

Beyond Static Features for Temporally Consistent 3D Human Pose and Shape from a Video

Hongsuk Choi¹

Gyeongsik Moon¹

Ju Yong Chang²

Kyoung Mu Lee¹

¹ECE & ASRI, Seoul National University, Korea

²ECE, Kwangwoon University, Korea

{redarknight, mks0601, kyoungmu}@snu.ac.kr, juyong.chang@gmail.com

input

TCMR (Ours)

VIBE

Figure 1: VIBE [15], the state-of-the-art video-based 3D human pose and shape estimation method, outputs very different 3D human poses per frame, although the frames have subtle differences. Our TCMR produces clearly more temporally consistent and smooth 3D human motion. *This is a video figure that is best viewed by Adobe Reader.*

Abstract

Despite the recent success of single image-based 3D human pose and shape estimation methods, recovering temporally consistent and smooth 3D human motion from a video is still challenging. Several video-based methods have been proposed; however, they fail to resolve the single image-based methods' temporal inconsistency issue due to a strong dependency on a static feature of the current frame. In this regard, we present a temporally consistent mesh recovery system (TCMR). It effectively focuses on the past and future frames' temporal information without being dominated by the current static feature. Our TCMR significantly outperforms previous video-based methods in temporal consistency with better per-frame 3D pose and shape

accuracy. We also release the codes¹.

1. Introduction

Various methods have been proposed to analyze humans from images, ranging from estimating a simplistic 2D skeleton to recovering 3D human pose and shape. Despite the recent improvements, estimating 3D human pose and shape from images is still a challenging task, especially in the monocular case due to depth ambiguity, limited training data, and complexity of human articulations.

Most of the previous methods [7, 12, 16, 17, 23, 27] attempt to recover 3D human pose and shape from a single image. They are generally based on parametric 3D hu-

¹https://github.com/hongsukchoi/TCMR_RELEASE

man mesh models, such as SMPL [19], and directly regress the model parameters from the input image. Although single image-based methods predict a reasonable output from a static image, they tend to produce temporally inconsistent and unsmooth 3D motion when applied to a video per frame. The temporal instability is from inconsistent 3D pose errors for consecutive frames. For example, the errors could occur in different 3D directions, or the following frames’ pose outputs could remain relatively the same, not reflecting the motion.

Several methods [13, 15, 20] have been proposed to extend the single image-based methods to the video case effectively. They feed a sequence of images to the pretrained single image-based 3D human pose and shape estimation networks [12, 16] to obtain a sequence of static features. All input frames’ static features are passed to a temporal encoder, which encodes a temporal feature for each input frame. Then, a body parameter regressor outputs SMPL parameters for each frame from the temporal feature of the corresponding time step.

Although the above works quantitatively improved the per-frame 3D pose accuracy and motion smoothness, their qualitative results still suffer from the temporal inconsistency aforementioned, as shown in Figure 1. We argue that the failure comes from a strong dependency on the static feature of the *current* frame. For terminological convenience, we use a word *current* to indicate the time step of a target frame where SMPL parameters to be estimated. The first reason for the strong dependency is a residual connection between the current frame’s static and temporal features. While the residual connection has been widely verified to facilitate a learning process, naively applying it to the temporal encoding can hinder the system from learning useful temporal information. Given that the static feature is extracted by the pretrained network [12, 16], it contains a strong cue for the SMPL parameters of the current frame. Thus, the residual connection’s identity mapping of the static feature can make the SMPL parameter regressor heavily depend on it and leverage the temporal feature marginally. This procedure can constrain the temporal encoder from encoding more meaningful temporal features. The second reason is the temporal encoding that takes static features from all frames, which include a current static feature. The current static feature has the largest potential to affect the current temporal feature, from which SMPL parameters are predicted. This phenomenon is caused by the current static feature having the most crucial information for 3D human pose and shape of a current frame. Although the dominance will increase the per-frame accuracy of 3D pose and shape estimation, it can prevent the temporal encoder from fully exploiting the past and future frames’ temporal information. Taken together, the existing video-based methods have a strong preference for the current static fea-

ture, and suffer from the temporal inconsistency issue as single image-based methods do.

In this work, we propose a temporally consistent mesh recovery system (TCMR). It is designed to resolve the strong dependency on the current static feature for temporally consistent and smooth 3D human motion output from a video. First, although we follow the previous video-based works [13, 15, 20] to encode a temporal feature of the current frame, we remove the residual connection between the static and temporal features. Moreover, we introduce PoseForecast, which consists of two temporal encoders, to *forecast* a current pose from the past and future frames without the current frame. The temporal features from PoseForecast are free from the current static feature; however, they contain essential temporal information of the past and future frames to *forecast* a current pose. The temporal features from PoseForecast are integrated with the current temporal feature, which is extracted from all input frames, to predict current SMPL parameters. The parameters estimated from the integrated temporal feature are the final output in inference time. By removing the strong dependency on the current static feature, our SMPL parameter regressor can have more chance to focus on the past and future frames without being dominated by the current frame.

Despite its simplicity, we observed that our newly designed temporal architecture is highly effective on obtaining the temporally consistent and smooth 3D human motion. It also improves the accuracy of the 3D pose and shape per frame by utilizing better temporal information. We show that the proposed TCMR outperforms the previous video-based methods [13, 15, 20] on various 3D video benchmarks, especially in temporal consistency.

Our contributions can be summarized as follows.

- We present a temporally consistent mesh recovery system (TCMR), which produces temporally consistent and smooth 3D human motion from a video. It effectively leverages temporal information from the past and future frames without being dominated by the static feature of the current frame.
- Despite its simplicity, TCMR not only improves the temporal consistency of 3D human motion but also increases per-frame 3D pose and shape accuracy compared to a baseline method.
- TCMR outperforms previous video-based methods in temporal consistency by a large margin while achieving better per-frame 3D pose and shape accuracy.

2. Related works

Single image-based 3D human pose and shape estimation. Most of the current single image-based 3D human

pose and shape estimation methods are based on the model-based approach, which predicts parameters of a predefined 3D human mesh model, SMPL [19]. Kanazawa *et al.* [12] proposed an end-to-end trainable human mesh recovery (HMR) system that uses adversarial loss to make their output 3D human mesh anatomically plausible. Pavlakos *et al.* [27] used 2D joint heatmaps and silhouette as cues for predicting accurate SMPL parameters. Omran *et al.* [24] proposed a similar system, which uses human part segmentation as a cue for regressing SMPL parameters. Pavlakos *et al.* [26] proposed a system that uses multi-view color consistency to supervise a network using multi-view geometry. Kolotouros *et al.* [16] introduced a self-improving system that consists of an SMPL parameter regressor and an iterative fitting framework [2]. Georgakis *et al.* [9] incorporated hierarchical kinematic prior on a human body to a network.

