

# Learning Local Recurrent Models for Human Mesh Recovery

Runze Li<sup>1,2</sup>, Srikrishna Karanam<sup>1</sup>, Ren Li<sup>1</sup>, Terrence Chen<sup>1</sup>, Bir Bhanu<sup>2</sup>, and Ziyang Wu<sup>1</sup>

<sup>1</sup>United Imaging Intelligence, Cambridge MA, USA

<sup>2</sup>University of California Riverside, Riverside CA, USA

{first.last}@uii-ai.com



Figure 1: We present **LMR**, a new method for video human mesh recovery. Unlike existing work, LMR captures local human part dynamics and interdependencies by learning multiple local recurrent models, resulting in notable performance improvement over the state of the art. Here, we show a few qualitative results on the 3DPW dataset.

## Abstract

We consider the problem of estimating frame-level full human body meshes given a video of a person with natural motion dynamics. While much progress in this field has been in single image-based mesh estimation, there has been a recent uptick in efforts to infer mesh dynamics from video given its role in alleviating issues such as depth ambiguity and occlusions. However, a key limitation of existing work is the assumption that all the observed motion dynamics can be modeled using one dynamical/recurrent model. While this may work well in cases with relatively simplistic dynamics, inference with in-the-wild videos presents many challenges. In particular, it is typically the case that different body parts of a person undergo different dynamics in the video, e.g., legs may move in a way that may be dynamically different from hands (e.g., a person dancing). To address these issues, we present a new method for video mesh recovery that divides the human mesh into several local parts following the standard skeletal model. We then model the dynamics of each local part with separate recurrent models, with each model conditioned appropriately based on

the known kinematic structure of the human body. This results in a structure-informed local recurrent learning architecture that can be trained in an end-to-end fashion with available annotations. We conduct a variety of experiments on standard video mesh recovery benchmark datasets such as Human3.6M, MPI-INF-3DHP, and 3DPW, demonstrating the efficacy of our design of modeling local dynamics as well as establishing state-of-the-art results based on standard evaluation metrics.

## 1. Introduction

We consider the problem of human mesh recovery in videos, i.e., fitting a parametric 3D human mesh model to each frame of the video. With many practical applications [2, 3], including in healthcare for COVID-19 [4–6], there has been much progress in this field in the last few years [1, 7, 8]. In particular, most research effort has been expended in single image-based mesh estimation where one

\*This work was done during the internships of Runze Li and Ren Li with United Imaging Intelligence. Corresponding author: Srikrishna Karanam.

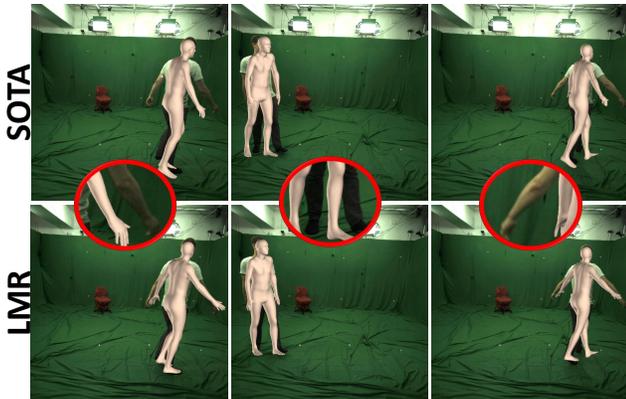


Figure 2: A qualitative comparison with VIBE [1], highlighting local regions (ellipses that show zoomed-in VIBE results) where LMR gives better performance.

seeks to fit the human mesh model to a single image. However, such 3D model estimation from only a single 2D projection (image) is a severely under-constrained problem since multiple 3D configurations (in this case poses and shapes of the mesh model) can project to the same image. Such ambiguities can be addressed by utilizing an extra dimension that is typically associated with images- the temporal dimension leading to video data and the problem of video mesh recovery.

The currently dominant paradigm for video mesh recovery involves the *feature-temporal-regressor* architecture. A deep convolutional neural network (CNN) is used to extract frame-level image feature vectors, which are then processed by a temporal encoder to learn the motion dynamics in the video. The representation from the temporal encoder is then processed by a parameter regressor module that outputs frame-level mesh parameter vectors. While methods vary in the specific implementation details, they mostly follow this pipeline. For instance, while Kanazawa *et al.* [9] implement the temporal encoder using a feed-forward fully convolutional model, Kocabas *et al.* [1] uses a recurrent model to encode motion dynamics. However, uniformly across all these methods, the parameter regressor is implemented using a “flat” regression architecture that takes in feature vectors as input and directly regresses all the model parameters, e.g., 85 values (pose, shape, and camera) for the popularly used skinned multi-person linear (SMPL) model [7, 10]. While this paradigm has produced impressive recent results as evidenced by the mean per-joint position errors on standard datasets (see Arnab *et al.* [11] and Kocabas *et al.* [1] for a fairly recent benchmark), a number of issues remain unaddressed that provide us with direction and scope for further research and performance improvement.

First, the above architectures implicitly assume that all motion dynamics can be captured using a single dynamical

system (e.g., a recurrent network). While this assumption may be reasonable for fairly simplistic human motions, it is not sufficient for more complex actions. For instance, while dancing, the motion dynamics of a person vary from one part of the body to the other. As a concrete example, the legs may remain static while the hands move vigorously, and these roles may be reversed after a certain period of time (static hands and moving legs several frames later), leading to more “locally” varying dynamics. Intuitively, this tells us that the motion of each local body part should in itself be modeled separately by a dynamical system, and that such a design should help capture this local “part-level” dynamical information more precisely as opposed to a single dynamical system for the entire video snippet.

Next, as noted above, the *regressor* in the *feature-temporal-regressor* architecture involves computing all the parameters of the SMPL model using a direct/flat regression design without due consideration given to the interdependent nature of these parameters (i.e., SMPL joint rotations are not independent but rather conditioned on other joints of other parts such as the root [10]). It has been noted in prior work [12] that such direct regression of rotation matrices, which form a predominant part of the SMPL parameter set, is challenging as is and only made further difficult due to these interdependencies in the SMPL model. In addition to direct rotation regression, the temporal module in the above *feature-temporal-regressor* also does not consider any joint and part interdependencies, i.e., modeling all motion dynamics using a single global dynamical system, thus only further exacerbating this problem.

