

# Learning Visibility for Robust Dense Human Body Estimation

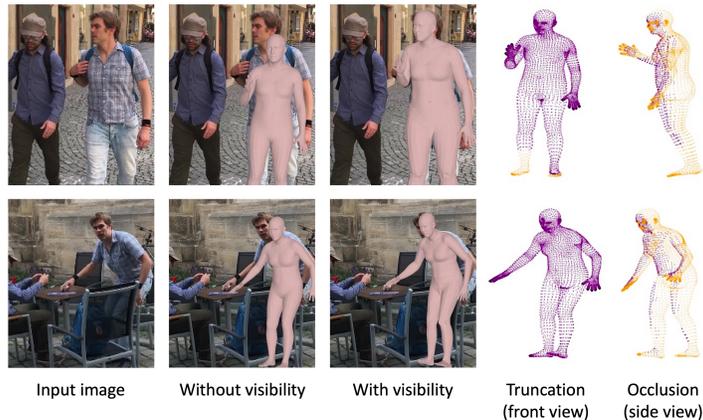
Chun-Han Yao<sup>1</sup> Jimei Yang<sup>2</sup> Duygu Ceylan<sup>2</sup> Yi Zhou<sup>2</sup>  
Yang Zhou<sup>2</sup> Ming-Hsuan Yang<sup>134</sup>

<sup>1</sup>UC Merced <sup>2</sup>Adobe <sup>3</sup>Google <sup>4</sup>Yonsei University

**Abstract.** Estimating 3D human pose and shape from 2D images is a crucial yet challenging task. While prior methods with model-based representations can perform reasonably well on whole-body images, they often fail when parts of the body are occluded or outside the frame. Moreover, these results usually do not faithfully capture the human silhouettes due to their limited representation power of deformable models (e.g., representing only the naked body). An alternative approach is to estimate dense vertices of a predefined template body in the image space. Such representations are effective in localizing vertices within an image but cannot handle out-of-frame body parts. In this work, we learn dense human body estimation that is robust to partial observations. We explicitly model the visibility of human joints and vertices in the x, y, and z axes separately. The visibility in x and y axes help distinguishing out-of-frame cases, and the visibility in depth axis corresponds to occlusions (either self-occlusions or occlusions by other objects). We obtain pseudo ground-truths of visibility labels from dense UV correspondences and train a neural network to predict visibility along with 3D coordinates. We show that visibility can serve as 1) an additional signal to resolve depth ordering ambiguities of self-occluded vertices and 2) a regularization term when fitting a human body model to the predictions. Extensive experiments on multiple 3D human datasets demonstrate that visibility modeling significantly improves the accuracy of human body estimation, especially for partial-body cases. Our project page with code is at: <https://github.com/chhankyao/visdb>.

## 1 Introduction

Estimating 3D human pose and shape from monocular images is a crucial task for various applications such as performance retargeting, virtual avatars, and human action recognition. It is a fundamentally challenging problem due to the depth ambiguity and the complex nature of human appearances that vary with articulation, clothing, lighting, viewpoint, and occlusions. To represent the complicated 3D human bodies via compact parameters, model-based methods like SMPL [24] have been widely used in the community. However, SMPL parameters represent human bodies in a holistic manner, causing their limited flexibility to fit real-world images faithfully via direct regression. More importantly,



**Fig. 1. Dense human body estimation with/without visibility modeling.** We propose to learn dense visibility to improve human body estimation in terms of faithfulness to the input image and robustness to truncation (top) or occlusions (bottom). We show the estimated meshes without/with visibility modeling in columns 2-3 and the vertex visibility labels in columns 4-5 (purple:visible, orange:invisible).

the regression-based methods tend to fail when a human body is not fully visible in the image, *e.g.*, occluded or out of frame [16]. In this work, we aim to learn human body estimation that is faithful to the input images and robust to partial-body cases.

Instead of directly regressing SMPL parameters, we train a neural network to predict the coordinate heatmaps in three dimensions for each human joint and mesh vertex. The dense heatmap-based representation can preserve the spatial relationship in the image domain and model the uncertainty of predictions. It is shown to be effective in localizing visible joints/vertices and flexible to fit an input image faithfully [39,27,28,29]. Nonetheless, the x and y-axis heatmaps are defined in the image coordinates, which cannot represent the out-of-frame (*i.e.*, truncated by image boundaries) body parts. In addition, occlusions by objects or the human body itself could cause ambiguity for depth-axis predictions. Without knowing which joints/vertices are visible, the network tends to produce erroneous outputs on partial-body images. To address this, we propose *Visibility-aware Dense Body (VisDB)*, a heatmap-based dense representation augmented by visibility. Specifically, we train a network to predict binary truncation and occlusion labels along with the heatmaps for each human joint and vertex. With the visibility modeling, the proposed network can learn to make more accurate predictions based on the observable cues. In addition, the vertex-level occlusion predictions can serve as a depth ordering signal to constrain depth predictions. Finally, by using visibility as the confidence of 3D mesh prediction, we demonstrate that VisDB is a powerful intermediate representation which allows us to regress and/or optimize SMPL parameters more effectively. In Figure 1, we

show examples of truncation and occlusions as well as the dense human body estimations with and without visibility modeling.

Considering that most existing 3D human datasets lack dense visibility annotations, we obtain pseudo ground-truths from dense UV estimations [8]. Given the estimated UV map of an image, we calculate the pixel-to-vertex correspondence by minimizing the distance of their UV coordinates. Each vertex mapped to a human pixel is considered visible, and vice versa. Note that this covers the cases of truncation, self-occlusions, and occlusions by other objects. We further show that the dense vertex-to-pixel correspondence provides a good supervisory signal to localize vertices in the image space. Since dense UV estimations are based on part-wise segmentation masks which are robust to partial-body images, the dense correspondence loss can mitigate the inaccurate pseudo ground-truth meshes and better align the outputs with human silhouettes. To demonstrate the effectiveness of our method, we conduct extensive experiments on multiple human datasets used by prior arts. Both qualitative and quantitative results on the Human3.6M [12], 3DPW [25], 3DPW-OCC [25,45], and 3DOH [45] datasets show that learning visibility significantly improves the accuracy of dense human body estimation, especially on images with truncated or occluded human bodies.

