

Task-Generic Hierarchical Human Motion Prior using VAEs

Jiaman Li^{1,2}, Ruben Villegas³, Duygu Ceylan³, Jimei Yang³,
Zhengfei Kuang^{1,2}, Hao Li^{4,5}, Yajie Zhao²

¹University of Southern California ²USC Institute for Creative Technologies
³Adobe Research ⁴Pinscreen ⁵UC Berkeley

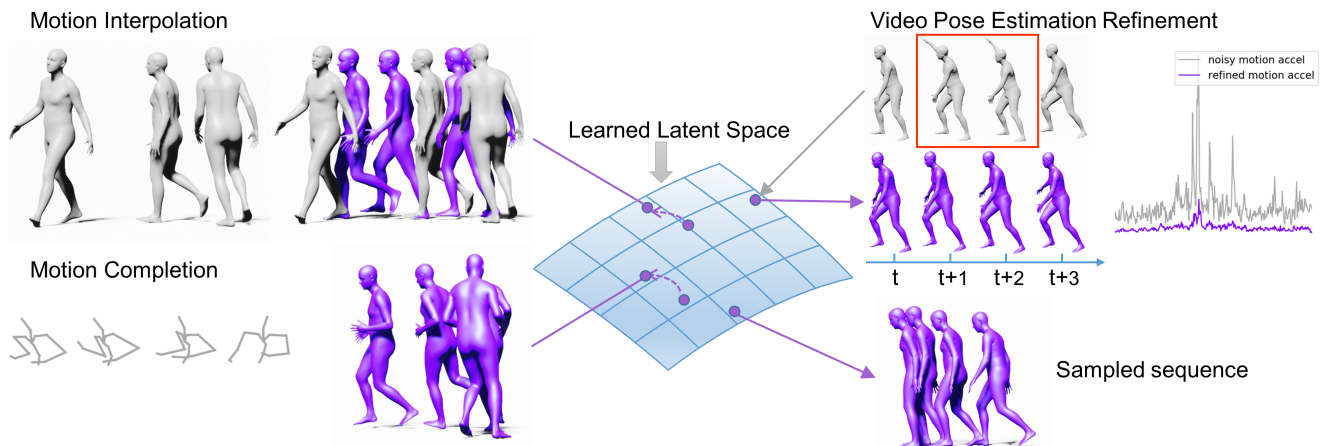


Figure 1. Our general purpose motion prior consists of a latent space of human motions and is learned using a hierarchical motion variational autoencoder (HM-VAE). Our approach is task-generic and can be directly adopted to a wide range of applications. *Left:* Motion interpolation and completion can be accomplished by traversing the latent space. *Right:* Noisy pose estimation can be refined by projecting noisy inputs into our latent space and decoding back. And a latent vector in the learned latent space is corresponding to a valid motion sequence.

Abstract

A deep generative model that describes human motions can benefit a wide range of fundamental computer vision and graphics tasks, such as providing robustness to video-based human pose estimation, predicting complete body movements for motion capture systems during occlusions, and assisting key frame animation with plausible movements. In this paper, we present a method for learning complex human motions independent of specific tasks using a combined global and local latent space to facilitate coarse and fine-grained modeling. Specifically, we propose a hierarchical motion variational autoencoder (HM-VAE) that consists of a 2-level hierarchical latent space. While the global latent space captures the overall global body motion, the local latent space enables to capture the refined poses of the different body parts. We demonstrate the effectiveness of our hierarchical motion variational autoencoder in a variety of tasks including video-based human pose estimation, motion completion from partial observations, and motion

synthesis from sparse key-frames. Even though, our model has not been trained for any of these tasks specifically, it provides superior performance than task-specific alternatives. Our general-purpose human motion prior model can fix corrupted human body animations and generate complete movements from incomplete observations.

1. Introduction

The modeling of human motions is a core component for many vision tasks, including pose estimation, action recognition, motion synthesis, and motion prediction. Several recent work have demonstrated new capabilities for generating complex body movements and capturing motion from unconstrained videos [23, 21, 55, 5, 8, 37, 20, 36]. While robustness and accuracy is constantly evolving for these methods, highly challenging scenes, occlusions, and body poses can still result in corrupted animations and noise.

Conventional techniques for reducing artifacts, include temporal filtering [19], inverse kinematics [51, 39, 10] and statistical human motion priors [48, 18, 25, 2]. While ef-

fective in reducing unwanted jitters and implausible poses, these methods do not generalize well to complex motions and the results are often inaccurate w.r.t. the ground truth.

To address this challenge, deep learning-based motion priors were proposed which are particularly effective in representing complex motion variations [16, 15, 23, 31]. These priors are generally designed for predetermined tasks, such as 3D pose estimation from a video, and a common problem is to be able to cover all possible input cases during training, such as occlusions, motion blur, etc. Ideally, we could build a prior model that describes the space of plausible human body movements, independently of the application and simply plug this model into any system. Training such model would simply consist of collecting high-quality motion capture data (task-generic), instead of fitting for example 3D models to an image (task-specific).

We introduce a generalized motion prior, that learns complex human body motions from high-fidelity motion capture data [32]. Similar to the work of [54] who developed a deep optimized prior for 3D modeling, our prior for motion is a general purpose one. We present a deep generative model based on a joint global and local latent space representation that can accurately capture the poses of different body parts while also modeling the global correlations across the body joints. Specifically, we adopt a two-level hierarchical motion variational autoencoder (HM-VAE) which maps the human motion to global and local latent spaces simultaneously. Our HM-VAE model adopts the recently proposed skeleton-aware architecture [1] and defines the global and local latent spaces via direct pooling and unpooling operations on the skeleton structure. While our HM-VAE successfully models the local human motion, we introduce an additional trajectory prediction component to model global motions. Taking local joint positions as input, our trajectory model estimates the root joint velocity at each timestep, enabling us to recover human motions in world space.

We show the generality and effectiveness of our human motion prior on various applications. First, we show that our task-generic model can refine human motions predicted from video [23, 21] by mapping noisy predictions into our motion prior latent space. We also demonstrate that our model can perform motion completion given partial observations (e.g., the upper body motion only) or motion synthesis given sparse keyframes. In both of these tasks, we optimize for both the global and local latent variables to match the partial observations and restore complete plausible motion sequences. While our model is not trained for any of these tasks specifically, it outperforms task-specific alternatives both qualitatively and quantitatively.

