

# Multi-initialization Optimization Network for Accurate 3D Human Pose and Shape Estimation

Zhiwei Liu<sup>1,2,5\*</sup>, Xiangyu Zhu<sup>1,2\*</sup>, Lu Yang<sup>3</sup>, Xiang Yan<sup>6</sup>, Ming Tang<sup>1,2</sup>, Zhen Lei<sup>1,2,4</sup>, Guibo Zhu<sup>1,2</sup>, Xuetao Feng<sup>7</sup>, Yan Wang<sup>7</sup>, and Jinqiao Wang<sup>1,2,5</sup>

<sup>1</sup>National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences

<sup>2</sup>School of Artificial Intelligence, University of Chinese Academy of Sciences

<sup>3</sup>Beijing University of Posts and Telecommunications

<sup>4</sup>Centre for Artificial Intelligence and Robotics, Hong Kong Institute of Science and Innovation, Chinese Academy of Sciences

<sup>5</sup>ObjectEye, Inc.

<sup>6</sup>Dilusense Technology Corporation.

<sup>7</sup>Alibaba Group

{zhiwei.liu,xiangyu.zhu}@nlpr.ia.ac.cn,soeaver@bupt.edu.cn,xiang9292@foxmail.com,{tangm,zlei,gbzhu}@nlpr.ia.ac.cn,{xuetao.fxt,wy843@inc.com,jqwang@nlpr.ia.ac.cn}

## ABSTRACT

3D human pose and shape recovery from a monocular RGB image is a challenging task. Existing learning based methods highly depend on weak supervision signals, e.g. 2D and 3D joint location, due to the lack of in-the-wild paired 3D supervision. However, considering the 2D-to-3D ambiguities existed in these weak supervision labels, the network is easy to get stuck in local optima when trained with such labels. In this paper, we reduce the ambiguity by optimizing multiple initializations. Specifically, we propose a three-stage framework named Multi-Initialization Optimization Network (MION). In the first stage, we strategically select different coarse 3D reconstruction candidates which are compatible with the 2D keypoints of input sample. Each coarse reconstruction can be regarded as an initialization leads to one optimization branch. In the second stage, we design a mesh refinement transformer (MRT) to respectively refine each coarse reconstruction result via a self-attention mechanism. Finally, a Consistency Estimation Network (CEN) is proposed to find the best result from multiple candidates by evaluating if the visual evidence in RGB image matches a given 3D reconstruction. Experiments demonstrate that our Multi-Initialization Optimization Network outperforms existing 3D mesh based methods on multiple public benchmarks.

## CCS CONCEPTS

• **Human-centered computing**; • **Computing methodologies** → **Shape inference**; **Reconstruction**;

\*Equal contribution.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

MM '21, October 20–24, 2021, Virtual Event, China.

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-8651-7/21/10...\$15.00

<https://doi.org/10.1145/3474085.3475355>

## KEYWORDS

3D human reconstruction, 3D pose estimation, deep learning

### ACM Reference Format:

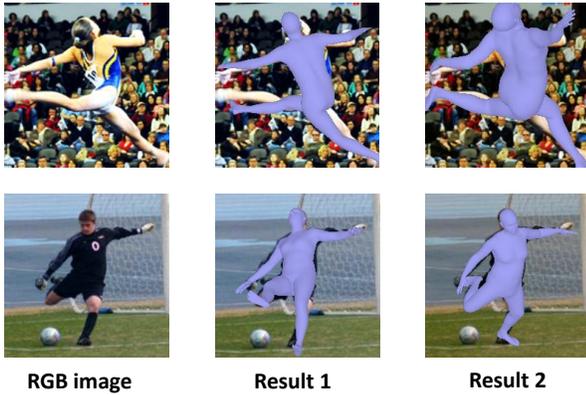
Zhiwei Liu<sup>1,2,5\*</sup>, Xiangyu Zhu<sup>1,2\*</sup>, Lu Yang<sup>3</sup>, Xiang Yan<sup>6</sup>, Ming Tang<sup>1,2</sup>, Zhen Lei<sup>1,2,4</sup>, Guibo Zhu<sup>1,2</sup>, Xuetao Feng<sup>7</sup>, Yan Wang<sup>7</sup>, and Jinqiao Wang<sup>1,2,5</sup>. 2021. Multi-initialization Optimization Network for Accurate 3D Human Pose and Shape Estimation. In *Proceedings of the 29th ACM Int'l Conference on Multimedia (MM '21)*, Oct. 20–24, 2021, Virtual Event, China. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3474085.3475355>

## 1 INTRODUCTION

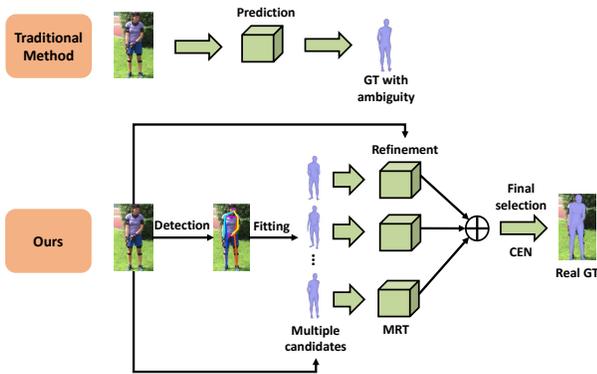
With the help of the recent developed parametric model of the human body, monocular image 3D human pose and shape reconstruction has achieved great advancements in recent years [6, 20, 21, 30]. Nowadays optimization-based and regression-based approaches are two representative research branches of this field. Optimization-based approach optimizes model parameters by fitting the human body model to 2D keypoints or other kinds of weak supervision labels, e.g. part segmentation [44], densepose [12, 45]. However, optimizing these weak labels inevitably suffers from 2D-to-3D ambiguity problem. For example, 2D observations lack depth information, and 3D keypoints lack the information of rotation angle on limb axis. Therefore, optimization-based approaches are sensitive to the choice of initialization and easy to converge into a local optimum or unrealistic result.

