

Anatomy-aware 3D Human Pose Estimation with Bone-based Pose Decomposition

Tianlang Chen, Chen Fang, Xiaohui Shen, Yiheng Zhu, Zhili Chen, and Jiebo Luo, *Fellow, IEEE*

Abstract—In this work, we propose a new solution to 3D human pose estimation in videos. Instead of directly regressing the 3D joint locations, we draw inspiration from the human skeleton anatomy and decompose the task into bone direction prediction and bone length prediction, from which the 3D joint locations can be completely derived. Our motivation is the fact that the bone lengths of a human skeleton remain consistent across time. This promotes us to develop effective techniques to utilize global information across *all* the frames in a video for high-accuracy bone length prediction. Moreover, for the bone direction prediction network, we propose a fully-convolutional propagating architecture with long skip connections. Essentially, it predicts the directions of different bones hierarchically without using any time-consuming memory units (e.g. LSTM). A novel joint shift loss is further introduced to bridge the training of the bone length and bone direction prediction networks. Finally, we employ an implicit attention mechanism to feed the 2D keypoint visibility scores into the model as extra guidance, which significantly mitigates the depth ambiguity in many challenging poses. Our full model outperforms the previous best results on Human3.6M and MPI-INF-3DHP datasets, where comprehensive evaluation validates the effectiveness of our model.

Index Terms—3D pose, Bone, Length, Direction, Long skip connections.

I. INTRODUCTION

3D human pose estimation in videos has been widely studied in recent years. It has extensive applications in action recognition, sports analysis and human-computer interaction. Current state-of-the-art approaches [1], [2], [3] typically decompose the task into 2D keypoint detection followed by 3D pose estimation. Given an input video, they first detect the 2D keypoints of each frame, and then predict the 3D joint locations of a frame based on the 2D keypoints.

When estimating the 3D joint locations from 2D keypoints, the challenge is to resolve depth ambiguity, as multiple 3D poses with different joint depths can be projected to the same 2D keypoints. Exploiting temporal information from the video has been demonstrated to be effective for reducing this ambiguity. Typically, to predict the 3D joint locations of a frame in a video, recent approaches [1], [4], [5] utilize temporal networks that additionally feed the adjacent frames'

2D keypoints as input. These approaches consider the adjacent local frames most associated with the current frame, and extract their information as extra guidance. However, such approaches are limited to exploiting information only from the neighboring frames. Given a 1-minute input video with a frame rate of 50, even though we choose the existing temporal network with largest temporal window size (i.e 243 frames) [1], it is limited to using a concentrated short segment (about one-twelfth length of the video) to predict a single frame. Such a design can easily make existing temporal networks fail when the current frame and its adjacent input frames correspond to a complex pose, because none of the input frames provide reliable and high-confidence information to the networks.

Considering this, our first contribution is proposing a novel approach that can effectively capture the knowledge from both local and distant frames to estimate the 3D joint locations of the current frame, by cleverly exploiting the anatomic properties of the human skeleton. We refer to it as *anatomy awareness*. Specifically, based on the anatomy of the human skeleton, we decompose the task of 3D joint location prediction into two sub-tasks – bone direction prediction and bone length prediction. We demonstrate that the combination of the two new tasks are essentially equivalent to the original task. The motivation is based on the fact that the bone lengths of a person remain consistent in a video over time (This can be verified by 3D human pose datasets such as Human3.6M and MPI-INF-3DHP). Hence, when we predict the bone lengths of a particular frame, we can leverage the frames distributed over the duration of the entire video for more accurate and smooth prediction. Note that although Sun et al. [6] transform the task into a generic bone-based representation, such a generic representation does not allow them to utilize that critical bone length consistency. In contrast, we decompose the task *explicitly into bone direction and bone length prediction*. We demonstrate that this explicit design leads to significant advantages over either the generic representation design in [6] or imposing a bone length consistency loss across frames.

However, it is nontrivial to implement this explicit design. One problem for training the proposed bone length prediction network is that the training dataset typically contains only a few skeletons. For example, the training set of Human3.6M contains 5 actors corresponding to 5 bone length settings. Directly training the network on the data from the 5 actors leads to serious overfitting. Therefore, we adopt the fully-connected residual network for bone length prediction and propose two effective mechanisms to prevent overfitting via a network design and data augmentation.

As for the bone directions, we adopt the temporal convolu-

T. Chen, and J. Luo are with the Department of Computer Science, University of Rochester, Rochester, NY, 14627 USA (e-mail: {zyang39, tusharku, tchen45, jluo}@cs.rochester.edu).

C. Fang, X. Shen, Y. Zhu and Z. Chen are with Bytedance AI Lab, Mountain View, CA, USA (e-mail: {fangchen, shenxiaohui, yiheng.zhu, zhili.chen}@bytedance.com).

Corresponding author: Jiebo Luo.

Copyright©2021 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending an email to pubs-permissions@ieee.org.

tional network in [1] to predict the direction of each bone in the 3D space for each frame. Motivated by [5], we believe it is beneficial to predict the directions of different bones hierarchically, instead of all at once as in [1]. Following the human skeleton anatomy, the directions of simple torso bones (*e.g.* lumbar vertebra) with less motion variation should be predicted first, and then guide the prediction of challenging limb bones (*e.g.* arms and legs). This strategy is applied straightforwardly by a recurrent neural network (RNN) with different joints predicted step by step in [5] for a single frame. However, the high computation complexity of RNN precludes the network from holding a large temporal window which has been shown to improve performance. To solve this issue, based on [1], we further propose a high-performance fully-convolutional propagating architecture, which contains multiple sub-networks with each predicting the directions of all the bones. The hierarchical prediction is implicitly performed via long skip connections between adjacent sub-networks.

Additionally, motivated by [6], we create an effective joint shift loss for the two sub-tasks (*i.e.*, bone direction prediction and bone length prediction) to learn jointly. The joint shift loss penalizes the relative joint shift between all long-range joint pairs, for example the left hand and right foot. Thus, it provides an extra strong supervision for the two networks to be trained to coordinate with each other and produce robust predictions.

Last but not least, we propose a simple yet effective approach to further reduce the depth ambiguity. Specifically, we incorporate 2D keypoint visibility scores into the model as a new feature, which indicates the probability of each 2D keypoint being visible in a frame and provides extra knowledge of the depth relation between specific joints. We argue that the scores are useful to those poses with body parts occluded or when the relative depth matters. For example, if a person keeps her/his hands in front of the chest in a frontal view, our model will be confused on whether the hands are in front of the chest (visible) or behind the back (occluded), since the occluded 2D keypoints can still be predicted sometimes. Furthermore, We adopt an implicit attention mechanism to dynamically adjust the importance of the visibility scores for better performance.

Our contributions are summarized as follows:

- We explicitly decompose the task of 3D joint estimation into bone direction prediction and bone length prediction. As such, the bone length prediction branch can fully utilize frames across the entire video.
- We propose a new fully-convolutional architecture for hierarchical bone direction prediction.
- We propose a high-performance bone length prediction network, two mechanisms are created to effectively prevent overfitting.
- We feed the visibility scores of 2D keypoint detection into the model to better resolve the depth ambiguity.
- Our model is inspired by the human skeleton anatomy and achieves the state-of-the-art performance on Human3.6M and MPI-INF-3DHP datasets.

II. RELATED WORK

3D human pose estimation has received much attention in recent years. To predict the 3D joint location from 2D image input, previous works of 3D pose estimation typically fall into two categories based on the training pipeline. For the approaches of the first category, they created an end-to-end convolutional neural network (CNN) model to directly predict the 3D joint location from the original input images. To establish a strong baseline, Pavlakos et al. [7] integrated the volumetric representation with a coarse-to-fine supervision scheme to figure out the 3D joint locations by the predicted 3D volumetric heat maps. Based on the ConvNet pose estimator and the volumetric heap map representation proposed by [7], recent approaches mainly made progress from two aspects. On the one hand, human-structure constraints such as the human shape constraints [8], body articulation constraints [9] and the joint angle constraints [10] were employed to prevent invalid pose prediction. On the other hand, effective training approaches were proposed, making the estimation process differentiable [11] and enabling the model to learn from weakly labeled data [12], [13]. To further enable the pose estimator to predict full 3D human mesh instead of the joint locations, Kanazawa et al. [14], [15] proposed end-to-end CNN frameworks for reconstructing the full 3D mesh of a human body from an image or a video. These approaches based on image-level input can directly capture rich knowledge contained in images. However, without intermediate feature and supervision, the model’s performance will also be affected by the image’s background, lighting and person’s clothing. More importantly, the large dimension of image-level input disables the 3D model from receiving a large number of images as input, bottlenecking the performance of 3D pose estimation in video.

