

# Delving Deep Into Hybrid Annotations for 3D Human Recovery in the Wild

Yu Rong<sup>1</sup> Ziwei Liu<sup>1</sup> Cheng Li<sup>2</sup> Kaidi Cao<sup>4</sup> Chen Change Loy<sup>3</sup>

<sup>1</sup>CUHK - SenseTime Joint Lab, The Chinese University of Hong Kong

<sup>2</sup>SenseTime Research <sup>3</sup>Nanyang Technological University <sup>4</sup>Stanford University

{ry017, zwliu}@ie.cuhk.edu.hk chengli@sensetime.com

kaidicao@cs.stanford.edu ccloy@ntu.edu.sg

## Abstract

Though much progress has been achieved in single-image 3D human recovery, estimating 3D model for in-the-wild images remains a formidable challenge. The reason lies in the fact that obtaining high-quality 3D annotations for in-the-wild images is an extremely hard task that consumes enormous amount of resources and manpower. To tackle this problem, previous methods adopt a hybrid training strategy that exploits multiple heterogeneous types of annotations including 3D and 2D while leaving the efficacy of each annotation not thoroughly investigated. In this work, we aim to perform a comprehensive study on cost and effectiveness trade-off between different annotations. Specifically, we focus on the challenging task of in-the-wild 3D human recovery from single images when paired 3D annotations are not fully available. Through extensive experiments, we obtain several observations: 1) 3D annotations are efficient, whereas traditional 2D annotations such as 2D keypoints and body part segmentation are less competent in guiding 3D human recovery. 2) Dense Correspondence such as DensePose [1] is effective. When there are no paired in-the-wild 3D annotations available, the model exploiting dense correspondence can achieve 92% of the performance compared to a model trained with paired 3D data. We show that incorporating dense correspondence into in-the-wild 3D human recovery is promising and competitive due to its high efficiency and relatively low annotating cost. Our model trained with dense correspondence can serve as a strong reference for future research<sup>1</sup>.

## 1. Introduction

Recovering 3D human model [23, 13] is essential in many applications such as augmented reality. Recent studies [24, 13, 21, 22] typically use a parametric model known



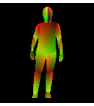


| Annotation         | Sparse<br>2D   | Dense<br>Labeling   | Dense<br>Correspondence   | Constrained<br>3D   | In-the-wild<br>3D   |
|--------------------|--|---|---|---|---|
| Examples           |  |  |  |  |  |
| Annotation<br>Cost | \$   | \$\$  | \$\$\$  | \$\$\$\$  | \$\$\$\$\$  |

Figure 1. **Annotations overview for 3D human recovery.** We study five kinds of annotations that are typically used in training deep networks for 3D human recovery. The number of ‘\$’ indicates the annotation cost of obtaining the corresponding annotations. A higher number of ‘\$’ suggests a higher cost.

as Skinned Multi-Person Linear Model (SMPL) [18] to represent 3D human models and estimate parameters of SMPL with a deep convolutional neural network (DCNN). Training such a deep network to handle 3D human recovery in the wild is challenging, as obtaining high-quality 3D annotations in unconstrained environments for training are both laborious and expensive. To circumvent this hurdle, one often has to adopt *hybrid annotations* for training, so as to leverage limited annotations from multiple datasets to avoid overfitting. For instance, Kanazawa *et al.* [13] train their models using both 3D joints of Human3.6M dataset [10] and 2D keypoints from COCO dataset [16]. Alternatively, apart from using an RGB image as an input to a network, one would introduce an auxiliary input as a prior to improve performance, *e.g.*, Omran *et al.* [21] use body part segmentation as an intermediate representation.

As summarized in Fig. 1, there are five common types of annotations: (a) Sparse 2D annotations such as 2D keypoints, (b) Dense labeling such as body part segmentation, (c) Dense correspondence such as the IUUV maps produced by DensePose [1, 19], (d) Constrained 3D annotations, *i.e.*, 3D annotations for images captured in constrained environments, such as Human3.6M [10], and (e) In-the-wild 3D annotations, *i.e.*, 3D annotations for in-the-wild images, such as UP-3D [15]. These annotations not only vary in their expressiveness but also their labeling cost. For instance, 3D annotations like SMPL are more expressive than the

<sup>1</sup>Code and models are available at the project page: [https://penincillin.github.io/dct\\_iccv2019](https://penincillin.github.io/dct_iccv2019)

dense correspondence, since the former encapsulates 3D deformable surface model while the latter only retains the UV fields. However, establishing 3D annotations requires a more complex annotation system than that required for annotating dense correspondence. Annotating dense correspondence such as DensePose [1] could be accomplished solely by human annotators while obtaining 3D annotations usually requires auxiliary facilities such as sparse markers [17] and IMUs [26].

In this study, we aim to perform a systematic study to investigate the cost and effectiveness trade-off between using different annotations in learning a deep network for 3D human recovery. We focus our study on the challenging task of recovering 3D human model from in-the-wild images, especially in the case when in-the-wild 3D annotations are insufficient, and how other annotation types could complement and bridge the gap. Our study is conducted using a unified and simple network, which could serve as a solid baseline for future study. Two aspects of using different annotations are investigated, *i.e.*, the effect of different annotations in serving as (a) a supervisory signal, (b) as an input to the network.

Our experiments reveal several observations:

**(1) 3D annotations are efficient for the in-the-wild scenario.** For in-the-wild images, models trained with paired 3D annotations achieve the best performance. Besides, excluding 80% paired in-the-wild 3D annotations only increases the reconstruction error by 5%. When there are no paired in-the-wild 3D annotations existing, incorporating constrained 3D annotations in the training phase can improve the performance and prevent a model from generating unnatural 3D human models.

**(2) Sparse 2D annotations and dense labeling alone are insufficient.** When there are no paired in-the-wild 3D annotations, using sparse 2D keypoints as the only supervision will decrease the models' performance by 60%. Besides, using dense labeling as input only brings marginal performance gain.

**(3) Dense correspondence such as IUUV map is an effective substitute for 3D annotations.** After a simple refinement step that removes noisy predictions, dense keypoints sampled from IUUV maps can serve as a strong supervision. IUUV map itself can also serve as a complementary input. Incorporating dense correspondence can further improve the models' performance by 2.9% or help the model trained with only 20% paired 3D annotations achieve similar performance of the model trained with full 3D annotations. Especially, when there are no paired 3D annotations available for in-the-wild images, the model using dense correspondence as supervision can achieve 92% of the performance of an upper-bound models that are trained with a full set of paired 3D in-the-wild annotations.