Regression-based approaches apply a deep CNN to regress the 3D human pose and shape parameters from a RGB input image. Due to the difficulty of in-the-wild 3D training data acquisition, existing regression based methods also heavily rely on the supervision of weak labels, including 2D projection loss [24], 3D keypoints loss [15], body silhouette loss [46] and densepose loss [34]. Since each weak label corresponds with many possible 3D bodies, network with a single output cannot always select the real 3D reconstruction from different possible results to match with the weak supervision label, leading to the unsatisfied training quality.

In general, both two kinds of approaches suffer from the ambiguity problem caused by the weak supervision. In this paper,



(a) 2D-to-3D ambiguity: different 3D results which are matched with the same 2D keypoints.



(b) Comparison between traditional regression based method and our Multi-initialization Optimization Network

**Figure 1: The problem caused by weak supervision label and our solution**

we aim to solve this issue and discover the most accurate result from multiple possible 3D reconstructions which are matched with the same weak label. Instead of using a single network that outputs a single reconstruction result in one propagation, we decompose the internal analysis process of network and propose a novel three-stage framework named Multi-Initialization Optimization Network (MION) to predict the most appropriate reconstruction result through multiple optimization branches.

Specially, in the first stage, we design a reconstruction candidate optimization strategy to optimize several different coarse reconstruction candidates for each sample. To achieve this, we first generate a human mesh candidate pool by clustering a big human motion capture dataset combination [27]. Then for each candidate in the pool, we optimize its camera parameters by fitting the 3D mesh model to a weak observation of the given sample. After this, the optimized candidate with a relative small fitting loss can be selected to get into the next inference stage and regarded as an initialization to start one optimization branch. Moreover, we accelerate the parameter optimization process by directly computing its

closed-form solution. Based on the multiple candidates, our method has more possibilities to avoid the local optimal and find the real 3D reconstruction.

In the second stage, we continue to push each optimization branch forward and design a mesh refinement transformer (MRT) to refine each coarse reconstruction candidate. This transformer framework has two advantages. Firstly, to achieve a refinement process, the initial 3D reconstruction information from the last stage can be effectively encoded into the transformer by a novel Projected Normalized Coordinate Code (PNCC) positional encoding. Secondly, transformer can adaptively learn the non-local relationships between different body joints during training stage, which is significant for enhancing the body structure prediction ability.

Finally, the refined 3D reconstruction candidates from different optimization branches need to be aggregated to one result. In the last stage, a Consistency Estimation Network (CEN) is proposed to distinguish if the form of each 3D reconstruction matches the visual evidence in the input image and select the best 3D reconstruction result. In order to get rid of the limitation of label ambiguity, CEN benefits from a specific data synthesis strategy which generates the accurate 3D ground-truth for each training sample. The overview of our Multi-Initialization Optimization Network (MION) is shown in Fig. 1.

Our contributions can be summarized as follows:

- In order to deal with the issue that weak supervision labels have ambiguities. We propose a novel three-stage framework named Multi-Initialization Optimization Network (MION) to predict appropriate human pose and shape reconstruction through multiple optimization branches.
- In our MION, instead of predicting a single result, for each sample, we calculate multiple coarse reconstruction candidates as the initializations to start different optimization branches. Compared with a single prediction, multiple candidates give more possibilities to avoid the local optimal and find the real 3D reconstruction.
- Given the different initial reconstruction candidates, we design a mesh refinement transformer (MRT) with a novel Projected Normalized Coordinate Code (PNCC) positional encoding to further refine each coarse reconstruction via a self-attention mechanism. Then a Consistency Estimation Network (CEN) is proposed to select the best 3D reconstruction result from all the optimization branches.
- Both qualitative and quantitative experiments show that our MION significantly improves the performance of monocular image 3D human reconstruction and achieves the state-of-the-art result among other methods.

## 2 RELATED WORKS

**Optimization based methods.** With the development of parametric 3D human body model, such as SCAPE [3], SMPL [25], optimization based 3D human reconstruction method becomes an important branch in the research community. This kind of method infers 3D reconstruction by fitting a parametric model to match the given 2D observation. Early optimization methods [1, 36] fit the human body model by the manually generated keypoints and silhouettes labels. These methods rely on manual intervention and

generalize badly to images in the wild. Federica *et al.* [5] propose first automatic 3D human model reconstruction method SMPLify which fits SMPL model to the 2D keypoints predicted by CNN detector [33]. Meanwhile, objection function contains different regularization terms to ensure the optimization can produce a plausible result. In order to further improve the performance, more supervision information are incorporated into the optimization target, such as silhouette [22], scene constrains [48] and multi-view [14].

Generally, these optimization based methods are sensitive to the choice of initialization and tends to have a slow optimization speed. Meanwhile, the optimization process only converges to one local optimal result by the 2D observation input without appearance information. Thus it suffers from the problem of label ambiguities.

**Learning based methods.** Another representative method is learning base reconstruction method, which learns a human model parameter regressor by a data-driven way. Because of the lack of in-the-wild 3D reconstruction paired training data. Existing methods focus on the weak supervision way to solve this issue. HMR [18] directly regresses SMPL parameters from images by a CNN and adds iterative regression to further improve the accuracy. It also proposes an adversarial prior in case the reconstruction is not realistic. SPIN [20] incorporates a optimization into the network learning process, where the predicted parameter is refined by a optimization process to further supervise the network. These two methods both apply 2D and 3D keypoints as weak supervision signal. Inspired by the dense correspondence representation used in DensePose [12], various learning based methods [11, 34] regard IUV map as intermediate representation or weak supervision label for improving the regression CNN. On the other hand, in order to remove the limitation of SMPL parameter space, many learning based methods do not rely on the parametric model and directly regress the coordinates of each vertices on the mesh. BodyNet [39] regresses a volumetric representation of 3D human by a Voxel-CNN. Densebody [46] and DaNet [49] use a UV position map to represent 3D human body.

All the aforementioned learning based methods rely on weak supervision label during training stage. It is flawed in that the label with ambiguities cannot always lead the network to predict the real reconstruction. In this work, with the same weak supervision labels, we propose a multi-path optimization based reconstruction framework to reduce the possibility of getting into local optimal. In order to deal with the occlusion cases with ambiguities, Biggs *et al.* [4] also proposes to predict a candidate set which contains different possible reconstruction results. However, their network is only trained on the single Human3.6m dataset [15] where each sample already has 3D reconstruction ground-truth. Thus it still cannot solve the problem raised by the weak label.