Conversely, the model-free approach estimates the shape directly instead of regressing the model parameters. Varol *et al.* [34] proposed BodyNet, which estimates 3D human shape in the 3D volumetric space. Kolotouros *et al.* [17] designed a graph convolutional human mesh regression system. Their graph convolutional network takes a template human mesh in a rest pose as input and predicts mesh vertex coordinates using image features from ResNet [10]. Moon and Lee [23] introduced a lixel-based 1D heatmap to localize mesh vertices in a fully convolutional manner. Choi *et al.* [7] proposed a graph convolutional network that recovers 3D human pose and mesh from a 2D human pose.

Despite moderate performance on a static image, the single image-based works suffer from temporal inconsistency (*e.g.*, sudden change of poses), when applied to a video.

Video-based 3D human pose and shape estimation. HMMR [13] extracts static features and encodes them to a temporal feature using a 1D fully convolutional temporal encoder. It learns temporal context representation to reduce the 3D prediction’s temporal inconsistency by predicting 3D poses in the nearby past and future frames. Dorrer *et al.* [8] trained their network on a sequence of optical flow and 2D poses to make their network generalize well to unseen videos. Sun *et al.* [33] proposed a skeleton-disentangling framework, which separates 3D human pose and shape estimation into multi-level spatial and temporal subproblems. They enforced the network to order shuffled frames to encourage temporal feature learning. VIBE [15] encodes static features from the input frames into a temporal feature by using a bi-directional gated recurrent unit (GRU) [6], and feeds it to an SMPL parameter regressor. A motion discriminator is introduced to encourage the regressor to produce plausible 3D human motion. MEVA [20] addresses the problem in a coarse-to-fine manner. Their system initially estimates the coarse 3D human motion using a variational motion estimator (VME), and predicts the residual motion with a motion residual regressor (MRR).

Temporally consistent 3D human motion from a video.

Although there have been many methods for video-based 3D human motion estimation [3, 8, 13, 15, 20, 22, 28, 29, 33], most of them showed their results only qualitatively, and did not report numerical evaluation on temporal consistency. After the HMMR [13] introduced the 3D pose acceleration error for the temporal consistency and smoothness of human motion, the following works [15, 20] have reported the error metric. HMMR and VIBE [15] lowered the acceleration error compared with the single image-based methods. However, they revealed a trade-off between per-frame accuracy and temporal consistency. The HMMR outputs smoother 3D human motion but provides low per-frame 3D pose accuracy. Conversely, the VIBE [15] shows high per-frame 3D pose accuracy; however, the output is temporally inconsistent in quantitative metrics and qualitative results compared with HMMR.

In this regard, MEVA [20] attempts to establish the balance between the per-frame 3D pose accuracy and the temporal smoothness. Although it provides better results in both metrics, the qualitative results still expose unsmooth 3D motion. The reason is that the system strongly depends on the current static feature to estimate the current 3D pose and shape. First, MEVA uses a residual connection between the current frames’ static and temporal features. In addition, the current temporal feature, which is used to refine initial 3D pose and shape by MRR, is encoded from static features of all frames, which include the current frame. This procedure can make the temporal feature dominated by the current static feature. As a result, the refinement is significantly driven by the current static feature, and the 3D errors from consecutive frames appear inconsistent. On the contrary, our TCMR is deliberately designed to reduce the strong dependency on the static feature. The residual connection is removed, and PoseForecast forecasts additional temporal features from past and future frames without a current frame. Our approach alleviates the dependency and provides temporally consistent and accurate 3D human motions in both qualitative and quantitative manners.

Forecasting 3D human poses from images.

Recently, [5, 13, 36, 37] proposed to predict a person’s future 3D human poses from RGB input. Chao *et al.* [5] leveraged a recurrent neural network (RNN) to forecast a sequence of 2D poses from a static image, and estimate 3D poses from the predicted 2D poses. The HMMR [13] predicts the current, future, and past 3D poses from a current input image using a hallucinator. It hallucinates the past and future 3D poses from a current frame and is self-supervised by the output of the 1D fully convolutional temporal encoder. Zhang *et al.* [37] proposed a neural autoregressive framework that takes past video frames as input to predict future 3D motion. Yuan *et al.* [36] adopted deep reinforcement learning to forecast future 3D human poses from egocentric videos.