Our contributions are as follows. First, we present an effective task-generic motion prior model, that can improve the performance of a wide range of applications. Second,

we propose a two-level hierarchical motion variational autoencoder (HM-VAE) that consists of a skeleton-aware architecture, allowing it to accurately capture the local motion of body parts and the global correlation between them. Finally, we introduce a trajectory prediction module to model the global trajectory conditioned on the local body motion.

2. Related Work

Deep Learning Based Priors. The ability of deep neural networks to model data priors has sparked research in a variety of domains. Deep Image Prior (DIP) [43] shows that a generator network without any learning is an effective prior for image restoration. Given randomly initialized weights, the neural network is able to perform image denoising or super resolution via optimization defined by a task-dependent energy term and a regularizer. A similar idea is proposed and validated in the video domain [26], by training a network to mimic specific image operators in a single test video, the learned video prior is able to eliminate temporal inconsistencies in various video processing tasks. Besides discovering priors in the 2D domain, the ability to capture 3D priors is also demonstrated in recent works [12, 54]. Point2Mesh [12] randomly samples a fixed vector and optimizes the network parameters to reconstruct a mesh with geometric details and showcases the effectiveness of self-prior. Deep Optimized Priors [54] propose to learn a pre-trained prior first which then serves as initialization for optimizing both the latent vector and the decoder parameters given a task-specific objective and regularization loss. In this work, we investigate data-driven priors in the human motion domain and validate the effectiveness of our method by applying it to various human motion tasks without explicitly training for any specific tasks.

Generative Motion Modeling. With the recent success of learning based methods, several works have focused on generative models for motion synthesis. Martinez *et al.* [33] propose a recurrent neural network model for generating future human motion by predicting future joint velocities and adding them to previous joint positions. MT-VAE [53] propose a probabilistic recurrent neural network method for generating multiple future human motions. Aksan *et al.* [3] propose to predict future human motion by exploiting the kinematic structure in human bodies. Following the autoregressive generative model formulation, Motion Transformers [27] are introduced to model the future pose distribution along with a discrete pose representation, leveraging the advantage of the Transformer [45] architecture. Motion-VAE [29] models the future pose distribution given previous pose using a variational autoencoder (VAE) [22] approach. Normalizing flows is another category of generative models recently applied to human motion modeling. MoGlow [14] uses normalizing flows for motion modeling and achieve re-

alistic motion generation taking root trajectory as the conditioning signal. Recent work also address the problem of motion in-betweening from a generative modeling perspective. Long-term motion in-betweening [57] uses a generative adversarial neural network (GAN) [9] approach to generate human motion given sparse key-frames. Robust Motion in-betweening [13] employs an LSTM to generate a motion sequence given initial frames and end frames while also allowing for motion variations. Here, we focus on extending VAEs to model long-term human motion. Our key difference is to embed multiple frames of motions into a hierarchical global and local latent spaces. Methods like MotionVAE model motion in a per-frame basis while our method maps a full motion sequence into a compact latent space.

Motion Estimation From Videos. A multitude of optimization-based [30, 6, 17], learning-based [41, 20], and hybrid methods [24, 35] have been proposed to tackle the problem of single-image 3D human pose estimation. Its rapid progress has stimulated research interest on the long-standing problem of extracting 3D human motion from videos [4, 44, 50, 49, 11, 42, 34, 36, 40]. VIBE [27] uses an LSTM to capture temporal information and introduce a discriminator training strategy to ensure the predicted poses lie in a valid manifold. MEVA [31] presents a coarse-to-fine strategy where a valid motion sequence is first extracted conditioned on a latent vector and then refined utilizing person-specific details. TCMR [7] focuses on avoiding the temporal jitters that exist in the VIBE results and proposes a strategy to explicitly leverage past and future frames to achieve smoother results. Texture-based tracking [52] is shown to improve the motion stability during optimization. Foot contacts and physically-based models [38] are also used for estimating realistic human motions from videos. In this work, we are not aiming to design a specific 3D video pose estimation method. Instead, we show that our human motion prior is capable of eliminating jitters and noises that exist in the results of current state-of-the-art methods. We demonstrate that our method can be applied to any pose estimation methods and in our experiments we outperform previous work both quantitatively and qualitatively.

3. Hierarchical Motion VAEs

The core of our method is a hierarchical motion variational autoencoder (HM-VAE) that models human motion by jointly learning a local and global latent space. Specifically, given a motion sequence $\mathbf{x} \in R^{T \times J \times D}$, represented as the D dimensional joint rotations in a fixed time window of size T^1 , we first learn an embedding of \mathbf{x} into local and

¹In our experiments, we use the SMPL [30] skeleton hence the number of joints J is 24 and we use the continuous 6D rotation representation [56],

global latent spaces represented by latent codes z_l and z_g respectively. Assuming the latent space in the local and global levels are independent [28], we then model the probability distribution of a motion sequence as:

$$p(\mathbf{x}, z) = p(x|z_l, z_g)p(z_l)p(z_g). \quad (1)$$

Our variational autoencoder adopts the recently proposed skeleton-aware architecture [1] to facilitate learning over the humanoid skeleton structure directly. Before we discuss the details of our model, we first provide a brief overview of the skeleton-aware architecture. We refer the reader to the original paper for more details.

3.1. Background

The skeleton-aware architecture consists of three critical components that we adopt in our model design: skeleton convolution, skeleton pooling, and skeleton unpooling.

Skeleton Convolution. Given a motion sequence $\mathbf{x}, \mathbf{x} \in R^{T \times J \times D}$, we denote $\mathbf{y}, \mathbf{y} \in R^{T' \times J \times D'}$ as the updated features after a skeleton convolution operation. For each bone i in the skeleton, the updated feature is calculated as $\mathbf{y}_i = \frac{1}{|N_i^d|} \sum_{j \in N_i^d} \mathbf{x}_j * W_j^i + b_j^i$, where the symbol $*$ denotes a one dimensional temporal convolution operation with the temporal filter $W_j^i \in R^{k \times D \times D'}$ and bias $b_j^i \in R^{D'}$. D' represents the number of temporal filters, k represents the temporal kernel size, and N_i^d represents the neighboring bones of bone i within distance d . The distance between two bones (j_1, j_2) is defined as the number of bones needed to cross to reach j_2 starting from j_1 along the kinematic chain. The skeleton convolution operation preserves the number of edges J while downsampling the temporal dimension to T' .