The contributions of our work are two-fold: 1) We sys-

tematically study the effectiveness of different annotations for in-the-wild 3D human recovery. We observe that while using paired 3D annotations leads to optimal results, it is not necessary for 3D human recovery, especially when considering its high annotating cost. 2) We reveal the effectiveness of incorporating dense correspondence into in-the-wild 3D human recovery. Our experiments show that when there are no in-the-wild annotations available, models trained with dense correspondence can still achieve the same performance as the models trained with 60% paired in-the-wild 3D annotations. The resulted model can serve as a strong and solid baseline for future studies.

## 2. Related Work

Recent studies on 3D human recovery mainly use a parametric model - SMPL [18] to represent human in 3D space. These studies can be divided into two groups: optimization-based methods and learning-based methods. Early works are mainly the optimization-based approach. Bogo *et al.* [3] propose to estimate parameters of SMPL through aligning the predicted models with 2D keypoints. Lassner *et al.* [15] extend the algorithm by adding the silhouettes matching loss and 91 landmarks. Tan *et al.* [24] propose an encoder-decoder architecture, in which the encoder predicts SMPL parameters from images and decoder predicts silhouettes from SMPL. The model is trained with heatmaps of silhouettes. BodyNet [25] proposes to predict volumetric 3D human first and then regress SMPL parameters from the predicted volumetric result.

Other recent works [13, 21, 22] share similar pipelines. They all design a CNN-based model to predict the parameters of SMPL. The models are trained with images that come with 2D annotations (2D keypoints) and 3D annotations (3D joints or ground-truth SMPL parameters). Kanazawa *et al.* [13] add adversarial loss [7] to judge whether the generated 3D human models are real or not. Pavlakos *et al.* [22] propose to first predict the silhouette and 2D keypoints heatmaps and then use them as the input for the SMPL parameters estimator. Omran *et al.* [21] argue that using body part segmentations to replace 2D images as input will enhance the performance of the model.

Most existing studies do not comprehensively investigate the efficiency of each annotation they use. The other works such as NBF [21] and HMR [13] have not completely evaluated the quality of generated 3D models. Their evaluation metrics are partial. Specifically, NBF [21] only evaluates the quality of predicted 3D poses, omitting the predicted shape. HMR [13] evaluates in-the-wild images using the accuracy of body part segmentation, which is only a 2D metric. In order to thoroughly evaluate how models' performance is affected by different factors, in this work, we conduct a series of experiments under a unified framework and training strategy. Besides, we use the Euclidean dis-

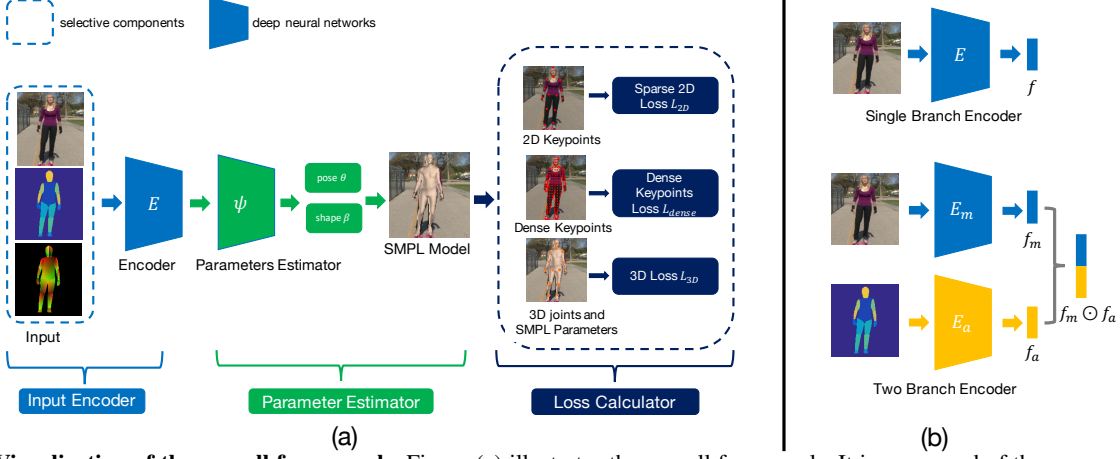


Figure 2. **Visualization of the overall framework.** Figure (a) illustrates the overall framework. It is composed of three components. 1) Input encoding part takes inputs and outputs encoded features. 2) Parameter estimator estimates the pose and shape parameters of the SMPL model given the outputs of the encoder. 3) Given estimated parameters, SMPL model generates predicted 3D joints, 2D keypoints and dense keypoints to calculate loss. Figure (b) shows two possible architectures of the input encoder. Input encoder could either be composed of a single branch that only takes one kind of inputs or two branches that takes original images and the other auxiliary inputs.

tance between predicted and ground-truth 3D meshes as the evaluation metric, which can faithfully reveal the quality of both pose and shape.

### 3. 3D Recovery with Hybrid Annotations

To evaluate the efficiency of different annotations for in-the-wild 3D human recovery, we conduct a series of experiments based on a unified framework and train-validation setting. In this section, we first introduce the framework used in the experiments. Then we describe five annotations investigated in this work. Finally, we discuss how to exploit the dense correspondence.

**3D Human Model.** Skinned Multi-Person Linear Model (SMPL) [18] is a 3D human body model parameterized by the pose and shape parameters. The shape parameters  $\beta \in \mathbb{R}^{10}$  are the first 10 coefficients of PCA components of shape space. Pose parameters  $\theta \in \mathbb{R}^{3 \times K}$  represent the 3D rotations for  $K = 23$  joints. In general, to specify a complete SMPL model,  $(23 + 1) \times 3 = 72$  pose parameters (three more parameters for global rotation) and 10 shape parameters are required.

**Framework.** The overall framework, as shown in Figure 2, is composed of three components: 1) input encoder 2) parameter estimator 3) loss calculator. The input encoder  $E$  has two variations of architectures: single branch and two branch. A two-branch encoder is composed of a main encoder  $E_m$  and an auxiliary encoder  $E_a$ . The main encoder takes images as input while the auxiliary encoder takes one auxiliary input that can either be body part segmentation or IUUV maps. The generated main features  $f_m$  and  $f_a$  are then concatenated to produce the final feature vector  $f = f_m \odot f_a$ . The single branch encoder has only one main branch  $E_m$  whose inputs are one category of original

Table 1. **Role of each annotation.** The role of different annotations in our experiments.

| Annotation  | Sparse 2D | Dense Labeling | Dense Correspondence | Constrained 3D | In-the-wild 3D |
|-------------|-----------|----------------|----------------------|----------------|----------------|
| Input       |           | ✓              | ✓                    |                |                |
| Supervision | ✓         |                | ✓                    | ✓              | ✓              |

images, body part segmentation and IUUV maps. It takes inputs and outputs encoded features  $f_m$ . For single branch encoder,  $f = f_m$ .