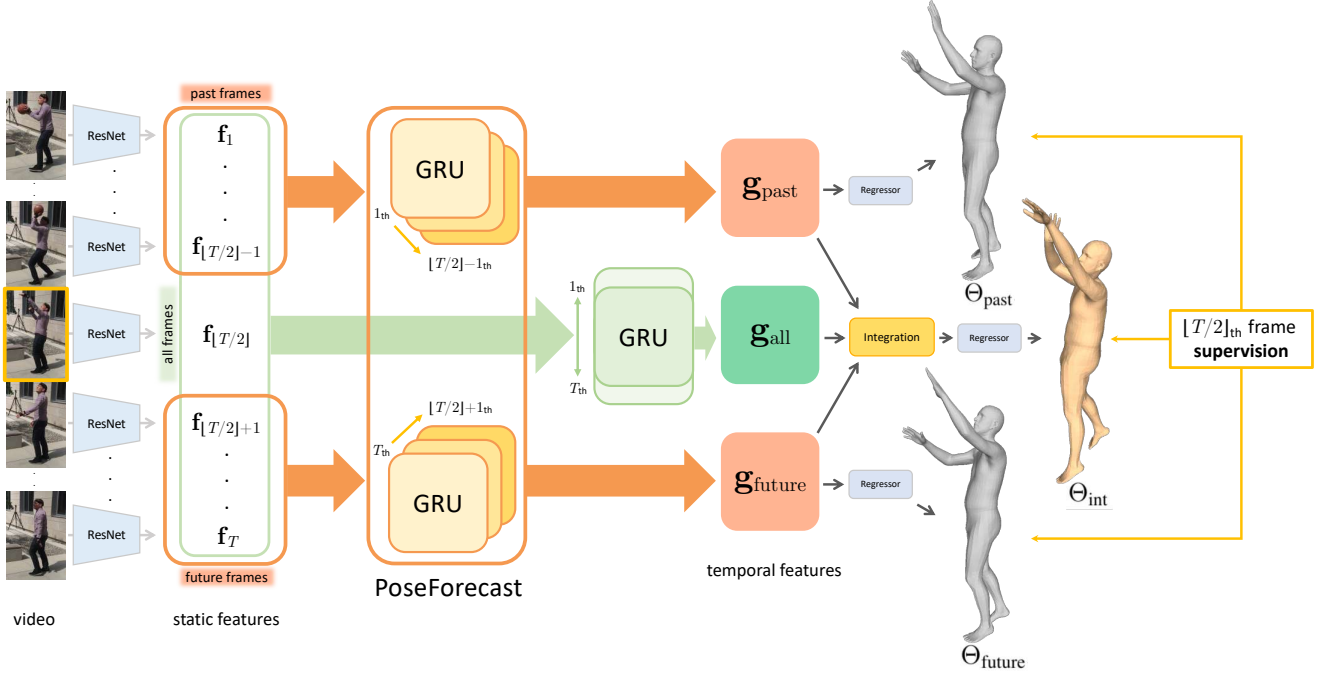


Figure 2: The overall pipeline of TCMR. The gold-colored output Θ_{int} is used in inference time, which is regressed from the integrated temporal feature.

Although the objective of the above methods is to forecast future 3D poses, our system aims to learn useful temporal features free from a current static feature by the forecasting.

3. TCMR

Figure 2 shows the overall pipeline of our TCMR. We provide descriptions of each part in the system as follows.

3.1. Temporal encoding from all frames

Given a sequence of T RGB frames $\mathbf{I}_1, \dots, \mathbf{I}_T$, ResNet [10], pretrained by Kolotouros *et al.* [16], extracts a static image feature per frame. Then, a global average pooling is applied on the ResNet outputs, which become $\mathbf{f}_1, \dots, \mathbf{f}_T$, where $\mathbf{f}_\bullet \in \mathbb{R}^{2048}$. The network weights of the ResNet are shared for all frames.

From the extracted static features of all input frames, we compute the current frame’s temporal feature using a bi-directional GRU, which consists of two uni-directional GRUs. We denote the bi-directional GRU as \mathcal{G}_{all} . The current frame is defined as a $\lfloor T/2 \rfloor$ th frame among T input frames. The two uni-directional GRUs extract temporal features from the input static features in the opposite time directions. The initial inputs of the two GRUs are \mathbf{f}_1 and \mathbf{f}_T , respectively, and the initial hidden states of them are initialized as zero tensors. Then, they recurrently updates their hidden states by aggregating the static features from the next frames $\mathbf{f}_2, \dots, \mathbf{f}_{\lfloor T/2 \rfloor}$ and $\mathbf{f}_{T-1}, \dots, \mathbf{f}_{\lfloor T/2 \rfloor}$, respectively. The concatenated hidden states of the GRUs at the current frame become the current temporal feature from all

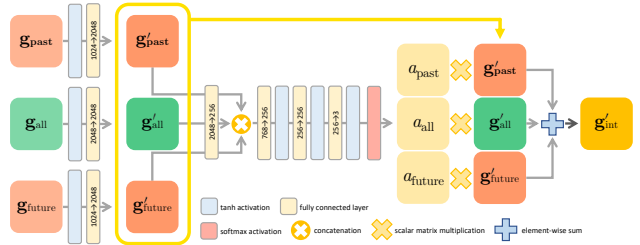


Figure 3: Temporal feature integration to estimate 3D human mesh for the current frame.

input frames $\mathbf{g}_{\text{all}} \in \mathbb{R}^{2048}$. Unlike VIBE [15], we do not add residual connection between $\mathbf{f}_{\lfloor T/2 \rfloor}$ and \mathbf{g}_{all} , such that the current temporal feature will not be dominated by $\mathbf{f}_{\lfloor T/2 \rfloor}$.

3.2. Temporal encoding by PoseForecast

PoseForecast *forecasts* additional temporal features for the current target pose from the past and future frames by employing two additional GRUs, denoted as $\mathcal{G}_{\text{past}}$ and $\mathcal{G}_{\text{future}}$, respectively. The past and future frames are defined as $1, \dots, (\lfloor T/2 \rfloor - 1)$ th frames and $(\lfloor T/2 \rfloor + 1), \dots, T$ th frames, respectively. The initial input of the $\mathcal{G}_{\text{past}}$ is \mathbf{f}_1 , and the initial hidden state is initialized as a zero tensor. Then, it recurrently updates its hidden state by aggregating the static features from the next frames $\mathbf{f}_2, \dots, \mathbf{f}_{\lfloor T/2 \rfloor - 1}$. The final hidden state of the $\mathcal{G}_{\text{past}}$ becomes the temporal feature from the past frames $\mathbf{g}_{\text{past}} \in \mathbb{R}^{1024}$. Similarly, $\mathcal{G}_{\text{future}}$ takes \mathbf{f}_T as an initial input with a zero-initialized hidden state, and re-

currently updates its hidden state by aggregating the static features from the next frames $\mathbf{f}_{T-1}, \dots, \mathbf{f}_{\lfloor T/2 \rfloor + 1}$. The final hidden state of the $\mathcal{G}_{\text{future}}$ becomes the temporal feature from the future frames $\mathbf{g}_{\text{future}} \in \mathbb{R}^{1024}$.

3.3. Temporal feature integration

We integrate the extracted temporal features from all frames \mathbf{g}_{all} , from the past frames \mathbf{g}_{past} , and from the future frames $\mathbf{g}_{\text{future}}$ for the final 3D mesh estimation, as illustrated in Figure 3. For the integration, we pass each temporal feature to ReLU activation function and a fully connected layer to change the size of the channel dimension to 2048. The outputs of the fully connected layer are denoted as \mathbf{g}'_{all} , $\mathbf{g}'_{\text{past}}$, and $\mathbf{g}'_{\text{future}}$. Then, the output features are resized to 256 by a shared fully connected layer and concatenated. The concatenated feature is passed to several fully connected layers, followed by the softmax activation function, which produces attention values $\mathbf{a} = (a_{\text{all}}, a_{\text{past}}, a_{\text{future}}) \in \mathbb{R}^3$. The attention values represent how much the system should give a weight for the feature integration. The final integrated temporal feature is obtained by $\mathbf{g}'_{\text{int}} = a_{\text{all}}\mathbf{g}'_{\text{all}} + a_{\text{past}}\mathbf{g}'_{\text{past}} + a_{\text{future}}\mathbf{g}'_{\text{future}}$.