Skeleton Pooling. Skeleton pooling merges the features of connected bones and extracts higher-level motion features by reducing the spatial resolution of the input. The pooling operation is applied to pairs of bones which are connected by a joint with degree of 2. For example, the thigh and calf which are connected by the knee. We recursively search such bone pairs starting from the root node (the hip), and merge their corresponding features using average pooling operation. As we perform pooling, the number of joints is reduced in subsequent layers of the network. Given disjoint sets of pooling bones denoted as $\{P(1), P(2), \dots, P(m)\}$, the pooling operation is defined as

$$F'_i = \text{Pool}(F_j | j \in P(i)).$$

Skeleton Unpooling. The unpooling operation mirrors skeleton pooling. Specifically, given the activation features F defined on a bone b obtained by merging the bones (i, j) ,

hence $D = 6$.

unpooling simply replaces b with the bones i and j where the new bone features are defined as $F_i = F, F_j = F$. The number of bones is increased after the unpooling layer.

3.2. Motion Prior Learning for Local Motion

Given a motion sequence $\mathbf{x}, \mathbf{x} \in R^{T \times J \times D}$, our HM-VAE consists of an encoder and a decoder as shown in Figure 2. The encoder learns the posterior distribution of the local, z_l , and global, z_g latent spaces given data \mathbf{x} :

$$q(z_l, z_g|x) = q(z_l|f_l(x))q(z_g|f_g(x)), \quad (2)$$

where $f_l(x), f_g(x)$ represent the motion features extracted from different layers in the encoder. Our VAE is then trained by maximizing the modified Evidence Lower Bound (ELBO) [22]:

$$\begin{aligned} \log p(x) \geq \mathbb{E}_{q(z_l, z_g|x)} [\log p(x|z_l, z_g)] - \\ \beta KL(q(z_l|f_l(x))||p(z_l)) - \\ \beta KL(q(z_g|f_g(x))||p(z_g)), \end{aligned} \quad (3)$$

where $q(z_l, z_g|x)$ is an encoder network that maps the input x into the local and global latent spaces, $p(x|z_l, z_g)$ is a decoder network that maps latent variables back into the input x , and $p(z_l)$ and $p(z_g)$ are assumed to be standard normal distributions $\mathcal{N}(0, I)$.

Encoder. The encoder consists of four building blocks B_1, B_2, B_3, B_4 where each building block is a combination of a skeleton convolution, skeleton pooling, and a LeakyReLU activation layer. As shown in Figure 2, we introduce a linear layer $W \in R^{T' \times d \times 2d_h}$ after B_1 and B_4 , mapping motion features of each corresponding block to a latent space. While the shallow layer features F_l after B_1 represent the local latent space, the deep layer features F_g after B_4 correspond to the global latent space. We enforce a normal distribution on each latent space:

$$z_l \sim \mathcal{N}(\mu_l(F_l), \sigma_l(F_l)), z_g \sim \mathcal{N}(\mu_g(F_g), \sigma_g(F_g)). \quad (4)$$

Decoder. The decoder has a symmetric architecture to the encoder. Each building block in decoder consists of temporal upsampling, skeleton unpooling, skeleton convolution and LeakyReLU activation layers. Given the latent codes z_l and z_g , the decoder first maps them to features through linear layers. Temporal upsampling and skeleton unpooling operations are used to increase the number of timesteps and bones gradually. The features obtained from the global latent code after a series of temporal upsampling, skeleton unpooling and convolution are concatenated with the features obtained from the local latent code. A final block of unpooling and convolution operations are used to reconstruct the

original motion sequence \mathbf{x} . We further add a forward kinematics layer proposed in [46] to convert \mathbf{x} into joint positions \mathbf{P} to define an additional joint position reconstruction loss. Also, we convert the 6D rotation representation to the rotation matrix \mathbf{R} and use an additional reconstruction loss defined on the rotation matrices. Overall, the reconstruction loss used to train the decoder is defined as:

$$L_{rec} = L_{6d} + L_{rot} + \lambda L_{joints}, \quad (5)$$

$$L_{6d} = \sum_{t=1}^T \|\mathbf{x}'_t - \mathbf{x}_t\|^2, \quad (6)$$

$$L_{rot} = \sum_{t=1}^T \|\mathbf{R}'_t - \mathbf{R}_t\|^2, \quad (7)$$

$$L_{joints} = \sum_{t=1}^T \|\mathbf{P}'_t - \mathbf{P}_t\|^2. \quad (8)$$

we experimentally set λ to 10 in our training process.

3.3. Trajectory Prediction

Given a motion sequence, we use the presented HM-VAE to model the local motion, i.e., the local joint rotations. In addition, we utilize a similar skeleton-aware architecture without reducing temporal dimension to model the global root joint trajectory. Specifically, given a sequence of joint positions denoted as $\mathbf{P} \in R^{T \times J \times 3}$, we apply four skeleton convolution layers with skeleton pooling layers to obtain motion features $F \in R^{T \times J' \times d}$. We use a linear layer that takes F as input and estimates the root velocity $V \in R^{T \times 3}$. By accumulating the root velocity in subsequent frames, we compute the global root trajectory $G \in R^{T \times 3}$. The root trajectory at any particular time t is defined as $G_t = \sum_{i=0}^t V_i$. We train the trajectory estimation module with a loss function that consists of both velocity and trajectory terms:

$$L_{traj} = \sum_{t=1}^T \|V'_t - V_t\|^2 + \|G'_t - G_t\|^2 \quad (9)$$

4. Application

Our HM-VAE provides a generalized motion prior that can be applicable in various tasks like 3D video pose estimation, motion interpolation and motion completion. In this section, we introduce the applications we consider and describe the effective strategy used for each application. We provide qualitative and quantitative results in the next section.