To address the aforementioned issues, we present a new architecture for capturing the human motion dynamics for estimating a parametric mesh model in videos. Please note that while we use the SMPL model [10] in this work, our method can be extensible to other kinds of hierarchical parametric human meshes as well. See Figure 1 for some qualitative results with our method on the 3DPW [13] dataset and Figure 2 for a comparison with a current state-of-the-art method. Our method, called *local recurrent models for mesh recovery (LMR)*, comprises several design considerations. First, to capture the need for modeling locally varying dynamics as noted above, LMR defines six local recurrent models (root, head, left/right arms, left/right legs), one each to capture the dynamics of each part. As we will describe later, each “part” here refers to a chain of several joints defined on the SMPL model. Note that such a part division is not ad hoc but grounded in the hierarchical and part-based design of the SMPL model itself, which divides the human body into the six parts above following the standard skeletal rigging procedure [10]. Next, to model the conditional interdependence of local part dynamics, LMR first infers root part dynamics (i.e., parameters of all joints in the root part). LMR then uses these root part parameters to subsequently

infer the parameters of all other parts, with the output of each part conditioned on the root output. For instance, the recurrent model responsible for producing the parameters of the left leg takes as input both frame-level feature vectors as well as frame-level root-part parameters from the root-part recurrent model.

Note the substantial differences between LMR’s design and those of prior work- (a) we use multiple local recurrent models instead of one global recurrent model to capture motion dynamics, and (b) such local recurrent modeling enables LMR to explicitly capture local part dependencies. Modeling these local dependencies enables LMR to infer motion dynamics and frame-level video meshes informed by the geometry of the problem, i.e., the SMPL model, which, as noted in prior work [12], is an important design consideration as we take a step towards accurate rotation parameter regression architectures. We conduct extensive experiments on a number of standard video mesh recovery benchmark datasets (Human3.6M [14], MPI-INF-3DHP [15], and 3DPW [13]), demonstrating the efficacy of such local dynamic modeling as well as establishing state-of-the-art performance with respect to standard evaluation metrics.

To summarize, the key contributions of our work are:

- We present LMR, the first local-dynamical-modeling approach to video mesh recovery where unlike prior work, we explicitly model the local dynamics of each body part with separate recurrent networks.
- Unlike prior work that regresses mesh parameters in a direct or “flat” fashion, our local recurrent design enables LMR to explicitly consider human mesh interdependencies in parameter inference, thereby resulting in a structure-informed local recurrent architecture.
- We conduct extensive experiments on standard benchmark datasets and report competitive performance, establishing state-of-the-art results in many cases.

## 2. Related Work

There is much recent work in human pose estimation, including estimating 2D keypoints [16–18], 3D keypoints [19–23], and a full mesh [1, 7–9, 11, 24, 25]. Here, we discuss methods that are relevant to our specific problem-fitting 3D meshes to image and video data.

**Single-image mesh fitting.** Most recent progress in human mesh estimation has been in fitting parametric meshes to single image inputs. In particular, following the availability of differentiable parametric models such as SMPL [10], there has been an explosion in interest and activity in this field. Kanazawa *et al.* [7] presented an end-to-end trainable regression architecture for this problem that could

in principle be trained with 2D-only keypoint data. Subsequently, many improved models have been proposed. Kolotourous *et al.* [25] and Georgakis *et al.* [8] extended this architecture to include more SMPL-structure-informed design considerations using either graph-based or parameter factorization-based approaches. There have also been attempts at SMPL-agnostic modeling of joint interdependencies, with Fang *et al.* [26] employing bidirectional recurrent networks and Isack *et al.* [27] learning priors between joints using a pre-defined joint connectivity scheme. While methods such as Georgakis *et al.* [8] and Zhou *et al.* [28] also take a local part-based kinematic approach, their focus is on capturing inter-joint spatial dependencies. On the other hand, LMR’s focus is on capturing inter-part temporal dependencies which LMR models using separate recurrent networks.

**Video mesh fitting.** Following the success of image-based mesh fitting methods, there has been a recent uptick in interest and published work in fitting human meshes to videos. Arnab *et al.* [11] presented a two-step approach that involved generating 2D keypoints and initial mesh fits using existing methods, and then using these initial estimates to further refine the results using temporal consistency constraints, e.g., temporal smoothness and 3D priors. However, such a two-step approach is susceptible to errors in either steps and our proposed LMR overcomes this issue with an end-to-end trainable method that provides deeper integration of the temporal data dimension both in training and inference. On the other hand, Kanazawa *et al.* [9] and Kocabas *et al.* [1] also presented end-to-end variants of the *feature-temporal-regressor* where frame-level feature vectors are first encoded using a temporal encoder (e.g., a single recurrent network) and finally processed by a parameter regressor to generate meshes. However, such a global approach to modeling motion dynamics (with only one RNN) does not capture the disparities in locally varying dynamics (e.g., hands vs. legs) which is typically the case in natural human motion. LMR addresses this issue by design with multiple local RNNs in its architecture, one for each pre-defined part of the human body. Such a design also makes mesh parameter regression more amenable by grounding this task in the geometry of the problem, i.e., the SMPL model itself.

## 3. Technical Approach

### 3.1. Parametric Mesh Representation

We use the Skinned Multi-Person Linear (SMPL) model [10] to parameterize the human body. SMPL uses two sets of parameter vectors to capture variations in the human body: shape and pose. The shape of the human body is represented using a 10-dimensional vector  $\beta \in \mathbb{R}^{10}$  whereas the pose of the body is represented using a 72-