In the training stage, we pass $\mathbf{g}'_{\text{past}}$, $\mathbf{g}'_{\text{future}}$, and \mathbf{g}'_{int} to the SMPL parameter regressor, which outputs Θ_{past} , Θ_{future} , and Θ_{int} from each input temporal feature, respectively. The regressor is shared for all outputs. Θ_{\bullet} denotes a union of SMPL parameter set $\{\theta_{\bullet}, \beta_{\bullet}\}$ and weak-perspective camera parameter set $\{s_{\bullet}, t_{\bullet}\}$. θ , β , s , and t represent SMPL pose parameter, identity parameter, scale, and translation, respectively. In the testing stage, we only pass \mathbf{g}'_{int} to the parameter regressor and use Θ_{int} as the final 3D human mesh.

3.4. Loss functions

For the training, we supervise all three outputs Θ_{past} , Θ_{future} , and Θ_{int} with current frame groundtruth. $L2$ loss between predicted and groundtruth SMPL parameters and 2D/3D joint coordinates are used, following VIBE [15]. The 3D joint coordinates are obtained by forwarding the SMPL parameters to the SMPL layer, and the 2D joint coordinates are obtained by projecting the 3D joint coordinates using the predicted camera parameters.

4. Implementation details

Following VIBE [15], we set the length of the input sequence T to 16 and the input video frame rate to 25-30 frames per second and initialize the backbone and regressor with the pretrained SPIN [16]. The weights are updated by the Adam optimizer [14] with a mini-batch size of 32. The human body region is cropped using a groundtruth box in both of training and testing stages following previous works [12, 15–17]. The cropped image is resized to 224×224 . Inspired by Sarandi *et al.* [32], we occlude the cropped image with various objects for data augmentation.

input -res +PF (Ours) +res -PF

Figure 4: Qualitative comparison between our TCMR (middle) and the baseline (right). TCMR learns more useful temporal features, and provides a more accurate 3D pose and temporally consistent 3D motion. res denotes the residual connection and PF is the abbreviation for PoseForecast. *This is a video figure that is best viewed by Adobe Reader.*

The occlusion augmentation reduces both pose and acceleration errors approximately by $1mm$. Following [13, 15], we precompute the static features from the cropped images by ResNet [10] to save training time and memory. All the 3D rotations of θ are initially predicted in the 6D rotational representation of Zhou *et al.* [39], and converted to the 3D axis-angle rotations. The initial learning rate is set to 5^{-5} and reduced by a factor of 10, when the 3D pose accuracy does not improve after every 5 epochs. We train the network for 30 epochs with one NVIDIA RTX 2080Ti GPU. PyTorch [25] is used for code implementation.

5. Experiment

5.1. Evaluation metrics and datasets.

Evaluation metrics. We report the per-frame and temporal evaluation metrics. For the per-frame evaluation, we use mean per joint position error (MPJPE), Procrustes-aligned MPJPE (PA-MPJPE), and mean per vertex position error (MPVPE). The position errors are measured in millimeter (mm) between the estimated and groundtruth 3D coordinates after aligning the root joint. Particularly, we use PA-MPJPE as the main metric for per-frame accuracy, since it excludes the effect of outputs’ scale ambiguity on errors. For the temporal evaluation, we use the acceleration error proposed in HMMR [13]. The acceleration error computes an average of the difference between the predicted and groundtruth acceleration of each joint in (mm/s^2).

Datasets. We use 3DPW [35], Human3.6M [11], MPI-INF-3DHP [21], InstaVariety [13], Penn Action [38], and PoseTrack [1] for training, following VIBE [15]. 3DPW is the only in-the-wild dataset that contains accurate groundtruth SMPL parameters. 3DPW, Human3.6M, MPI-INF-3DHP are also used for evaluation. More details are in the supplementary material.

Table 1: Comparison between different temporal architectures. All networks estimate only on the middle frame of the input sequence.

remove residual	PoseForecast	PA-MPJPE↓	Accel↓
✗	✗	55.6	29.2
✗	✓	55.0	24.9
✓	✗	54.2	8.7
✓ (Ours)	✓	53.9	7.7

Table 2: Comparison between PoseForecast that takes a current frame and that does not take a current frame.

PoseForecast input	PA-MPJPE↓	Accel↓
w. current frame	53.8	10.3
wo. current frame (Ours)	53.9	7.7

5.2. Ablation study

In this study, we show how each component of our temporal architecture reduces the dependency of the model on a current static feature, and make it focus on temporal features from the past and future. We take the same baseline used in VIBE [15]. The baseline has a single bi-directional GRU that encodes temporal features from all input frames and a residual connection between the static and temporal features as VIBE. It also predicts each 3D pose and shape for all input frames in a single feed-forward, but does not use the motion discriminator. We use 3DPW [35], MPI-INF-3DHP [21], InstaVariety [13], and Penn Action [38] for training, and 3DPW for evaluation.

Effectiveness of residual connection removal. To analyze the effect of the residual connection between the static and temporal features, we compare the models with and without it. As shown in Table 1, removing the residual connection decreases the acceleration error significantly, which indicates a considerable improvement in temporal consistency and smoothness of 3D human motion. This finding verifies that the identity mapping of the current static feature inside the residual connection hinders a model from learning meaningful temporal features. Moreover, the increased temporal consistency of 3D motion improves the per-frame 3D pose accuracy. Figure 4 illustrates how the enhanced temporal consistency contributes to better per-frame 3D pose estimation. The sudden change of poses, caused by the inaccurate 3D pose estimation on specific frames, is disappeared. The above comparisons clearly validate the effectiveness of removing the residual connection in terms of both per-frame and temporal metrics.

Effectiveness of PoseForecast We compare the models with and without PoseForecast to verify the effectiveness of forecasting current temporal features only from the past and future frames. On the basis of the results in Table 1, PoseForecast consistently improves per-frame and temporal

Table 3: Comparison between different supervision on estimated SMPL parameters from the PoseForecast.

PoseForecast supervision target	PA-MPJPE↓	Accel↓
none	55.1	8.3
GT of past and future frames	54.1	8.5
GT of current frame (Ours)	53.9	7.7

metrics regardless of the residual connection. Particularly, the acceleration error consistently decreases by over 11%. Thus, the temporal encoding that takes all frames with the current frame may be suboptimal, and forecasting the current temporal features from the past and future frames is beneficial for temporally consistent 3D human motion.