For the approaches of the second category, they built a 3D joint estimation model on top of a high-performance 2D keypoint detector. Given an input image, these approaches first utilized the 2D keypoint detector to predict the image’s 2D keypoints. The predicted 2D keypoints were then lifted as the 3D joint estimation model’s input to predict the final 3D joint locations. As an earlier work, Chen et al. [16] regarded the 3D pose estimation as a matching problem. They found the best matching 3D pose of the 2D keypoint input from the 3D pose pool by a nearest-neighbor (NN) model. Considering that the ground-truth 3D pose of the input may be non-corresponding to all the 3D poses in the pool, Martinez et al. [17] proposed an effective fully-connected residual network to regress the 3D joint locations from 2D keypoint input. In addition to utilizing effective human-structure information as the approaches of the first category, based on [17], recent approaches of this category further improved the pose estimation performance by hierarchical joint prediction [5], 2D keypoint refinement [18] and view-invariant constraint [3], [19]. Overall, the approaches in such a “image-2D-3D” pipeline outperform the end-to-end counterparts. One important reason is that the 2D detector can be trained by large-scale indoor/outdoor images. It provides the 3D model a strong intermediate feature to build upon.

When estimating the 3D poses in a video, recent approaches

exploited temporal information into the model to alleviate incoherent predictions. As an earlier work, Mehta et al. [20] applied simple temporal filtering across 2D and 3D poses from previous frames to predict a temporally consistent 3D pose. As Long Short Term Memory networks (LSTM) were created to adaptively capture information from temporal input by the well-designed input gate, output gate and forget gate, Lin et al. [21] presented the LSTM-based Recurrent 3D Pose Sequence Machine. It automatically learns the image-dependent structural constraint and sequence-dependent temporal context by a multi-stage sequential refinement. Similar to [21], Rayat et al. [4] predicted temporally consistent 3D poses by learning the temporal context of a sequence using sequence-to-sequence LSTM-based network. Considering the high computational complexity of LSTM, Pavllo et al. [1] further introduced a temporal fully-convolutional model which enables parallel processing of multiple frames and supports very long 2D keypoint sequence as input. All these approaches essentially leverage the adjacent frames to benefit the current frame’s prediction. Compared with them, we are the first to make all the frames in a video contribute to the 3D prediction. Motivated by [22], [23], [14] that created effective sub-tasks for human pose estimation and mesh recovery, we propose a novel solution to decompose the 3D pose estimation task into two bone-based sub-tasks. It should be noticed that Sun et al. [6] also transformed the 3D joint into a bone-based representation. They trained the model to regress short and long range relative shifts between different joints. We demonstrate that completely decomposing the task into the bone length and bone direction prediction achieves the best performance and makes better use of the relative joint shift supervision.

III. OUR MODEL

In this section, we formally present our 3D pose estimation model. In section III-A, we first describe the overall anatomy-aware framework that decomposes the 3D joint location prediction task into bone length and direction prediction. In section III-B, we present the fully-convolutional propagating network for hierarchical bone direction prediction. In Section III-C, the architecture and training details of bone length prediction network are presented. In Section III-D, we describe the framework’s overall training strategy. In section III-E, an implicit attention mechanism is introduced to feed the keypoint visibility scores into the model as extra guidance. Our framework’s overall architecture is shown as Fig. 1.

A. Anatomy-aware Framework

As in [1], [4], [3], given the predicted 2D keypoints of each frame in a video, we aim at predicting the normalized 3D locations of j pre-defined joints for each frame. The 3D location of joint “Pelvis” is commonly defined as the origin of the 3D coordinates. Given a human joint set that contains j joints as in Fig. 2, they correspond to $(j - 1)$ directed bones with each joint being the vertex of at least one bone. This enables us to transform the 3D joint coordinates to the presentation of bone lengths and bone directions.

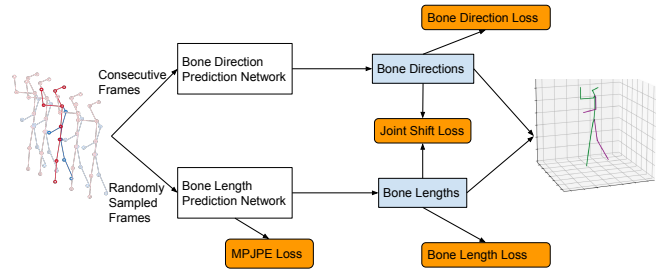


Fig. 1. The overview of the proposed anatomy-aware framework. It predicts the bone directions and bone lengths of the current frame using consecutive local frames and randomly sampled frames across the entire video, respectively.

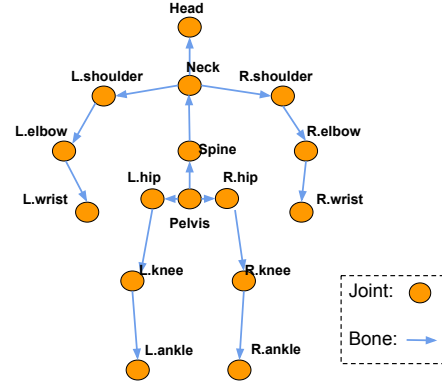


Fig. 2. The joint and bone representation of a human pose.

Formally, to predict the 3D joint locations of a specific (*i.e.* current) frame, we decompose the task to predict the length and direction of each bone. For the k -th joint, its 3D location \vec{J}_k can be derived as:

$$\vec{J}_k = \sum_{b \in B^k} \vec{D}_b \cdot L_b \quad (1)$$

Here \vec{D}_b and L_b are the direction and length of bone b , respectively. B^k contains all the bones in the path from “Pelvis” to the k -th joint.

We use two separate sub-networks to predict the bone lengths and directions of the current frame, respectively, as bone length prediction needs global input to ensure consistency across all the frames, whereas bone directions should be estimated within a local temporal window. Meanwhile, to ensure consistency between predicted bone lengths and directions, motivated by [6], we add a joint shift loss between the two predictions in addition to their own losses, as shown in Fig. 1. Specifically, the joint shift loss is defined as follows:

$$\mathcal{L}_{JS} = \sum_{k_1, k_2 \in \mathcal{P}, k_1 < k_2} \left\| X_{JS}^{k_1, k_2} - Y_{JS}^{k_1, k_2} \right\|_2^2 \quad (2)$$

Here $Y_{JS}^{k_1, k_2}$ is the 3-dimensional ground-truth relative joint shift of the current frame from the k_1 -th joint to the k_2 -th joint, $X_{JS}^{k_1, k_2}$ is the corresponding predicted relative joint shift derived from the predicted bone lengths and bone directions of the current frame. \mathcal{P} contains all the joint pairs that are not directly connected as a bone. With the joint shift loss, the two sub-networks are connected and enforced to learn from each other jointly. We describe the details of the two sub-networks in the following two sections.

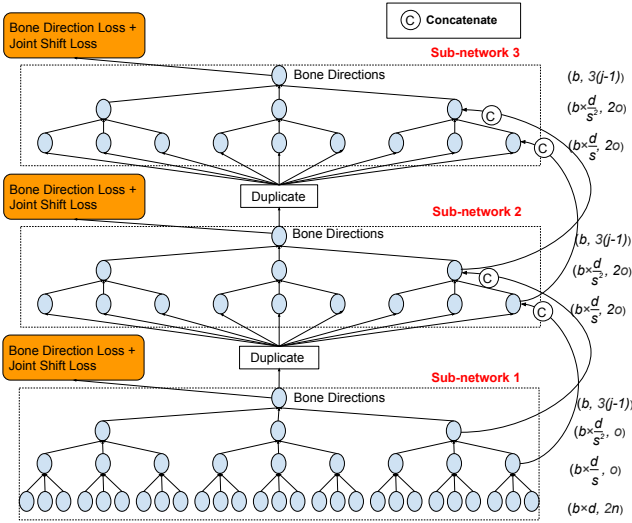


Fig. 3. The architecture of the bone direction prediction network. Long skip connections are added between adjacent sub-networks. We illustrate the dimension of each input/output. b is the batch size. o is the output channel number of fully-connected layer. s is the stride of 1D convolution layer in the network ($s = 3$). n is the size of the 2D keypoint set. d is the input frame number of the bottom sub-network.