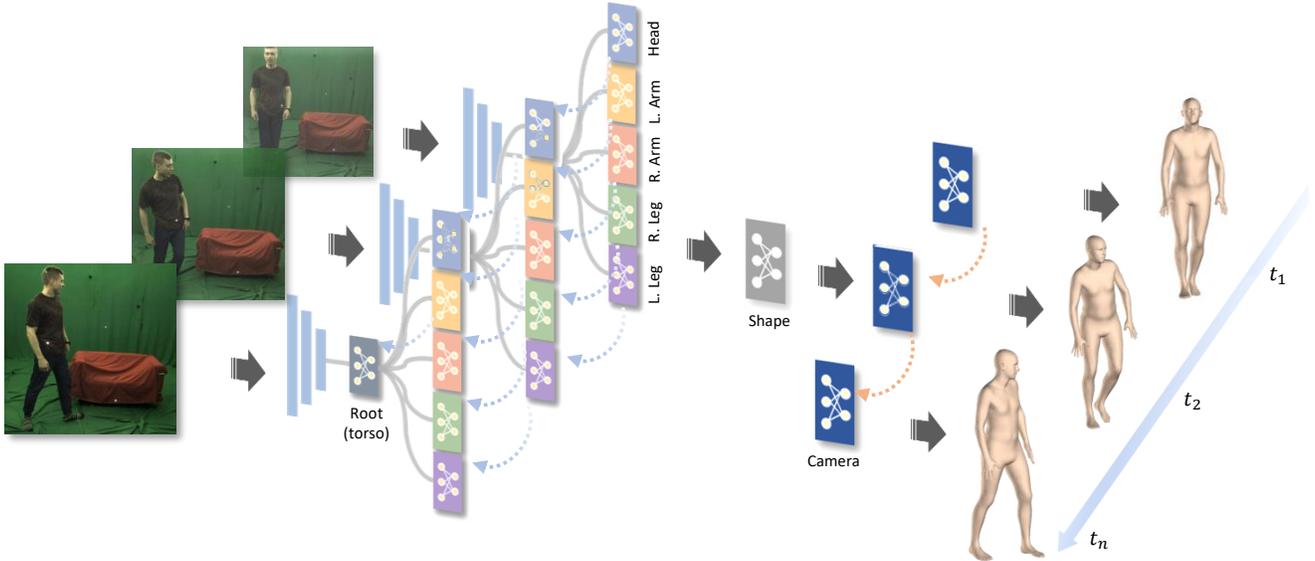


Figure 3: The proposed local recurrent modeling approach to human mesh recovery.

dimensional vector  $\theta \in \mathbb{R}^{72}$ . While  $\beta$  corresponds to the first ten dimensions of the PCA projection of a shape space,  $\theta$  captures, in axis-angle format [29], the global rotation of the root joint (3 values) and relative (to the root) rotations of 23 other body joints (69 values). Given  $\beta$ ,  $\theta$ , and a learned model parameter set  $\psi$ , SMPL defines the mapping  $M(\beta, \theta, \psi) : \mathbb{R}^{82} \rightarrow \mathbb{R}^{3 \times N}$  from the 82-dimensional parametric space to a vertex space of  $N = 6890$  3D mesh vertices. One can then infer the 24 3D joints of interest (e.g., hips, legs, etc.)  $\mathbf{X} \in \mathbb{R}^{3 \times K}$ ,  $K = 24$  using a pre-learned joint regression matrix  $\mathbf{W}$  as  $\mathbf{X} = \mathbf{W}\mathbf{J}$ . Using a known camera model, e.g., a weak-perspective model as in prior work [7], one can then obtain the corresponding 24 2D image points  $\mathbf{x} \in \mathbb{R}^{2 \times K}$  as:

$$\mathbf{x} = s\Pi(\mathbf{X}(\beta, \theta)) + \mathbf{t}, \quad (1)$$

where the scale  $s \in \mathbb{R}$  and translation  $\mathbf{t} \in \mathbb{R}^2$  represent the camera model, and  $\Pi$  is an orthographic projection. Therefore, fitting 3D SMPL mesh to a single image involves estimating the parameter set  $\Theta = \{\beta, \theta, s, \mathbf{t}\}$ . In video mesh recovery, we take this a step forward by estimating  $\Theta$  for every frame in the video.

### 3.2. Learning Local Recurrent Models

As noted in Section 1, existing video mesh fitting methods formulate the problem in the *feature-temporal-regressor* design where all motion dynamics in the video are captured using a single RNN. We argue that this is insufficient for mesh estimation due to the inherently complex nature of human actions/motion, more so in challenging in-the-wild scenarios. Our key insight is that natural human motion dynamics has a more locally varying characteristic that can

more precisely be captured using locally learned recurrent networks. We then translate this idea into a conditional local recurrent architecture, called **LMR** and visually summarized in Figure 3, where we define multiple recurrent models, one each to capture the dynamics of the corresponding local region in the human body. During training and inference, LMR takes as input a segment of an input video  $\mathbf{V} = \{\mathbf{I}_1, \mathbf{I}_2, \dots, \mathbf{I}_t, t = 1, 2, \dots, T\}$ , where  $T$  is a design parameter corresponding to the length of the input sequence. LMR first processes each frame with its feature extraction module to produce frame-level feature vectors  $\Phi = \{\phi_1, \phi_2, \dots, \phi_t\}$  for each of the  $T$  frames. LMR then processes  $\Phi$  with its local part-level recurrent models and associated parameter regressors, and aggregates all part-level outputs to obtain the mesh and camera parameters  $\Theta_t, t = 1, 2, \dots, T$  for each frame, finally producing the output video mesh.

#### 3.2.1 LMR Architecture

As shown in Figure 3(a), our architecture comprises a feature extractor followed by our proposed LMR module. The LMR module is responsible for processing the frame-level representation  $\Phi$  to output the per-frame parameter vectors  $\Theta_t$ . Following the design of the SMPL model and prior work [8, 10], we divide the human body into six local parts- *root* (4 joints in the root region), *head* (2 joints in the head region), *left arm* (5 joints on left arm), *right arm* (5 joints on right arm), *left leg* (4 joints on left leg), and *right leg* (4 joints on right leg). Given this division, the pose of local part  $p_i, i = 1, \dots, 6$  can be expressed as  $\theta^i = [\mathbf{r}_1, \dots, \mathbf{r}_{n_i}], i = 1, \dots, 6$ , where  $\mathbf{r}_q (q = 1, \dots, n_i)$

is a rotation parameterization (e.g.,  $r_q \in \mathbb{R}^3$  in case of axis angle) of joint  $q$  and  $n_i$  is the number of joints defined in part  $i$ . The overall pose parameter vector  $\theta$  can then be aggregated as  $\theta = [\theta^1, \dots, \theta^6]$ .