3D Video Pose Estimation. Our learned motion prior provides an effective strategy to refine video based pose estimations. Concretely, we take potentially noisy pose estimates as input to the encoder, then decode refined poses

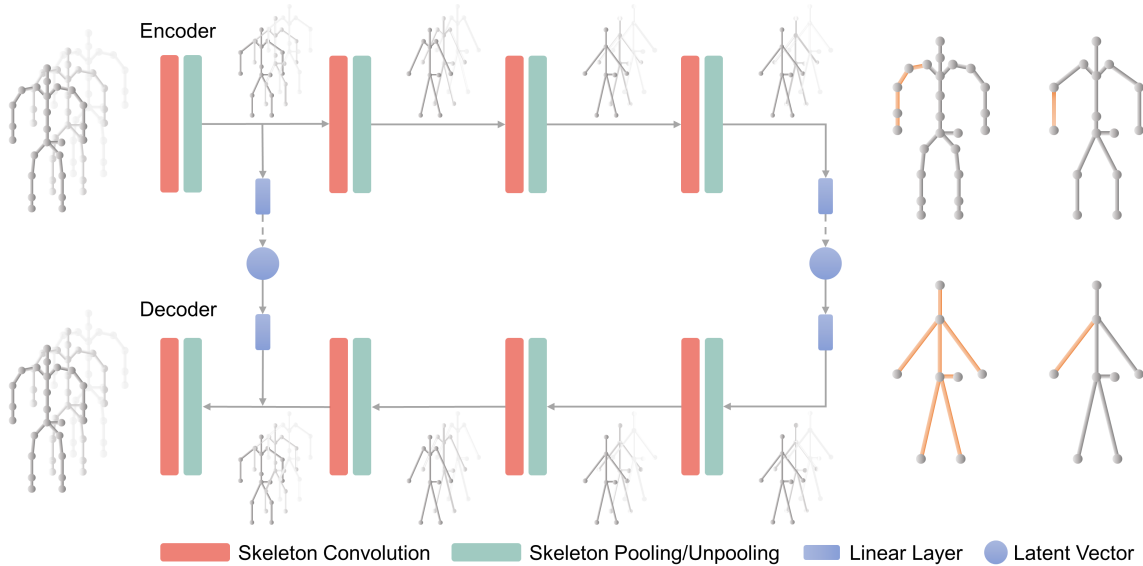


Figure 2. Model Overview. *Left*: Model architecture. (We omit activation layers and temporal upsampling layers here for simplicity.) *Right*: Illustration of receptive field in shallow (B_1) and deep layers (B_4).

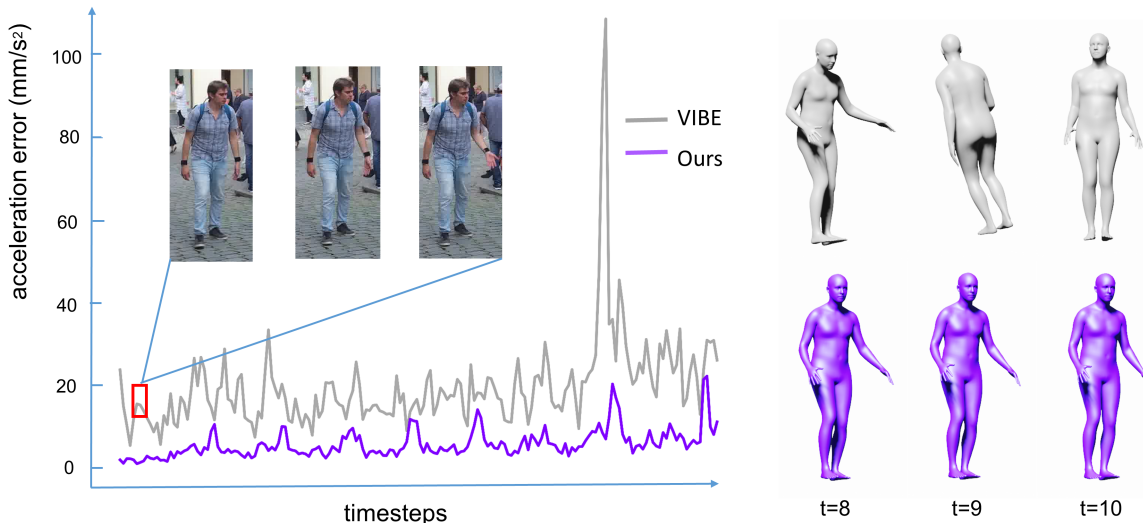


Figure 3. Acceleration error curves for VIBE [23] results and our refined results. The right figure shows poses in consecutive timesteps corresponding to the reference images on the left.

using the encoded latent vector. Our HM-VAE is designed for a fixed window size of T frames. In order to have our method process sequences of any length, we could simply partition the input sequence into windows of T frames. However, with no overlap across the time windows, we observe that this may result in discontinuities. Therefore, we propose a sliding window strategy using center frames to process arbitrarily long sequences. Specifically, for each time window we process, we only update the pose of the center frame with the refined result and shift the time window one step. We take the $\frac{T}{2}$ th frame $M_{\frac{T}{2}}$ as the refined final result, added to our final refined sequence S . And the window is shifted by one timestep along the input motion sequence for processing the next window. For each window

size of pose sequences, we formally define the process as follows, where W_2 represents next window.

$$z = Enc(N_1, \dots, N_T), \quad (10)$$

$$M_1, \dots, M_{\frac{T}{2}}, \dots, M_T = Dec(z), \quad (11)$$

$$S = S \cup M_{\frac{T}{2}} \quad (12)$$

$$W_2 = N_2, \dots, N_{T+1} \quad (13)$$

Motion Interpolation and Completion. A common setup in motion synthesis is to generate motion sequences given a sparse set of keyframes, which we refer to as motion interpolation. Motion completion, on the other hand,

focuses on synthesizing complete body motion from partial observations, e.g., completing the motion of the lower body by observing the upper body. For both motion interpolation and completion tasks, we simply utilize the decoder of HM-VAE to synthesize motion while searching for an optimal latent code to match the given observations (i.e., sparse keyframes or partial body motion). The optimization objective is to minimize the reconstruction error between the given observations and the corresponding decoded poses. We define the reconstruction objective as a combination of three terms including matching the joint rotations using both 6D rotation and rotation matrix representations and matching the joint positions after forward kinematics:

$$L_{rec} = L_{6d} + L_{rot} + \lambda_1 L_{joints} \quad (14)$$

Concretely, we perform optimization in two phases. Starting with randomly sampled latent vectors z_l and z_g , in the first phase, we optimize for the latent vectors that minimize L_{rec} as the only objective. The decoder parameters θ are fixed during this phase. In the next phase, we optimize for the decoder parameters θ' [54] while keeping the latent vectors fixed. In this second phase, we introduce a regularization loss to constrain θ' and prevent it from deviating too much from the pre-trained parameters θ . Thus, the optimization objective in the second phase becomes:

$$L_{opt} = L_{rec} + \lambda_2 \|\theta' - \theta\|^2 \quad (15)$$

5. Experiments

In this section, we first describe the dataset we use for training and evaluation. Then we showcase the results of applying our HM-VAE in the applications we introduced in the previous section. Finally, we perform an ablation study to validate the effectiveness of our overall approach.