B. Bone Direction Prediction Network

We adopt the temporal fully-convolutional network proposed by Pavlo et al. [1] as the backbone architecture of our bone direction prediction network. Specifically, the 2D keypoints of d consecutive frames are concatenated to form the input to the network, with the 2D keypoints of the current frame in the center. In essence, to predict the bone directions of the current frame, the temporal network captures the information of the current frame and the context from its adjacent frames as well. A bone direction loss based on mean squared error is applied to train the network:

$$\mathcal{L}_D = \|X_D - Y_D\|_2^2 \quad (3)$$

Here X_D and Y_D represent the predicted and ground-truth $3(j-1)$ -dimensional bone direction vector of the current frame, respectively.

It should be noted that the joint shift loss introduced in Section III-A makes the predicted directions of different bones mutually relevant. For example, if the predicted direction of the left lower arm is inaccurate, the predicted direction of the left upper arm will also be affected, since the model is encouraged to regress a long range shift from left shoulder to left wrist. Intuitively, it would benefit the overall prediction if we could first predict those easy and high-confident cases, and let them guide the subsequent prediction of other joints. As poses may vary significantly, it is difficult to pre-determine the hierarchy of the prediction. Motivated by [5], here we propose a fully-convolutional propagating architecture with long skip connections, and let the network itself to learn the prediction hierarchy instead, as in Fig. 3.

Specifically, the architecture is a stack of several sub-networks, with each sub-network being a temporal fully-convolutional network with residual blocks proposed by [1]. The output of each sub-network is the predicted bone directions of the current frame. Except the top sub-network, we temporally duplicate the output of each sub-network $\frac{d}{s}$ times as

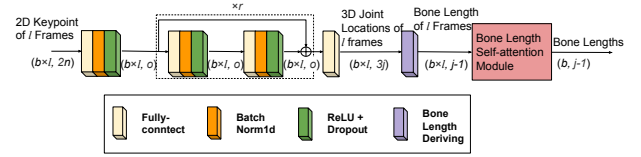


Fig. 4. Detailed structure of the bone length prediction network. r is the number of residual block.

the input to the next sub-network. For each residual block of a specific sub-network, we concatenate its output with the output of the corresponding residual block in the adjacent upper sub-network on channel level. This forms the long skip connections between adjacent sub-networks. We adopt an independent training strategy for each sub-network, that is, we train each sub-network by the loss of the bone direction prediction network, the back propagation is blocked between different sub-networks. By doing that, the bottom networks would not be affected by the upper ones, and instead would propagate high-confident predictions to guide subsequent predictions. In the process, the model automatically learns the hierarchical order of the prediction. In Section IV, we demonstrate the effectiveness of the proposed architecture.

C. Bone Length Prediction Network

As discussed in Section III-A, the prediction of bone lengths requires global inputs from the entire video. However, taking too many frames as the input would make the computation prohibitively expensive. To capture the global context efficiently, we choose to randomly sample l frames across the entire video as the input to the network. The detailed structure of the network is shown as Fig. 4.

We adopt the fully-connected residual network for bone length prediction. Specifically, it has the same structure and layer number as the bottom sub-network of the bone direction prediction network. However, since the randomly sampled frames do not have temporal connections, we replace each 1D convolution layer by the fully-connected layer in the network. This adapts the network for single-frame input instead of multi-frame consecutive inputs. The fully-connected network predicts the $(j-1)$ bone lengths of each sampled frame.

Intuitively, we can average the predicted bone lengths of each sampled frame as the predicted bone lengths of the current frame. In such a way, a similar bone length loss can be applied to train the fully-connected network:

$$\mathcal{L}_L = \|X_L - Y_L\|_2^2 \quad (4)$$

Here X_L and Y_L are the predicted and ground-truth $(j-1)$ -dimensional bone length vector of the current frame.

However, since the training datasets usually only contain very limited number of actors, and the bone lengths in the videos performed by the same actor are identical. Such a training loss would lead to severe overfitting. To solve this problem, instead of directly predicting the bone lengths, the fully-connected residual network is modified to predict the 3D joint locations of each sampled frame, supervised by the mean per joint position error (MPJPE) loss as in [1]:

$$\mathcal{L}_J = \frac{1}{j} \sum_{k=1}^j \|X_J^k - Y_J^k\|_2^2 \quad (5)$$

Here X_j^k and Y_j^k are the predicted and ground-truth 3-dimensional 3D joint locations of the k -th joint. Since each frame would predict a set of 3D joint locations, and the number of input frames l is usually large, minimizing the MPJPE for the predictions of all the frames would make the convergence speed of the bone length prediction network much faster than the bone direction prediction network. This decreases the performance of jointly training them by Equation 2. We choose to randomly sample one frame from l input frames and calculate its corresponding MPJPE loss. We find that such a training strategy works quite well as shown in the experiments.

Since the 3D joint locations would vary in each frame, the overfitting problem would be largely avoided. The bone lengths of each of the l frame can then be derived from the 3D joint locations accordingly.

When averaging the bone length predictions from the input frames, the prediction accuracy for a specific bone depends on the poses. For example, some of the bones in a certain pose are hardly visible due to occlusion or foreshortening, making the prediction unreliable. Motivated by the self-attention mechanism that is widely used for vision and language tasks [24], [25], [26], [27], [28], we further incorporate a self-attention module at the top of the fully-connected network to predict the bone length vector X_L of the current frame:

$$\begin{aligned} A_i &= \frac{\exp(\gamma W X_{iJ})}{\sum_{i=1}^l \exp(\gamma W X_{iJ})} \\ X_L &= \sum_{i=1}^l A_i \odot X_{iL} \end{aligned} \quad (6)$$

Here X_{iJ} and X_{iL} are the predicted 3D-dimensional 3D joint location vector and the corresponding derived bone length vector of the i -th input frame. W is a learnable matrix of the self-attention module. \odot indicates element-wise multiplication. A_i is the $(j-1)$ -dimensional attention weights that indicate the bone-specific importance of the i -th frame's predicted bone lengths for X_L . γ is a hyper-parameter to control the degree of attention. During training, the fully-connected residual network and the bone length self-attention module are optimized independently by L_J (Equation 5) and L_L (Equation 4).

To further solve the overfitting problem of the self-attention module, we augment the training data by generating samples with variant bone lengths. In particular, for each training video, we randomly create a new group of $(j-1)$ bone lengths. We modify the ground-truth 3D joint locations of each frame to make them accordant with the new bone lengths. Because the camera parameter is available, we can reconstruct the 2D keypoints of each frame from its modified ground-truth 3D joint locations. For each training iteration, we additionally feed a batch of randomly sampled l frames from the augmented videos and use the corresponding 2D keypoints and bone lengths to optimize the self-attention module by L_L . As the self-attention module is only used for bone length re-weighting, we consider it valid to train this module by a combination of predicted 2D keypoints and reconstructed clean 2D keypoints.

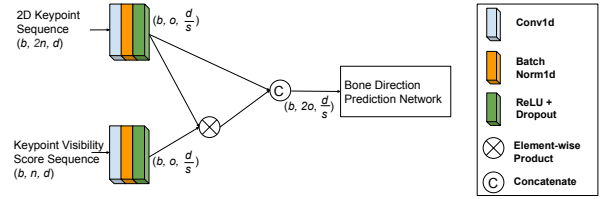


Fig. 5. Detailed network structure for feeding the visibility scores. d is the input frame number of the bone direction prediction network, s and o are the stride and output channel number of 1D convolution layer, respectively.

D. Overall Loss Function

By combining the losses in each sub-network and the joint shift loss, the overall loss function for our framework is given as:

$$\mathcal{L} = \lambda_D \mathcal{L}_D + \lambda_L \mathcal{L}_L + \lambda_J \mathcal{L}_J + \lambda_{JS} \mathcal{L}_{JS} \quad (7)$$

Here λ_D , λ_L , λ_J and λ_{JS} are hyper-parameters regulating the importance of each term.

During training, only the parameters of the bone direction prediction network are updated by the joint shift loss. Essentially, the joint shift loss supervises the model to predict robust bone directions that match with the predicted bone lengths for long range objectives. In Section IV, we prove that the proposed anatomy-aware framework better exerts the joint shift loss's potential than [6].

During inference, to estimate the 3D joint locations of a specific frame, we adopt the same strategy as the training process. We still randomly sample l frames of the video to predict the bone lengths of this frame. We find that taking all the frames of the video as input for bone length prediction does not lead to better performance. In Section IV, we provided more details regarding the frame sampling strategy for bone length prediction.

