 2021. [2](#)
- [56] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. Synthetic training for accurate 3d human pose and shape estimation in the wild. *arXiv preprint arXiv:2009.10013*, 2020. [7](#)
- [57] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. Hierarchical kinematic probability distributions for 3D human shape and pose estimation from images in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11219–11229, 2021. [5](#)
- [58] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. Probabilistic 3D human shape and pose estimation from multiple unconstrained images in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16094–16104, 2021. [5](#)
- [59] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [6](#)
- [60] Cristian Sminchisescu and Alexandru C Telea. Human pose estimation from silhouettes. a consistent approach using distance level sets. In *10th International Conference on Computer Graphics, Visualization and Computer Vision (WSCG'02)*, volume 10, 2002. [5](#)
- [61] Irwin Sobel. An Isotropic 3x3 Image Gradient Operator. *Presentation at Stanford A.I. Project 1968*, 02 2014. [5](#)
- [62] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5693–5703, 2019. [1](#)
- [63] Neerja Thakkar, Georgios Pavlakos, and Hany Farid. The Reliability of Forensic Body-Shape Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 44–52, 2022. [3](#), [5](#)
- [64] Yating Tian, Hongwen Zhang, Yebin Liu, and Limin Wang. Recovering 3d human mesh from monocular images: A survey. *arXiv preprint arXiv:2203.01923*, 2022. [2](#)
- [65] Garvita Tiwari, Dimitrije Antić, Jan Eric Lenssen, Nikolaos Sarafianos, Tony Tung, and Gerard Pons-Moll. PoseNDF: Modeling Human Pose Manifolds with Neural Distance Fields. In *European Conference on Computer Vision*, pages 572–589. Springer, 2022. [3](#), [7](#)
- [66] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for StyleGAN image manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–14, 2021. [6](#)
- [67] Jiahang Wang, Sheng Jin, Wentao Liu, Weizhong Liu, Chen Qian, and Ping Luo. When human pose estimation meets robustness: Adversarial algorithms and benchmarks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11855–11864, 2021. [2](#)
- [68] Zhenzhen Weng, Kuan-Chieh Wang, Angjoo Kanazawa, and Serena Yeung. Domain Adaptive 3D Pose Augmentation for In-the-wild Human Mesh Recovery. In *International Conference on 3D Vision*, 2022. [2](#)
- [69] Stephen Wright, Jorge Nocedal, et al. Numerical optimization. *Springer Science*, 35(67-68):7, 1999. [6](#)
- [70] Xian Wu, Rui-Long Li, Fang-Lue Zhang, Jian-Cheng Liu, Jue Wang, Ariel Shamir, and Shi-Min Hu. Deep portrait image completion and extrapolation. *IEEE Transactions on Image Processing*, 29:2344–2355, 2019. [5](#)
- [71] Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. GHUM & GHUML: Generative 3d human shape and articulated pose models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6184–6193, 2020. [1](#), [2](#)
- [72] Andrei Zanfir, Eduard Gabriel Bazavan, Mihai Zanfir, William T Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Neural descent for visual 3d human pose and shape. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14484–14493, 2021. [3](#)
- [73] 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. [1](#), [2](#), [3](#), [7](#), [8](#)
- [74] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. [6](#)
- [75] Zibo Zhao, Wen Liu, Yanyu Xu, Xianing Chen, Weixin Luo, Lei Jin, Bohui Zhu, Tong Liu, Binqiang Zhao, and Shenghua Gao. Prior based human completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7951–7961, 2021. [5](#)