Given encoded feature vectors, the parameter estimator  $\psi$ , which is composed of two fully-connected layers, predicts the pose and shape parameters of SMPL. The SMPL model then generates the final 3D meshes. Follow the practice in previous works [20, 5, 13], the parameter estimator outputs the residual of parameters  $\Delta\theta$ . The final parameters are then obtained by adding the residual with the mean parameters  $\bar{\theta}$ . This strategy helps the model to focus on the variance of different images and thus leads to faster convergence. The parameter estimation process is formulated as follows:  $\theta = \bar{\theta} + \psi(E(I))$ , where  $I$  denotes inputs.

In the training phase, the loss calculator further regresses predicted 3D joints, 2D keypoints and dense keypoints obtained from SMPL vertices. The corresponding losses are then calculated using the ground-truth annotations.

#### 3.1. Hybrid Annotations

In this section, we discuss different annotations investigated in this work. The annotations include constrained and in-the-wild 3D annotations, sparse 2D annotations, dense labeling and dense correspondence. Depending on the nature of each annotation, they can serve as either input or supervision or both. The role of each annotation in our experiments is listed in Table 1.

**3D Annotations.** 3D annotations can be divided into two

categories according to whether the images are captured in constrained environments or in the wild. Since this paper mainly focuses on in-the-wild scenarios, constrained annotations are mainly used for pre-training. It will also take part in training when there are no paired in-the-wild 3D annotations available. In the loss calculating phase, for images with ground-truth SMPL parameters, we minimize the distance between predicted and ground-truth parameters. For numerical stability, each pose parameter  $\theta_i$  is converted into a  $3 \times 3$  rotation matrix using the Rodrigues formula [18]. For images with 3D joints annotation, we further minimize the distance between predicted and ground-truth 3D joints. 3D Loss  $L_{3D}$  is defined as follows:

$$\begin{aligned} L_{3D-joints} &= \sum_{i=1}^M \| (J_i^{3D} - \hat{J}_i^{3D}) \|_2, \\ L_{SMPL} &= \sum_{i=1}^O \| R(\theta_i) - R(\hat{\theta}_i) \|_2 + \| \beta_i - \hat{\beta}_i \|_2, \\ L_{3D} &= L_{3D-joints} + L_{SMPL}, \end{aligned} \quad (1)$$

where  $[\theta_i, \beta_i]$  and  $[\hat{\theta}_i, \hat{\beta}_i]$  are the predicted and ground-truth SMPL parameters, respectively.  $M$  and  $O$  represent the number of images with 3D joints annotation and ground-truth SMPL parameters.  $R : \mathbb{R}^3 \rightarrow \mathbb{R}^{3 \times 3}$  represents the Rodrigues formula.

**Sparse 2D Annotations.** To estimate 2D keypoints, the parameter estimator predicts three additional parameters to model the camera  $C \in \mathbb{R}^3$ , two parameters for the camera translation and one parameter for the focal length.  $C$  is then used to project the predicted 3D joints  $\hat{J}^{3D}$  to 2D keypoints  $\hat{J}^{2D}$ . The sparse 2D loss  $L_{2D}$  can then be defined as:

$$L_{2D} = \sum_{i=1}^S \| (J_i^{2D} - \hat{J}_i^{2D}) \times \mu_i \|_1, \quad (2)$$

where  $S$  is the number of training data with 2D keypoints annotation.  $J_i^{2D}$  and  $\hat{J}_i^{2D}$  denote the predicted and ground-truth 2D keypoints for the  $i$ th data sample, respectively.  $\mu_i$  represent the visibility vectors, where  $\mu_{ij} = 1$  means the  $j$ -th joint of  $i$ -th sample is visible, otherwise  $\mu_{ij} = 0$ .

**Dense Labeling.** Dense labeling investigated in this work is body part segmentation. In this work, dense labeling is only used as input. It can either be the sole input or serve as the auxiliary input. In our experiments, body part segmentation is not used as supervision, since the process of obtaining body part segmentation from SMPL predictions is not differentiable.

**Dense Correspondence.** Our work is in parallel with HoloPose [8] to incorporate dense correspondence into 3D human reconstruction. We exploit DensePose [1, 19], which establishes dense correspondence between RGB images and human bodies. Each pixel on a given image can be assigned with a  $(I, U, V)$  coordinate, which indicates a specific position on the surface-based human body.  $I \in \mathbb{Z}$  indicates which body part this point belongs to and  $(U, V) \in \mathbb{R}^2$  is

the coordinate of the precise location on the unrolled surface of the body part specified by  $I$ .

There is a close connection between SMPL and IUUV in that each vertex of the SMPL model can be assigned an  $(I, U, V)$  coordinate. In this way, for each point annotated with  $(I, U, V)$ , we calculate which triangle face of SMPL this point belongs to and the distances from this point to each vertex of the triangle face. These distances form the barycentric coordinates specific to this triangle face. Consequently, we have a mapping function  $\phi$  that can map the points annotated with  $(I, U, V)$  to the vertices of SMPL model. The mapping is provided in the following equation:

$$[v_1, v_2, v_3], [b_1, b_2, b_3] = \phi(I, U, V), \quad (3)$$

where  $v_i$  denotes the index of selected vertices and  $b_i$  represent the barycentric coordinate. We show some examples in Figure 3 to demonstrate the relationship between DensePose model and SMPL.

In the training phase, IUUV maps generated by DensePose can either be used as inputs or used for providing supervision. When serving as supervision, dense keypoints are sampled from IUUV maps and used to calculate dense keypoint loss. Each dense keypoint is composed of two parts: the coordinate  $(x, y)$  on the RGB images and the coordinate  $(I, U, V)$ . For simplicity of notation, we denote  $(I, U, V)$  coordinate as  $D$ . Given  $D$ , Equation (3) is used to calculate which vertices  $\mathbf{f} = [v_1, v_2, v_3]$  this point is closest to and the corresponding barycentric coordinates  $\mathbf{b} = [b_1, b_2, b_3]$ . After obtaining  $\mathbf{f}$  and  $\mathbf{b}$ , we project predicted SMPL vertices  $\hat{P} \in \mathbb{R}^{3 \times N}$  to 2D space  $\hat{P}^{2D} \in \mathbb{R}^{2 \times N}$  using the similar method of projecting 3D joints to 2D keypoints. Finally, we can obtain the predicted dense keypoints by weighted averaging the selected 2D vertices using barycentric coordinates and calculate the dense keypoint loss between the pixel coordinates of predicted and ground-truth dense keypoints. The whole process is formulated as:

$$\begin{aligned} [v_{i1}, v_{i2}, v_{i3}], [b_{i1}, b_{i2}, b_{i3}] &= \phi(D_i), \\ \hat{X}_i &= \sum_{j=1}^3 \hat{P}_i^{2D}[v_{ij}] \times b_{ij}, \\ L_{dense} &= \sum_{i=1}^T \| (X_i - \hat{X}_i) \|_1, \end{aligned} \quad (4)$$

where  $T$  is the number of images with dense keypoints annotations,  $\phi : \mathbb{Z} \times \mathbb{R}^2 \rightarrow \mathbb{Z}^3 \times \mathbb{R}^3$  is the mapping function defined in Equation (3).