E. Incorporating the Visibility Scores

Our 3D joint prediction network takes the predicted 2D keypoints as the input, which sometimes have ambiguities. For example, when a person has his/her legs crossed in a front view, the corresponding 2D keypoints cannot provide information about which leg is in the front. A more common situation happens when the person put his/her hands in front of the chest or behind the back, the 3D model will be confused by the relative depth between the hands and chest. Even though temporal information is exploited, the above problems still exist.

We provide a simple yet effective approach to solve the problem without feeding the original images into the model. Specifically, we predict the visibility score of each 2D keypoint in a frame, and incorporate it into the model for 3D joint estimation. The visibility score indicates the confidence of each keypoint being visible in the frame, which can be extracted from most 2D keypoint detectors. Compared with the position, a keypoint's visibility is easy to be predicted, this score is thus reliable.

We argue that the importance of a specific keypoint's visibility score is related to the corresponding pose. For example, the visibility scores of the hands become useless if the hands are stretched far away from the body. Therefore,

we adopt an implicit attention mechanism to adaptively adjust the importance of the visibility scores.

Given the 2D keypoint sequence of d consecutive frames as the input of the network in Section III-B, we feed the keypoint visibility score sequence of the d frames as in Fig. 5. An 1D convolutional block is first applied, which maps the visibility score sequence into a temporal hidden feature that has the same dimension as the hidden feature of the 2D keypoint sequence. After that, we do element-wise multiplication for the two temporal hidden features as the weighted visibility score feature. In the end, we concatenate the weighted visibility score feature and the hidden feature of the 2D keypoint sequence on channel level and feed the concatenated feature to the next 1D convolution layer of the network. The whole process can be expressed as follows:

$$O_I = [C_{1D}(X_I) \odot C_{1D}(V_I); C_{1D}(X_I)] \quad (8)$$

Here X_I , V_I and O_I represent the input 2D keypoint sequence, the input keypoint visibility score sequence and the outputted concatenated features. C_{1D} represent the operations performed by a combination of 1D convolution layer, batch normalization layer and ReLU unit. $[\cdot; \cdot]$ indicates concatenation. Similar to [29], the hidden feature of the 2D keypoint sequence can be regarded as implicit attention weights to adjust the visibility score importance.

IV. EXPERIMENTS

A. Datasets and Evaluation

We evaluate the proposed model on two well established 3D human pose estimation datasets: Human3.6M [30] and MPI-INF-3DHP [31].

- **Human3.6M** contains 3.6 million video frames with the corresponding annotated 3D and 2D human joint positions, from 11 actors. Each actor performs 15 different activities captured from 4 camera views. Following previous works [1], [5], [4], [6], [17], the model is trained on five subjects (S1, S5, S6, S7, S8) and evaluated on two subjects (S9 and S11) on a 17-joint skeleton. We follow the standard protocols to evaluate the models on Human3.6M. The first one (*i.e.* Protocol 1) is the mean per-joint position error (MPJPE) in millimeters that measures the mean Euclidean distance between the predicted and ground-truth joint positions without any transformation. The second one (*i.e.* Protocol 2) is the normalized variant P-MPJPE after aligning the predicted 3D pose with the ground-truth using a similarity transformation. In addition, to measure the smoothness of predictions over time, which is important for video, we also report the joint velocity errors (MPJVE) created by [1] corresponding to the MPJPE of the first derivative of the 3D pose sequences.
- **MPI-INF-3DHP** is a recently proposed 3D dataset consisting of both constrained indoor and complex outdoor scenes. It records 8 actors performing 8 activities from 14 camera views. Following [20], [31], on a 14-joint skeleton, we consider all the 8 actors in the training set and select sequences from 8 camera views in total (5 chest-high cameras, 2 head-high cameras and 1 knee-high camera) for

training. Evaluation is performed on the independent MPI-INF-3DHP test set that has different scenes, camera views and relatively different actions from the training set. This design implicitly covers the cross-dataset evaluation. We report the Percentage of Correct Keypoints (PCK) within 150mm range, Area Under Curve (AUC), and MPJPE.

B. Implementation details

For Human3.6M, we use the predicted 2D keypoints released by [1] from the Cascaded Pyramid Network (CPN) as the input of our 3D pose model. For MPI-INF-3DHP, the predicted 2D keypoints are acquired from the pretrained AlphaPose model [32]. In addition to the 2D keypoints, the keypoint visibility scores for both datasets are also extracted from the pretrained AlphaPose model.

We use the Adam optimizer to train our model in an end-to-end manner. For each training iteration, the mini-batch size is set to 1024 for both original samples and augmented samples. We set $\lambda_D = 0.02$, $\lambda_L = 0.05$, $\lambda_J = 1$ and $\lambda_{JS} = 0.1$ for the loss terms in Equation 7. For the bone length self-attention module, we set $\gamma = 10$ in Equation 8. The sampled frame number of the bone length prediction network l is set to 50 for both the training and inference process. For the proposed architecture in Section III-B, the number of sub-networks is set to 2. As in [1], the output channel number of each 1D convolution layer and fully-connected layer is set to 1024. For actual implementation, instead of manually deriving the 3D joint locations and relative joint shifts from the predicted bone lengths and bone directions, we regress the two objectives by feeding the concatenation of the predicted bone length vector and bone direction vector into two fully-connected layers, respectively. The fully-connected layers are trained together with the whole network. This achieves slightly better performance.

C. Experiment results

Table I shows the quantitative results of our proposed full model and other baselines on Human3.6M. Following [1], we present the performance of our 81-frame and 243-frame models which receive 81 and 243 consecutive frames, respectively, as the input of the bone direction prediction network. We also experiment with a causal version of our model to enable real-time prediction. During the training/inference process, the causal model only receives d consecutive and l randomly sampled frames from the past/current frames for the current frame’s estimation. Overall, our model has low average error on both Protocol 1, Protocol 2 and MPJVE. On a great number of actions, we achieve the best performance. Compared with the baseline model [1] that shares the same 2D keypoint detector, our model achieves more smooth prediction with lower MPJVE and achieves significantly better performance on complex activities such as “Sitting” (-3.4mm in Protocol 1) and “Sitting down” (-5.6mm in Protocol 1). We attribute it to the accurate prediction of the bone lengths for these activities. Even though the person bends his/her body, based on the predicted bone lengths, the joint shift loss can effectively guide the model to predict high-quality bone directions. Fig. 6

TABLE I
QUANTITATIVE COMPARISONS BETWEEN THE ESTIMATED POSE AND THE GROUND-TRUTH ON HUMAN3.6M UNDER PROTOCOLS 1,2 AND MPJVE. (*)
WE REPORT THE RESULT WITHOUT DATA AUGMENTATION USING VIRTUAL CAMERAS.