To further validate the forecasting, we compare our PoseForecast with its variations. First, we show the effectiveness of taking past and future frames without a current frame in Table 2. As the table shows, additionally taking current frames increases the acceleration error by 33%. Thus, maintaining the temporal features free from the current static feature is important for temporally consistent and smooth 3D human motion. Second, we validate the effectiveness of supervising the predicted SMPL parameters from PoseForecast (*i.e.*, Θ_{past} and Θ_{future}) with groundtruth of the current frame in Table 3. As shown in the table, supervising the predicted parameters with the current groundtruth provides better per-frame 3D pose accuracy and temporal consistency than the other supervisions. When we supervise the predicted parameters with the groundtruth of $\lfloor T/2 \rfloor - 1$ th and $\lfloor T/2 \rfloor + 1$ th frames (the second row), the acceleration error increases by 10%. The performance degrades, because the temporal features of PoseForecast are encoded from the input including the static features of the target frames (*i.e.*, $\lfloor T/2 \rfloor - 1$ th and $\lfloor T/2 \rfloor + 1$ th frames). As verified in Table 2, including the target static feature hinders PoseForecast from learning useful temporal information for temporally consistent and smooth 3D human motion. The encoded temporal feature is likely to be dominated by the target static feature and marginally leverage temporal information from other frames. When no supervision is observed (the first row), both 3D pose accuracy and temporal consistency decrease compared with ours. Hence, designing our PoseForecast to forecast the current SMPL parameters by supervising it with the current target (the third row) facilitates the network to learn more useful temporal features.

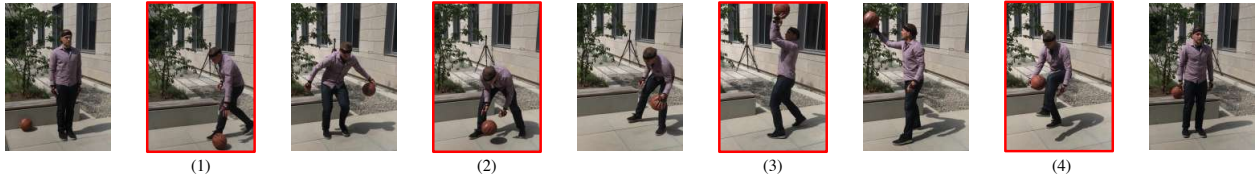
In summary, the above comparisons show that forecasting current temporal features from past and future frames is effective for temporally consistent 3D human motion by reducing the strong dependency on a current static feature.

5.3. Comparison with state-of-the-art methods

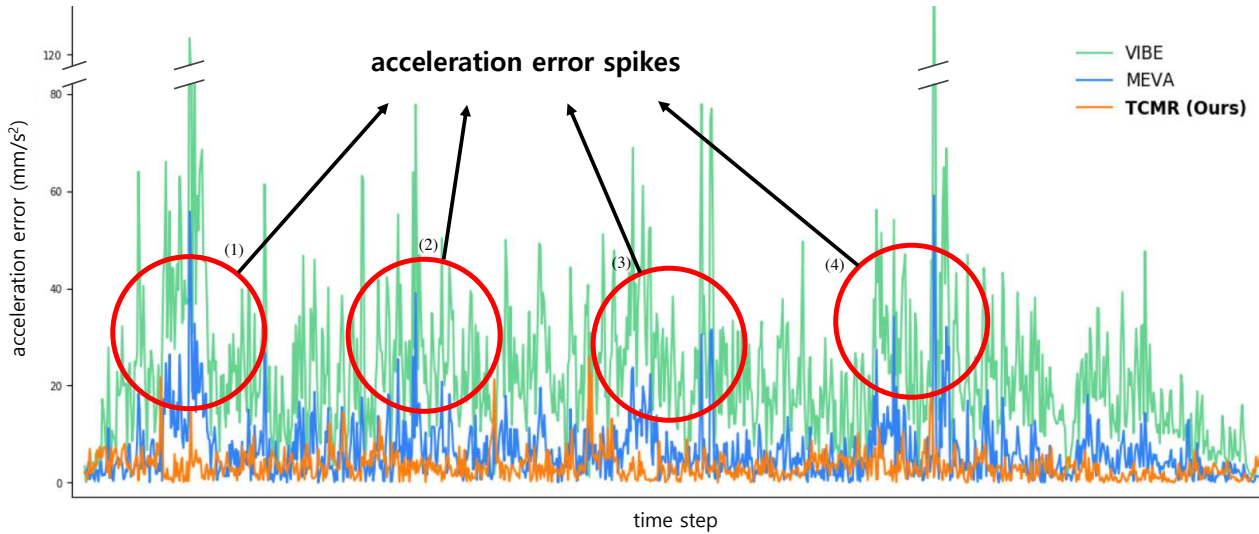
Comparison with video-based methods. We compare our TCMR with previous state-of-the-art video-based methods

Table 4: Evaluation of state-of-the-art methods on on 3DPW [35], MPI-INF-3DHP [21], and Human3.6M [11]. All methods except HMMR [13] do not use Human3.6M SMPL parameters from Mosh [18], but use 3DPW train set for training following MEVA [20]. The number of input frames are following the protocols of the papers.

method	3DPW				MPI-INF-3DHP			Human3.6M			number of input frames
	PA-MPJPE ↓	MPJPE ↓	MPVPE ↓	Accel ↓	PA-MPJPE ↓	MPJPE ↓	Accel ↓	PA-MPJPE ↓	MPJPE ↓	Accel ↓	
HMMR [13]	72.6	116.5	139.3	15.2	-	-	-	56.9	-	-	20
VIBE [15]	57.6	91.9	-	25.4	68.9	103.9	27.3	53.3	78.0	27.3	16
MEVA [20]	54.7	86.9	-	11.6	65.4	96.4	11.1	53.2	76.0	15.3	90
TCMR (Ours)	52.7	86.5	103.2	6.8	63.5	97.6	8.5	52.0	73.6	3.9	16



(a) sampled frames of 'courtyard_basketball_01' in order



(b) acceleration error plot of 'courtyard_basketball_01'

Figure 5: Comparison between the acceleration errors of the proposed TCMR, MEVA [20], and VIBE [15]. Our TCMR shows clearly lower acceleration errors along the time step than previous methods, which indicates temporally consistent 3D motion output. The previous methods reveal extreme acceleration error spikes compared to our TCMR.