To capture locally varying dynamics across the video sequence, LMR defines one recurrent model for each of the six parts defined above (see Figure 3(b)). The recurrent model for part  $i$  is responsible for predicting its corresponding  $\theta^i$ . To capture the conditional dependence between parts, the information propagation during training and inference is defined as follows. Given the frame-level feature representation  $\Phi$ , the mean pose vector  $\theta_{\text{mean}}$ , and the mean shape vector  $\beta_{\text{mean}}$  (note that it is common [1, 7, 9] to initialize mesh fitting with these mean values), the recurrent model responsible for the root part (number 1) first predicts its corresponding pose vector  $\theta_t^1, t = 1, \dots, T$  for each of the  $t$  frames using the concatenated vector  $[\Phi_t, \theta_{\text{mean}}^1, \beta_{\text{mean}}]$  as input for the current frame  $t$ . Note that  $\Phi_t$  is the feature vector for frame  $t$  and  $\theta_{\text{mean}}^1$  represents the mean pose parameters of part  $p_1$ . All other recurrent models (parts 2 through 6) then take in as input the concatenated vector  $[\Phi_t, \theta_{\text{mean}}^k, \beta_{\text{mean}}, \theta_t^1]$  in predicting their corresponding pose vectors  $\theta_t^k, k = 2, \dots, 6$  and  $t = 1, \dots, T$ , where  $\theta_{\text{mean}}^k$  represents the mean pose parameters of part  $p_k$ . Note this explicit dependence of part  $k$  on the root (part 1) prediction  $\theta^1$ . Given the aggregated (over all 6 parts) pose vector  $\theta_t$ , LMR has a fully-connected module that takes as input the concatenated vector  $[\Phi_t, \theta_t, \beta_{\text{mean}}]$  for each frame  $t$  to predict the per-frame shape vectors  $\beta_t, t = 1, \dots, T$ . Finally, given an initialization for the camera model  $c_{\text{init}} = [s_{\text{init}}, t_{\text{init}}]$ , LMR uses the concatenated vector  $[\Phi_t, \theta_t, \beta_t, c_{\text{init}}]$  as part of its camera recurrent model to predict the camera model  $c_t, t = 1, \dots, T$  for each frame. Note that while we have simplified the discussion and notation here for clarity of exposition, LMR actually processes each batch of input in an iterative fashion, which we next describe in more mathematical detail.

### 3.2.2 Training an LMR model

As noted above and in Figure 3, the proposed LMR module takes as input the video feature set  $\Phi$  and the mean pose and shape parameters  $\theta_{\text{mean}}$  and  $\beta_{\text{mean}}$  and produces the set of parameter vectors  $\Theta_t = [\theta_t, \beta_t, c_t]$  for each frame  $t$ . The LMR block processes each input set in an iterative fashion, with the output after each iteration being used as a new initialization point to further refine the result. The final output  $\Theta_t$  is then obtained at the end of  $L$  such iterations. Here, we provide further details of this training strategy.

Let each iteration step above be denoted by the letter  $v$ . At step  $v = 0$ , the initial pose and shape values for frame  $t$  will then be  $\theta_{t,v} = \theta_{\text{mean}}$  and  $\beta_{t,v} = \beta_{\text{mean}}$ . The  $t, v$  notation refers to the  $v^{\text{th}}$  iterative step of LMR for frame

number  $t$ . So, given  $\Phi, \beta_{t,v}$ , and the root pose  $\theta_{t,v}^1$  (recall root is part number 1 from above), the input to the root RNN will be the set of  $t$  vectors  $[\Phi_t, \theta_{t,v}^1, \beta_{t,v}]$  for each of the  $t$  frames. The root RNN then estimates an intermediate residual pose  $\Delta\theta_{t,v}^1$ , which is added to the input  $\theta_{t,v}^1$  to give the root RNN output  $\theta_{t,v}^1 = \theta_{t,v}^1 + \Delta\theta_{t,v}^1$ .

Given the root prediction  $\theta_{t,v}^1$  at iteration  $v$ , each of the other dependent part RNNs then use this information to produce their corresponding pose outputs. Specifically, for part RNN  $k$ , the input vector set (across the  $t$  frames) will be  $[\Phi_t, \theta_{t,v}^k, \beta_{t,v}, \theta_{t,v}^1]$  for  $k = 2, \dots, 6$ . Each part RNN first gives its corresponding intermediate residual pose  $\Delta\theta_{t,v}^k$ . This is then added to its corresponding input part pose, giving the outputs  $\theta_{t,v}^k = \theta_{t,v}^k + \Delta\theta_{t,v}^k$  for  $k = 2, \dots, 6$ .

After producing all the updated pose values at iteration  $v = 0$ , LMR then updates the shape values. Recall that the shape initialization used at  $v = 0$  is  $\beta_{t,v} = \beta_{\text{mean}}$ . Given  $\Phi$ , the updated and aggregated pose vector set  $\theta_{t,v} = [\theta_{t,v}^1, \dots, \theta_{t,v}^6]$ , and the shape vector set  $\beta_{\text{mean}}$ , LMR then uses the input vector set  $[\Phi_t, \theta_{t,v}, \beta_{\text{mean}}]$  as part of the shape update module to produce the new shape vector set  $\beta_{t,v}$  for each frame  $t$  during the iteration  $v$ .

Given these updated  $\theta_{t,v}$  and  $\beta_{t,v}$ , LMR then updates the camera model parameters (used for image projection) with a camera model RNN. We use an RNN to model the camera dynamics to cover scenarios where the camera might be moving, although a non-dynamical fully-connected neural network can also be used in cases where the camera is known to be static. Given an initialization for the camera model  $c_{t,v} = c_{\text{init}}$  at iteration  $v = 0$ , the camera RNN processes the input vector set  $[\Phi_t, \theta_{t,v}, \beta_{t,v}, c_{\text{init}}]$  to produce the new camera model set  $c_{t,v}$  for each frame  $t$ .