The main contributions of our work are:

- We propose VisDB, a heatmap-based human body representation augmented with dense visibility. We train a neural network to predict the 3D coordinates of human joints and vertices as well as their truncation and occlusion labels. We obtain pseudo ground-truths of visibility labels from image-based dense UV estimates, which are also used as additional supervision signal to better align our predictions with the input image.
- We show how the dense visibility predictions can be used for robust human body estimation. First, we exploit occlusion labels to supervise vertex depth predictions. Second, we regress and optimize SMPL parameters to fit VisDB (partial-body) outputs by using visibility as confidence weighting.

## 2 Related Work

**Model-based human body estimation.** Most existing methods on human body estimation adopt a model-based representation. For instance, SMPL [24] is a widely-used statistical human body model that maps a set of pose  $\theta \in \mathbb{R}^{72}$  and shape  $\beta \in \mathbb{R}^{10}$  parameters to a 3D human mesh  $V \in \mathbb{R}^{6890 \times 3}$ . In SMPL,  $\theta$  represents the axis-angle 3D rotations of 24 joints, and  $\beta$  is the top-10 PCA coefficients of a statistical human shape space. Early methods iteratively optimize the SMPL parameters to fit the estimated 2D keypoints [2] or silhouettes [20]. Several recent works [13,35,31,34,17,19] train a deep neural network to directly regress SMPL parameters from an input image. However, the SMPL representation is not always informative enough for a network to learn as it embeds the articulated body shapes in a low dimensional space. The regression-based methods often fail on truncation and occlusion cases since the networks tend to make holistic predictions based on certain body parts only [16]. Instead, we show that

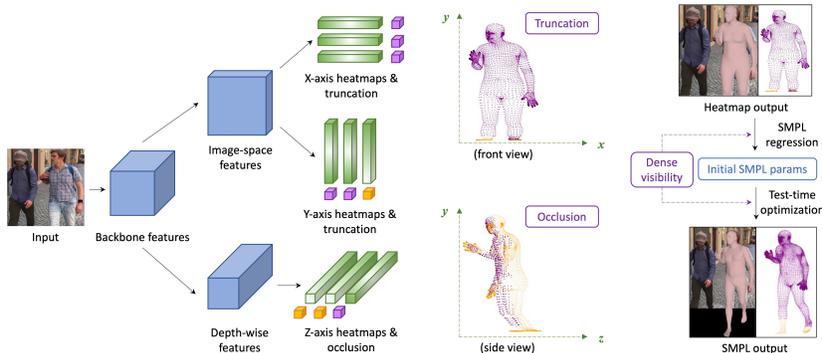
localizing 3D vertices is a more suitable task to learn for such scenarios. The network needs to learn the relationship between the parameters and the shape as well in order to estimate accurate SMPL parameters.

**Dense human body representations.** To fit the complicated shapes more faithfully, dense human body representations have been proposed, including volumetric space [41], occupancy field [37,38], dense UV correspondence [1,44], and 3D mesh [18,4,29,21,22]. Among these methods, I2L-MeshNet [29] proposes an efficient heatmap representation to estimate human joints and vertices in the image space and root-relative depth axis. It can fit the input images accurately since heatmaps preserve the spatial relationship in image features extracted by a convolutional neural network (CNN). Nonetheless, even when certain body parts are not visible in the image, this model is designed to localize all the joints and vertices within the image frame. We show that it can negatively affect the model performance and emphasize the importance of additional visibility information.

**Occlusion-aware methods.** Several methods have been proposed to deal with the challenging scenarios where human bodies are partially truncated or occluded. Muller *et al.* [30] and Hassan *et al.* [9] introduce explicit modeling of human body self-contact and human-scene interactions, respectively. These methods require ground-truth annotations which are hard to obtain. Other methods leverage human-centric heatmaps, part segmentation masks, or dense UV estimations [8], to increase the model robustness on truncated images [36], crowded scenes (occluded by other people) [40] or general occlusions [43,7,45,16]. Although effective in particular scenarios, most of them directly regress SMPL parameters which still suffer from the limited representation strength. To the best of our knowledge, the proposed VisDB representation is the first to explicitly model dense human body visibility (including truncation and all occlusion scenarios), which is trained with pseudo ground-truth visibility labels from dense UV estimates.

### 3 Approach

We illustrate our overall framework in Figure 2. In Section 3.1, we describe a heatmap-based representation which we build our method upon. Then, we introduce the proposed Visibility-aware Dense Body (VisDB) in Section 3.2. Each human joint and mesh vertex is represented by 1) three 1D heatmaps (x, y, z dimensions) which define its 3D coordinate and 2) three binary labels indicating its visibility in three dimensions. We train a network model to predict the dense heatmaps and visibility, which represents a partial body faithful to the input image. The visibility estimations can be interpreted as depth ordering signals or prediction confidence. In Section 3.3, we design a visibility-guided depth ordering loss to self-supervise depth estimation. In Section 3.4, we show that VisDB outputs can be used to fit SMPL models accurately and efficiently. We train a regression network to estimate SMPL parameters based on the joint and vertex coordinates as well as their visibility labels. During inference, we initialize the SMPL parameters by the regressor and further optimize them to



**Fig. 2. VisDB framework overview** (best viewed in color). Given an input image, the network extracts features in the image and depth coordinates, from where we predict the x, y, z heatmaps for each human joint and vertex. In addition, we predict a binary visibility label (purple:visible, orange:invisible) of each axis, *i.e.*, x-truncation, y-truncation, To obtain a more regularized and complete human body, we train a regression network to estimate SMPL parameters based on the dense 3D coordinates and visibility. At test time, we can further optimize the regressed SMPL parameters to fit the partial-body predictions from heatmaps.

align with the VisDB predictions. Finally, in Section 3.5, we exploit dense UV correspondence to obtain robust pseudo labels of visibility and weakly supervise vertex localization in the image space.