**Synthetic data.** Since it is difficult to collect in-the-wild samples with 3D ground-truth, synthetic data plays an important role in 3D human pose and shape estimation task. Pavlakos *et al.* [32] adopt joint heatmap and silhouette as intermediate representation and design two decoders to respectively predict the pose and shape parameter from the above two representations. Without introducing appearance information, the decoders can be trained by synthetic data. Xu *et al.* [43] propose a network to decode human body from a synthetic IUV map. Although the above methods benefit from the abundant synthetic training data, the inference process from

weak label to 3D ground-truth is an ambiguity task which cannot be solved by even human.

Existing synthetic data based methods usually ignore the appearance information. Different from them, we utilize synthetic data to help the network discriminate if the visual evidence in RGB image matches a given 3D reconstruction, which can be used to select the best result from multiple candidates.

**Human structure dependence.** Human body structure in natural world has a strong prior, which builds the dependencies between different parts of the whole human body. Some existing methods explicit model such dependencies in their learning framework. CMR [21] apply a Graph-CNN [19] to model the interactions between different vertices in the inference framework. METRO [23] replaces Graph-CNN with a Transformer encoder [41] to model the interactions. However the query in its Transformer only contains the global feature of the input image and abandons detailed local information. In this work, we apply a full encoder-decoder Transformer structure which maintains the local image appearance information to refine the body structure.

### 3 METHOD

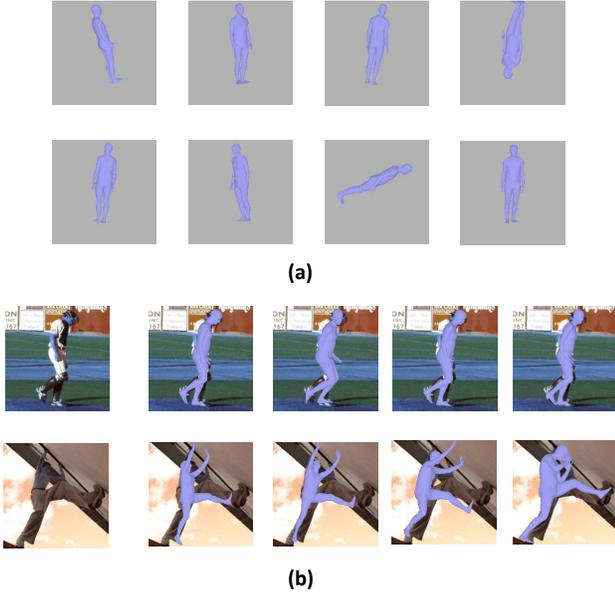
As above mentioned, the whole framework of our Multi-Initialization Optimization Network (MION) has three stages. In this section, we explain the details of each stage in sequence.

#### 3.1 reconstruction candidate selection

For common traditional CNN regression based 3D human reconstruction frameworks, training with weak supervision labels might make the network fall into local optimum. As long as the predicted 3D reconstruction matches with the provided weak supervision labels, for e.g., the 2D landmarks or masks. the network training will stop no matter if the current prediction is correct. In this work, we aim at designing a new three-stage framework begins with multiple initializations.

At our first stage, considering that nowadays 2D body keypoints telenology already achieves a high performance, For each training sample, we first utilize HRNet [38] to detect its 2D body keypoints which represents the weak supervision label. Then we expect to coarsely locate all the representative possible solutions according to this 2D keypoints in the whole solution space. To this end, a candidate selection strategy is proposed to calculate multiple possible 3D reconstruction candidates for each sample. Each candidate can be regard as an initial optimization point.

Specifically, we regard a huge human motion capture dataset named AMASS [27] as the source of our candidate and assume it contains all the possible human poses shown in natural images. Each 3D human mesh sample in AMASS is expressed by a set of SMPL [25] parameters. SMPL is a parametric model for human body mesh representation, which maps the shape parameters  $\beta \in R^{10}$  and the pose parameters  $\theta \in R^{72}$  to the human body mesh  $M \in R^{V \times 3}$  by a linear model. Based on this, to generate the reconstruction candidates of a specific sample, we first optimize a perspective projection matrix for each candidate in AMASS by fitting the 3D keypoints from candidate meshes to the 2D keypoints from the given sample. Then the candidate with a low fitting loss



**Figure 2: (a) Visualization of camera orientation cluster result. (b) Visualization of the selected candidates**

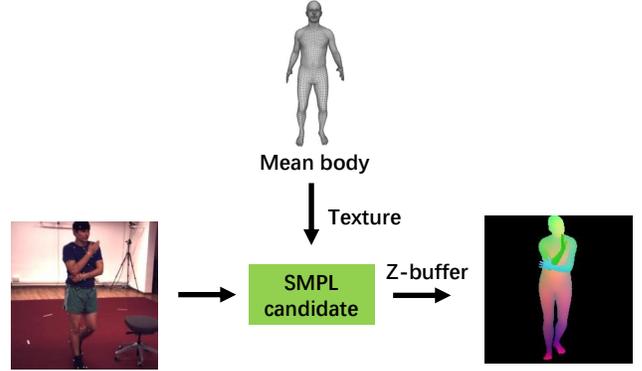
value can be selected as an possible solution to get into the next inference stage.

However, AMASS is a huge candidate dataset. It is unrealistic to optimize all the perspective projection matrix parameters for each candidate because of the high computational cost. In order to reduce the cost, we calculate the representative candidate bodies from the whole AMASS dataset by k-means clustering. Then the approximate optimal solution can be directly selected from the cluster centroids. In this work, we set the number of pose parameter clusters to 10000 and the number of camera orientation parameter clusters to 30, making the whole candidate pool contain 300000 members. Fig 2(a) visualizes different cluster centroids of camera orientation parameter. The optimization target function can be defined as:

$$\max_T \mathcal{L} = \|\Pi_K(J_{cand}^{3d}, T) - J_{gt}^{2d}\|_2^2 \quad (1)$$

$$\Pi_T(J_{cand}^{3d}, T) = \begin{bmatrix} f & 0 & c_1 \\ 0 & f & c_2 \end{bmatrix} \begin{pmatrix} \frac{J_x + T_x}{J_z + T_z} \\ \frac{J_y + T_y}{J_z + T_z} \\ 1 \end{pmatrix} \quad (2)$$

where  $J_{cand}^{3d} \in R^{N \times 3}$  denotes the coordinates of a set of 3D keypoints regressed from a candidate mesh, and  $J_{gt}^{2d} \in R^{N \times 2}$  is the corresponding 2D keypoint.  $\Pi_K$  is the perspective projection function from 3D to 2D.  $T$  is camera translation parameters for optimizing.  $f$  is the pre-defined focal length of camera,  $c_1$  and  $c_2$  are the pre-defined camera center parameters. Since the camera orientation parameters and pose parameters are already known, only the three camera translation parameters are required to be optimized. To further accelerate the computing speed, we directly compute the