Dataset. We use the AMASS dataset [32] for training HM-VAE. AMASS dataset is a large collection of 15 motion capture datasets with a unified data representation. The dataset has more than 40 hours of motion data and serves as a great testbed for motion modeling. We use the same validation and testing split introduced in VIBE [23]. For refining video based pose estimates, we use 3DPW [47], a 3D motion in the wild dataset, as our test set. For the motion interpolation task, we train our HM-VAE on the LAFAN1 dataset [13] to provide quantitative comparisons to the baseline methods. LAFAN1 consists of high-quality motion capture data with specific action types. We follow the data split proposed in [13] and use subjects 1, 2, 3 and 4 as training, and subject 5 for testing.

Implementation Details. We use a batch size of 8 for training. The KL divergence weight β is set to 0.003. Un-

	PA-MPJPE	MPJPE	ACCEL	ACCER
HD [21]	72.17	115.97	14.96	14.73
HD [21] w Prior	71.39	113.90	5.21	8.36
VIBE [23]	56.56	93.59	27.12	27.99
VIBE [23] w Prior	55.84	92.43	6.03	9.15

Table 1. 3D Video Human Pose Estimation Results in 3DPW Testing Dataset.

less noted otherwise, we train HM-VAE with motion sequences of length 64. While training our HM-VAE, to prevent the learning dominated by either shallow or deep latent vector, we use similar strategy proposed in [28]. We first only train our model with deep latent vector, then start training both shallow and deep latent vectors after 50000 iterations. For the motion interpolation experiments, we found our method converged at around 150 iterations of optimization, with 50 iterations for the first phase and 100 iterations for the second phase. For the motion completion experiments, we found our optimization converged at around 300 iterations with 100 iterations belonging to the first phase.

5.1. Results

3D Video Pose Estimation. In this section, we show that our model can be used to refine the results of off-the-shelf 3D video pose estimation methods. In order to adapt HM-VAE to different global rotations and frame rate among different datasets, we train our HM-VAE with data augmentation. Our data augmentation consists of different frame rates and random global rotations. Also, we use the HM-VAE model trained with a window size 8 in this application which we observe has a better reconstruction quality.

We show quantitative results in Table 1 where we test our method with inputs obtained by both VIBE [23] and HumanDynamics (HD) [21]. We report errors using the same metrics as VIBE [23]. Specifically, we report the mean per joint position error with (PA-MPJPE) and without (MPJPE) the Procrustes-alignment, as well as the mean per joint acceleration and acceleration error. In Figure 3, we show the acceleration error curves as well as example poses obtained for consecutive timesteps. Compared to current state-of-the-art methods, our refined motions have smaller acceleration errors. While previous approaches are prone to abrupt changes across consecutive poses as shown in Figure 3, our model smooths out these noisy estimates. We refer the readers to the supplementary video for a detailed comparison.

Motion Interpolation. In order to demonstrate the effectiveness of HM-VAE for the motion interpolation task, we compare it to appropriate baseline methods. Specifically, in order to interpolate local joint rotations, we use the standard spherical linear interpolation (Slerp). Since interpolation quality is directly related to the number of missing frames, we perform our evaluations in four settings where 5, 15, 30, 45 frames are missing in each setting. Following the same setting as [13], given the starting 10 frames and

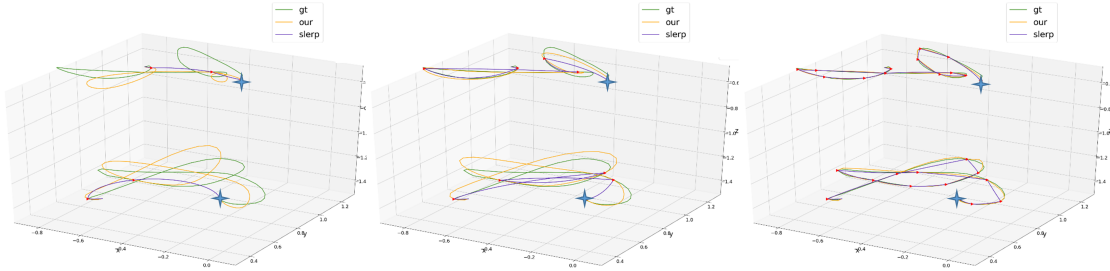


Figure 4. Local trajectory comparison for motion interpolation in AMASS data. From left to right is the trajectory when key frame interval is 30, 15, 5 respectively. The upper curves represent the left wrist, the lower curves represent the right ankle. The star symbol represents the starting point, the arrow symbol represent the position of key frames. Our results show similar moving patterns to ground truth, while Slerp differs a lot when key frame interval is large.

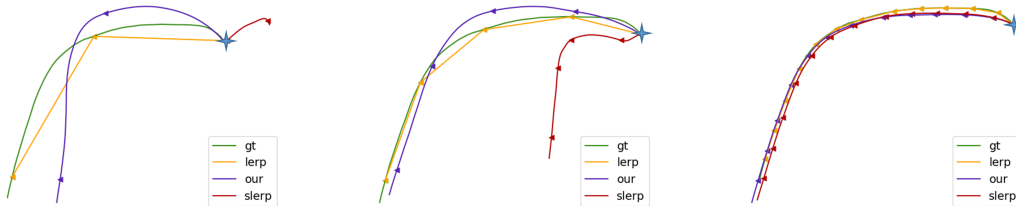


Figure 5. Global root trajectory comparison for motion interpolation in AMASS data. From left to right is the trajectory (in xy plane) when key frame interval is 30, 15, 5 respectively. The star symbol represents the starting point, the arrow symbol represents the position of key frames.