### 3.1 Preliminaries: Heatmap-based Representation

Given an input image, a prior heatmap-based method [29] estimates three 1D heatmaps  $H = \{H^x, H^y, H^z\}$  for each human joint and mesh vertex. The x and y-axis heatmaps  $H^x, H^y$  are defined in the image space, and the z-axis heatmaps  $H^z$  are defined in the depth space relative to root joint. We denote the joint heatmaps as  $H_J \in \mathbb{R}^{N_J \times D \times 3}$  and vertex heatmaps as  $H_V \in \mathbb{R}^{N_V \times D \times 3}$ , where  $N_J$  is the number of joints,  $N_V$  is the number of vertices, and  $D$  is the heatmap resolution. The heatmaps are predicted based on image features  $F \in \mathbb{R}^{c \times h \times w}$  extracted by a backbone network as follows:

$$\begin{aligned}
 H^x &= f^{1D,x}(\text{avg}^y(f^{\text{up}}(F))), \\
 H^y &= f^{1D,y}(\text{avg}^x(f^{\text{up}}(F))), \\
 H^z &= f^{1D,z}(\psi(\text{avg}^{x,y}(F))),
 \end{aligned} \tag{1}$$

where  $f^{1D,i}$  is 1-by-1 1D convolution for the  $i$ -th axis heatmaps,  $\text{avg}^i$  is  $i$ -axis marginalization by averaging,  $f^{\text{up}}$  denotes up-sampling by deconvolution, and  $\psi$  is a 1D convolution layer followed by reshaping operation. Finally, the continuous 3D coordinates of joints  $J \in \mathbb{R}^{N_J \times 3}$  and vertices  $V \in \mathbb{R}^{N_V \times 3}$  can be obtained by

applying soft-argmax on the discrete heatmaps  $H_J$  and  $H_V$ , respectively. More details can be found in [29] and supplementary material.

### 3.2 Visibility-aware Dense Body

Heatmap-based representations are shown effective in estimating human bodies in the image space. However, they often fail when the human bodies are occluded or truncated since the predictions are based on spatial image features and limited by the image boundaries. Without knowing which joints/vertices are invisible, fitting a SMPL model on the entire body tends to generate erroneous outputs. To deal with more practical scenarios where only partial bodies are visible, we make the following adaptations to a heatmap-based representation: 1) To augment the x and y-axis heatmaps, we predict binary truncation labels  $S^x, S^y$ , indicating whether a joint or vertex is within the image frame, 2) For the z-axis heatmaps, we predict a binary occlusion label  $S^z$  which specifies the depth-wise visibility. The visibility labels are predicted in a similar fashion as the heatmaps in Eq. (1):

$$\begin{aligned} S^x &= \sigma(\text{avg}^x(g^{1\text{D},x}(\text{avg}^y(f^{\text{up}}(F))))), \\ S^y &= \sigma(\text{avg}^y(g^{1\text{D},y}(\text{avg}^x(f^{\text{up}}(F))))), \\ S^z &= \sigma(\text{avg}^z(g^{1\text{D},z}(\psi(\text{avg}^{x,y}(F))))), \end{aligned} \quad (2)$$

where  $g^{1\text{D}}$  is a 1-by-1 1D convolutional layer similar to  $f^{1\text{D}}$  and  $\sigma$  is a sigmoid operator. We then concatenate the  $\{S^x, S^y, S^z\}$  predictions to obtain joint visibility  $S_J \in \mathbb{R}^{N_J \times 3}$  and vertex visibility  $S_V \in \mathbb{R}^{N_V \times 3}$ . By applying the soft-argmax operators to the predicted 1D heatmaps, the final output of our network becomes  $\{J, V, S_J, S_V\}$ , referred to as Visibility-aware Dense Body (VisDB). With the visibility information, the network model can learn to focus on the visible body parts and push the invisible parts towards the image boundaries. In our experiments (Table 3), we demonstrate that visibility modeling significantly reduces the errors of visible vertices. Moreover, the visibility labels can be seen as the confidence of coordinate predictions, which are essential to mesh regularization and completion via SMPL model fitting as described in Section 3.4.

We denote the ground-truth VisDB as  $\{J^*, V^*, S_J^*, S_V^*\}$  and train the network by using the following losses. The joint coordinate loss  $\mathcal{L}_{joint}$  is defined as:

$$\mathcal{L}_{joint} = \|J - J^*\|_1. \quad (3)$$

The vertex coordinate loss  $\mathcal{L}_{vert}$  is defined as:

$$\mathcal{L}_{vert} = \|V - V^*\|_1. \quad (4)$$

We also regress the joints from vertices using a pre-defined regressor  $W \in \mathbb{R}^{N_V \times N_J}$  and calculate a regressed-joint loss  $\mathcal{L}_{r-joint}$ :

$$\mathcal{L}_{r-joint} = \|WV - J^*\|_1. \quad (5)$$

Similar to [29], we apply losses on the mesh surface normal and edge length as shape regularization. The normal loss  $\mathcal{L}_{norm}$  and edge loss  $\mathcal{L}_{edge}$  are:

$$\mathcal{L}_{norm} = \sum_f \sum_{\{v_i, v_j\} \subset f} \left| \left\langle \frac{v_i - v_j}{\|v_i - v_j\|_2}, n_f^* \right\rangle \right|, \quad (6)$$

$$\mathcal{L}_{edge} = \sum_f \sum_{\{v_i, v_j\} \subset f} \left| \|v_i - v_j\|_2 - \|v_i^* - v_j^*\|_2 \right|, \quad (7)$$

where  $f$  is a mesh surface,  $n_f$  is the unit normal vector of  $f$ , and  $v_i, v_j$  are the coordinates of vertex  $i$  and  $j$ , respectively. Finally, we define the joint and vertex visibility loss  $\mathcal{L}_{vis}$  with binary cross entropy (BCE):

$$\mathcal{L}_{vis} = \text{BCE}(S_J, S_J^*) + \text{BCE}(S_V, S_V^*). \quad (8)$$

The VisDB prediction is illustrated in Figure 2 (left).