**Figure 3: The process of PNCC generating**

closed-form solution of the ternary homogeneous linear equations. This implementation omits the time cost on Gaussian elimination compared with the traditional least square method. Based on this, our candidate selection process only takes 18 ms on GPU.

After fitting all the 3D candidates to the given 2D keypoints, we select a set of candidates with low fitting loss and a large pose variance, although the selected candidates match with the same 2D keypoints, they still might distribute dispersedly in the pose solution space to give more initializations for avoiding the local optimal. The selected reconstruction candidates of two samples is shown in Fig. 2(b).

### 3.2 Mesh Refinement Transformer

Most selected candidates from the first stage are coarsely optimized and need further refinement. In order to alleviate local optimal solutions and increase the possibility of finding the real 3D reconstruction, in this stage, we provide an individual optimization branch to refine each coarse candidate. Each branch only focuses on refining the details to make the initialization more compatible with the weak supervision label.

Inspired by a series of recent vision transformer (ViT) works [6, 7, 10, 50], we determine to use the ViT architecture which has two advantages for the body reconstruction refinement task. Firstly, ViT follows a sequence prediction format by regarding the input image as a sequence of different local patches. This manner allows the whole inference framework pay attention on the details of each local patch, which is suitable for the refinement task. Secondly, the self-attention mechanism of transformers, which explicitly models all pairwise interactions between elements in a sequence makes our architectures particularly suitable for learning the relationship of different body parts.

To construct our transformer network, we first use a backbone network (e.g. resnet) to extract image feature [7]. Then three deconvolution layers are added to the top layer of backbone to upsample the feature map and recover more spatial information. Finally we flatten feature map to make a feature vector sequence and input it into a transformer encoder, as shown in Fig. 4. Thus each element in the sequence represents a local patch of input image.

In order to refine the SMPL parameter of a given candidate, we need to encode the candidate as a 2D map for CNN, which

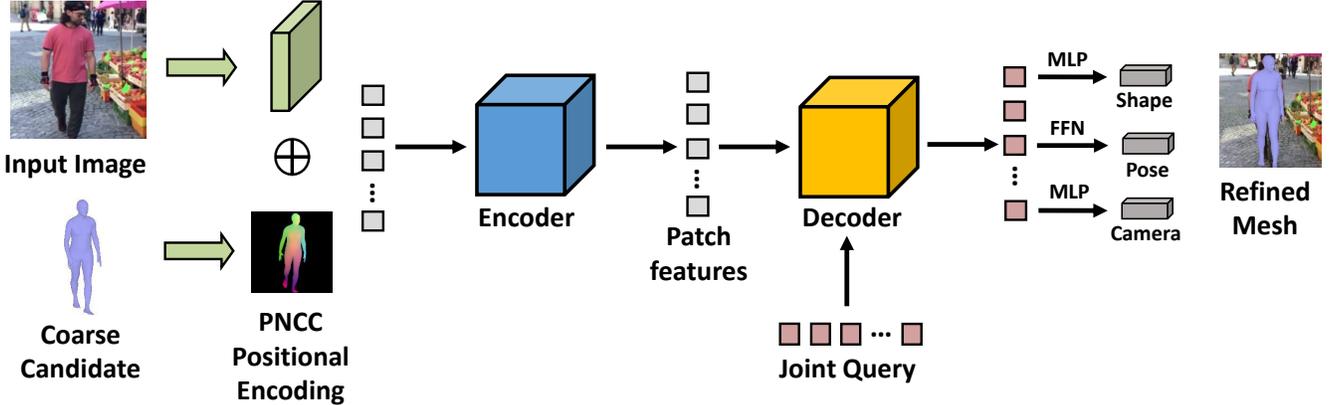


Figure 4: The framework of Mesh Refinement Transformer

corresponds to the positional encoding in the common transformer. To this end, inspired by Projected Normalized Coordinate Code (PNCC) used in 3D face alignment work 3DDFA [51], we propose a Projected Normalized Coordinate Code based positional encoding (PNCC-PE) which encodes the initial SMPL parameter information into the transformer input. Specifically, we first normalize the mean body point cloud to 0-1 in  $x, y, z$  axis as Eq. (3)

$$NCC_d = \frac{\bar{S}_d - \min(\bar{S}_d)}{\max(\bar{S}_d) - \min(\bar{S}_d)} \quad (d = x, y, z), \quad (3)$$

where  $\bar{S} \in R^{V \times 3}$  is the point cloud tensor of mean SMPL model. We regard the three channel normalized coordinate code (NCC) of each vertice as its texture. Then PNCC is generated by adopting Z-Buffer to render the body mesh with initial SMPL parameter on a zero value background, as shown in Eq. (4)

$$PNCC = Z\text{-Buffer}(\mathcal{M}(\theta_{cand}, \beta), \gamma_{cand}, NCC) \quad (4)$$

where  $\theta_{cand}$  and  $\gamma_{cand}$  are the pose and camera parameters of a candidate predicted by RCG. All the candidates apply a mean SMPL shape parameter  $\beta$ . The PNCC calculation process is shown in Fig. 3 After we get PNCC map, the sine and cosine functions used in original transformer positional encoding [41] are applied to transfer the 3-channel PNCC into the final positional encoding feature (PNCC-PE). Each value on the PNCC-PE can be represented by:

$$\begin{aligned} PNCC - PE_{(pos, 2i)} &= \sin(PNCC_{(pos)} / 10000^{2i/d_{model}}) \\ PNCC - PE_{(pos, 2i+1)} &= \cos(PNCC_{(pos)} / 10000^{2i/d_{model}}) \end{aligned} \quad (5)$$

where  $pos$  is denotes the spatial position index on PNCC map,  $i$  is the channel index of PNCC-PE. PNCC-PE and backbone feature have the same channel dimension number. Thus they can be added and input to the transformer encoder. Based on this, PNCC-PE builds the relationship between each local image patch and its corresponding body part by the provided prior information of candidate. This design helps the tranformer pay more attention to the useful local patch details when refining each body parts.