After going through one round of pose update, shape update, and camera update as noted above, LMR then re-initializes this prediction process with the updated pose and shape vectors from the previous iteration. Specifically, given the updated  $\theta_{t,v}$  and  $\beta_{t,v}$  at the end of iteration  $v = 0$ , the root RNN at iteration  $v = 1$  then takes as input the set  $[\Phi_t, \theta_{t,v}^1, \beta_{t,v}]$ , where the pose and shape values are not the mean vectors (as in iteration  $v = 0$ ) but the updated vectors from iteration  $v = 0$ . LMR repeats this process for a total of  $V$  iterations, finally producing the parameter set  $\Theta_t = [\theta_t, \beta_t, c_t]$  for each frame  $t$ . Note that this iterative strategy is similar in spirit to the iterative error feedback strategies commonly used in pose estimators [7, 30–32].

All the predictions above are supervised using several cost functions. First, if ground-truth SMPL model parameters  $\Theta_t^{gt}$  are available, we enforce a Euclidean loss between the predicted and the ground-truth set:

$$L_{\text{smp}} = \frac{1}{T} \sum_{t=1}^T \|\Theta_t^{gt} - \Theta_t\|_2 \quad (2)$$

where the summation is over the  $t = T$  input frames in the

current batch of data.

Next, if ground-truth 3D joints  $\mathbf{X}_t^{gt} \in \mathbb{R}^{3 \times K}$  (recall  $K=24$  from Section 3.1) are available, we enforce a mean per-joint L1 loss between the prediction 3D joints  $\mathbf{X}_t \in \mathbb{R}^{3 \times K}$  and  $\mathbf{X}_t^{gt}$ . To compute  $\mathbf{X}_t$ , we use the predicted parameter set  $\Theta_t$  and the SMPL vertex mapping function  $M(\beta, \theta, \psi) : \mathbb{R}^{82} \rightarrow \mathbb{R}^{3 \times N}$  and the joint regression matrix  $W$  (see Section 3.1). The loss then is:

$$L_{3D} = \frac{1}{T} \frac{1}{K} \sum_{t=1}^T \sum_{k=1}^K \|\mathbf{X}_{k,t}^{gt} - \mathbf{X}_{k,t}\|_1 \quad (3)$$

where each column of  $\mathbf{X}_{k,t}^{gt} \in \mathbb{R}^3$  and  $\mathbf{X}_{k,t} \in \mathbb{R}^3$  is one of  $K$  joints in three dimensions and the outer summation is over  $t = T$  frames as above.

Finally, to provide supervision for camera prediction, we also enforce a mean per-joint L1 loss between the prediction 2D joints  $\mathbf{x}_t \in \mathbb{R}^{2 \times K}$  and the ground-truth 2D joints  $\mathbf{x}_t^{gt}$ . To compute  $\mathbf{x}_t$ , we use the 3D joints prediction  $\mathbf{X}_t$  and the camera prediction  $\mathbf{c}_t$  to perform an orthographic projection following Equation 1. The loss then is:

$$L_{2D} = \frac{1}{T} \frac{1}{K} \sum_{t=1}^T \sum_{k=1}^K \|\mathbf{x}_{k,t}^{gt} - \mathbf{x}_{k,t}\|_1 \quad (4)$$

where each column  $\mathbf{x}_{k,t}^{gt} \in \mathbb{R}^2$  and  $\mathbf{x}_{k,t} \in \mathbb{R}^2$  of  $\mathbf{x}_t^{gt}$  and  $\mathbf{x}_t$  respectively is one of  $K$  joints on the image and the outer summation is over  $t = T$  frames as above.

The overall LMR training objective then is:

$$L_{LMR} = w_{\text{smp}} L_{\text{smp}} + w_{3D} L_{3D} + w_{2D} L_{2D} \quad (5)$$

where  $w_{\text{smp}}$ ,  $w_{3D}$ , and  $w_{2D}$  are the corresponding loss weights.

## 4. Experiments and Results

### 4.1. Datasets and Evaluation

Following Kocabas *et al.* [1], we use a mixture of both datasets with both 2D (e.g., keypoints) as well as 3D (e.g., mesh parameters) annotations. For 2D datasets, we use PennAction [36], PoseTrack [37], and InstaVariety [9], whereas for 3D datasets, we use Human3.6M [14], MPI-INF-3DHP [15], and 3DPW [13]. In all our experiments, we use exactly the same settings as Kocabas *et al.* [1] for a fair benchmarking of the results. To report quantitative performance, we use evaluation metrics that are now standard in the human mesh research community. On all the test datasets, we report both mean-per-joint position error (MPJPE) as well as Procrustes-aligned mean-per-joint position error (PAMPJPE). Additionally, following Kanazawa *et al.* [9] and Kocabas *et al.* [1], on the 3DPW test set, we also report the acceleration error (“Accel.”), which is the average (across

all keypoints) difference between the ground truth and predicted acceleration of keypoints, and the per-vertex error (PVE).

### 4.2. Ablation Results

We first present results of an ablation experiment conducted to study the efficacy of the proposed design of LMR, i.e., the use of multiple local recurrent models as opposed to a single recurrent model as is done in prior work [1]. Here, we follow the same pipeline as Figure 3 in spirit, with the only difference being the use of only one RNN to infer all the pose parameters  $\theta$  instead of the six RNNs depicted in Figure 3(b). All other design choices, e.g., for the shape model or the camera model, remain the same as LMR. We show qualitative results of this experiment in Figure 4 and quantitative results in Table 1. In Figure 4, we show two frames from two different video sequences in (a) and (b). The first row shows results with this single RNN baseline and the second row shows corresponding results with our full model, i.e., LMR. One can note that LMR results in better mesh fits, with more accurate  $\Theta$ -inference in regions such as hands and legs. We further substantiate this performance gap quantitatively in Table 1, where one can note the proposed LMR gives consistently better performance than its baseline single RNN counterpart across all datasets as well as evaluation metrics.

### 4.3. Comparison with the state-of-the-art results

We compare the performance of LMR with a wide variety of state-of-the-art image-based and video-based methods. We first begin with a discussion on relative qualitative performance. In Figure 5, we show three frames from two different video sequences in (a) and (b) comparing the performance of the image-based HMR method [7] (first row) and our proposed LMR. Since LMR is a video-based method, one would expect substantially better performance, including in cases where there are self-occlusions. From Figure 5, one can note this is indeed the case. In the first column of Figure 5, HMR is unable to infer the correct head pose (it infers front facing when the person is actually back facing), whereas LMR is able to use the video information from prior to this frame to infer the head pose correctly. Note also HMR’s incorrect inference in other local regions, e.g., legs, in the subsequent frames in Figure 5(a). This aspect of self-occlusions (i.e., invisible face keypoints) is further demonstrated in Figure 5(b), where HMR is unstable (front facing on a few and back facing on a few frames), whereas LMR consistently infers the correct pose.