### 3.3 Resolving Depth Ambiguity via Visibility

Vertex-level visibility can not only be seen as model confidence for SMPL fitting but also provide depth ordering information. Intuitively, visible vertices should have lower depth value compared to the invisible vertices projected to the same pixel. We observe that VisDB network generally predicts accurate 2D coordinates and visibility, but sometimes fails at depth predictions when the human body occludes itself and the pose is less common in the training datasets. To resolve the depth ambiguity in self-occlusion cases, we propose a depth ordering loss  $\mathcal{L}_{depth}$  based on vertex visibility as follows:

$$\mathcal{L}_{depth} = \sum_x \sum_y \text{ReLU} \left( \max_{v \in Q(x,y)} v^z - \min_{\bar{v} \in \bar{Q}(x,y)} \bar{v}^z \right), \quad (9)$$

where  $Q(x, y)$  is the set of vertices projected to a discretized image coordinate  $(x, y)$  which belong to the front (occluding) part, and  $\bar{Q}$  contains the vertices of the back (occluded) part(s). The definition can be written as:

$$\begin{aligned} Q(x, y) &= \left\{ v \mid v \mapsto (x, y) \wedge P(v) = p^*(x, y) \right\} \\ \bar{Q}(x, y) &= \left\{ v \mid v \mapsto (x, y) \wedge P(v) \neq p^*(x, y) \right\}, \end{aligned} \quad (10)$$

where  $\mapsto$  denotes the discrete projection and  $P(v)$  is the part label of vertex  $v$  defined in DensePose [8]. We define the front part  $p^*(x, y)$  by finding the vertex with highest z-axis visibility score  $s^z$  as:

$$p^*(x, y) = P \left( \arg \max_{v \mapsto (x,y)} s_v^z \right). \quad (11)$$

$\mathcal{L}_{depth}$  is designed to push the self-occluded part(s)  $\bar{Q}$  to the back and non-occluded part  $Q$  to the front, where the occlusion information is given by the

z-axis visibility. Note that we compare the maximum depth (back side) of  $Q$  and the minimum depth (front side) of  $\bar{Q}$ , and thus  $\mathcal{L}_{depth}$  will be nonzero if the depth ordering disagrees with occlusion prediction and zero if the parts do not overlap anymore. Since this loss depends on accurate visibility estimations, we only apply it during the fine-tuning stage.

### 3.4 SMPL Fitting from Visible Dense Body

From the VisDB predictions, we can obtain the 3D coordinates and visibility of human joints and vertices. While the partial-body outputs are faithful to the input image from the front view, they sometimes look abnormal from a side view or contain rough surfaces. To regularize the body shape and complete the truncated parts, we perform model fitting on the visible dense body predictions. Given the coordinates and visibility of joints and vertices, we train a regression network to estimate SMPL pose  $\theta \in \mathbb{R}^{72}$  and shape  $\beta \in \mathbb{R}^{10}$  parameters. The regressed parameters are then forwarded to the SMPL model to generate the mesh coordinates denoted as  $\text{SMPL}(\theta, \beta) \in \mathbb{R}^{N_v \times 3}$ . Unlike prior art [29] which regresses a SMPL model from all the joints regardless of their visibility, our VisDB representation allows us to fit the visible partial body only. The training objectives of the SMPL regressor include SMPL parameter error, vertex error, joint error, and the negative log-likelihood of a pose prior distribution. The SMPL parameter loss  $\mathcal{L}_{smpl}$  is defined as:

$$\mathcal{L}_{smpl} = \|\theta - \theta^*\|_1 + \|\beta - \beta^*\|_1, \quad (12)$$

where  $\theta^*$  and  $\beta^*$  are the ground-truth pose and shape parameters. The SMPL vertex loss  $\mathcal{L}_{smpl-vert}$  and joint loss  $\mathcal{L}_{smpl-joint}$  are defined similarly as in Eq. (4) and (5) but weighted by visibility  $S_V, S_J$  as:

$$\mathcal{L}_{smpl-vert} = S_V \odot \|\text{SMPL}(\theta, \beta) - V_c^*\|_1, \quad (13)$$

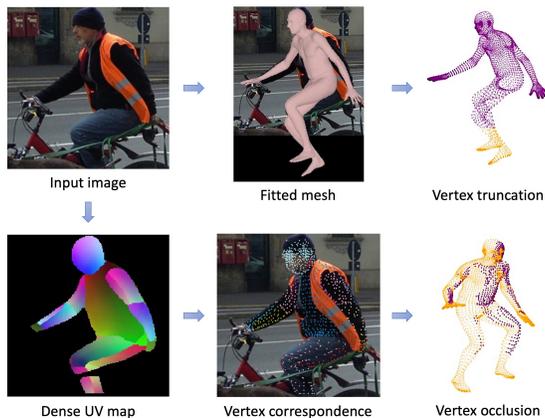
$$\mathcal{L}_{smpl-joint} = S_J \odot \|\text{WSMPL}(\theta, \beta) - J_c^*\|_1, \quad (14)$$

where  $\odot$  denotes element-wise multiplication, and  $(V_c^*, J_c^*)$  are the ground-truth root-relative coordinates of vertices and joints in the camera space. Ideally, the VisDB network makes more confident predictions on the clearly visible joints and vertices. Hence, we see the visibility labels as prediction confidence and use them to weight the coordinate losses. In addition, we apply a pose prior loss  $\mathcal{L}_{prior}$  using a fitted Gaussian Mixture Model (GMM) provided by [33]:

$$\mathcal{L}_{prior} = -\log\left(\sum_i G_i(\theta)\right), \quad (15)$$

where  $G_i$  is the  $i$ -th component of GMM.