The encoder of Mesh Refinement Transformer (MRT) follows the standard multi-head attention and feed-forward networks architecture. Since the pose parameter of a SMPL model is composed of a set of joint rotation vectors, we decompose the pose representation and make each input query of our decoder represent the embedding of one joint rotation vector. Then the output sequence of decoder is the refined pose parameter. Meanwhile, we add two MLP networks which respectively predict the camera parameter and shape parameter from the top sequence feature. Therefore, all the weak supervision labels (e.g. 2D keypoints) can be involved into the training by a perspective projection. The overview of our MRT is shown in Fig. 4

The total loss of our MRT is as follows:

$$L = w_1 * L_{smpl} + w_2 * (L_J^{3d} + L_J^{2d}) \quad (6)$$

$$L_{smpl} = \|[\theta_{reg}, \beta_{reg}] - [\theta_{gt}, \beta_{gt}]\|_2 \quad (7)$$

$$L_J^{3d} = \|J_{gt}^{3d} - R(V_{reg})\|_2 \quad (8)$$

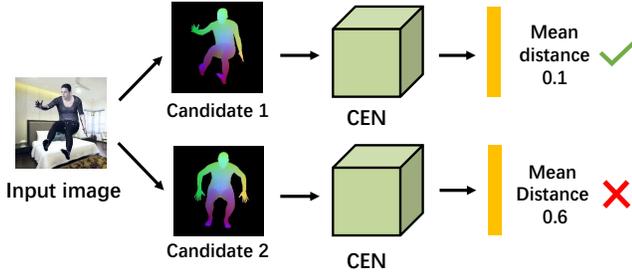
$$L_J^{2d} = \|J_{gt}^{2d} - \Pi_K(R(V_{reg}))\|_2 \quad (9)$$

where  $L_v$ ,  $L_J^{3d}$  and  $L_J^{2d}$  are the vertice loss, 3D joint loss and 2D joint projection loss.  $w_1$  and  $w_2$  are the weights of different loss functions.

### 3.3 Consistency Estimation Network

Applying multiple initalizations usually leads to different optimized results. Some results might fall into local optimum and some results might be close to the real 3D ground-truth. Therefore, it is necessary to find the best reconstruction candidate among all the refined candidates. To this end, in the last stage, we consider the optimum path selection as a scoring problem and propose a Consistency Estimation Network (CEN) to solve this problem by utilizing the synthesis training data SURREAL [40].

The target of the network is to identify if one 3D reconstruction matches the visual evidence of input human image. For this purpose, we need a dataset with ground-truth 3D body mesh, which can be achieved by data synthesis. Specifically, when generating a synthesis sample, we select one training sample and randomly pick two candidates from its candidate collection generated in the



**Figure 5: The inference process of Consistency Estimation Network.**

first stage. Then we use the SMPL parameter of one candidate to render a RGB synthesis image. Following the synthesis method of SURREAL [40] dataset, the body texture are selected from its own texture set and the background image is selected from a subset of LSUN dataset [47].

The other picked candidate is used to simulate a predicted result from the last stage. Following the same way as MRT, we use its SMPL parameters to render a PNCC map to represent the reconstruction information. Finally, as shown in Fig. 5, PNCC is concatenated with the input image and sent to a CNN which regresses the vertice wise distance of the two picked candidates. The loss function of our Consistency Estimation Network (CEN) is as follows:

$$L_i = \|P_{reg}^i - \|V_{gt}^i - V_{cand}^i\|_2\|_2 \quad (10)$$

where  $i$  is the index of one vertice on the body point cloud,  $P_{reg}^i$  is the  $i$ th output of our CEN, meaning the regressed score of  $i$ th vertice.  $\|V_{gt}^i - V_{cand}^i\|_2$  denotes the distance between the PNCC encoding body and RGB image encoding body in a normalized point cloud space. During the inference stage, we compute the average distance of all the output vertice distance and choose the reconstruction result in the optimization branch with the lowest distance as our final result.

### 3.4 Implementation details

For our candidate selection strategy in the first stage, after we fit all the candidates to the 2D keypoints of current sample, all the candidates with a fitting loss lower than 2000 are chosen to be an available candidate. Then we iteratively select the candidate which has the largest pose parameter distance with the selected candidates and put it into the final candidate collection.

For the Mesh Refinement Transformer (MRT), the backbone network adopts the architecture of ResNet-50 [13]. Note that we remove the last fully connection layers in original ResNet-50 and add three deconvolution layers to make a fully convolution network (FCN) as our backbone. The FCN receives the  $224 \times 224$  input image and produces  $56 \times 56$  feature maps with 384 channels. In order to match with image feature, the rendered PNCC has the same resolution of  $56 \times 56$  and each channel of PNCC is transferred into a position encoding map with 128 channel. During MRT training stage, the loss weight of SMPL parameter regression is set to 1 and the loss weight of 2D keypoints and 3D keypoints regression is set to 5. The data augmentation techniques includes rotation  $[-60^\circ, 60^\circ]$

color jittering  $[0.6, 1.4]$  and flipping, are applied randomly to input images. We adopt the AdamW [26] optimizer with an initial learning rate of  $5 \times 10^{-5}$  to train the MRT model, and reduce the learning rate to  $5 \times 10^{-6}$  after 20 epochs. The training process stops after 60 epochs. MRT is trained on 4 Titan X GPUs with a batch size of 64. During training and inference stage, all the 2D keypoints used in the first stage are predicted by the HRNet-W48 network in MMPose toolbox [9].

For training the Consistency Estimation Network (CEN), we adopts a ResNet-34 network as backbone. The initial learning rate is set to 0.01. We train our model for 100 epochs and lower the learning rate by a factor of 10 after 50 epochs.