Protocol 1		Dir.	Disc.	Eat	Greet	Phone	Photo	Pose	Purch.	Sit	SitD.	Smoke	Wait	WalkD.	Walk	WalkT.	Avg
Martinez et al. [17] ICCV'17	51.8	56.2	58.1	59.0	69.5	78.4	55.2	58.1	74.0	94.6	62.3	59.1	65.1	65.1	49.5	52.4	62.9
Sun et al. [6] ICCV'17	52.8	54.8	54.2	54.3	61.8	67.2	53.1	53.6	71.7	86.7	61.5	53.4	61.6	61.6	47.1	53.4	59.1
Pavlakos et al. [12] CVPR'18	48.5	54.4	54.4	52.0	59.4	65.3	49.9	52.9	65.8	71.1	56.6	52.9	60.9	44.7	47.8	56.2	56.2
Yang et al. [9] CVPR'18	51.5	58.9	50.4	57.0	62.1	65.4	49.8	52.7	69.2	85.2	57.4	58.4	43.6	60.1	47.7	58.6	58.6
Luvizon et al. [33] CVPR'18	49.2	51.6	47.6	50.5	51.8	60.3	48.5	51.7	61.5	70.9	53.7	48.9	57.9	44.4	48.9	53.2	53.2
Hossain & Little [4] ECCV'18	48.4	50.7	57.2	55.2	63.1	72.6	53.0	51.7	66.1	80.9	59.0	57.3	62.4	46.6	49.6	58.3	58.3
Lee et al. [5] ECCV'18	40.2	49.2	47.8	52.6	50.1	75.0	50.2	43.0	55.8	73.9	54.1	55.6	58.2	43.3	43.3	52.8	52.8
Chen et al. [3] CVPR'19	41.1	44.2	44.9	45.9	46.5	39.3	41.6	54.8	73.2	46.2	48.7	42.1	35.8	46.6	38.5	46.3	46.3
Pavlo et al. [1] (243 frames, Causal) CVPR'19	45.9	48.5	44.3	47.8	51.9	57.8	46.2	45.6	59.9	68.5	50.6	46.4	51.0	34.5	35.4	49.0	49.0
Pavlo et al. [1] (243 frames) CVPR'19	45.2	46.7	43.3	45.6	48.1	55.1	44.6	44.3	57.3	65.8	47.1	44.0	49.0	32.8	33.9	46.8	46.8
Lin et al. [2] BMVC'19	42.5	44.8	42.6	44.2	48.5	57.1	42.6	41.4	56.5	64.5	47.4	43.0	48.1	33.0	35.1	46.6	46.6
Cai et al. [34] ICCV'19	44.6	47.4	45.6	48.8	50.8	59.0	47.2	43.9	57.9	61.9	49.7	46.6	51.3	37.1	39.4	48.8	48.8
Cheng et al. [35] ICCV'19 (*)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	44.8
Yeh et al. [36] NIPS'19	44.8	46.1	43.3	46.4	49.0	55.2	44.6	44.0	58.3	62.7	47.1	43.9	48.6	32.7	33.3	46.7	46.7
Xu et al. [18] CVPR'20	37.4	43.5	42.7	42.7	46.6	59.7	41.3	45.1	52.7	60.2	45.8	43.1	47.7	33.7	37.1	45.6	45.6
Ours (243 frames, Causal)	42.5	45.4	42.3	45.2	49.1	56.1	43.8	44.9	56.3	64.3	47.9	43.6	48.1	34.3	35.2	46.6	46.6
Ours (81 frames)	42.1	43.8	41.0	43.8	46.1	53.5	42.4	43.1	53.9	60.5	45.7	42.1	46.2	32.2	33.8	44.6	44.6
Ours (243 frames)	41.4	43.5	40.1	42.9	46.6	51.9	41.7	42.3	53.9	60.2	45.4	41.7	46.0	31.5	32.7	44.1	44.1
Protocol 2		Dir.	Disc.	Eat	Greet	Phone	Photo	Pose	Purch.	Sit	SitD.	Smoke	Wait	WalkD.	Walk	WalkT.	Avg
Martinez et al. [17] ICCV'17	39.5	43.2	46.4	47.0	51.0	56.0	41.4	40.6	56.5	69.4	49.2	45.0	49.5	38.0	43.1	47.7	47.7
Sun et al. [6] ICCV'17	42.1	44.3	45.0	45.4	51.5	53.0	43.2	41.3	59.3	73.3	51.0	44.0	48.0	38.3	44.8	48.3	48.3
Pavlakos et al. [12] CVPR'18	34.7	39.8	41.8	38.6	42.5	47.5	38.0	36.6	50.7	56.8	42.6	39.6	43.9	32.1	36.5	41.8	41.8
Yang et al. [9] CVPR'18	26.9	30.9	36.3	39.9	43.9	47.4	28.8	29.4	36.9	58.4	41.5	30.5	29.5	42.5	32.2	37.7	37.7
Hossain & Little [4] ECCV'18	35.7	39.3	44.6	43.0	47.2	54.0	38.3	37.5	51.6	61.3	46.5	41.4	47.3	34.2	39.4	44.1	44.1
Chen et al. [3] CVPR'19	36.9	39.3	40.5	41.2	42.0	34.9	38.0	51.2	67.5	42.1	42.5	37.5	30.6	40.2	34.2	41.6	41.6
Pavlo et al. [1] (243 frames, Causal) CVPR'19	35.1	37.7	36.1	38.8	38.5	44.7	35.4	34.7	46.7	53.9	39.6	35.4	39.4	27.3	28.6	38.1	38.1
Pavlo et al. [1] (243 frames) CVPR'19	34.1	36.1	34.4	37.2	36.4	42.2	34.4	33.6	45.0	52.5	37.4	33.8	37.8	25.6	27.3	36.5	36.5
Lin et al. [2] BMVC'19	32.5	35.3	34.3	36.2	37.8	43.0	33.0	32.2	45.7	51.8	38.4	32.8	37.5	25.8	28.8	36.8	36.8
Cai et al. [34] ICCV'19	35.7	37.8	36.9	40.7	39.6	45.2	37.4	34.5	46.9	50.1	40.5	36.1	41.0	29.6	33.2	39.0	39.0
Cheng et al. [35] ICCV'19 (*)	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	34.1
Xu et al. [18] CVPR'20	31.0	34.8	34.7	34.4	36.2	43.9	31.6	33.5	42.3	49.0	37.1	33.0	39.1	26.9	31.9	36.2	36.2
Ours (243 frames, Causal)	33.6	36.0	34.4	36.6	37.5	42.6	33.5	33.8	44.4	51.0	38.3	33.6	37.7	26.7	28.2	36.5	36.5
Ours (81 frames)	33.1	35.3	33.4	35.9	36.1	41.7	32.8	33.3	42.6	49.4	37.0	32.7	36.5	25.5	27.9	35.6	35.6
Ours (243 frames)	32.6	35.1	32.8	35.4	36.3	40.4	32.4	32.3	42.7	49.0	36.8	32.4	36.0	24.9	26.5	35.0	35.0
MPJVE		Dir.	Disc.	Eat	Greet	Phone	Photo	Pose	Purch.	Sit	SitD.	Smoke	Wait	WalkD.	Walk	WalkT.	Avg
Pavlo et al. [1] (243 frames) CVPR'19	3.0	3.1	2.2	3.4	2.3	2.7	2.7	2.7	3.1	2.1	2.9	2.3	2.4	3.7	3.1	2.8	2.8
Ours (243 frames)	2.7	2.8	2.0	3.1	2.0	2.4	2.4	2.8	1.8	2.4	2.0	2.1	3.4	2.7	2.4	2.5	2.5

shows the visualized qualitative results from the baseline and our full model on ‘‘Sitting’’ and ‘‘Sitting down’’ poses.

TABLE II
QUANTITATIVE COMPARISONS OF MODELS TRAINED/EVALUATED ON HUMAN3.6M USING THE GROUND-TRUTH 2D INPUT.

	Protocol 1	Protocol 2
Martinez et al. [17]	45.5	37.1
Hossain & Little [4]	41.6	31.7
Lee et al. [5]	38.4	-
Pavlo et al. (243 frames) [1]	37.2	27.2
Ours (243 frames)	32.3	25.2

Moreover, our model sharply improves the lower bound of 3D pose estimation when using the ground-truth 2D keypoints as input. For this experiment, data augmentation is not applied as it can be regarded as using extra ground truth 2D keypoints. From Table II, the gap between our model and the baseline is nearly 5mm on Protocol 1. It indicates that if the performance of the bottom 2D keypoint detector is improved, our model can further boost the improvement.

Following [1], we report the model parameter number and an estimate of the floating-point operations (FLOPs) per frame at inference time to compare different models’ computation complexity. The FLOPs are estimated in the same way as [1]. Besides, to evaluate the models’ inference efficiency, we report the inference frame per second (FPS) of different models by letting them estimate the 3D poses of a 10,000-frame test video frame by frame on a single GeForce GTX 2080 Ti GPU. As shown in Table III, our 9-frame model holds similar computation complexity to the 243-frame baseline model [1]. It achieves lower MPJPE with sharply fewer input frames, demonstrating the effectiveness of our proposed approach. On the other hand, even though the proposed model’ inference speed is about 3.5 times lower than [1], they still hold an acceptable FPS, which is significantly higher than common 2D

keypoint detection models [32], [37]. As the complete 3D pose estimation process is a combination of 2D keypoint detection and 3D pose estimation, the inference speed of the proposed model will not be the bottleneck.

We further investigate the sensitivity of different hyper-parameters mentioned in Section IV-B. Overall, as shown in Table IV, the model’s performance is insensitive to the setting of different hyper-parameters, indicating its strong robustness. Indeed, keeping increasing the sub-network number of the bone direction prediction network can still reduce the MPJPE. We set the final sub-network number to 2 as a good trade-off between the performance and the computational complexity.