### 3.2. Learning

**Sampling Strategy for Dense Correspondence.** The dense points drawn from IUUV maps cannot be employed directly since they frequently contain wrong predictions. For example, the left foot could be wrongly predicted as the right foot. To avoid erroneous points corrupting our model,



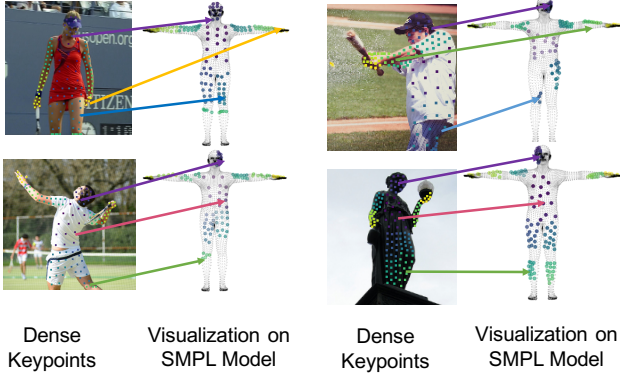


Figure 3. **Relationship between DensePose and SMPL.** Corresponding keypoints are annotated with same color.

Table 2. **FLOPs and model size of different architectures.**

| Encoder               | FLOPs $\times 10^9$ | Model Size (mb) |
|-----------------------|---------------------|-----------------|
| ResNet-101            | 7.803               | 174.97          |
| ResNet-50             | 4.090               | 102.27          |
| ResNet-50 & ResNet-18 | 5.905               | 150.97          |
| ResNet-18 & ResNet-18 | 3.630               | 97.783          |

we perform refinement by using accurate sparse keypoints as the reference. For each visible 2D keypoint, we check the values of IUUV map in the  $3 \times 3$  grid centering at it and select the value of ‘I’ (which indicates body part) that appears most frequently as the body part prediction of IUUV map surrounding this keypoint. Then we check whether the body part prediction matches the 2D keypoint or not.

After finding the erroneous region, our sampling scheme sets the IUUV map of this sub-area to be background in a recursive manner: We first set the IUUV value of the keypoint to be background, then we check the  $3 \times 3$  grid around it and determine the pixels whose value of ‘I’ equals to the *surrounding IUUV* and set their IUUV values to be background. Further, we check the  $3 \times 3$  grids centering at these pixels and determine more pixels using the same condition. The process is conducted recursively until there are no more pixels found. The above process is conducted on each keypoint to refine the whole IUUV map before we use the map as input and for sampling dense keypoints. A more detailed description along with an illustration figure can be found in the appendix A.

**Overall Loss Function.** The overall loss  $L$  is defined as:

$$L = \lambda_1 L_{3D} + \lambda_2 L_{2D} + \lambda_3 L_{dense}. \quad (5)$$

Detail values of  $\lambda$  used in the experiments is listed in the appendix C.

## 4. Experiments

We first introduce the datasets and evaluation metrics used in this work. In our experiments, we employ four

datasets: Human3.6M [10], COCO-DensePose [1], UP-3D [15] and 3DPW [26]. Experiments are mainly conducted on UP-3D dataset since it is the only in-the-wild dataset with SMPL annotations. We compare our methods with previous state-of-the-arts on UP-3D, 3DPW and COCO-DensePose datasets.

**Human3.6M.** Human3.6M [10] is an indoor dataset. Following HMR [22], we use Mosh [17] to collect ground-truth SMPL parameters from raw 3D Mocap markers. In our experiment, the data of Human3.6M is used in pre-training. It is also used in training when there are no paired in-the-wild 3D annotations available.

**COCO-DensePose.** COCO-DensePose dataset [1] is a newly released dataset that builds dense correspondence between images and body part surface. Images in this dataset are all selected from the keypoints MS-COCO dataset [16]. Researchers in [1] re-annotate each selected image with about 100 ~ 150 dense keypoints. We train our model on the training set and test the models on the evaluation set.

**UP-3D.** This dataset is built by Lassner *et al.* [15]. They pick images from four pose estimation datasets including: LSP [11], LSP-extended [12], MPII [2] and FashionPose [6]. The researchers extend SMPLify [3] and fit the model to those images. Then they ask human annotators to pick the samples with good fitness.

**3DPW.** This dataset is built by Von *et al.* [26]. They estimate 3D poses using a single hand-held camera and a set of IMUs attached at body limbs. 3D body shapes are obtained through 3D scans. This dataset cannot be counted as a totally in-the-wild dataset since the data are collected by several actors performing different actions. We compare our methods with previous state-of-the-arts, *e.g.*, HMR [13].

**Evaluation Metrics.** For COCO-DensePose dataset, the evaluation metric is the dense keypoints distance introduced in Equation (4). It is abbreviated as DKD in the following sections. For other datasets with SMPL annotations, we use the mean per-vertex error (PVE) proposed by Pavlakos *et al.* [22] as the metric, which computes the Euclidean distance between ground-truth SMPL vertices and the predicted SMPL vertices. We also report mean per joint position error (MPJPE) on SMPL joints to reveal the quality of pose recovery and PVE between SMPL vertices whose shape parameters come from ground-truth and prediction while pose parameters are set to be the same (in the experiment, pose parameters are all set to be zero). We use this metric to reveal the quality of shape recovery and abbreviate it as PVE-T, where ‘T’ refers to T-pose.

**Implementation Details.** All images are cropped according to the bounding boxes of humans. These images are further padded and scaled to  $224 \times 224$ . During training, images are randomly flipped and scaled for data augmentation. As depicted in Figure 2, the input encoder has two architectures. In most experiments, the single branch encoder is

Table 3. **Influence of different annotations.** The evaluation metrics are PVE, MPJPE and PVE-T, separately. For all metrics, lower is better. “3D” refers to paired in-the-wild 3D annotations. “20% 3D” refers to 20% randomly selected 3D annotations. “Sparse 2D” refers to sparse 2D keypoints. “Dense” refers to dense correspondence, namely, IUUV maps generated by DensePose [1, 19].