Next, we compare the performance of LMR with the state-of-the-art video-based VIBE method [1]. In Figure 6, we show three frames from two different video sequences in (a) and (b). One can note substantial performance improvement in several local regions from these results. In

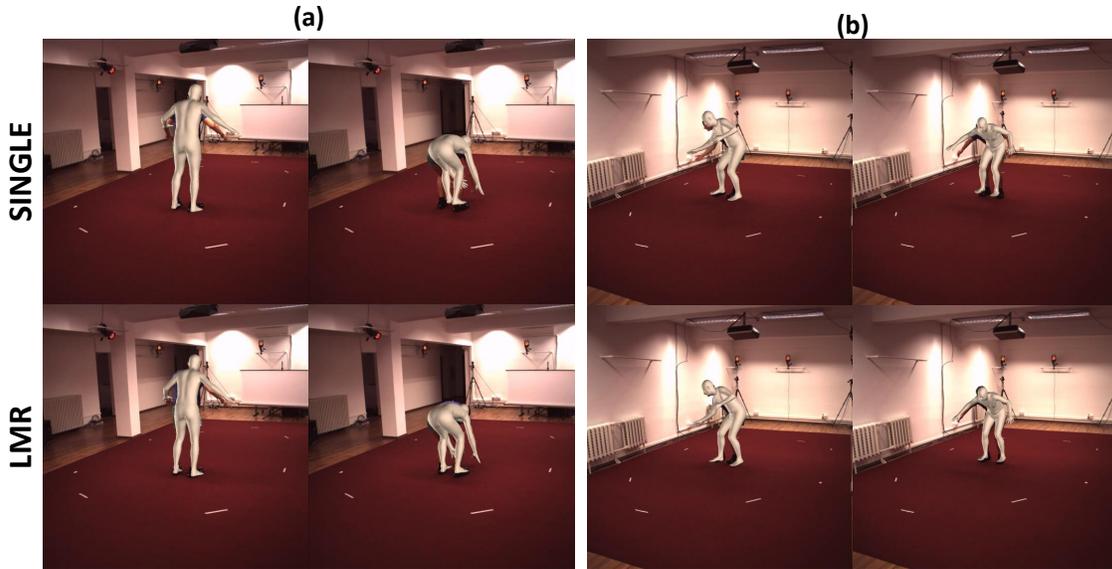


Figure 4: Two sets of qualitative results comparing LMR with a single-RNN baseline model.

Methods	Human3.6M		MPI-INF-3DHP		3DPW			
	MPJPE↓	Rec. Error↓	MPJPE↓	Rec. Error↓	MPJPE↓	Rec. Error↓	PVE↓	Accel↓
Single RNN	69.2	45.6	100.0	66.7	87.7	55.3	101.0	19.0
LMR no root dependencies	66.7	43.5	97.1	64	86.3	55.1	98.9	17.6
<b>LMR</b>	<b>61.9</b>	<b>42.5</b>	<b>94.6</b>	<b>62.4</b>	<b>81.7</b>	<b>51.2</b>	<b>93.6</b>	<b>15.6</b>

Table 1: Results of an ablation study comparing LMR with a single RNN baseline.

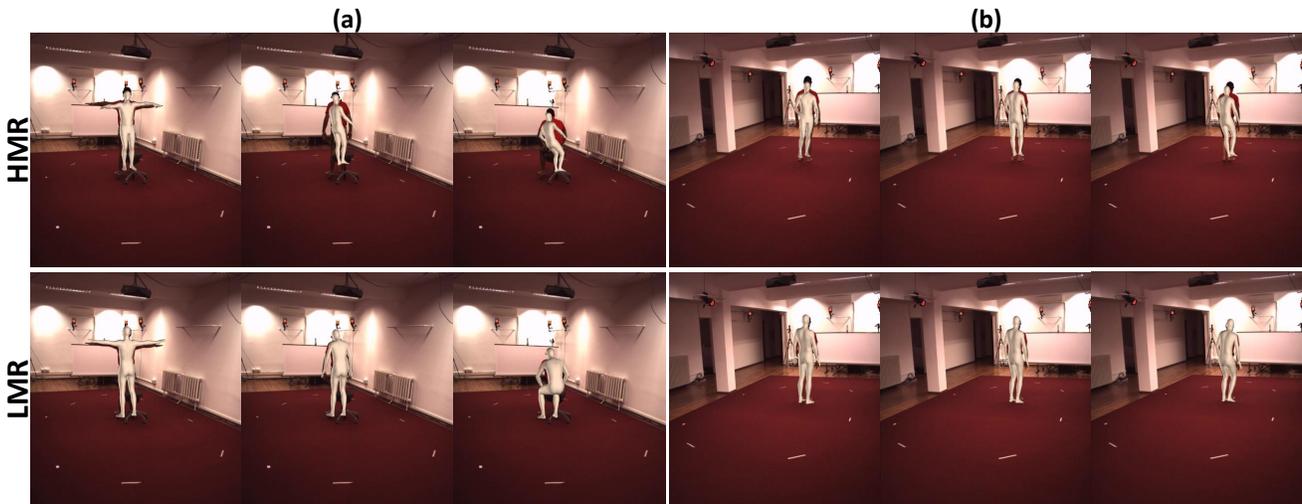


Figure 5: Two sets of qualitative results comparing the performance of LMR with the image-based HMR [7] method.

particular, LMR infers more accurate hand pose and camera model parameters in Figure 6(a) when compared to VIBE. The results in Figure 6(b), a more challenging scenario, best illustrates the benefits offered by proposed local design of LMR. Given the variety of body movements in this set of

frames, one can note the improved performance of LMR in several regions- hands and legs in the first column, head in the second column, and hands and legs again in the third column. These results are further substantiated in the quantitative comparison we discuss next.