We observe that the regressed SMPL meshes roughly capture the human pose and shape but do not always align with the VisDB predictions in details.



**Fig. 3. Dense UV correspondence and visibility labels.** Given an input image, we obtain a fitted SMPL mesh and dense UV estimation from off-the-shelf algorithms. To acquire the dense visibility labels for training, we identify the truncated vertices from the fitted mesh. From the dense UV map, we calculate the pixel-to-vertex correspondence to obtain pseudo ground-truths of vertex occlusions as well as image-space coordinates for weak supervision.

Therefore, we use the regressed parameters as initialization and propose efficient test-time optimization to further optimize the SMPL parameters against VisDB predictions. For this optimization, we apply similar losses as in Eq. (13)-(15), except that the ground-truths  $\{V_c^*, J_c^*\}$  are replaced by the VisDB predictions converted into root-relative coordinates in the camera space. Please refer to the supplemental material for details on estimating the root joint coordinate in the camera space. Since we initialize the SMPL parameters by the regression network and the use strong supervisory signal, *i.e.*, 3D joint and vertex coordinates, the test-time optimization only takes around 100 iterations to converge using an Adam optimizer [14]. We illustrate the process of SMPL regression and optimization in Figure 2 (right).

### 3.5 Exploiting Dense UV Correspondence

Most existing 3D human datasets do not provide joint visibility labels, and none annotates vertex visibility. To train our VisDB network, we obtain pseudo ground-truths from the fitted SMPL meshes and dense UV estimations. For x and y-axis truncation, we can simply identify the truncated joints/vertices by projecting the fitted mesh onto the image plane. Occlusion, however, cannot be easily inferred from the input image or fitted mesh alone. One can estimate self-occlusion by rendering a fitted mesh, but this does not capture occlusions by other objects. More importantly, the fitting algorithm used to get the pseudo ground-truth meshes is not robust to partial-body cases. To address this, we

propose to exploit dense UV correspondence between the input image and a SMPL mesh. Dense UV estimation provides the part-based segmentation mask of a human body as well as continuous UV coordinates of each human pixel, which are robust to truncation and occlusions. We calculate the UV coordinate of each pixel by applying an off-the-shelf dense UV estimation method [8]. For each human pixel  $p$ , we then find the corresponding mesh vertex  $v$  whose UV coordinate is closest to the pixel. The pixel-to-vertex  $M_P$  and vertex-to-pixel  $M_V$  mappings can be expressed as:

$$\begin{aligned} M_P &= \{p \rightarrow v \mid v = \operatorname{argmin}_{v'} \|\operatorname{UV}(v') - \operatorname{UV}(p)\|_2 \forall p\} \\ M_V &= \{v \rightarrow \{p'\} \mid M_P(p') = v \forall v\}. \end{aligned} \quad (16)$$

A vertex mapped to at least one pixel is labeled as visible or occluded otherwise.

Similar to [7,43,44], we also utilize the dense vertex-pixel correspondence as weak supervision for better alignment with the human silhouettes. For each vertex  $v$ , we calculate the center of its corresponding pixels  $M_V(v)$  and define a UV correspondence loss  $\mathcal{L}_{uv}$  as:

$$\mathcal{L}_{uv} = \sum_v s_v^z \left\| v^{x,y} - \sum_{p \in M_V(v)} \frac{p}{|M_V(v)|} \right\|_1, \quad (17)$$

where  $v^{x,y}$  is the 2D projection of vertex  $v$  and  $s_v^z$  is the binary occlusion label with  $s_v^z = 1$  indicating that the vertex  $v$  is visible. The UV correspondence loss can not only mitigate the inaccurate pseudo ground-truth meshes, but improve the faithfulness to human silhouettes since it is based on segmentation mask predictions. We empirically discover that this direct vertex-level supervision is more efficient and effective for VisDB training compared to rendering-based losses [43,6]. The proposed vertex-pixel correspondence and visibility labeling are illustrated in Figure 3.

### 3.6 Model Training and Inference

We first train the VisDB network on 3D data with mesh annotations, then fine-tune it on all training data by adding the depth ordering and UV correspondence losses. The regressor network is trained to estimate the SMPL parameters based on the estimated coordinates and visibility of joints and vertices. During inference, we apply optional optimization on the regressed SMPL parameters to best align with the VisDB predicted mesh. For the VisDB network backbone, we use a ResNet50 [11] model pre-trained on the ImageNet dataset [5]. The weights are updated by the Adam optimizer [14] with a mini-batch size of 64. We represent a human body by  $N_J = 30$  joints and  $N_V = 6890$  vertices, and the heatmap resolution  $D = 64$ . In addition, we use the ground-truth bounding boxes to crop the human region from an input image and resize it to  $256 \times 256$ . The bounding boxes of testing data are estimated by a pre-trained Mask R-CNN [10] model if not available in the dataset. We apply common data augmentations such as random scaling ( $\pm 25\%$ ), rotation ( $\pm 45^\circ$ ), horizontal flip, and color jittering ( $\pm 20\%$ )

**Table 1. Quantitative evaluations on Human3.6M [12] and 3DPW [25].** To align the settings, we train our baseline, I2L-MeshNet [29], on the same datasets, and denote it by I2L-MeshNet<sup>†</sup>. Both our mesh and SMPL parameter outputs perform favorably against the prior state-of-the-arts.