Table 5: Comparison between ours and previous methods applied with average filtering on 3DPW [35].

method	PA-MPJPE ↓	MPJPE ↓	Accel ↓
VIBE [15]	57.6	91.9	25.4
+ Avg. filter	57.8	91.6	13.5
MEVA [20]	54.7	86.9	11.6
+ Avg. filter	55.5	87.7	8.2
TCMR (Ours)	52.7	86.5	6.8
+ Avg. filter	55.0	88.7	6.2

[13, 15, 20] that report the acceleration error in Table 4. On the basis of the study of Luo *et al.* [20], all methods, except HMMR [13] are trained on the train set including

3DPW [35], but do not leverage Human3.6M [11] SMPL parameters obtained from Mosh [18] for supervision. The numbers of VIBE [15] are from MEVA [20], but we validated them independently. As shown in the table, our proposed system outperforms the previous video-based methods on all benchmarks both in per-frame 3D pose accuracy and temporal consistency. These results prove that our system effectively leverages temporal information of the past and future by resolving the system's strong dependency on a current static feature. Although MEVA [20] also improves the per-frame and temporal metrics, the model consumes nearly 6 times more input frames during training and testing, and provides worse results than ours. In addition,

Table 6: Evaluation of state-of-the-art methods on 3DPW [35], MPI-INF-3DHP [21], and Human3.6M [11]. All methods do not use 3DPW [35] on training. ‘single image’ or ‘video’ denotes whether the input of a method is a single image or a video.

method	3DPW				MPI-INF-3DHP			Human3.6M			
	PA-MPJPE ↓	MPJPE ↓	MPVPE ↓	Accel ↓	PA-MPJPE ↓	MPJPE ↓	Accel ↓	PA-MPJPE ↓	MPJPE ↓	Accel ↓	
single image	HMR [12]	76.7	130.0	-	37.4	89.8	124.2	-	56.8	88.0	-
	GraphCMR [17]	70.2	-	-	-	-	-	-	50.1	-	-
	SPIN [16]	59.2	96.9	116.4	29.8	67.5	105.2	-	41.1	-	18.3
	I2L-MeshNet [23]	57.7	93.2	110.1	30.9	-	-	-	41.1	55.7	13.4
	Pose2Mesh [7]	58.3	88.9	106.3	22.6	-	-	-	46.3	64.9	23.9
	HKMR [9]	-	-	-	-	-	-	-	-	59.6	-
video	HMMR [13]	72.6	116.5	139.3	15.2	-	-	-	56.9	-	-
	Doersch <i>et al.</i> [8]	74.7	-	-	-	-	-	-	-	-	-
	Sun <i>et al.</i> [33]	69.5	-	-	-	-	-	-	42.4	59.1	-
	VIBE [15]	56.5	93.5	113.4	27.1	63.4	97.7	29.0	41.5	65.9	18.3
	TCMR (Ours)	55.8	95.0	111.3	6.7	62.8	96.5	9.5	41.1	62.3	5.3

MEVA requires at least 90 input frames, which means that it can not be trained and tested on short videos. Figure 5 describes the clear advantage of our TCMR on the temporal consistency among video-based methods. The previous methods expose numerous spikes, which represent unstable and unsmooth 3D motion estimation. Our TCMR provides relatively low acceleration errors along the time step, which indicates temporally consistent 3D motion output. The figure’s acceleration errors are measured on a sequence of the 3DPW validation set that has a diverse motion.

To further confirm the effectiveness of the proposed system on temporal consistency, we compare our TCMR with VIBE [15] and MEVA [20] with an average filter applied as post-processing in Table 5. Average filtering is performed by spherical linear interpolation in the quaternions of estimated SMPL [19] pose parameters following MEVA. The numbers of other methods are from MEVA. As shown in the table, our system outperforms other methods even when they are applied with the average filtering. Moreover, the results imply that the average filtering can decrease the per-frame 3D pose accuracy by smoothing out the details of 3D human motion. However, each component of our TCMR decreases the acceleration error while improving the per-frame 3D pose accuracy, as shown in Table 1.

In summary, our newly designed system significantly outperforms the previous state-of-the-art methods in temporal consistency and smoothness of 3D human motion without any post-processing while also increasing the per-frame 3D pose accuracy. Note that the comparison in Table 4 and 5 is the fairest comparison between the video-based methods, since all methods, except HMMR [13], used the same training datasets.

Comparison with single image-based and video-based methods. We compare our system with previous 3D pose and shape estimation methods, including single image-based methods in Table 6. None of the methods are trained on 3DPW [35]. For evaluation on Human3.6M [11], we use the frontal view images following [13, 16], whereas all views are tested in Table 4 and 5. In addition, to confirm

the acceleration error of VIBE [15] on MPI-INF-3DHP [21] and Human3.6M, we re-evaluate the model using the pre-trained weights provided in the official code repository.

As shown in the table, our method outperforms all the previous methods on 3DPW, a challenging in-the-wild benchmark, and MPI-INF-3DHP in per-frame 3D pose accuracy (PA-MPJPE) and temporal consistency. Especially the temporal consistency is largely improved compared with single image-based methods. While VIBE decreases the acceleration error of SPIN [16] by 9% and is defeated by Pose2Mesh [7] in the temporal consistency, our system provides over 3 times better performance than both SPIN and Pose2Mesh in 3DPW. Moreover, VIBE gives a higher acceleration error than I2L-MeshNet [23] but our TCMR outperforms it by a wide margin in Human3.6M.

We provide qualitative comparison with VIBE [15] and MEVA [20] on 3DPW, qualitative results of TCMR on Internet videos, and failure cases in this link ².

6. Conclusion

We present TCMR, a novel and powerful system that estimates a 3D human mesh from a RGB video. Previous video-based methods suffer from the temporal inconsistency issue because of the strong dependency on the static feature of the current frame. We resolve the issue by removing the residual connection between the static and temporal features, and employing PoseForecast that forecasts the current temporal feature from the past and future frames. In comparison with the previous video-based methods, the proposed TCMR provides highly temporally consistent 3D motion and a more accurate 3D pose per frame.

Acknowledgements. This work was supported by IITP grant funded by the Ministry of Science and ICT of Korea (No. 2017-0-01780) and AIRS Company in Hyundai Motor Company & Kia Motors Corporation through HKMC-SNU AI Consortium Fund.