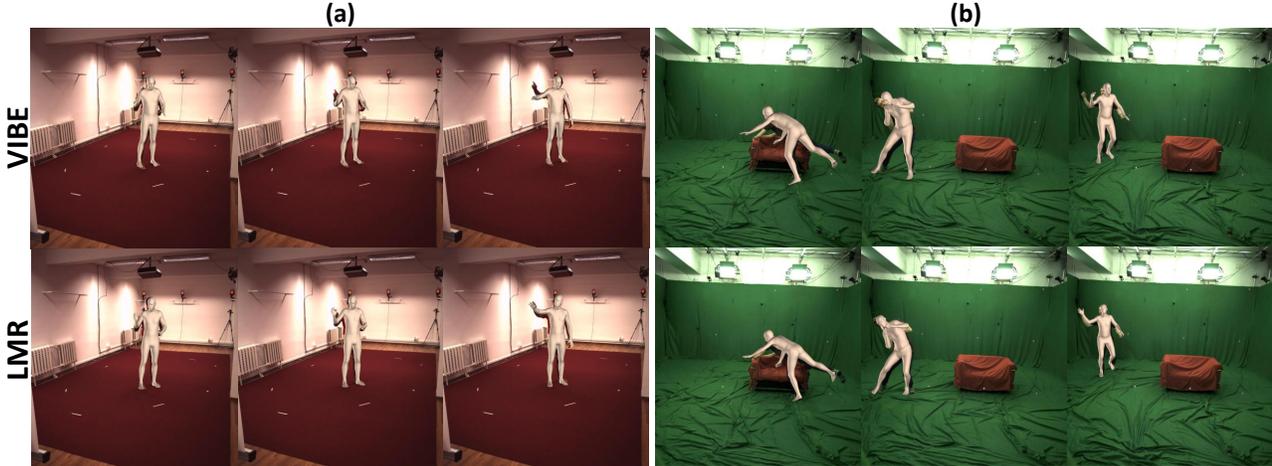


Figure 6: Two sets of qualitative results comparing the performance of LMR with the video-based VIBE [1] method.

Methods		Human3.6M		MPI-INF-3DHP		3DPW			
		MPJPE ↓	Rec. Error ↓	MPJPE ↓	Rec. Error ↓	MPJPE ↓	Rec. Error ↓	PVE ↓	Accel ↓
Image-based	Kanazawa <i>et al.</i> [7]	88.0	56.8	124.2	89.8	130	76.7	-	37.4
	Omran <i>et al.</i> [33]	-	59.9	-	-	-	-	-	-
	Pavlakos <i>et al.</i> [24]	-	75.9	-	-	-	-	-	-
	Kolotouros <i>et al.</i> [25]	-	50.1	-	-	-	70.2	-	-
	Georgakis <i>et al.</i> [8]	67.7	50.1	-	-	-	-	-	-
Extra-fitting	Kolotouros <i>et al.</i> [34]	62.2	<b>41.1</b>	105.2	67.5	96.9	59.2	116.4	29.8
Video-based	Kanazawa <i>et al.</i> [9]	-	56.9	-	-	116.5	72.6	139.3	<b>15.2</b>
	Arnab <i>et al.</i> [11]	77.8	54.3	-	-	-	72.2	-	-
	Doersch <i>et al.</i> [35]	-	-	-	-	-	74.7	-	-
	Kocabas <i>et al.</i> [1]	65.6	41.4	96.6	64.6	82.9	51.9	99.1	23.4
	<b>LMR</b>	<b>61.9</b>	42.5	<b>94.6</b>	<b>62.4</b>	<b>81.7</b>	<b>51.2</b>	<b>93.6</b>	15.6

Table 2: Comparing LMR to the state of the art (“-”: unavailable result in the corresponding paper).

We provide a quantitative comparison of the performance of LMR to various state-of-the-art image- and video-based methods in Table 2. We make several observations. First, as expected, LMR gives substantially better performance when compared to the image-based method of Kanazawa *et al.* [7] (MPJPE of 61.9 mm for LMR vs. 88.0 mm for HMR on Human3.6M, 94.6 mm for LMR vs. 124.2 mm for HMR on MPI-INF-3DHP, and 81.7 mm for LMR vs. 130.0 mm for HMR on 3DPW). This holds with other image-based methods as well (first half of Table 2). Next, LMR gives competitive performance when compared to state-of-the-art video-based methods as well. In particular, further substantiating the discussion above, LMR generally outperforms Kocabas *et al.* [1] with margins that are higher on the “in-the-wild” datasets (MPJPE of 94.6 mm for LMR vs. 96.6 mm for Kocabas *et al.* [1] on MPI-INF-3DHP, Accel. of 15.6 mm/s<sup>2</sup> for LMR vs. 23.4 mm/s<sup>2</sup> for Kocabas *et al.* [1] on 3DPW), further highlighting the efficacy of LMR’s local dynamic modeling.

Finally, in Table 2, we also compare our results with

those of Kolotouros *et al.* [34] that uses an additional step of in-the-loop model fitting. Note that despite our proposed LMR **not** doing this extra model fitting, it outperforms Kolotouros *et al.* [34] in most cases, with particularly substantial performance improvements on MPI-INF-3DHP (MPJPE of 94.6 mm for LMR vs. 105.2 mm for Kolotouros *et al.* [34]) and 3DPW (MPJPE of 81.7 mm for LMR vs. 96.9 mm for Kolotouros *et al.* [34]).

## 5. Conclusions

We considered the problem of video human mesh recovery and noted that the currently dominant design paradigm of using a single dynamical system to model all motion dynamics, in conjunction with a “flat” parameter regressor is insufficient to tackle challenging in-the-wild scenarios. We presented an alternative design based on local recurrent modeling, resulting in a structure-informed learning architecture where the output of each local recurrent model (representing the corresponding body part) is appropriately conditioned based on the known human kinematic structure.

We presented results of an extensive set of experiments on various challenging benchmark datasets to demonstrate the efficacy of the proposed local recurrent modeling approach to video human mesh recovery.