Method	Output	Human3.6M		3DPW		
		MPJPE↓	PA-MPJPE↓	MPJPE↓	PA-MPJPE↓	MPVE↓
GraphCMR [18]	Mesh	-	50.1	-	70.2	-
Pose2Mesh [4]	Mesh	64.9	47.0	89.2	58.9	109.3
I2L-MeshNet [29]	Mesh	55.7	41.1	93.2	57.7	109.2
I2L-MeshNet <sup>†</sup> [29]	Mesh	-	-	84.5	51.1	98.2
METRO [21]	Mesh	54.0	36.7	77.1	47.9	88.2
Mesh Graphormer [22]	Mesh	51.2	<b>34.5</b>	74.7	45.6	87.7
VisDB (mesh)	Mesh	<b>51.0</b>	<b>34.5</b>	<b>73.5</b>	<b>44.9</b>	<b>85.5</b>
NBF [31]	Param	-	59.9	-	-	-
HMR [13]	Param	88.0	56.8	-	81.3	-
DenseRaC [43]	Param	76.8	48.0	-	-	-
I2L-MeshNet [29]	Param	-	-	100.0	60.0	121.5
OOH [45]	Param	-	41.7	-	-	-
SPIN [17]	Param	-	41.1	-	59.2	116.4
I2L-MeshNet <sup>†</sup> [29]	Param	-	-	88.0	55.5	102.3
DSR [6]	Param	60.9	40.3	85.7	51.7	99.5
VIBE [15]	Param	65.6	41.4	82.0	51.9	99.1
TCMR [3]	Param	62.3	41.1	-	-	-
DecoMR [44]	Param	60.6	39.3	-	-	-
PARE [16]	Param	-	-	79.1	46.4	94.2
VisDB (param)	Param	<b>50.0</b>	<b>33.8</b>	<b>72.1</b>	<b>44.1</b>	<b>83.5</b>

during training. Considering that truncation and occlusion examples are rare in most 3D human datasets, we include random occlusion masks and bounding box shifting ( $\pm 25\%$ ) as additional augmentations to increase the partial-body/whole-body ratio. Our models are implemented with PyTorch [32] and trained with NVIDIA Tesla V100 GPUs. More implementation details are presented in the supplemental material.

## 4 Experiments

### 4.1 Datasets and Metrics

Following most prior arts, we adopt mixed 2D-3D training on the MSCOCO [23], Human3.6M [12], MuCo-3DHP [26], and 3DPW [25] datasets. The pseudo ground-truth meshes of Human3.6M and MSCOCO are obtained by applying SMPLify-X [33] to fit the joint annotations. We evaluate our models on the Human3.6M, 3DPW, 3DPW-OCC [25,45], and 3DOH [45] testing sets. Note that 3DOH is composed of images with object occlusions and 3DPW-OCC contains a subset of 3DPW sequences where the human bodies are partially occluded. For quantitative evaluation, we calculate the common joint and vertex error metrics in the camera space and report them in millimeters (mm), including MPJPE (mean per-joint position error) [12], PA-MPJPE (Procrustes-aligned mean per-joint position error) [46], and MPVE (mean per-vertex error) [35].

**Table 2. Quantitative evaluations on 3DOH [45] and 3DPW-OCC [25,45].** We compare VisDB with prior occlusion-aware methods to demonstrate its robustness on partial-body cases. For VisDB and I2L-MeshNet<sup>†</sup> [29], We report both the mesh and SMPL parameter (mesh/param) results.

Method	3DOH		3DPW-OCC		
	MPJPE↓	PA-MPJPE↓	MPJPE↓	PA-MPJPE↓	MPVE↓
OOH [45]	-	58.5	-	72.2	-
I2L-MeshNet <sup>†</sup> [29]	67.0/69.3	46.3/47.9	96.5/98.0	61.0/62.6	120.2/127.0
PARE [16]	63.3	44.3	91.4	57.4	115.3
VisDB	<b>62.1/60.9</b>	<b>43.2/42.7</b>	<b>90.3/87.3</b>	<b>57.1/56.0</b>	<b>114.0/110.5</b>

## 4.2 Quantitative Comparisons

**Human3.6M and 3DPW.** In Table 1, we compare the performance of our method and prior arts on the Human3.6M [12] and 3DPW [25] datasets. For VisDB and I2L-MeshNet [29], we report the results of both heatmap-based mesh outputs (mesh) and SMPL parameters (param). Our SMPL parameters are obtained from regression and test-time optimization. Note that each method uses different network backbone, human body representation, training datasets, and inference strategy. For instance, METRO [21] and Mesh Graphormer [22] adopt a transformer-based [42] network while the others use CNN backbones. VIBE [15] and TCMR [3] are video-based approaches whereas the others only take images as input. Despite these differences, VisDB performs favorably against prior methods in term of most evaluation metrics. Particularly, our method achieves larger performance gains on the 3DPW dataset since it contains more truncation and occlusion cases. The VisDB performance is most directly comparable with I2L-MeshNet [29] as we adopt similar training settings. For fair comparisons, we re-train its model on the same datasets and denote it as I2L-MeshNet<sup>†</sup>. The results demonstrate that our visibility learning improves both the mesh and SMPL outputs significantly. In prior literature, SMPL parameters generally lead to higher errors compared to dense mesh outputs, which we conjecture is caused by the difficulty to directly regress low-dimensional parameters. On the contrary, VisDB is a powerful intermediate representation that provides dense 3D information of visible partial body, allowing us to regress and optimize SMPL parameters more accurately. In our experiments, we observe that VisDB (mesh) captures the human silhouettes better but VisDB (param) produces lower errors since the ground-truth meshes are also regularized by SMPL representation.

**3DPW-OCC and 3DOH.** To emphasize the robustness on partial-body images, we further evaluate on two occlusion datasets: 3DPW-OCC [25,45] and 3DOH [45]. As shown in Table 2, VisDB produces lower errors on both datasets compared to prior occlusion-aware methods. While I2L-MeshNet<sup>†</sup> performs considerably worse on these images, the errors by our model remain relatively low.