²<https://www.youtube.com/watch?v=WB3nTnSQDII>

Supplementary Material for

Beyond Static Features for Temporally Consistent 3D Human Pose and Shape from a Video

7. More qualitative results

We provide more qualitative results in the online video ³, which consists of three parts. The first part shows the qualitative results of our TCMR on in-the-wild videos that have fast and diverse motions from 3DPW [35]. We also provide the outputs rendered from the opposite view. The second part compares the proposed TCMR with VIBE [15] and MEVA [20]. The results are rendered on a plain background with a fixed camera to clearly compare the temporal consistency and smoothness of 3D human motion following MEVA [20]. The fixed camera has the fixed weak-perspective camera parameters s and t , which are set to one and zero, respectively. The last part provides the results of TCMR on Internet videos. The bounding boxes of people in the videos are tracked by a multi-person tracker that uses YOLOv3 [30]. With the cropped images from the bounding boxes, our TCMR processes 41 frames per second (fps) for the video ⁴ with 5 people. A single NVIDIA RTX 2080Ti GPU is used for the test.

8. Human evaluation.

We surveyed 50 people to pick the most realistic motion from TCMR, MEVA, and VIBE outputs on 20 sequences of 3DPW [35] validation and test sets. TCMR, MEVA, and VIBE got 69%, 26%, and 5% votes, respectively. The result is coherent with the acceleration error results of the three methods in the main manuscript.

9. Attention values in feature integration.

During the temporal feature integration, the past and future temporal features are weighted more than the current temporal feature, and the variation range of each attention value is $\pm 20\%$. The past and future temporal features' attention values tend to become larger when the current pose is difficult or the motion is fast. The attached videos^{5,6} plot the attention values of the past, future, and current temporal features on two sequences of 3DPW [35]. The values are written at the top-right of frames, and the sum is always 1. As the video shows, the attention value of the current temporal feature does not drop below 0.4 when a subject is

walking in slow motion, whereas the value overall stays below 0.4 when a subject is playing basketball with fast movement and complex poses.

10. Datasets

3DPW. 3DPW [35] is captured from in-the-wild and contains 3D human pose and shape annotations. It consists of 60 videos and 51K video frames in total, which are captured with a phone at 30 fps. IMU sensors are leveraged to acquire the groundtruth 3D human pose and shape. We follow the official split protocol to train and test our model, where train, validation, test sets consist of 24, 12, 24 videos, respectively. Also, we report MPVPE on 3DPW because it only has groundtruth 3D shape among the datasets we used. We use 14 joints defined by Human3.6M [11] for evaluating PA-MPJPE and MPJPE following the previous works [12, 13, 15, 16].

Human3.6M. Human3.6M [11] is a large-scale indoor 3D human pose benchmark, which consists of 15 action categories and 3.6M video frames. Following [15], our TCMR is trained on 5 subjects (S1, S5, S6, S7, S8) and tested on 2 subjects (S9, S11). We subsampled the dataset to 25 fps (originally 50 fps) for training and evaluation on the acceleration error. 14 joints defined by Human3.6M are used for computing PA-MPJPE and MPJPE.

MPI-INF-3DHP. MPI-INF-3DHP [21] is a 3D benchmark mostly captured from indoor environment. The train set has 8 subjects, 16 videos per subject, and 1.3M video frames captured at 25 fps in total. It exploits a marker-less motion capture system and provides 3D human pose annotations.

The test set contains 6 subjects performing 7 actions in both the indoor and outdoor environment. The positional errors (*i.e.*, PA-MPJPE and MPJPE) of TCMR are measured on the valid frames, which are composed of every 10th frame approximately, using 17 joints defined by MPI-INF-3DHP. The acceleration error is computed using all frames.

InstaVariety. InstaVariety is a 2D human dataset curated by HMMR [13], whose videos are collected from Instagram using 84 motion-related hashtags. There are 28K videos with an average length of 6 seconds, and OpenPose [4] is used to acquire pseudo-groundtruth 2D pose annotations.

Penn Action. Penn Action [38] contains 2.3K video sequences of 15 different sports actions. It has a total of 77K video frames annotations for 2D human poses, bounding boxes, and action categories.

³<https://www.youtube.com/watch?v=WB3nTnSQDII>

⁴<https://www.youtube.com/watch?v=Opry3F6aB1I>

⁵<https://youtu.be/dFQ6hkfkWz0>

⁶<https://youtu.be/otdL5WVjwPg>

Table 7: Comparison between different models using ResNet with different initialization to extract static features. All models use the same SMPL parameter regressor pretrained by SPIN [16].

ResNet initialization	remove residual	PoseForecast	PA-MPJPE↓	Accel↓
ResNet with random initialization	✗	✗	126.5	24.3
ResNet pretrained on ImageNet [31]	✗	✗	103.7	65.5
ResNet from SPIN [16]	✗	✗	55.6	29.2
ResNet from SPIN [16] (TCMR. Ours.)	✓	✓	53.9	7.7

PoseTrack. PoseTrack [1] is a 2D benchmark for multi-person pose estimation and tracking in videos. It contains 1.3K videos and 46K annotated frames in total. The videos are captured at different fps, varying around 25 fps. We use 792 videos from the official train set, which has 2D pose annotations for 30 frames in the middle of the video.

11. Effect of pretrained ResNet

Due to lack of video data, our TCMR and previous video-based methods [13, 15, 20] employ ResNet [10] pretrained by the single image-based 3D human pose and shape estimation methods [12, 16] to extract static features from input frames. The pretrained ResNet is trained on large-scale in-the-wild 2D human pose datasets and provides reliable static features. However, it is also one reason for the strong dependency of the system on the current static feature. The current static feature extracted by the pretrained ResNet already contains a strong cue on the current 3D human pose and shape, leading the system to leverage temporal information marginally.

In this regard, an alternative to our TCMR, one could train models from scratch without using the ResNet pretrained by [12, 16] to extract static features to reduce the strong dependency. Table 7 compares our TCMR, the baseline (the third row), and the models that do not use the ResNet pretrained by SPIN [16]. As the table shows, the models that do not use the ResNet pretrained by SPIN [16] reveal very high per-frame 3D pose errors. This indicates that training the models with only video data in the current literature is not sufficient for accurate 3D human pose estimation. The interesting part is that the model using ResNet with random initialization provides the highest 3D pose error but the lowest acceleration error among the models without our TCMR. While the high pose error attributes to the lack of train data, the low acceleration error implies that the strong cue of the current static feature adversely affects the temporal consistency of 3D human motion.

In summary, with the insufficient video data in the current literature, the proposed TCMR significantly improves the temporal consistency of 3D human motion by reducing the strong dependency on the current static feature. It also preserves the per-frame 3D pose accuracy by leveraging the ResNet pretrained on large-scale in-the-wild 2D hu-

Table 8: Performance comparison between two networks taking different input fps on 3DPW [35]. The numbers in the second row are from Table 4 of the main manuscript.

input fps	PA-MPJPE↓	Accel↓
15	53.5	15.3
30	52.7	7.1

man pose datasets to extract useful static features.