## 4 EXPERIMENTS

This section focuses on the empirical evaluation of the proposed method. First, we present the datasets and evaluation metrics that we employed for quantitative and qualitative evaluation. Then, we conduct extensive ablation experiments and comparison experiments to verify the effectiveness of our method.

### 4.1 Datasets and Evaluation Metrics

The networks mentioned in this work are trained on various datasets. Specifically, the training sets of our Mesh Refinement Transformer (MRT) is dominated by weak supervision datasets, including Human3.6M [15], LSP [16], MPII [2], COCO [24], LSP-Extended [17], MPI-INF-3DHP [28]. Our Consistency Estimation Network (CEN) is trained on UP-3D [22] and SURREAL [40] which have 3D ground-truth. We conduct the evaluations on the test set of Human3.6M and 3DPW [42]. To get a pair comparison with earlier state-of-the-art method, we use the same evaluation metric with SPIN [20].

**Human3.6M:** It is an indoor benchmark for 3D human pose estimation. It includes multiple subjects performing actions like Eating, Sitting and Walking. Following typical protocols, e.g., [20], we use subjects S1, S5, S6, S7, S8 for training and we evaluate on subjects S9 and S11.

**LSP:** LSP with its extension is a standard 2D human pose estimation dataset which is collected by the images from sports activities. This dataset has large variance in terms of appearance and especially articulations. Each person in this dataset is labeled with total 14 joints which are used for weak supervision by the perspective projection.

**MPII:** MPII is a 2D human pose estimation dataset which covers a wide range of human activities with 25k images containing over 40k people. We use it for weak supervision during training stage.

**MPI-INF-3DHP:** It is a dataset captured with a multi-view setup mostly in indoor environments. No markers are used for the capture, so 3D pose data tend to be less accurate compared to other datasets. We use the provided training set (subjects S1 to S8) for training.

**UP-3D:** It is a recent dataset that collects color images from 2D human pose benchmarks. SMPLify [5] is utilized to generate 3D human shape candidates for each sample. The candidates were evaluated by human annotators to select only the images with good

Method	MPJPE	PA-MPJPE
1-path (MRT)	62.31	45.48
2-path (MION)	61.59	44.75
3-path (MION)	59.98	42.81
4-path (MION)	58.78	41.86
5-path (MION)	<b>56.88</b>	<b>41.59</b>
6-path (MION)	58.71	42.02

**Table 1: Analysis the difference of using different number of optimization paths. MPJPE and PA-MPJPE are used as evaluation metric.**

3D shape fits. It comprises 8515 images, where 7818 are used for training.

**SURREAL:** It provide a tool to generate synthetic image examples with 3D shape ground truth. In this work, we select the SMPL pose parameters from AMASS and select the background images from LSUN dataset [47]. Each training sample is generated by a common rendering pipeline and has an accurate 3D ground-truth.

**Evaluation Metrics:** Following previous method [11, 31, 49], the evaluation is conducted by two popular protocols: Mean Per Joint Position Error (MPJPE) and the MPJPE after rigid alignment of the prediction with ground truth using Procrustes Analysis (MPJPE-PA). Both two metrics measures the Euclidean distances between the ground truth joints and the predicted joints.

## 4.2 Ablation experiment

To evaluate the effectiveness of each component proposed in our method, we conduct ablation experiments on Human3.6M under various settings.

**Effect of different number of initializations** In order to study the effect of using different number of initializations in our 3D human reconstruction framework, we respectively test our Multi-Initialization Optimization Network (MION) under different numbers of optimization initializations. As shown in Tab. 1, when only using one initialization to inference the reconstruction, our MPTN degrades to a single Mesh Refinement Transformer (MRT). Not surprisingly, one initialization method leads to the worst result. When applying more candidates into the whole framework, the performance is continuously improved until the number of candidates reaches 5. This phenomenon indicates that applying multiple different optimization initializations is an effective way to alleviate falling into local optimum and improve the overall performance in 3D human reconstruction task

**Effectiveness of PNCC position encoding** To demonstrate the superiority of our Projected Normalized Coordinate Code based positional encoding (PNCC-PE) over traditional position encoding method, we further conduct experiment to investigate the impact of PNCC-PE. Tab. 2 shows the ablation study on Human3.6M. MION (w/o PNCC-PE) replaces the PNCC-PE by the traditional sinusoidal position encoding [41]. As we can see, compared with traditional position encoding, the proposed PNCC-PE significantly improve the performance on 3D human reconstruction task: 56.88%

Method	MPJPE	PA-MPJPE
MION w/o PNCC-PE	61.14	43.01
MION	<b>56.88</b>	<b>41.59</b>

**Table 2: Analysis of the effective of using the PNCC position encoding in a transformer framework. MPJPE and PA-MPJPE are used as evaluation metric.**

Method	MPJPE	PA-MPJPE
MION w/o CEN	62.47	44.71
MION with CEN	<b>56.88</b>	<b>41.59</b>
MION best path	52.17	39.82

**Table 3: Analysis of the effective of using Consistency Estimation Network (CEN) in a transformer framework. MPJPE and PA-MPJPE are used as evaluation metric.**

vs 61.14%. We attribute this phenomenon to the initialization information brought by the PNCC-PE makes the joint queries of decoder easily focus on the corresponding local patch feature.

**Effectiveness of Consistency Estimation Network** Consistency Estimation Network (CEN) is one of the key steps in our algorithm. To demonstrate the effectiveness of our CEN on selecting the best optimized result, we make a performance comparison between different 3D reconstruction selection methods in the third stage. First, we apply the random selection strategy instead of CEN. As indicated in Tab. 3, compared with random selection strategy MION w/o CEN, our CEN reduces the MPJPE error from 62.47% to 56.88%, which proves the effectiveness of CEN on predicting the consistency between RGB image and given human body parameter. We also evaluate the upper bound of CEN by always selecting the result with lowest error. As we can see, the performance gets further improvement: from 56.88% to 52.17%, meaning that the multiple initialization method has the potential to achieve a better result.