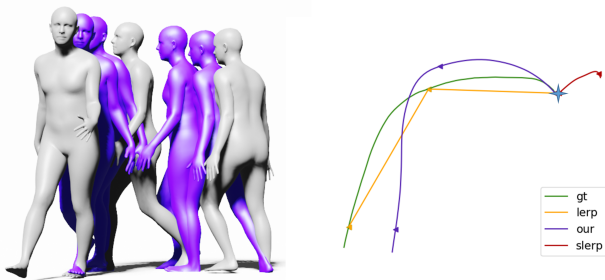


Figure 6. Motion interpolation results in AMASS test data. The gray mesh shows key frame poses, the purple mesh show the generated poses. The interval between two key frames is 30 frames. The right figure shows the global trajectory comparison for this motion sequence.

ending 1 frame as key frames, we aim to generate the frames in-between. We show quantitative comparisons in terms of local pose estimation in Table 2. In addition to the metrics introduced before, we also report the global quaternion loss proposed by the original benchmark [13]. We show that our method outperforms the Slerp baseline quantitatively. We also show that the performance achieved by our human motion prior is competitive against the in-betweening specific method from [13] in the global quaternion loss metric. Please note that we use the LAFAN1 dataset for this evaluation to compare against the global quaternion errors directly reported by [13] since their code is not published and the au-

thors were not able to run their model on our dataset. We also provide additional qualitative results in the AMASS dataset. We visualize local joint trajectories in Figure 4 for a walking motion sequence. Our results preserve the original motion patterns while Slerp fails to model the local motion when the interval between two key frames becomes large. We further demonstrate global trajectory interpolation with our method and the alternatives. Specifically, we use our global trajectory estimation module by providing the local motion predicted by our method as well as Slerp. In addition, we also define a simple baseline where we linearly interpolate the global root position of the sparse keyframes (lerp). As shown in Figure 5, the trajectory estimated by our method more closely resembles the ground truth. We also show a mesh visualization result for motion interpolation in Figure 6. For more qualitative results, we encourage readers to check our accompanying video.

Motion Completion. Given only upper body joint rotations as target, we aim to recover the complete body motion sequences. For this experiment, we use motion sequences from the testing and validation split of the AMASS dataset. As shown in Figure 7, our approach is able to restore complete motions since the global latent space capture the correlations among different joints. Therefore, the missing lower legs movement that matches the given upper body is retrieved from the learned latent space for human motion.

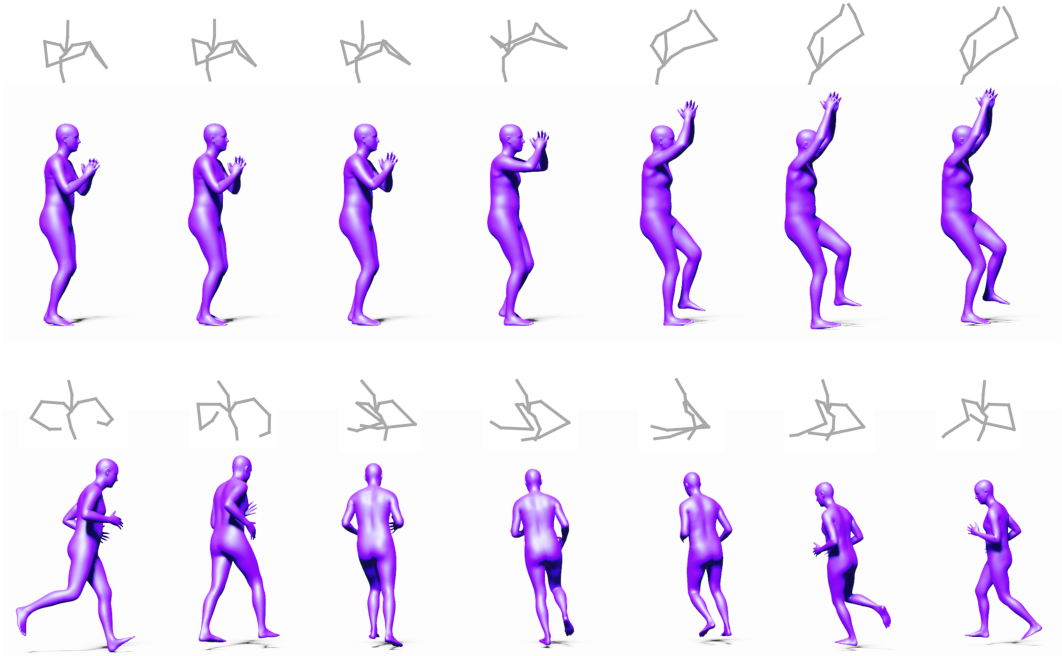


Figure 7. Motion Completion Results. Given upper body joint rotation as optimization objective, the prior model can complete whole motion sequences.

	5	15	30	45
MPJPE-Slerp	16.02	57.13	96.54	118.96
MPJPE-Ours	14.08	45.09	90.41	117.93
PAMPJPE-Slerp	15.82	54.11	83.66	92.4
PAMPJPE-Ours	12.03	38.37	72.21	86.06
ACCEL-Slerp	1.75	0.78	0.35	0.23
ACCEL-Ours	4.79	4.48	4.08	3.61
ACCER-Slerp	5.98	5.97	6.05	6.06
ACCER-Ours	5.31	5.83	6.54	6.75
Global Quat-Slerp	0.22	0.62	0.98	1.25
Global Quat-[13]	0.17	0.42	0.69	0.94
Global Quat-Ours	0.24	0.54	0.94	1.25

Table 2. Quantitative Evaluation for Motion Interpolation in LAFAN1 Dataset.

5.2. Ablation Study

In order to motivate the design choices we made, we perform an ablation study where we compare our HM-VAE with a non-hierarchical motion VAE (M-VAE) and a VAE with only temporal convolution layers (TCN-VAE). The temporal convolution layers were used in training an autoencoder for motion processing [16, 15]. We compare our model to the alternatives in the task of motion reconstruction using the AMASS dataset. Specifically, for each testing sequence, we take the local joint rotations as input to the encoder and then decode the motion from the mean vector. We measure the mean joint reconstruction error as shown in Table 3. Our HM-VAE model outperforms the M-VAE by a large margin in motion reconstruction evaluation. And the model with skeleton-aware architecture has superior performance than its temporal convolution counterpart. Therefore, we show that skeleton operations from the skeleton-aware architecture are important for

	PA-MPJPE	MPJPE	ACCEL	ACCER
TCN-VAE	87.27	103.60	1.66	6.46
M-VAE	59.71	74.34	2.36	6.15
HM-VAE	45.82	58.46	2.29	5.98

Table 3. Motion Reconstruction Results in AMASS test data.

modeling the human body structure in comparison to using standard temporal convolution. Moreover, modeling a global and local motion latent spaces further improve the human motion modeling power of the skeleton-aware architecture.