12. Effect of input fps

Table 8 shows the effect of input fps. The acceleration error doubles when input fps reduces by half, whereas the accuracy remains relatively the same. The result indicates that TCMR can still fix invalid poses using relatively sparse temporal information. The result also implies that temporally dense information is critical for temporal consistency of outputs, which is intuitive.

13. Pose2Mesh with temporal smoothing

We performed temporal smoothing on Pose2Mesh [7], the state-of-the-art single image-based 3D human pose and shape estimation method. Pose2Mesh wins the first in MPJPE, MPVPE, and acceleration error and the second in PA-MPJPE among single image-based methods according to Table 6 of the main manuscript. Pose2Mesh with euro-filter achieves PA-MPJPE 58.6, MPJPE 89.6, acceleration error 12.9 on 3DPW. TCMR still outperforms the smoothed Pose2Mesh by nearly twice in temporal consistency without any post-processing.

References

- [1] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. Posetrack: A benchmark for human pose estimation and tracking. *CVPR*, 2018. 5, 10
- [2] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it SMPL: Automatic estimation of 3D human pose and shape from a single image. *ECCV*, 2016. 3
- [3] Yujun Cai, Lihao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann. Exploiting spatial-temporal relationships for 3D pose estimation via graph convolutional networks. *ICCV*, 2019. 3
- [4] Zhe Cao, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Realtime multi-person 2d pose estimation using part affinity fields. 2017. 9
- [5] Yu-Wei Chao, Jimei Yang, Brian Price, Scott Cohen, and Jia Deng. Forecasting human dynamics from static images. *CVPR*, 2017. 3
- [6] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *EMNLP*, 2014. 3
- [7] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Pose2Mesh: Graph convolutional network for 3D human pose and mesh recovery from a 2D human pose. *ECCV*, 2020. 1, 3, 8, 10
- [8] Carl Doersch and Andrew Zisserman. Sim2real transfer learning for 3D human pose estimation: motion to the rescue. *NeurIPS*, 2019. 3, 8
- [9] Georgios Georgakis, Ren Li, Srikrishna Karanam, Terrence Chen, Jana Košecká, and Ziyang Wu. Hierarchical kinematic human mesh recovery. *ECCV*, 2020. 3, 8
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CVPR*, 2016. 3, 4, 5, 10
- [11] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI*, 2014. 5, 7, 8, 9
- [12] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. *CVPR*, 2018. 1, 2, 3, 5, 8, 9, 10
- [13] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3D human dynamics from video. *CVPR*, 2019. 2, 3, 5, 6, 7, 8, 9, 10
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2014. 5
- [15] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. VIBE: Video inference for human body pose and shape estimation. *CVPR*, 2020. 1, 2, 3, 4, 5, 6, 7, 8, 9, 10
- [16] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. *ICCV*, 2019. 1, 2, 3, 4, 5, 8, 9, 10
- [17] Nikos Kolotouros, Georgios Pavlakos, and Kostas Daniilidis. Convolutional mesh regression for single-image human shape reconstruction. *CVPR*, 2019. 1, 3, 5, 8
- [18] Matthew Loper, Naureen Mahmood, and Michael J. Black. Mosh: Motion and shape capture from sparse markers. *ACM TOG*, 2014. 7
- [19] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. SMPL: A skinned multi-person linear model. *ACM TOG*, 2015. 2, 3, 8
- [20] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 3D human motion estimation via motion compression and refinement. *ACCV*, 2020. 2, 3, 7, 8, 9, 10
- [21] 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. *3DV*, 2017. 5, 6, 7, 8, 9
- [22] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. Vnect: Real-time 3d human pose estimation with a single rgb camera. *ACM TOG*, 2017. 3
- [23] Gyeongsik Moon and Kyoung Mu Lee. I2L-MeshNet: Image-to-Lixel prediction network for accurate 3D human pose and mesh estimation from a single RGB image. *ECCV*, 2020. 1, 3, 8
- [24] Mohamed Omran, Christoph Lassner, Gerard Pons-Moll, Peter Gehler, and Bernt Schiele. Neural Body Fitting: Unifying deep learning and model based human pose and shape estimation. *3DV*, 2018. 3
- [25] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 5
- [26] Georgios Pavlakos, Nikos Kolotouros, and Kostas Daniilidis. TexturePose: Supervising human mesh estimation with texture consistency. *ICCV*, 2019. 3
- [27] Georgios Pavlakos, Luyang Zhu, Xiaowei Zhou, and Kostas Daniilidis. Learning to estimate 3D human pose and shape from a single color image. *CVPR*, 2018. 1, 3
- [28] Dario Pavullo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3D human pose estimation in video with temporal convolutions and semi-supervised training. *CVPR*, 2019. 3
- [29] Mir Rayat Imtiaz Hossain and James J Little. Exploiting temporal information for 3D human pose estimation. *ECCV*, 2018. 3
- [30] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 9
- [31] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. *IJCV*, 2015. 10
- [32] István Sárádi, Timm Linder, Kai O Arras, and Bastian Leibe. How robust is 3d human pose estimation to occlusion? *IROS workshop*, 2018. 5
- [33] Yu Sun, Yun Ye, Wu Liu, Wenpeng Gao, YiLi Fu, and Tao Mei. Human mesh recovery from monocular images via a skeleton-disentangled representation. *ICCV*, 2019. 3, 8

- [34] Gul Varol, Duygu Ceylan, Bryan Russell, Jimei Yang, Ersin Yumer, Ivan Laptev, and Cordelia Schmid. BodyNet: Volumetric inference of 3D human body shapes. *ECCV*, 2018. 3
- [35] Timo von Marcard, Roberto Henschel, Michael Black, Bodo Rosenhahn, and Gerard Pons-Moll. Recovering accurate 3d human pose in the wild using imus and a moving camera. *ECCV*, 2018. 5, 6, 7, 8, 9, 10
- [36] Ye Yuan and Kris Kitani. Ego-pose estimation and forecasting as real-time pd control. *ICCV*, 2019. 3
- [37] Jason Y Zhang, Panna Felsen, Angjoo Kanazawa, and Jitendra Malik. Predicting 3d human dynamics from video. *ICCV*, 2019. 3
- [38] Weiyu Zhang, Menglong Zhu, and Konstantinos G Derpanis. From actemes to action: A strongly-supervised representation for detailed action understanding. *ICCV*, 2013. 5, 6, 9
- [39] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. *CVPR*, 2019. 5