Table V shows the quantitative results of our full model and other baselines on MPI-INF-3DHP. Overall, MPI-INF-3DHP contains fewer training samples than human3.6M. This leads to better performance for the 81-frame models than the 243-frame models. Still, our model outperforms the baselines by a large margin.

Overall, the comparison with the baseline model [1] on Table I, II, V demonstrates that our superior performance is not only from the strong 2D keypoint detector, but also highly related to the proposed high-performance 3D model.

D. Ablation Study

We next perform ablation experiments on Human3.6M under Protocol 1. For all the comparisons, we use the 81-frame models for the baseline [1] and ours. They receive 81 consecutive frames as the input to predict the 3D joint locations and the bone directions of the current frame, respectively.

We first show how each proposed module improves the model’s performance in Table VI. For the very naive anatomy-aware framework, we adopt **Baseline** as the bone direction prediction network. We train the bone length prediction network and bone direction prediction network by the bone

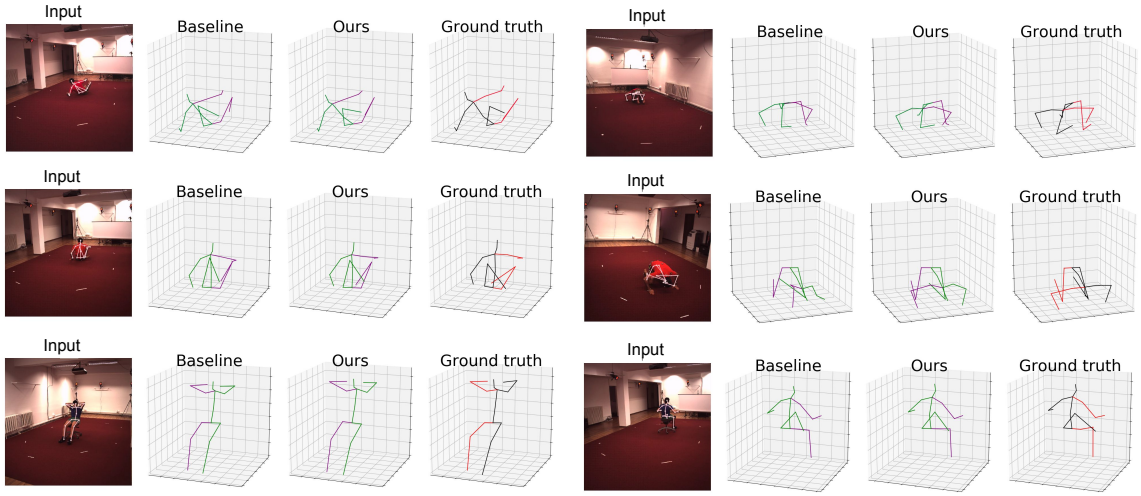


Fig. 6. Qualitative comparison between the proposed 243-frame model and the baseline 243-frame model [1] on typical poses.

TABLE III

COMPUTATIONAL COMPLEXITY, MPJPE, AND FRAME PER SECOND (FPS) OF DIFFERENT MODELS UNDER PROTOCOL 1 ON HUMAN3.6. THE COMPUTATIONAL COMPLEXITY IS COMPUTED WITHOUT TEST-TIME AUGMENTATION USED BY ALL THE MODELS. THE TWO NUMBERS OF “ \approx FLOPs” FOR OUR MODELS REPRESENT THE ESTIMATED FLOATING-POINT OPERATIONS OF THE BONE DIRECTION PREDICTION NETWORK AND THE BONE LENGTH PREDICTION NETWORK, RESPECTIVELY.

	Parameters	\approx FLOPs	MPJPE	FPS
Hossain & Little [4] ECCV’18	16.96M	33.88M	58.3	-
Pavlo et al. [1] (27 frames) CVPR’19 w/o dilation	29.53M	59.03M	49.3	486
Pavlo et al. [1] (27 frames) CVPR’19	8.56M	17.09M	48.8	1492
Pavlo et al. [1] (81 frames) CVPR’19	12.75M	25.48M	47.7	1121
Pavlo et al. [1] (243 frames) CVPR’19	16.95M	33.87M	46.8	863
Ours (9 frames)	18.24M	29.97M + 4.58M	46.3	759
Ours (27 frames)	31.88M	53.03M + 8.67M	45.3	410
Ours (81 frames)	45.53M	76.1M + 12.76M	44.6	315
Ours (243 frames)	59.18M	99.17M + 16.95M	44.1	264

TABLE IV

PARAMETER SENSITIVITY TEST FOR OUR 81-FRAME MODEL UNDER PROTOCOL 1 ON HUMAN3.6.

λ_D	λ_L	λ_{JS}	γ	l	Num. of sub-network	MPJPE
0.002	0.05	0.1	10	50	2	44.9
0.02	0.05	0.1	10	50	2	44.6
0.2	0.05	0.1	10	50	2	44.6
0.02	0.005	0.1	10	50	2	45.0
0.02	0.5	0.1	10	50	2	44.6
0.02	0.05	0.01	10	50	2	45.0
0.02	0.05	1	10	50	2	44.9
0.02	0.05	0.1	1	50	2	44.8
0.02	0.05	0.1	10	50	2	44.6
0.02	0.05	0.1	100	50	2	45.3
0.02	0.05	0.1	10	10	2	45.0
0.02	0.05	0.1	10	50	2	44.6
0.02	0.05	0.1	10	100	2	44.6
0.02	0.05	0.1	10	50	1	45.3
0.02	0.05	0.1	10	50	2	44.6
0.02	0.05	0.1	10	50	3	44.4

TABLE V

QUANTITATIVE COMPARISONS OF DIFFERENT MODELS ON MPI-INF-3DHP.

	PCK	AUC	MPJPE
Mehta et al. [31] 3DV’17	75.7	39.3	117.6
Mehta et al. [20] ACM ToG’17	76.6	40.4	124.7
Pavlo et al. [1] (81 frames) CVPR’19	86.0	51.9	84.0
Pavlo et al. [1] (243 frames) CVPR’19	85.5	51.5	84.8
Lin et al. [2] BMVC’19	82.4	49.6	81.9
Ours (81 frames)	87.9	54.0	78.8
Ours (243 frames)	87.8	53.8	79.1

TABLE VI

COMPARISON OF DIFFERENT MODELS UNDER PROTOCOL 1 ON HUMAN3.6M. **BASILINE** REPRESENTS THE BASELINE 81-FRAME MODEL [1]. OTHER ROWS REPRESENTS THE PROPOSED ANATOMY-AWARE FRAMEWORK THAT DECOMPOSES THE TASKS INTO BONE LENGTH PREDICTION AND BONE DIRECTION PREDICTION. **ML** REFERS TO THE USE OF MPJPE LOSS TO SOLVE THE OVERFITTING PROBLEM OF THE FULLY-CONNECTED RESIDUAL NETWORK AS SECTION III-C. **BONEATT** REFERS TO THE FEEDING OF THE BONE LENGTH SELF-ATTENTION MODULE FOR BONE LENGTH RE-WEIGHTING, INSTEAD OF DIRECTLY AVERAGING THE PREDICTED BONE LENGTHS. **AUG** REFERS TO THE APPLYING OF DATA AUGMENTATION TO SOLVE THE OVERFITTING PROBLEM OF THE SELF-ATTENTION MODULE. **JSL**, **LSC** AND **SCOREATT** REFER TO THE APPLYING OF THE JOINT SHIFT LOSS, THE INCORPORATING OF THE FULLY-CONVOLUTIONAL PROPAGATING ARCHITECTURE AND THE FEEDING OF VISIBILITY SCORES BY AN IMPLICIT ATTENTION MECHANISM, RESPECTIVELY.

	ML	BoneAtt	AUG	JSL	LSC	ScoreAtt	MPJPE
Baseline	X	X	X	X	X	X	47.7
Baseline	X	X	X	✓	✓	✓	46.4
AF	✓	✓	✓	✓	✓	✓	48.4
	✓	✓	✓	✓	✓	✓	46.6
	✓	✓	✓	✓	✓	✓	46.5
	✓	✓	✓	✓	✓	✓	46.3
	✓	✓	✓	✓	✓	✓	45.8
	✓	✓	✓	✓	✓	✓	45.1
	✓	✓	✓	✓	✓	✓	44.6

length loss and bone direction loss, respectively. We observe a drop of performance from 47.7mm to 48.4mm caused by the overfitting of the bone length prediction network. Once we handle the overfitting by applying the MPJPE loss and doing data augmentation which can be regarded as intrinsic parts of our anatomy-aware framework, MPJPE sharply drops to

46.3mm that is 1.4mm lower than **Baseline**. This demonstrates the framework’s effectiveness. Moreover, we find that the joint shift loss, the fully-convolutional propagating architecture and the attentive feeding of visibility score reduce the error about 0.5mm, 0.7mm and 0.5mm, respectively. In the end, if we apply them to the baseline 81-frame work without the proposed anatomy-aware structure, the MPJPE increases to 46.4mm.