**Table 3. Ablation studies of VisDB.** We compare the joint/vertex errors of VisDB mesh outputs on 3DPW [25] with/without individual components. The results show that truncation modeling ( $\mathcal{L}_{vis}^{x,y}$ ), occlusion modeling ( $\mathcal{L}_{vis}^z$ ), depth ordering loss  $\mathcal{L}_{depth}$ , and UV correspondence loss  $\mathcal{L}_{uv}$  each reduces the errors by a clear margin.

$\mathcal{L}_{vis}^{x,y}$	$\mathcal{L}_{vis}^z$	$\mathcal{L}_{depth}$	$\mathcal{L}_{uv}$	MPJPE	PA-MPJPE	MPVE
				84.5	51.1	98.2
✓	✓	✓	✓	79.4	47.8	91.1
✓		✓	✓	75.8	45.5	88.0
✓	✓		✓	77.3	46.3	88.9
✓		✓		74.9	45.6	87.1
✓	✓	✓	✓	<b>73.5</b>	<b>44.9</b>	<b>85.5</b>

**Table 4. Ablation studies of SMPL models.** We report the performance of SMPL outputs on the 3DPW dataset [25], which shows the effectiveness of our optimization and the importance of visibility in both regression and optimization work flows.

Regression	Optimization	MPJPE	PA-MPJPE	MPVE
-	-	73.5	44.9	85.5
w/o vis	-	79.0	48.8	96.2
w/o vis	w/o vis	77.6	47.0	93.9
w/ vis	-	74.9	45.3	87.3
w/ vis	w/ vis	<b>72.1</b>	<b>44.1</b>	<b>83.5</b>

### 4.3 Ablation Studies

**VisDB network training.** To evaluate the contribution of individual components in our method, we perform ablation studies on the 3DPW dataset [25]. Table 3 shows the performance of VisDB mesh outputs with/without truncation modeling  $\mathcal{L}_{vis}^{x,y}$ , occlusion modeling  $\mathcal{L}_{vis}^z$ , depth ordering loss  $\mathcal{L}_{depth}$ , and dense UV correspondence loss  $\mathcal{L}_{uv}$ . Without  $\mathcal{L}_{vis}^{x,y}$ ,  $\mathcal{L}_{vis}^z$ ,  $\mathcal{L}_{depth}$ , and  $\mathcal{L}_{uv}$ , the vertex error increases by 6.3mm, 3.1mm, 3.9mm, and 1.9mm, respectively. These results show that both visibility modeling and depth ordering loss play a crucial role in VisDB training.

**SMPL parameter fitting.** In Table 4, we quantitatively compare the SMPL models obtained from different methods. Given an estimated VisDB mesh, we can regress the SMPL parameters and/or optimize them during inference, and each process can be done with/without dense visibility weighting (Eq. (13) and (14)). By using visibility, the mean vertex error of regressed SMPL models drops by 8.7mm. With the proposed test-time optimization, we can further reduce the error by 3.8mm.

### 4.4 Qualitative Results

Figure 4 shows sample results by VisDB and I2L-MeshNet [29] on the 3DPW dataset [25]. I2L-MeshNet [29] regresses SMPL parameters from the entire heatmap-based mesh output, which leads to erroneous output meshes on truncated or occluded examples. VisDB predicts accurate vertex visibility labels, improving both the image-space dense body estimation and SMPL parameter optimization.



## References

1. Alldieck, T., Pons-Moll, G., Theobalt, C., Magnor, M.: Tex2shape: Detailed full human body geometry from a single image. In: ICCV (2019) [4](#)
2. Bogo, F., Kanazawa, A., Lassner, C., Gehler, P., Romero, J., Black, M.J.: Keep it smpl: Automatic estimation of 3d human pose and shape from a single image. In: ECCV. pp. 561–578 (2016) [3](#)
3. Choi, H., Moon, G., Chang, J.Y., Lee, K.M.: Beyond static features for temporally consistent 3d human pose and shape from a video. In: CVPR. pp. 1964–1973 (2021) [11](#), [12](#)
4. Choi, H., Moon, G., Lee, K.M.: Pose2mesh: Graph convolutional network for 3d human pose and mesh recovery from a 2d human pose. In: ECCV. pp. 769–787 (2020) [4](#), [11](#)
5. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR. pp. 248–255 (2009) [10](#)
6. Dwivedi, S.K., Athanasiou, N., Kocabas, M., Black, M.J.: Learning to regress bodies from images using differentiable semantic rendering. In: ICCV. pp. 11250–11259 (2021) [10](#), [11](#)
7. Guler, R.A., Kokkinos, I.: Holopose: Holistic 3d human reconstruction in-the-wild. In: CVPR. pp. 10884–10894 (2019) [4](#), [10](#)
8. Güler, R.A., Neverova, N., Kokkinos, I.: Densepose: Dense human pose estimation in the wild. In: CVPR. pp. 7297–7306 (2018) [3](#), [4](#), [7](#), [10](#)
9. Hassan, M., Choutas, V., Tzionas, D., Black, M.J.: Resolving 3d human pose ambiguities with 3d scene constraints. In: ICCV. pp. 2282–2292 (2019) [4](#)
10. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: ICCV. pp. 2961–2969 (2017) [10](#)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016) [10](#)
12. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. PAMI **36**(7), 1325–1339 (2013) [3](#), [11](#), [12](#)
13. Kanazawa, A., Black, M.J., Jacobs, D.W., Malik, J.: End-to-end recovery of human shape and pose. In: CVPR. pp. 7122–7131 (2018) [3](#), [11](#)
14. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) [9](#), [10](#)
15. Kocabas, M., Athanasiou, N., Black, M.J.: Vibe: Video inference for human body pose and shape estimation. In: CVPR. pp. 5253–5263 (2020) [11](#), [12](#)
16. Kocabas, M., Huang, C.H.P., Hilliges, O., Black, M.J.: PARE: Part attention regressor for 3D human body estimation. In: ICCV. pp. 11127–11137 (2021) [2](#), [3](#), [4](#), [11](#), [12](#)
17. Kolotouros, N., Pavlakos, G., Black, M.J., Daniilidis, K.: Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In: ICCV. pp. 2252–2261 (2019) [3](#), [11](#)
18. Kolotouros, N., Pavlakos, G., Daniilidis, K.: Convolutional mesh regression for single-image human shape reconstruction. In: CVPR. pp. 4501–4510 (2019) [4](#), [11](#)
19. Kolotouros, N., Pavlakos, G., Jayaraman, D., Daniilidis, K.: Probabilistic modeling for human mesh recovery. In: ICCV (2021) [3](#)
20. Lassner, C., Romero, J., Kiefel, M., Bogu, F., Black, M.J., Gehler, P.V.: Unite the people: Closing the loop between 3d and 2d human representations. In: CVPR. pp. 6050–6059 (2017) [3](#)