6. Conclusion

We propose a task-generic motion prior using a hierarchical motion VAE. We demonstrate the effectiveness of the prior in various applications including 3D video pose estimation, motion interpolation, and motion completion. By learning a global and local embedding, our prior can faithfully model human motion. While our prior enables to refine video-based human motion estimation results by reducing jitters, it also performs on-par with task specific methods for motion interpolation and completion. There are some limitations of our method we would like to address in future work. We observe that there are accumulation of errors when predicting the global trajectory for a long sequence. Exploring more constraints like foot contact during both training and inference might be a potential approach to address this. While we show that our prior is effective in different applications, using few-shot learning to better adapt to specific tasks is another interesting direction. Finally, incorporating certain physical properties and action conditions are also promising directions.

Acknowledgements. This research was conducted at Adobe, at University of Southern California and USC Institute for Creative Technologies. Research was sponsored by the Army Research Office and was supported under Cooperative Agreement Number W911NF-20-2-0053, and sponsored by the U.S. Army Research Laboratory (ARL) under contract number W911NF-14-D-0005, the CONIX Research Center, one of six centers in JUMP, a Semiconductor Research Corporation (SRC) program sponsored by DARPA; and in part by the ONR YIP grant N00014-17-S-FO14. Affiliation with Pinscreen and the University of California at Berkeley was supported by DARPA under cooperative agreement HR00112020054. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation.

References

- [1] Kfir Aberman, Peizhuo Li, Dani Lischinski, Olga Sorkine-Hornung, Daniel Cohen-Or, and Baoquan Chen. Skeleton-aware networks for deep motion retargeting. *SIGGRAPH*, 2020. [2](#), [3](#)
- [2] Ijaz Akhter and Michael J Black. Pose-conditioned joint angle limits for 3d human pose reconstruction. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1446–1455, 2015. [1](#)
- [3] Emre Aksan, Manuel Kaufmann, and Otmar Hilliges. Structured prediction helps 3d human motion modelling. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2019. First two authors contributed equally. [2](#)
- [4] Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. Scape: shape completion and animation of people. In *ACM SIGGRAPH 2005 Papers*, pages 408–416. 2005. [3](#)
- [5] Anurag Arnab, Carl Doersch, and Andrew Zisserman. Exploiting temporal context for 3d human pose estimation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3395–3404, 2019. [1](#)
- [6] Federica Bogo, Angjoo Kanazawa, Christoph Lassner, Peter Gehler, Javier Romero, and Michael J Black. Keep it simple: Automatic estimation of 3d human pose and shape from a single image. In *European Conference on Computer Vision*, pages 561–578. Springer, 2016. [3](#)
- [7] Hongsuk Choi, Gyeongsik Moon, and Kyoung Mu Lee. Beyond static features for temporally consistent 3d human pose and shape from a video. *arXiv e-prints*, pages arXiv:2011.2020. [3](#)
- [8] Rishabh Dabral, Anurag Mundhada, Uday Kusupati, Safeer Afaque, Abhishhek Sharma, and Arjun Jain. Learning 3d human pose from structure and motion. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 668–683, 2018. [1](#)
- [9] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *arXiv preprint arXiv:1406.2661*, 2014. [3](#)
- [10] Keith Grochow, Steven L Martin, Aaron Hertzmann, and Zoran Popović. Style-based inverse kinematics. In *ACM SIGGRAPH 2004 Papers*, pages 522–531. 2004. [1](#)
- [11] Riza Alp Guler and Iasonas Kokkinos. Holopose: Holistic 3d human reconstruction in-the-wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10884–10894, 2019. [3](#)
- [12] Rana Hanocka, Gal Metzger, Raja Giryes, and Daniel Cohen-Or. Point2mesh: A self-prior for deformable meshes. *arXiv preprint arXiv:2005.11084*, 2020. [2](#)
- [13] Félix G Harvey, Mike Yurick, Derek Nowrouzezahrai, and Christopher Pal. Robust motion in-betweening. *ACM Transactions on Graphics (TOG)*, 39(4):60–1, 2020. [3](#), [6](#), [7](#), [8](#)
- [14] Gustav Eje Henter, Simon Alexanderson, and Jonas Beskow. Moglow: Probabilistic and controllable motion synthesis using normalising flows. *arXiv preprint arXiv:1905.06598*, 2019. [2](#)
- [15] Daniel Holden, Jun Saito, and Taku Komura. A deep learning framework for character motion synthesis and editing. *ACM Transactions on Graphics (TOG)*, 35(4):1–11, 2016. [2](#), [8](#)
- [16] Daniel Holden, Jun Saito, Taku Komura, and Thomas Joyce. Learning motion manifolds with convolutional autoencoders. In *SIGGRAPH Asia 2015 Technical Briefs*, pages 1–4. 2015. [2](#), [8](#)
- [17] Yinghao Huang, Federica Bogo, Christoph Lassner, Angjoo Kanazawa, Peter V Gehler, Javier Romero, Ijaz Akhter, and Michael J Black. Towards accurate marker-less human shape and pose estimation over time. In *2017 international conference on 3D vision (3DV)*, pages 421–430. IEEE, 2017. [3](#)
- [18] Leslie Ikemoto, Okan Arıkan, and David Forsyth. Generalizing motion edits with gaussian processes. *ACM Transactions on Graphics (TOG)*, 28(1):1–12, 2009. [1](#)
- [19] Rudolph Emil Kalman. A new approach to linear filtering and prediction problems. *Transactions of the ASME—Journal of Basic Engineering*, 82(Series D):35–45, 1960. [1](#)
- [20] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. End-to-end recovery of human shape and pose. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7122–7131, 2018. [1](#), [3](#)
- [21] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. Learning 3d human dynamics from video. In *CVPR*, 2019. [1](#), [2](#), [6](#)
- [22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *ICLR*, 2014. [2](#), [4](#)
- [23] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. Vibe: Video inference for human body pose and shape estimation. In *CVPR*, 2020. [1](#), [2](#), [5](#), [6](#)
- [24] 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 *Proceedings of the IEEE International Conference on Computer Vision*, pages 2252–2261, 2019. [3](#)
- [25] Andreas M Lehrmann, Peter V Gehler, and Sebastian Nowozin. Efficient nonlinear markov models for human motion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1314–1321, 2014. [1](#)
- [26] Chenyang Lei, Yazhou Xing, and Qifeng Chen. Blind video temporal consistency via deep video prior. *Advances in Neural Information Processing Systems*, 33, 2020. [2](#)
- [27] Jiaman Li, Yihang Yin, Hang Chu, Yi Zhou, Tingwu Wang, Sanja Fidler, and Hao Li. Learning to generate diverse dance motions with transformer. *arXiv preprint arXiv:2008.08171*, 2020. [2](#), [3](#)