| Supervision →<br>Input ↓ | 3D & Dense &<br>Sparse 2D   | 20% 3D & Dense &<br>Sparse 2D | 3D & Sparse 2D       | Dense & Sparse 2D    | Sparse 2D Only        |
|--------------------------|-----------------------------|-------------------------------|----------------------|----------------------|-----------------------|
| IUV Only                 | <b>120.0 / 103.1</b> / 31.8 | 125.0 / 107.2 / 32.6          | 125.2 / 106.4 / 32.1 | 138.7 / 121.2 / 54.7 | 204.3 / 177.0 / 92.1  |
| Segment Only             | 123.0 / 105.1 / 32.7        | 126.7 / 110.0 / 33.2          | 124.8 / 107.8 / 31.7 | 147.4 / 130.1 / 55.9 | 203.8 / 176.7 / 93.3  |
| Image Only               | 123.7 / 105.9 / 30.9        | 127.5 / 110.6 / 32.2          | 127.4 / 108.5 / 30.7 | 137.7 / 120.3 / 51.7 | 203.2 / 178.5 / 106.2 |
| Image & IUV              | 122.4 / 105.1 / <b>30.2</b> | 125.0 / 107.6 / 32.1          | 125.5 / 107.3 / 30.7 | 133.8 / 117.2 / 52.5 | 197.3 / 172.8 / 107.9 |
| Image & Segment          | 121.5 / 104.3 / 31.0        | 126.4 / 107.0 / 31.6          | 125.8 / 106.8 / 31.5 | 142.2 / 124.2 / 56.6 | 201.2 / 177.5 / 101.7 |

based on ResNet-101 [9] while the main encoder and auxiliary encoder of the two branch architecture are based on ResNet-50 and ResNet-18, separately. In this way, models with different architectures have comparable FLOPs and model size. The overall FLOPs and size of models adopting different input encoders are listed in Table 2.

We assign additional fully-connected layers at the top of the input encoder to map the feature vectors to 85 dimensions. The final output vectors contain pose parameters  $\theta$  (72 dimensions), shape parameters  $\beta$  (10 dimensions) and camera model  $C$  (3 dimensions).

#### 4.1. The Effectiveness of Hybrid Annotations

In this subsection, we study the efficiency of different annotations when serving as inputs or supervisions. In all the experiments, sparse 2D keypoints are always assumed to be available, as annotating 2D keypoints is quite cheap. Alternatively, precise results can be obtained using state-of-the-arts 2D pose estimation algorithms [27, 4]. For each input type, we adopt five different supervision combinations, including 3D annotations, 3D annotations plus dense correspondence, randomly selected 20% 3D annotations plus dense correspondence, dense correspondence only and sparse 2D keypoints only. The results are listed in Table 3.

**Influence of Supervision.** Detailed numbers in this subsection are calculated by comparing the models that take images as the only input (the fourth row of Table 3). Same conclusions can be drawn from other models that use different inputs. It is not surprising that 3D annotations can provide the best guidance for in-the-wild 3D human recovery while sparse 2D keypoints are not as efficient. Dense correspondence, namely, IUUV maps generated by DensePose [1, 19], is an effective annotations for in-the-wild 3D human recovery. The model trained with sampled dense keypoints and sparse 2D keypoints can achieve the 92% performance of the model trained with full set 3D annotations. Furthermore, the model trained with hybrid of only 20% 3D annotations and dense correspondence achieve comparable performance with the model trained with full 3D annotations. Besides, the performance of the model trained with full 3D annotations can be improved by 2.9% through incorporating dense correspondence into training.

**Influence of Input.** Five input combinations are exploited

Table 4. **Influence of pose and shape parameters.** The evaluation metrics are: PVE, MPJPE and PVE-T, separately.

| 3D Loss →<br>Other Supervision ↓ | 3D Pose Only          | Shape parameters Only |
|----------------------------------|-----------------------|-----------------------|
| DC & Sparse 2D                   | 131.3 / 116.6 / 59.0  | 148.5 / 127.3 / 30.6  |
| Sparse 2D Only                   | 164.0 / 148.2 / 117.0 | 220.0 / 180.6 / 31.4  |

in our experiments, including 1) images only, 2) IUUV maps only, 3) body part segmentation only, 4) images plus IUUV maps, 5) images plus body part segmentation. The first three categories adopt a single branch architecture and the last two use the two-branch architecture. For a fair comparison, IUUV maps and body part segmentations are both generated by DensePose [1] model. Experimental results in Table 3 show that when sparse 2D keypoints serve as the only supervision, incorporating auxiliary inputs including body part segmentation or IUUV maps can only improve the models performance by 1.5% in average. It is marginal when compared with 32% improvement brought by incorporating sampled dense keypoints from IUUV maps into supervision while still using the images as the only input.

#### 4.2. Exploit 3D Annotations

**Influence of Separate Parameters.** We separately evaluate the influence of SMPL pose and shape parameters by using only one of them during training. The results shown in Table 4 suggest that: (1) 3D poses and SMPL parameters explicitly affect MPJPE and PVE-T, respectively. (2) 3D poses have more influence on the model’s overall performance. Besides, the results in Table 3 show that when both pose and shape parameters are used in training, MPJPE and PVE-T are nearly consistent with the PVE. Therefore, we only report PVE in the following experiments.

**Efficiency of 3D Annotations.** We then evaluate the effectiveness of in-the-wild 3D annotations. Models in this section are all trained with 3D annotations and sparse 2D annotations. In these experiments, the number of paired 3D annotations is reduced gradually from 100% to 0% (0% means only using sparse 2D annotations in training). The results are shown in Figure 4. We only show detailed results of models taking images as the only inputs. Detailed experiment results of all the models can be found in the appendix B. From the Figure 4, we find that 3D annotations are efficient. For instance, the reconstruction error only in-

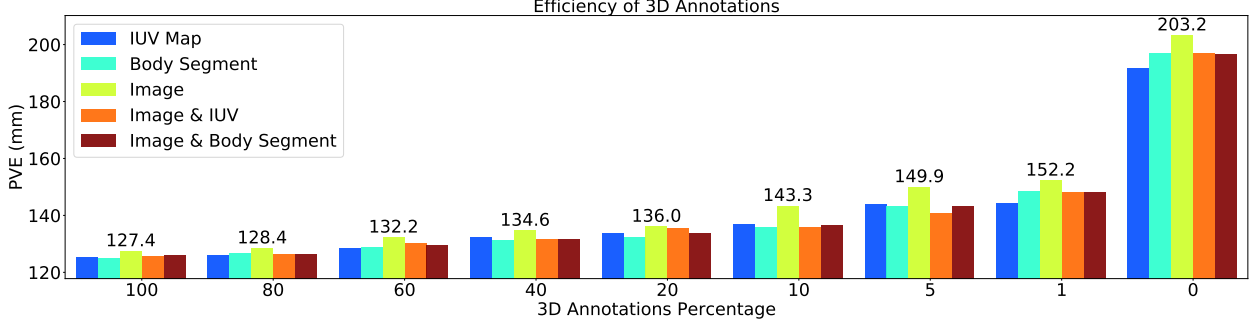


Figure 4. **Influence of 3D annotations.** We test different models on the test set of UP-3D [15] using the per-vertex error (abbreviated as PVE, the unit is mm.) as the metric. The figure shows that 3D annotations are very efficient.