21. Lin, K., Wang, L., Liu, Z.: End-to-end human pose and mesh reconstruction with transformers. In: CVPR. pp. 1954–1963 (2021) [4](#), [11](#), [12](#)
22. Lin, K., Wang, L., Liu, Z.: Mesh graphormer. In: ICCV (2021) [4](#), [11](#), [12](#)
23. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV. pp. 740–755 (2014) [11](#)
24. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: Smpl: A skinned multi-person linear model. TOG **34**(6), 1–16 (2015) [1](#), [3](#)
25. von Marcard, T., Henschel, R., Black, M.J., Rosenhahn, B., Pons-Moll, G.: Recovering accurate 3d human pose in the wild using imus and a moving camera. In: ECCV. pp. 601–617 (2018) [3](#), [11](#), [12](#), [13](#), [14](#)
26. Mehta, D., Sotnychenko, O., Mueller, F., Xu, W., Sridhar, S., Pons-Moll, G., Theobalt, C.: Single-shot multi-person 3d pose estimation from monocular rgb. In: 3DV. pp. 120–130 (2018) [11](#)
27. Moon, G., Chang, J.Y., Lee, K.M.: V2v-posenet: Voxel-to-voxel prediction network for accurate 3d hand and human pose estimation from a single depth map. In: CVPR. pp. 5079–5088 (2018) [2](#)
28. Moon, G., Chang, J.Y., Lee, K.M.: Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In: ICCV. pp. 10133–10142 (2019) [2](#)
29. Moon, G., Lee, K.M.: I2l-meshnet: Image-to-lixel prediction network for accurate 3d human pose and mesh estimation from a single rgb image. In: ECCV. pp. 752–768 (2020) [2](#), [4](#), [5](#), [6](#), [7](#), [8](#), [11](#), [12](#), [13](#), [14](#)
30. Muller, L., Osman, A.A., Tang, S., Huang, C.H.P., Black, M.J.: On self-contact and human pose. In: CVPR. pp. 9990–9999 (2021) [4](#)
31. Omran, M., Lassner, C., Pons-Moll, G., Gehler, P., Schiele, B.: Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In: 3DV. pp. 484–494 (2018) [3](#), [11](#)
32. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. NeurIPS **32**, 8026–8037 (2019) [11](#)
33. Pavlakos, G., Choutas, V., Ghorbani, N., Bolkart, T., Osman, A.A., Tzionas, D., Black, M.J.: Expressive body capture: 3d hands, face, and body from a single image. In: CVPR. pp. 10975–10985 (2019) [8](#), [11](#)
34. Pavlakos, G., Kolotouros, N., Daniilidis, K.: Texturepose: Supervising human mesh estimation with texture consistency. In: ICCV. pp. 803–812 (2019) [3](#)
35. Pavlakos, G., Zhu, L., Zhou, X., Daniilidis, K.: Learning to estimate 3d human pose and shape from a single color image. In: CVPR. pp. 459–468 (2018) [3](#), [11](#)
36. Rockwell, C., Fouhey, D.F.: Full-body awareness from partial observations. In: ECCV. pp. 522–539 (2020) [4](#)
37. Saito, S., Huang, Z., Natsume, R., Morishima, S., Kanazawa, A., Li, H.: Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In: ICCV. pp. 2304–2314 (2019) [4](#)
38. Saito, S., Simon, T., Saragih, J., Joo, H.: Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In: CVPR. pp. 84–93 (2020) [4](#)
39. Sun, X., Xiao, B., Wei, F., Liang, S., Wei, Y.: Integral human pose regression. In: ECCV. pp. 529–545 (2018) [2](#)
40. Sun, Y., Bao, Q., Liu, W., Fu, Y., Black, M.J., Mei, T.: Monocular, one-stage, regression of multiple 3d people. In: ICCV. pp. 11179–11188 (2021) [4](#)

41. Varol, G., Ceylan, D., Russell, B., Yang, J., Yumer, E., Laptev, I., Schmid, C.: Bodynet: Volumetric inference of 3d human body shapes. In: ECCV. pp. 20–36 (2018) [4](#)
42. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: NeurIPS. pp. 5998–6008 (2017) [12](#)
43. Xu, Y., Zhu, S.C., Tung, T.: Denserac: Joint 3d pose and shape estimation by dense render-and-compare. In: ICCV. pp. 7760–7770 (2019) [4](#), [10](#), [11](#)
44. Zeng, W., Ouyang, W., Luo, P., Liu, W., Wang, X.: 3d human mesh regression with dense correspondence. In: CVPR. pp. 7054–7063 (2020) [4](#), [10](#), [11](#)
45. Zhang, T., Huang, B., Wang, Y.: Object-occluded human shape and pose estimation from a single color image. In: CVPR. pp. 7376–7385 (2020) [3](#), [4](#), [11](#), [12](#)
46. Zhou, X., Zhu, M., Pavlakos, G., Leonardos, S., Derpanis, K.G., Daniilidis, K.: Monocap: Monocular human motion capture using a cnn coupled with a geometric prior. PAMI **41**(4), 901–914 (2018) [11](#)