- [28] Zhiyuan Li, Jaideep Vitthal Murkute, Prashna Kumar Gyawali, and Linwei Wang. Progressive learning and disentanglement of hierarchical representations. *arXiv preprint arXiv:2002.10549*, 2020. 3, 6
- [29] Hung Yu Ling, Fabio Zinno, George Cheng, and Michiel Van De Panne. Character controllers using motion vaes. *ACM Transactions on Graphics (TOG)*, 39(4):40–1, 2020. 2
- [30] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. *ACM transactions on graphics (TOG)*, 34(6):1–16, 2015. 3
- [31] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 3d human motion estimation via motion compression and refinement. In *Proceedings of the Asian Conference on Computer Vision*, 2020. 2, 3
- [32] Naureen Mahmood, Nima Ghorbani, Nikolaus F Troje, Gerard Pons-Moll, and Michael J Black. Amass: Archive of motion capture as surface shapes. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5442–5451, 2019. 2, 6
- [33] Julieta Martinez, Michael J. Black, and Javier Romero. On human motion prediction using recurrent neural networks. In *CVPR*, 2017. 2
- [34] 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 Transactions on Graphics (TOG)*, 36(4):1–14, 2017. 3
- [35] Georgios Pavlakos, Vasileios Choutas, Nima Ghorbani, Timo Bolkart, Ahmed AA Osman, Dimitrios Tzionas, and Michael J Black. Expressive body capture: 3d hands, face, and body from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 10975–10985, 2019. 3
- [36] Dario Pavlo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7753–7762, 2019. 1, 3
- [37] Mir Rayat Imtiaz Hossain and James J Little. Exploiting temporal information for 3d human pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 68–84, 2018. 1
- [38] Davis Rempe, Leonidas J Guibas, Aaron Hertzmann, Bryan Russell, Ruben Villegas, and Jimei Yang. Contact and human dynamics from monocular video. In *European Conference on Computer Vision*, pages 71–87. Springer, 2020. 3
- [39] Charles F Rose III, Peter-Pike J Sloan, and Michael F Cohen. Artist-directed inverse-kinematics using radial basis function interpolation. In *Computer Graphics Forum*, volume 20, pages 239–250. Wiley Online Library, 2001. 1
- [40] Mingyi Shi, Kfir Aberman, Andreas Aristidou, Taku Komura, Dani Lischinski, Daniel Cohen-Or, and Baoquan Chen. Motionet: 3d human motion reconstruction from monocular video with skeleton consistency. *ACM Transactions on Graphics (TOG)*, 40(1):1–15, 2020. 3
- [41] Jun Kai Vince Tan, Ignas Budvytis, and Roberto Cipolla. Indirect deep structured learning for 3d human body shape and pose prediction. 2017. 3
- [42] Bugra Tekin, Isinsu Katircioglu, Mathieu Salzmann, Vincent Lepetit, and Pascal Fua. Structured prediction of 3d human pose with deep neural networks. *arXiv preprint arXiv:1605.05180*, 2016. 3
- [43] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9446–9454, 2018. 2
- [44] Raquel Urtasun, David J Fleet, and Pascal Fua. Temporal motion models for monocular and multiview 3d human body tracking. *Computer vision and image understanding*, 104(2-3):157–177, 2006. 3
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 2
- [46] Ruben Villegas, Jimei Yang, Duygu Ceylan, and Honglak Lee. Neural kinematic networks for unsupervised motion retargetting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8639–8648, 2018. 4
- [47] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 601–617, 2018. 6
- [48] Jack M Wang, David J Fleet, and Aaron Hertzmann. Gaussian process dynamical models for human motion. *IEEE transactions on pattern analysis and machine intelligence*, 30(2):283–298, 2007. 1
- [49] Xiaolin Wei and Jinxiang Chai. Videomocap: Modeling physically realistic human motion from monocular video sequences. In *ACM SIGGRAPH 2010 papers*, pages 1–10. 2010. 3
- [50] Xiaolin K Wei and Jinxiang Chai. Modeling 3d human poses from uncalibrated monocular images. In *2009 IEEE 12th International Conference on Computer Vision*, pages 1873–1880. IEEE, 2009. 3
- [51] Chris Welman. *Inverse kinematics and geometric constraints for articulated figure manipulation*. PhD thesis, Theses (School of Computing Science)/Simon Fraser University, 1993. 1
- [52] Donglai Xiang, Hanbyul Joo, and Yaser Sheikh. Monocular total capture: Posing face, body, and hands in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10965–10974, 2019. 3
- [53] Xinchen Yan, Akash Rastogi, Ruben Villegas, Kalyan Sunkavalli, Eli Shechtman, Sunil Hadap, Ersin Yumer, and Honglak Lee. Mt-vae: Learning motion transformations to generate multimodal human dynamics. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 2
- [54] Mingyue Yang, Yuxin Wen, Weikai Chen, Yongwei Chen, and Kui Jia. Deep optimized priors for 3d shape modeling and reconstruction. *arXiv preprint arXiv:2012.07241*, 2020. 2, 6
- [55] Andrei Zanfir, Eduard Gabriel Bazavan, Hongyi Xu, Bill Freeman, Rahul Sukthankar, and Cristian Sminchisescu. Weakly supervised 3d human pose and shape reconstruction with normalizing flows. *ECCV*, 2020. 1
- [56] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5745–5753, 2019. 3

- [57] Yi Zhou, Jingwan Lu, Connelly Barnes, Jimei Yang, Sitao Xiang, et al. Generative tweening: Long-term inbetweening of 3d human motions. *arXiv preprint arXiv:2005.08891*, 2020. [3](#)