To further validate the proposed modules, the following models are compared to answer the following questions:

- Q: Comparison between the proposed anatomy-aware network and the generic bone based representation.

M: **Baseline + Composition**

A: Sun et al. [6] propose compositional pose regression that transforms the task into a generic bone based representation. To demonstrate the effectiveness of the our anatomy-aware framework that decomposes the task explicitly into bone length and bone direction, we apply compositional pose regression on **Baseline**. Specifically, *exactly* as [6], we modify **Baseline** to predict the bones and train it by the long range and short range joint shift loss (*i.e.* compositional loss). The MPJPE is 47.4mm. Compared with our comparable anatomy-aware framework (without feeding **LSC** and **ScoreATT**) whose MPJPE is 45.8mm, this suggests that the explicit bone length and bone direction representation is more effective than the generic bone based representation because it can utilize global information across all the frames. It also makes better use of the relative joint shift supervision as our model obtains larger improvement from **JSL**.

- Q: Comparison between our anatomy-aware network and directly imposing a bone length consistency loss.

M: **Baseline + ConsistencyLoss**

A: A bone length consistency loss (as opposed to an explicit hard design) can be imposed in a straightforward manner across frames. To evaluate this idea, we further add a training loss term on **Baseline** to reduce the predicted bone length difference between randomly sampled frame pairs of a same video. The best MPJPE is 47.7mm. This indicates the uselessness of utilizing the bone length consistency by a form of training loss for this supervised learning task, compared with our solution.

- Q: Whether the distant frames indeed help the prediction of bone length.

M: **AF(consecutive) + ML + BoneAtt + AUG + JSL**

A: To validate this, we investigate the model that receives l consecutive local frames as the input of the bone length prediction network for training/inference, with the current one in the center. As Section IV-B, l is still set to 50 for both the training and inference process. The error increases from 45.8mm to 46.7mm. This demonstrates that randomly sampling frames from the entire video for bone length prediction indeed improves the model’s performance, which is consistent with our motivation to decompose the task.

- Q: Whether the long skip connections and the independent training strategy indeed help the prediction of bone direction.

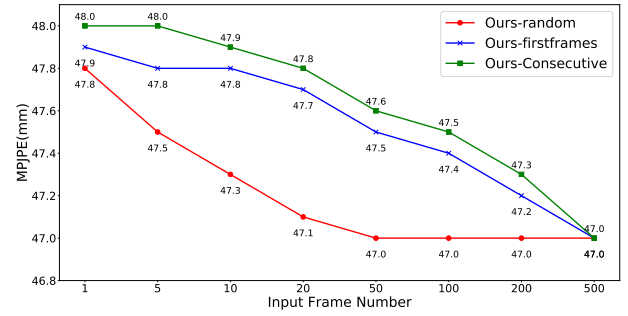


Fig. 7. The influence of the input frame number M on different real-time models.

M: **AF + ML + BoneAtt + AUG + JSL + Baseline-D**

A: It should be noticed that for the bone direction prediction network, **LSC** is two times deeper than **Baseline**. For fair comparison, we further design **Baseline-D** as the bone direction prediction network. Its structure is same as **LSC**, but with the long skip connections removed. Also, same as **Baseline**, the loss is only applied at the top and the back-propagation is not blocked between different sub-networks. The MPJPE is still 45.8mm. This indicates the uselessness of simply increasing the layer/parameter number of the temporal network.

- Q: Whether utilizing implicit attention mechanism is an effective way to feed the visibility score feature.

M: **AF + ML + BoneAtt + AUG + JSL + LSC + Fusion**

A: We create a model for which we directly concatenate the visibility scores with the input 2D keypoints before feeding them into the network instead of utilizing implicit attention. The MPJPE is 44.8mm. It proves that the implicit attention mechanism provides a more effective way to feed the visibility score feature.

E. Real-time 3D Pose Estimation

As mentioned in Section IV-C, our causal version model support real-time 3D human pose estimation. For the real-time 3D pose estimation, the inference speed is important and highly related to the frame selection strategy of the bone length prediction network. In this section, we compare different frame selection strategies for bone length prediction during the inference process. For all the comparisons, same as Section IV-D, we choose the 81-frame models for the baseline [1] and ours. It should be noticed that the 3D pose estimation model can only leverage the information of the current and past frames.

- **BS-causal** refers to the baseline model with causal convolutions [1]. To predict the 3D joint locations of the current frame, the input of the baseline model contains the 2D keypoint sequence of 81 consecutive frames with the current frame in the rightmost. The MPJPE is 49.8mm.
- **Ours-random** refers to our real-time model that adopts the same frame selection strategy as our standard model presented in the main paper. For the bone direction prediction network, we input the 2D keypoint sequence of 81 consecutive frames as **BS-causal**. To predict the bone lengths of the t -th frame, we randomly sample M

frames before the $(t + 1)$ -th frame and input their 2D keypoints to the bone length prediction network. If M is larger than t , we input the 2D keypoints of all the frames before the $(t + 1)$ -th frame.

- **Ours-firstframe** refers to our real-time model that doesn't need to iteratively compute the bone lengths. Specifically, during the inference process, to predict the bone lengths of the t -th frame, we always select the 2D keypoints of the first M frames of the video as input. In this situation, the model does not need to recalculate the bone lengths from the $(M + 1)$ -th frame. If M is larger than t , we input the 2D keypoints of all the frames before the $(t + 1)$ -th frame.
- **Ours-consecutive** refers to our real-time model that receives consecutive local frames for bone length prediction. To predict the bone length of the t -th frame, we input the 2D keypoint sequence of M consecutive frames with the t -th frame in the rightmost. Still, if M is larger than t , we input the 2D keypoints of all the frames before the $(t + 1)$ -th frame.

TABLE VII
RELATIVE AVERAGE REAL-TIME INFERENCE SPEED ON ALL THE VIDEOS
OF HUMAN3.6M TEST SET.

	Relative Inference Speed
BS-causal	1.0
Ours-random	0.28
Ours-firstframes	0.46

During the training process, for all the three models we propose (i.e. **Ours-random**, **Ours-firstframe**, **Ours-consecutive**), we still randomly sample 50 frames and input their 2D keypoints to the bone length prediction network. We don't observe further improvement of the model's performance when inputting the 2D keypoints of the first 50/50 consecutive local frames to train **Ours-firstframe/Ours-consecutive** as their inference process.

Fig. 7 shows how different settings of the input frame number M for inference influences the models' performance. As we expect, **Ours-random** achieves best performance when M is small. Moreover, randomly sampling 50 frames is sufficient for **Ours-random** to reach the best performance. However, it needs to predict the bone lengths frame by frame. On the other hand, **Ours-firstframe** achieves slightly better performance than **Ours-consecutive**. It indicates that predicting the bone lengths from the first frames of the video is more accurate than from consecutive local frames. We attribute it to the fact that the person's poses at the beginning of a video are commonly simpler on Human3.6M, this benefits the bone length prediction. To our surprise, even though we just update the bone lengths for the first 50 frames of each video and fix them from the 51-th frame, **Ours-firstframes** still achieves great performance. The MPJPE is 47.5mm, which is 2.3mm lower than the baseline. More importantly, it is more efficient than **Ours-random**. Tables VII shows different models' relative average inference speed on Human3.6M test set, M is set to 50 for all the models. The inference speed of **Ours-firstframes** is about two times lower than **BS-causal**, because its bone direction prediction network that proposed

in Section 3.2 of the main paper is nearly two times deeper than the baseline. However, **Ours-firstframes** is more efficient than **Ours-random** without updating the bone lengths frame by frame.

V. CONCLUSION

We present a new solution to estimating the human 3D pose. Instead of directly regressing the 3D joint locations, we transform the task into predicting the bone lengths and directions. For bone length prediction, we make use of the frames across the entire video and propose an effective fully-connected residual network with a bone length re-weighting mechanism. For bone direction prediction, we add along skip connections into a fully-convolutional architecture for hierarchical prediction. Extensive experiments have demonstrated that the combination of bone length and bone direction is an effective intermediate representation to bridge the 2D keypoints and 3D joint locations.