## 4.3 Comparison experiment

### Comparison on the In-door Dataset.

We evaluate the performance of our methods on the in-door Dataset Human3.6M in terms of 3D pose estimation accuracy. We train our model following the setting of SPIN [20] and utilize Human3.6M, LSP, MPII, COCO and MPI-INF-3DHP as the training set. Quantitative results are reported in Tab. 4. It shows the results of our approach against other the state of the art methods which output a full mesh of the human body (SMPL, in particular). As we can see, training with the same amount of weak supervision label, our method significantly outperforms other methods on MPJPE metric and achieves a competitive performance on PA-MPJPE metric.

**Comparison on In-the-wild Dataset.** Since the lack of in-the-wild 3D supervision labels, reconstructing 3D human model on in-the-wild outdoor images is much more challenging due to factors such as extreme poses, appearance variations and heavy occlusions. We conduct evaluation experiments on 3DPW datasets to compare

Method	MPJPE	PA-MPJPE
SMPLify [5]	-	82.3
NBF [30]	-	59.9
HMR [18]	88.0	56.8
GraphCMR [21]	-	50.1
HoloPose [11]	64.3	50.6
TexturePose [31]	-	49.7
DaNet [49]	61.5	48.6
DenseRaC [43]	76.8	48.0
Pose2Mesh [8]	64.9	47.0
SPIN [20]	62.3	<b>41.1</b>
MION	<b>56.88</b>	<b>41.59</b>

**Table 4: Comparison with state of the art on Human3.6M dataset. MPJPE and PA-MPJPE are used as evaluation metric.**

Method	MPJPE	PA-MPJPE
HMR [18]	-	81.3
GraphCMR [21]	-	70.2
STRAPS [35]	-	66.8
SPIN [20]	-	59.2
Pose2Mesh [8]	89.2	58.9
I2LMeshNet [29]	93.2	57.7
Song <i>et al.</i> [37]	-	55.9
MION	<b>81.98</b>	<b>52.34</b>

**Table 5: Comparison with state of the art on 3DPW dataset. MPJPE and PA-MPJPE are used as evaluation metric.**

our MION with previous 3D human pose and shape estimation methods. As indicated in Tab. 5, we can see our method outperforms all the other methods under the challenging scenarios, which proves the robustness and generalization of our framework.

#### 4.4 Analysis experiment

**Running speed analysis.** We evaluate the inference speed of our method and the state of the art method SPIN [20] on the same hardware platform (one Titan X Pascal GPU). Both two methods apply the same ResNet50 as the backbone. The whole running time of our MION is 128 ms. Note that our candidate selection process from 300k candidate pool only takes 18 ms, which is much less than the following CNN inference stages spend (63 ms for the second stage, 47 ms for the third stage). In general, since the running time of SPIN is 59ms and we perform much better than SPIN on 3DPW dataset, 52.34 vs 59.2 in terms of PA-MPJPE, we believe it is worth taking more time for this better solution.

**The performance on shape recovery.** In order to verify the effectiveness of our method on shape recovery, we evaluate the shape accuracy of three different methods on 3DPW dataset, including baseline (1-path), MION (5-path) and the state of art method SPIN. We apply two evaluation metrics in the experiment, the first

Method	Vertex L2 error	Shape Param L2 error
Baseline (1-path)	0.142	5.383
SPIN	0.116	4.398
MION (5-path)	<b>0.094</b>	<b>3.227</b>

**Table 6: The comparison of different methods in terms of shape recovery accuracy. Mean vertex L2 error and Shape Parameter L2 error are used as evaluation metric.**

one is mean vertex L2 error of two body clouds, the second one is the L2 distance between two SMPL shape parameter vectors. The results are shown in Tab. 6

## 5 CONCLUSION

This work aims to solve the 2D-to-3D ambiguity problem of training with weak supervision labels. Instead of training a regression network with one output, we propose to apply multiple initializations and different optimization branches to avoid the network easily get stuck in local optimum. Specifically, we propose a three-stage framework named Multi-Initialization Optimization Network (MION). In the first stage, we strategically select different coarse 3D reconstruction candidates which are compatible with the 2D keypoints of input sample. Regarding each candidate as an initialization, in the second stage, we design a mesh refinement transformer (MRT) to respectively refine each coarse reconstruction result via a self-attention mechanism. Finally, a Consistency Estimation Network (CEN) is proposed to find the best result from multiple candidates by evaluating if the visual evidence in RGB image matches a given 3D reconstruction. Experiments demonstrate that our framework outperforms existing 3D mesh based methods on multiple public benchmarks. Future work can focus on improving the efficiency of this framework.

## ACKNOWLEDGMENTS

Acknowledgement: This work was supported by the Research and Development Projects in the Key Areas of Guangdong Province (No.2019B010153001) and National Natural Science Foundation of China under Grants No.61772527, No.61976210, No.62076235, No.61806200, No.62002356, No.62002357, No.62006230.

## REFERENCES

- [1] Ankur Agarwal and Bill Triggs. 2005. Recovering 3D human pose from monocular images. *IEEE transactions on pattern analysis and machine intelligence* 28, 1 (2005), 44–58.
- [2] Mykhaylo Andriluka, Leonid Pishchulin, Peter Gehler, and Bernt Schiele. 2014. 2d human pose estimation: New benchmark and state of the art analysis. In *Proceedings of the IEEE Conference on computer vision and Pattern Recognition*. 3686–3693.
- [3] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. 2005. Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*. 408–416.
- [4] Benjamin Biggs, Sébastien Ehrhardt, Hanbyul Joo, Benjamin Graham, Andrea Vedaldi, and David Novotny. 2020. 3D Multi-bodies: Fitting Sets of Plausible 3D Human Models to Ambiguous Image Data. *arXiv preprint arXiv:2011.00980* (2020).
- [5] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. 2016. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. In *European conference on computer vision*. Springer, 561–578.