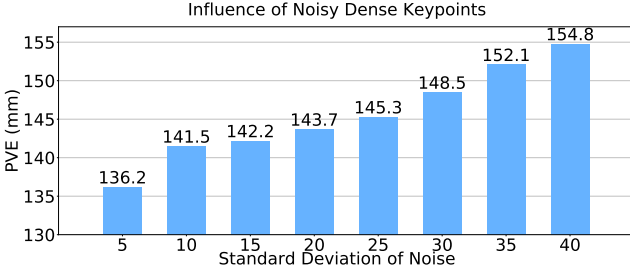


Figure 5. **Influence of noisy dense correspondence.** In this experiment, we add Gaussian noise to IUUV maps. The mean ( $\mu$ ) is fixed to be 0 and standard deviation varies from 5 to 40.

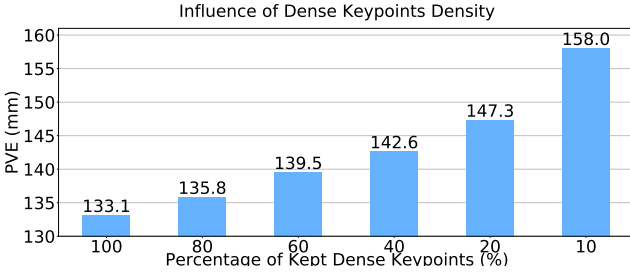


Figure 6. **Influence of dense keyoints density** In this experiment, sampled dense keyoints are randomly discarded. Performance of the model drops gracefully when more than 60% keyoints are retained. Even with only 10 ~ 15 dense keyoints kept, they are still significantly more efficient than sparse 2D keyoints.

creases by 6% when 80% 3D annotations are excluded from training. On the contrary, sparse 2D annotations are incompetent in guiding 3D human recovery. When there are no paired 3D data available, the performance drops drastically. The reconstruction error is 34% larger than models trained with only 1% of paired 3D data.

### 4.3. Exploit Dense Correspondence

Inspired by surprising efficiency of dense correspondence as observed in Table 3, we further investigate its effectiveness in this subsection. Models in this subsection all take images and IUUV maps as inputs.

**Influence of Noisy Dense Correspondence.** As stated before, dense keyoints used as supervision are sampled from IUUV maps, which might contain errors. We refine IUUV

maps as described in section 3.2. If we directly use raw IUUV maps the performance drops by 20.1%. We further study how noise in  $U$  and  $V$  influence the models’ performance, since the refinement process only removes potential errors in  $I$ . We add Gaussian noise to  $U$  and  $V$ , whose values lie in  $[0, 255]$ . The mean ( $\mu$ ) of Gaussian noise is fixed to be 0 and the standard deviation ( $\sigma$ ) varies from 5 to 40. The result is illustrated in Figure 5. The results show that our method is robust to the noise. The performance of the model drops gracefully while the variance of noise is less than 10. Even the variance of noise increases to 40, using noisy dense keyoints could still enhance the performance of the model considerably.

**Influence of Dense Keyoints Density.** Each image in the COCO-DensePose dataset is annotated with 100 ~ 150 dense keyoints. We sample same amount of dense keyoints on UP-3D. In this subsection, we study the influence of dense keyoints density by randomly discarding part of dense keyoints and train the model using the remaining ones. The number of dense keyoints is reduced gradually from 100% to 0% (0 means using only 2D keyoints.). The results are shown in Figure 6. The performance drops gracefully when more than 60% dense keyoints are retained. Besides, models trained with only 10 ~ 15 dense keyoints still have significantly higher performance than models trained with only sparse 2D keyoints. Experiment results in this subsection is useful for real-life application in that lots of efforts in annotating dense keyoints could be saved with a little sacrifice in the final performance.

### 4.4. Comparison with State-of-the-arts

**Quantitative Results.** For UP-3D, we compare our model with both the optimization-based methods [15] and the learning-based methods [13, 22, 21]. For COCO-DensePose, we mainly compare our method with HMR [13], since HMR is the only method that has been trained on COCO [16] dataset, which covers all the images in CODP dataset. For 3DPW, we train HMR on the training set and compare our methods with it on the testing set. The results are shown in Table 5. “Ours-3D”



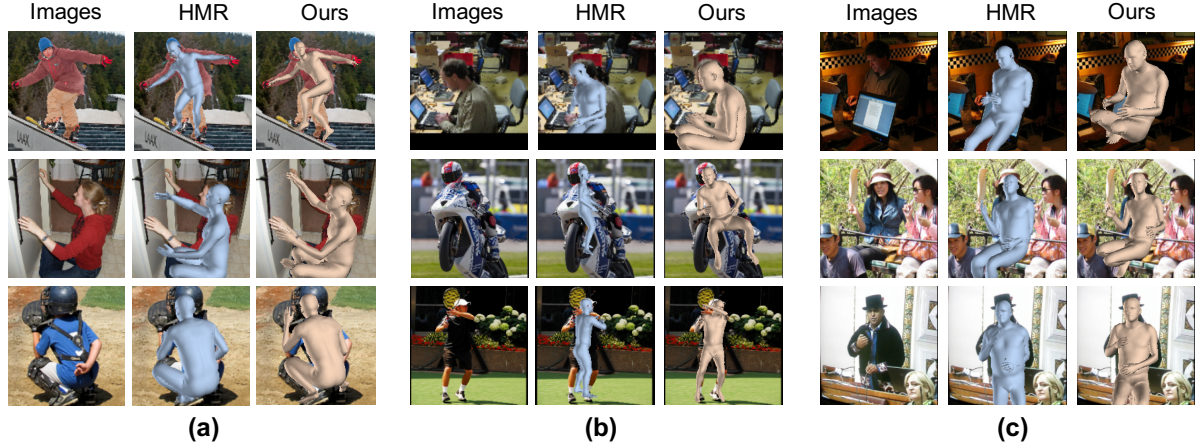


Figure 7. **A comparison between our model and HMR [13].** “Our model” refers to the model that adopts the framework in Figure 2. It uses images and IUUV maps as input and it is trained with dense correspondence and sparse 2D keypoints. (a) shows that our model can generate better-aligned results. (b) shows that our model still works well on some tough samples. (c) shows that our model is capable of generating natural results when HMR fails. The images all come from COCO-DensePose dataset [1].