## References

- [1] Muhammed Kocabas, Nikos Athanasiou, and Michael J. Black. VIBE: Video inference for human body pose and shape estimation. In *CVPR*, 2020. 1, 2, 3, 5, 6, 8
- [2] Vivek Singh, Kai Ma, Birgi Tamersoy, Yao-Jen Chang, Andreas Wimmer, Thomas O'Donnell, and Terrence Chen. DARWIN: Deformable patient avatar representation with deep image network. In *MICCAI*, 2017. 1
- [3] Angel Martínez-González, Michael Villamizar, Olivier Canévet, and Jean-Marc Odobez. Real-time convolutional networks for depth-based human pose estimation. In *IROS*, 2018. 1
- [4] Jianhai Li, Unni K Udayasankar, Thomas L Toth, John Seamans, William C Small, and Mannudeep K Kalra. Automatic patient centering for MDCT: effect on radiation dose. *American journal of roentgenology*, 188(2):547–552, 2007. 1
- [5] William Ching, John Robinson, and Mark McEntee. Patient-based radiographic exposure factor selection: a systematic review. *Journal of medical radiation sciences*, 61(3):176–190, 2014.
- [6] Srikrishna Karanam, Ren Li, Fan Yang, Wei Hu, Terrence Chen, and Ziyang Wu. Towards contactless patient positioning. *IEEE Transactions on Medical Imaging*, 2020. 1
- [7] Angjoo Kanazawa, Michael J. Black, David W. Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *CVPR*, 2018. 1, 2, 3, 4, 5, 6, 7, 8
- [8] Georgios Georgakis, Ren Li, Srikrishna Karanam, Terrence Chen, Jana Kosecka, and Ziyang Wu. Hierarchical kinematic human mesh recovery. In *ECCV*, 2020. 1, 3, 4, 8
- [9] Angjoo Kanazawa, Jason Y. Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. 2, 3, 5, 6, 8
- [10] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. SMPL: A skinned multi-person linear model. *ACM Transactions on Graphics*, 34(6):1–16, 2015. 2, 3, 4
- [11] Anurag Arnab, Carl Doersch, and Andrew Zisserman. Exploiting temporal context for 3d human pose estimation in the wild. In *CVPR*, 2019. 2, 3, 8
- [12] Alex Kendall and Roberto Cipolla. Geometric loss functions for camera pose regression with deep learning. In *CVPR*, 2017. 2, 3
- [13] Timo von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*, 2018. 2, 3, 6
- [14] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6M: Large scale datasets and predictive methods for 3d human sensing in natural environments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(7):1325–1339, 2013. 3, 6
- [15] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. Monocular 3d human pose estimation in the wild using improved cnn supervision. In *3DV*, 2017. 3, 6
- [16] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-glass networks for human pose estimation. In *ECCV*, 2016. 3
- [17] Z Cao, T Simon, SE Wei, YA Sheikh, et al. Openpose: Realtime multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019.
- [18] Feng Zhang, Xiatian Zhu, Hanbin Dai, Mao Ye, and Ce Zhu. Distribution-aware coordinate representation for human pose estimation. In *CVPR*, 2020. 3
- [19] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. A simple yet effective baseline for 3d human pose estimation. In *ICCV*, 2017. 3
- [20] Georgios Pavlakos, Xiaowei Zhou, and Kostas Daniilidis. Ordinal depth supervision for 3d human pose estimation. In *CVPR*, 2018.
- [21] Ikhsanul Habibie, Weipeng Xu, Dushyant Mehta, Gerard Pons-Moll, and Christian Theobalt. In the wild human pose estimation using explicit 2d features and intermediate 3d representations. In *CVPR*, 2019.
- [22] Umar Iqbal, Pavlo Molchanov, and Jan Kautz. Weakly-supervised 3d human pose learning via multi-view images in the wild. In *CVPR*, 2020.
- [23] Zhe Wang, Liyan Chen, Shauray Rathore, Daeyun Shin, and Charless Fowlkes. Geometric pose affordance: 3d human pose with scene constraints. In *Arxiv*, 2019. 3
- [24] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3d human pose and shape from a single color image. In *CVPR*, 2018. 3, 8
- [25] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. In *CVPR*, 2019. 3, 8
- [26] Haoshu Fang, Yuanlu Xu, Wenguan Wang, Xiaobai Liu, and Song-Chun Zhu. Learning pose grammar to encode human body configuration for 3d pose estimation. In *AAAI*, 2018. 3
- [27] Hossam Isack, Christian Haene, Cem Keskin, Sofien Bouaziz, Yuri Boykov, Shahram Izadi, and Sameh Khamis. Repose: Learning deep kinematic priors for fast human pose estimation. *arXiv preprint arXiv:2002.03933*, 2020. 3
- [28] Xingyi Zhou, Xiao Sun, Wei Zhang, Shuang Liang, and Yichen Wei. Deep kinematic pose regression. In *ECCV*, 2016. 3
- [29] Fabrizio Caccavale, Ciro Natale, Bruno Siciliano, and Luigi Villani. Six-dof impedance control based on angle/axis representations. *IEEE Transactions on Robotics and Automation*, 15(2):289–300, 1999. 4
- [30] Piotr Dollár, Peter Welinder, and Pietro Perona. Cascaded pose regression. In *CVPR*, 2010. 5
- [31] Markus Oberweger, Paul Wohlhart, and Vincent Lepetit. Training a feedback loop for hand pose estimation. In *ICCV*, 2015.
- [32] Joao Carreira, Pulkit Agrawal, Katerina Fragkiadaki, and Jitendra Malik. Human pose estimation with iterative error feedback. In *CVPR*, 2016. 5
- [33] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Pe-

ter V. Gehler, and Bernt Schiele. Neural body fitting: Unifying deep learning and model-based human pose and shape estimation. In *3DV*, 2018. 8

- [34] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3d human pose and shape via model-fitting in the loop. In *ICCV*, 2019. 8
- [35] Carl Doersch and Andrew Zisserman. Sim2real transfer learning for 3d human pose estimation: motion to the rescue. In *NeurIPS*, 2019. 8
- [36] Weiyu Zhang, Menglong Zhu, and Konstantinos G. Derpanis. From actemes to action: A strongly-supervised representation for detailed action understanding. In *ICCV*, December 2013. 6
- [37] M. Andriluka, U. Iqbal, E. Ensafutdinov, L. Pishchulin, A. Milan, J. Gall, and Schiele B. PoseTrack: A benchmark for human pose estimation and tracking. In *CVPR*, 2018. 6