In recent years, as 3D human pose estimation has become a significant research topic for researchers to study, multiple directions are demonstrated to be promising for exploration. First, effective data augmentation algorithms [35], [38] are continuously proposed to guide the model to handle occluded or complex pose inputs. Moreover, the creation of generative adversarial network (GAN) [39] enables a number of approaches [9], [40] to utilize GAN for realistic and reasonable pose prediction, even in a weakly supervised setting. In addition, high-performance temporal models [1], [4] are created, which support very long 2D keypoint sequence as input and can adaptively capture significant information from keyframes. These directions are regarded as general directions since the proposed temporal models, adversarial training, and data augmentation algorithms can be generally applied to different research tasks other than 3D human pose estimation. In this paper, we focus on a more fundamental aspect of human pose estimation and create an effective learning representation for this task. We believe that exploring the human pose's learning representation is promising as the human body is a special Kinematic Tree-based structure different from other objects. The motion of the human body is drove by joint rotation with fixed bone lengths. We are delighted to see the human pose's learning representation evolved from the joint level to the bone vector level to the bone length/direction level that constantly improves human pose estimation. Based on our work, it may be illuminating for future works to keep exploring the relationship between the tasks of bone length prediction and bone direction prediction. Currently, we adopt two independent networks to predict bone directions and bone lengths. It is valuable to study whether the model can further improve the performance of each task by utilizing the knowledge captured from the other task. Applying the joint shift loss is one useful way. However, we believe that capturing this relationship at the network level as [41], [42] for visual question answering and image-text matching will make extra improvement for accurate and smooth 3D human pose estimation in video.

REFERENCES

- [1] D. Pavlo, C. Feichtenhofer, D. Grangier, and M. 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*, 2019, pp. 7753–7762.
- [2] J. Lin and G. H. Lee, "Trajectory space factorization for deep video-based 3d human pose estimation," *arXiv preprint arXiv:1908.08289*, 2019.
- [3] X. Chen, K.-Y. Lin, W. Liu, C. Qian, and L. Lin, "Weakly-supervised discovery of geometry-aware representation for 3d human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10895–10904.
- [4] M. Rayat Imtiaz Hossain and J. J. Little, "Exploiting temporal information for 3d human pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 68–84.
- [5] K. Lee, I. Lee, and S. Lee, "Propagating lstm: 3d pose estimation based on joint interdependency," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 119–135.
- [6] X. Sun, J. Shang, S. Liang, and Y. Wei, "Compositional human pose regression," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2602–2611.
- [7] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Coarse-to-fine volumetric prediction for single-image 3d human pose," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7025–7034.
- [8] X. Zhou, Q. Huang, X. Sun, X. Xue, and Y. Wei, "Towards 3d human pose estimation in the wild: a weakly-supervised approach," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 398–407.
- [9] W. Yang, W. Ouyang, X. Wang, J. Ren, H. Li, and X. Wang, "3d human pose estimation in the wild by adversarial learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5255–5264.
- [10] R. Dabral, A. Mundhada, U. Kusupati, S. Afaq, A. Sharma, and A. Jain, "Learning 3d human pose from structure and motion," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 668–683.
- [11] X. Sun, B. Xiao, F. Wei, S. Liang, and Y. Wei, "Integral human pose regression," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 529–545.
- [12] G. Pavlakos, X. Zhou, and K. Daniilidis, "Ordinal depth supervision for 3d human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7307–7316.
- [13] M. Kocabas, S. Karagoz, and E. Akbas, "Self-supervised learning of 3d human pose using multi-view geometry," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1077–1086.
- [14] A. Kanazawa, M. J. Black, D. W. Jacobs, and J. Malik, "End-to-end recovery of human shape and pose," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7122–7131.
- [15] A. Kanazawa, J. Y. Zhang, P. Felsen, and J. Malik, "Learning 3d human dynamics from video," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5614–5623.
- [16] C.-H. Chen and D. Ramanan, "3d human pose estimation= 2d pose estimation+ matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7035–7043.
- [17] J. Martinez, R. Hossain, J. Romero, and J. J. Little, "A simple yet effective baseline for 3d human pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2640–2649.
- [18] J. Xu, Z. Yu, B. Ni, J. Yang, X. Yang, and W. Zhang, "Deep kinematics analysis for monocular 3d human pose estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 899–908.
- [19] G. Wei, C. Lan, W. Zeng, and Z. Chen, "View invariant 3d human pose estimation," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2019.
- [20] D. Mehta, S. Sridhar, O. Sotnychenko, H. Rhodin, M. Shafiei, H.-P. Seidel, W. Xu, D. Casas, and C. Theobalt, "Vnect: Real-time 3d human pose estimation with a single rgb camera," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, p. 44, 2017.
- [21] M. Lin, L. Lin, X. Liang, K. Wang, and H. Cheng, "Recurrent 3d pose sequence machines," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 810–819.
- [22] J. Yu, C. Hong, Y. Rui, and D. Tao, "Multitask autoencoder model for recovering human poses," *IEEE Transactions on Industrial Electronics*, vol. 65, no. 6, pp. 5060–5068, 2017.
- [23] X. Zhou, X. Sun, W. Zhang, S. Liang, and Y. Wei, "Deep kinematic pose regression," in *European Conference on Computer Vision*. Springer, 2016, pp. 186–201.
- [24] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, pp. 5998–6008, 2017.
- [25] J. Yu, J. Li, Z. Yu, and Q. Huang, "Multimodal transformer with multi-view visual representation for image captioning," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [26] Z. Yu, Y. Cui, J. Yu, M. Wang, D. Tao, and Q. Tian, "Deep multimodal neural architecture search," *arXiv preprint arXiv:2004.12070*, 2020.
- [27] X. Li, J. Song, L. Gao, X. Liu, W. Huang, X. He, and C. Gan, "Beyond rnns: Positional self-attention with co-attention for video question answering," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8658–8665.
- [28] M. Gu, Z. Zhao, W. Jin, D. Cai, and F. Wu, "Video dialog via multi-grained convolutional self-attention context multi-modal networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 12, pp. 4453–4466, 2020.
- [29] J.-H. Kim, S.-W. Lee, D. Kwak, M.-O. Heo, J. Kim, J.-W. Ha, and B.-T. Zhang, "Multimodal residual learning for visual qa," in *Advances in neural information processing systems*, 2016, pp. 361–369.
- [30] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu, "Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 7, pp. 1325–1339, 2013.
- [31] D. Mehta, H. Rhodin, D. Casas, P. Fua, O. Sotnychenko, W. Xu, and C. Theobalt, "Monocular 3d human pose estimation in the wild using improved cnn supervision," in *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, pp. 506–516.
- [32] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "RMPE: Regional multi-person pose estimation," in *ICCV*, 2017.
- [33] D. C. Luvizon, D. Picard, and H. Tabia, "2d/3d pose estimation and action recognition using multitask deep learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5137–5146.
- [34] Y. Cai, L. Ge, J. Liu, J. Cai, T.-J. Cham, J. Yuan, and N. M. Thalmann, "Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2272–2281.
- [35] Y. Cheng, B. Yang, B. Wang, W. Yan, and R. T. Tan, "Occlusion-aware networks for 3d human pose estimation in video," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 723–732.
- [36] R. A. Yeh, Y.-T. Hu, and A. G. Schwing, "Chirality nets for human pose regression," *arXiv preprint arXiv:1911.00029*, 2019.
- [37] Z. Cao, G. H. Martinez, T. Simon, S.-E. Wei, and Y. A. Sheikh, "Openpose: realtime multi-person 2d pose estimation using part affinity fields," *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [38] S. Li, L. Ke, K. Pratama, Y.-W. Tai, C.-K. Tang, and K.-T. Cheng, "Cascaded deep monocular 3d human pose estimation with evolutionary training data," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6173–6183.
- [39] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [40] B. Wandt and B. Rosenhahn, "Repnet: Weakly supervised training of an adversarial reprojection network for 3d human pose estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 7782–7791.
- [41] H. Nam, J.-W. Ha, and J. Kim, "Dual attention networks for multimodal reasoning and matching," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 299–307.
- [42] K.-H. Lee, X. Chen, G. Hua, H. Hu, and X. He, "Stacked cross attention for image-text matching," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 201–216.