Table 5. **Comparison with state-of-the-art methods.** This table presents the evaluation results on COCO-DensePose dataset [1] (CODP is used for the simplicity of notation.) using DKD (Dense Keypoints Distance), the unit is mm. It also presents evaluation results on UP-3D dataset [15] and 3DPW dataset [26] using PVE (Per-Vertex Error), the unit is mm. For all the metrics, lower is better. “Ours-3D” refers to the proposed model trained using paired 3D annotations. “Ours-DC” refers to the proposed model trained using only dense correspondence and sparse 2D annotations.

| Dataset →                   | CODP [1]    | UP-3D [15]   | 3DPW [26]    |
|-----------------------------|-------------|--------------|--------------|
| Metric →                    | DKD         | PVE          | PVE          |
| Methods ↓                   | (mm)        | (mm)         | (mm)         |
| Lassner <i>et al.</i> [15]  | –           | 169.8        | –            |
| NBF [21]                    | –           | 134.6        | –            |
| HMR [13]                    | 102.7       | 149.2        | 161.0        |
| Pavlakos <i>et al.</i> [22] | –           | <b>117.7</b> | –            |
| Ours-3D                     | –           | 122.2        | <b>152.9</b> |
| Ours-DC                     | <b>51.8</b> | 137.5        | 165.3        |

refers to the proposed model trained using paired 3D annotations. “Ours-DC” refers to the proposed model trained with only dense correspondence and sparse 2D annotations. These two models both adopt the two-branch encoder that takes images and IUUV maps as input. We use ResNet-18 as the backbone for “Ours-3D” and “Ours-DC” to guarantee a fair comparison, since models of previous works such as HMR [22] and NBF [21] are all based on ResNet-50 [9].

When 3D data is available, our method surpasses or performs comparably with previous state-of-the-arts, demonstrating that our model is simple yet efficient. On UP-3D dataset, it is noteworthy that our model trained using dense correspondence are comparable with most of the previous methods despite no paired in-the-wild 3D annotations are used in training.

**Qualitative Results.** We show some qualitative results of our model and HMR [13] in Figure 7. “Our model” refers to the model that adopts the framework in Figure 2, which uses images and IUUV maps as input and it is trained with

dense correspondence and sparse 2D keypoints. The observations for each subfigure are given as follows: (a) shows that our model generates better-aligned and more precise 3D human models than HMR does. (b) shows that when HMR fails on images with extreme poses or scales, our model can still generate plausible results. (c) shows that in some cases HMR generates erroneous 3D models while our method generates more natural results.

## 5. Conclusion

We have performed a systematic study of the cost and efficiency trade-off of hybrid annotations used in in-the-wild 3D human recovery. Through extensive experiments, we find that paired in-the-wild 3D annotations are not irreplaceable as commonly believed. Interestingly, in the absence of paired 3D data, the models that exploits dense correspondence can achieve 92% of the performance compared to the models trained with paired 3D data. We further benchmark against previous state-of-the-art methods on UP-3D [15] and 3DPW [26] dataset. Without paired in-the-wild 3D annotations, the model achieves comparable performance with most of the previous state-of-the-arts methods trained with paired 3D annotations. We demonstrate that dense correspondence is a new supervision form that is promising and competitive for in-the-wild 3D human recovery. Considering its high efficiency and relatively low annotating cost, our models can serve as a strong reference for future research.

**Acknowledgements.** This work is partially supported by the Collaborative Research grant from SenseTime Group (CUHK Agreement No. TS1610626 & No. TS1712093), the General Research Fund of the Hong Kong (CUHK 14209217), Singapore MOE AcRF Tier 1 (M4012082.020), NTU SUG, and NTU NAP.



## References

- [1] Rza Alp Gler, Natalia Neverova, and Iasonas Kokkinos. Densepose: Dense human pose estimation in the wild. In *CVPR*, 2018. 1, 2, 4, 5, 6, 8, 10
- [2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *CVPR*, 2014. 5
- [3] 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 *ECCV*, 2016. 2, 5
- [4] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. In *CVPR*, 2017. 6
- [5] Joao Carreira, Pulkrit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4733–4742, 2016. 3
- [6] Matthias Dantone, Juergen Gall, Christian Leistner, and Luc Van Gool. Body parts dependent joint regressors for human pose estimation in still images. *TPAMI*, 36(11):2131–2143, 2014. 5
- [7] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014. 2
- [8] Riza Alp Guler and Iasonas Kokkinos. Holopose: Holistic 3d human reconstruction in-the-wild. In *CVPR*, pages 10884–10894, 2019. 4
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5, 8
- [10] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI*, 36(7):1325–1339, 2014. 1, 5, 10
- [11] Sam Johnson and Mark Everingham. Clustered pose and nonlinear appearance models for human pose estimation. In *BMVC*, 2010. 5
- [12] Sam Johnson and Mark Everingham. Learning effective human pose estimation from inaccurate annotation. In *CVPR*, 2011. 5
- [13] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 1, 2, 3, 5, 7, 8
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015. 10
- [15] 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 *CVPR*, 2017. 1, 2, 5, 7, 8, 10
- [16] 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 *ECCV*, 2014. 1, 5, 7
- [17] Matthew Loper, Naureen Mahmood, and Michael J Black. Mosh: Motion and shape capture from sparse markers. *TOG*, 33(6):220, 2014. 2, 5
- [18] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *TOG*, 34(6):248, 2015. 1, 2, 3, 4
- [19] Natalia Neverova, James Thewlis, Riza Alp Guler, Iasonas Kokkinos, and Andrea Vedaldi. Slim densepose: Thrifty learning from sparse annotations and motion cues. In *CVPR*, pages 10915–10923, 2019. 1, 4, 6
- [20] Markus Oberweger, Paul Wohlhart, and Vincent Lepetit. Training a feedback loop for hand pose estimation. In *ICCV*, pages 3316–3324, 2015. 3
- [21] 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 *3DV*, 2018. 1, 2, 7, 8
- [22] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *CVPR*, 2018. 1, 2, 5, 7, 8
- [23] Yongbin Sun, Ziwei Liu, Yue Wang, and Sanjay E Sarma. Im2avatar: Colorful 3d reconstruction from a single image. *arXiv preprint arXiv:1804.06375*, 2018. 1
- [24] J Tan, Ignas Budvytis, and Roberto Cipolla. Indirect deep structured learning for 3d human body shape and pose prediction. In *BMVC*, 2017. 1, 2
- [25] Gul Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. Bodynet: Volumetric inference of 3d human body shapes. In *ECCV*, pages 20–36, 2018. 2
- [26] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 2, 5, 8, 10
- [27] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *ECCV*, 2018. 6

## A. Sampling Dense Keypoints

Since dense keypoint annotations are only available in COCO-DensePose dataset and training models purely using sparse 2D keypoints will lead to suboptimal results, we present an effective method for generating dense keypoints for other in-the-wild images that only annotated with sparse 2D keypoints. An effective way is to directly sample points from the IUV maps produced by the DensePose model.

The dense points drawn from IUV maps cannot be employed directly since the maps frequently contain wrong predictions. As Figure 8 (a) shows, the left foot is wrongly predicted as the right foot while the right foot is predicted as the opposite. To avoid erroneous points corrupting the learning of our model, we perform sampling of dense points by using accurate sparse keypoints as reference. Specifically, for each visible 2D keypoint, we check the values of IUV map in the  $3 \times 3$  grid centering at it and select the value of ‘I’ (which indicates body part) that appears most frequently as the body part prediction of IUV map surrounding this keypoint. Then we check whether the *surrounding IUV* is consistent with the 2D keypoint. For example, if a keypoint is labeled as “right ankle” but the *surrounding IUV* is “left foot”, then this sub-area is assigned as erroneous region.