- [6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *European Conference on Computer Vision*. Springer, 213–229.
- [7] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. 2020. Pre-trained image processing transformer. *arXiv preprint arXiv:2012.00364* (2020).
- [8] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. 2020. Pose2Mesh: Graph convolutional network for 3D human pose and mesh recovery from a 2D human pose. In *European Conference on Computer Vision*. Springer, 769–787.
- [9] MMPose Contributors. 2020. OpenMMLab Pose Estimation Toolbox and Benchmark. <https://github.com/open-mmlab/mmpose>.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiuhua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [11] Riza Alp Guler and Iasonas Kokkinos. 2019. Holopose: Holistic 3d human reconstruction in-the-wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10884–10894.
- [12] Riza Alp Güler, Natalia Neverova, and Iasonas Kokkinos. 2018. Densepose: Dense human pose estimation in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7297–7306.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [14] Yinghao Huang, Federica Bogo, Christoph Lassner, Angjoo Kanazawa, Peter V Gehler, Javier Romero, Ijaz Akhter, and Michael J Black. 2017. Towards accurate marker-less human shape and pose estimation over time. In *2017 international conference on 3D vision (3DV)*. IEEE, 421–430.
- [15] Catalin Ionescu, Dragos Papaya, Vlad Olaru, and Cristian Sminchisescu. 2013. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE transactions on pattern analysis and machine intelligence* 36, 7 (2013), 1325–1339.
- [16] Sam Johnson and Mark Everingham. 2010. Clustered Pose and Nonlinear Appearance Models for Human Pose Estimation.. In *bmvc*, Vol. 2. Citeseer, 5.
- [17] Sam Johnson and Mark Everingham. 2011. Learning effective human pose estimation from inaccurate annotation. In *CVPR 2011*. IEEE, 1465–1472.
- [18] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. 2018. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 7122–7131.
- [19] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [20] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. 2019. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 2252–2261.
- [21] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. 2019. Convolutional mesh regression for single-image human shape reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4501–4510.
- [22] Christoph Lassner, Javier Romero, Martin Kiefel, Federica Bogo, Michael J Black, and Peter V Gehler. 2017. Unite the people: Closing the loop between 3d and 2d human representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 6050–6059.
- [23] Kevin Lin, Lijuan Wang, and Zicheng Liu. 2020. End-to-End Human Pose and Mesh Reconstruction with Transformers. *arXiv preprint arXiv:2012.09760* (2020).
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [25] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. 2015. SMPL: A skinned multi-person linear model. *ACM transactions on graphics (TOG)* 34, 6 (2015), 1–16.
- [26] Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).
- [27] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. 2019. AMASS: Archive of motion capture as surface shapes. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5442–5451.
- [28] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. 2017. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *2017 international conference on 3D vision (3DV)*. IEEE, 506–516.
- [29] Gyeongsik Moon and Kyoung Mu Lee. 2020. I2L-MeshNet: Image-to-lixel prediction network for accurate 3D human pose and mesh estimation from a single RGB image. *arXiv preprint arXiv:2008.03713* (2020).
- [30] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. 2018. Neural body fitting: Unifying deep learning and model based human pose and shape estimation. In *2018 international conference on 3D vision (3DV)*. IEEE, 484–494.
- [31] Georgios Pavlakos, Nikos Kolotouros, and Kostas Daniilidis. 2019. Texturepose: Supervising human mesh estimation with texture consistency. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 803–812.
- [32] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. 2018. Learning to estimate 3D human pose and shape from a single color image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 459–468.
- [33] Leonid Pishchulin, Eldar Insafutdinov, Siyu Tang, Bjoern Andres, Mykhaylo Andriluka, Peter V Gehler, and Bernt Schiele. 2016. Deepcut: Joint subset partition and labeling for multi person pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4929–4937.
- [34] Yu Rong, Ziwei Liu, Cheng Li, Kaidi Cao, and Chen Change Loy. 2019. Delving deep into hybrid annotations for 3d human recovery in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 5340–5348.
- [35] Akash Sengupta, Ignas Budvytis, and Roberto Cipolla. 2020. Synthetic Training for Accurate 3D Human Pose and Shape Estimation in the Wild. *arXiv preprint arXiv:2009.10013* (2020).
- [36] Leonid Sigal, Alexandru Balan, and Michael Black. 2007. Combined discriminative and generative articulated pose and non-rigid shape estimation. *Advances in neural information processing systems* 20 (2007), 1337–1344.
- [37] Jie Song, Xu Chen, and Otmar Hilliges. 2020. Human Body Model Fitting by Learned Gradient Descent. *arXiv preprint arXiv:2008.08474* (2020).
- [38] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. 2019. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5693–5703.
- [39] Gul Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. 2018. Bodynet: Volumetric inference of 3d human body shapes. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 20–36.
- [40] Gul Varol, Javier Romero, Xavier Martin, Naureen Mahmood, Michael J Black, Ivan Laptev, and Cordelia Schmid. 2017. Learning from synthetic humans. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 109–117.
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762* (2017).
- [42] Timo von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. 2018. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 601–617.
- [43] Yuanlu Xu, Song-Chun Zhu, and Tony Tung. 2019. Denserac: Joint 3d pose and shape estimation by dense render-and-compare. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 7760–7770.
- [44] Lu Yang, Qing Song, Zhihui Wang, Mengjie Hu, Chun Liu, Xueshi Xin, Wenhe Jia, and Songcen Xu. 2020. Renovating parsing R-CNN for accurate multiple human parsing. In *European Conference on Computer Vision*. Springer, 421–437.
- [45] Lu Yang, Qing Song, Zhihui Wang, and Ming Jiang. 2019. Parsing r-cnn for instance-level human analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 364–373.
- [46] Pengfei Yao, Zheng Fang, Fan Wu, Yao Feng, and Jiwei Li. 2019. Densebody: Directly regressing dense 3d human pose and shape from a single color image. *arXiv preprint arXiv:1903.10153* (2019).
- [47] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. 2015. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365* (2015).
- [48] Andrei Zanfir, Elisabeta Marinou, and Cristian Sminchisescu. 2018. Monocular 3d pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2148–2157.
- [49] Hongwen Zhang, Jie Cao, Guo Lu, Wanli Ouyang, and Zhenan Sun. 2019. Danet: Decompose-and-aggregate network for 3d human shape and pose estimation. In *Proceedings of the 27th ACM International Conference on Multimedia*. 935–944.
- [50] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip HS Torr, et al. 2020. Re-thinking Semantic Segmentation from a Sequence-to-Sequence Perspective with Transformers. *arXiv preprint arXiv:2012.15840* (2020).
- [51] Xiangyu Zhu, Zhen Lei, Xiaoming Liu, Hailin Shi, and Stan Z Li. 2016. Face alignment across large poses: A 3d solution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 146–155.