After finding the erroneous region, our sampling scheme will set the IUV map of this sub-area to be background in a recursive manner: We first set the IUV value of the keypoint to be background, then we check the  $3 \times 3$  grid around it and determine the pixels whose value of ‘I’ equals to the *surrounding IUV* and set their IUV values to be background. Further, we check the  $3 \times 3$  grids centering at these pixels and determine more pixels using the same condition. The process is conducted recursively until there are no more pixels found. The above process is conducted on each keypoint to refine the whole IUV map before we use the map as the complementary input and for sampling dense keypoints. The sampling process is depicted in Figure 8 (b).

## B. Efficiency of 3D Annotations.

**Detailed experiment results.** Detailed experiment results in Figure 4 is listed in Table 6. In experiments, the amount of paired 3D annotations used in the training phase is reduced gradually from 100% to 0% (0% means only using sparse 2D annotations in training). From the table, we find that 3D annotations are quite efficient. The reconstruction error only increases by 6% when 80% 3D annotations are excluded from training.

**Influence of constrained 3D.** We also investigate constrained annotations. The experiment results are listed in Table 7. When paired in-the-wild 3D annotations exist, using constrained 3D annotations barely brings improvement. However, when there are no paired in-the-wild 3D annotations

exist, incorporating constrained 3D annotations into training improves the performance of models by 30%.

## C. Implementation Details

In this section, we discuss more implementation details. In the training phase, the whole model is first pretrained using 3D data from Human3.6M dataset [10], then it is finetuned on the COCO-DensePose [1], UP-3D [15] and 3DPW [26]. For COCO-DensePose dataset, we train our model with ground truth dense keypoints and 2D keypoints. For UP-3D and 3DPW dataset, our model is trained with the combination of 3D annotations, 2D keypoints and sampled dense keypoints. The sampled dense keypoints are obtained based on the method described in Section A.

In the training phase, the batch size is set to 128. Adam optimizer [14] with  $1e-4$  is adopted in the whole training phase. The model gets converged after  $40 \sim 50$  epochs. Especially, if all the losses including 3D, dense and 2D are used in training, their balance weights are 10, 1, 10, respectively. If only two losses are used, their balance weights are set to be both 10.

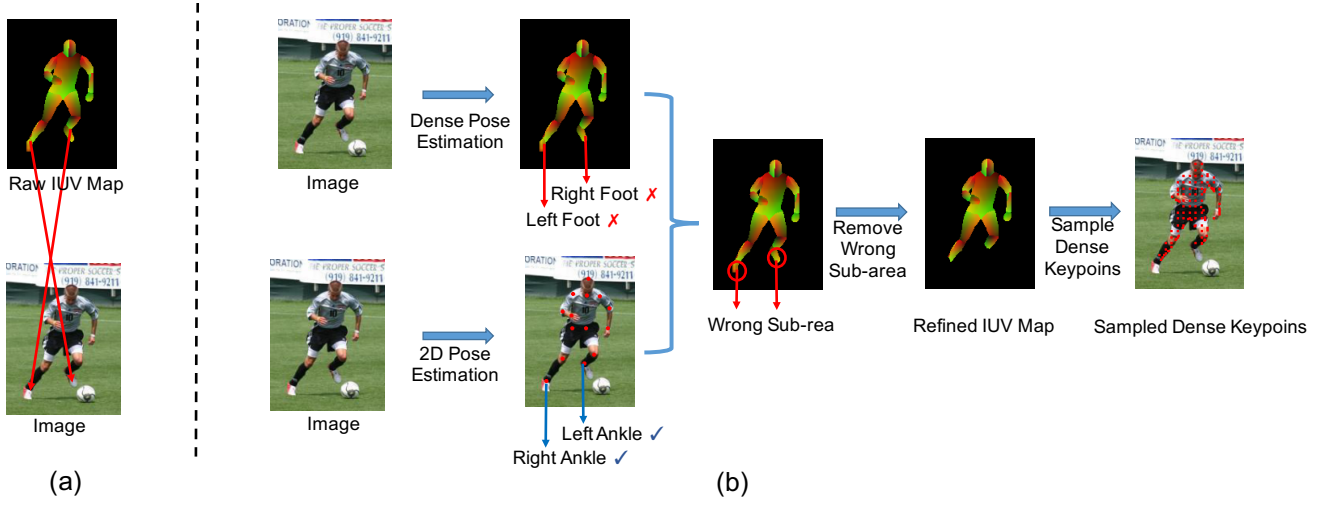


Figure 8. Figure (a) demonstrates that the raw IUUV map might contain errors. Figure (b) shows the process of refining the IUUV maps. The generated IUUV map is compared with the 2D keypoints. If they are not consistent, *e.g.*, the sub-area around “right ankle” is predicted as “left foot”, then we discard this sub-area by assigning it as background. We compare each keypoint with the predicted IUUV maps surrounding it and remove the inconsistent part.

Table 6. **Influence of 3D annotations.** This table lists detailed experiment results of Figure 4.

| Kept 3D Annotations (%) →<br>Input ↓ | 100   | 80    | 60    | 40    | 20    | 10    | 5     | 1     | 0     |
|--------------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| IUV Map                              | 125.2 | 125.9 | 128.3 | 132.3 | 133.6 | 136.8 | 144.0 | 144.3 | 191.5 |
| Body Segment                         | 124.8 | 126.7 | 128.9 | 131.3 | 132.3 | 135.9 | 143.0 | 148.5 | 196.7 |
| Image                                | 127.4 | 128.4 | 132.2 | 134.6 | 136.0 | 143.3 | 149.9 | 152.2 | 203.2 |
| Image & IUV                          | 125.5 | 126.2 | 130.1 | 131.6 | 135.3 | 135.9 | 140.6 | 148.0 | 197.0 |
| Image & Body Segment                 | 125.8 | 126.1 | 129.5 | 131.4 | 133.7 | 136.5 | 143.3 | 148.0 | 196.4 |

Table 7. **Influence of constrained 3D annotations.** The inputs of the models are all single images.

| Other Supervisions → | 100% 3D &<br>Sparse 2D | 20% 3D &<br>Sparse 2D | Dense &<br>Sparse 2D | Sparse 2D<br>Only |
|----------------------|------------------------|-----------------------|----------------------|-------------------|
| with Constrained 3D  | 127.4                  | 137.7                 | 137.3                | 203.2             |
| w/o Constrained 3D   | 128.9                  | 138.1                 | 173.4                | 230.9             |