# Exploiting Temporal Contexts with Strided Transformer for 3D Human Pose Estimation

Wenhao Li, Hong Liu†, Runwei Ding, Mengyuan Liu, Pichao Wang, and Wenming Yang

*Abstract*—Despite the great progress in 3D human pose estimation from videos, it is still an open problem to take full advantage of a redundant 2D pose sequence to learn representative representations for generating one 3D pose. To this end, we propose an improved Transformer-based architecture, called Strided Transformer, which simply and effectively lifts a long sequence of 2D joint locations to a single 3D pose. Specifically, a Vanilla Transformer Encoder (VTE) is adopted to model long-range dependencies of 2D pose sequences. To reduce the redundancy of the sequence, fully-connected layers in the feed-forward network of VTE are replaced with strided convolutions to progressively shrink the sequence length and aggregate information from local contexts. The modified VTE is termed as Strided Transformer Encoder (STE), which is built upon the outputs of VTE. STE not only effectively aggregates long-range information to a single-vector representation in a hierarchical global and local fashion, but also significantly reduces the computation cost. Furthermore, a full-to-single supervision scheme is designed at both full sequence and single target frame scales applied to the outputs of VTE and STE, respectively. This scheme imposes extra temporal smoothness constraints in conjunction with the single target frame supervision and hence helps produce smoother and more accurate 3D poses. The proposed Strided Transformer is evaluated on two challenging benchmark datasets, Human3.6M and HumanEva-I, and achieves state-of-the-art results with fewer parameters. Code and models are available at https://github.com/Vegetebird/StridedTransformer-Pose3D.

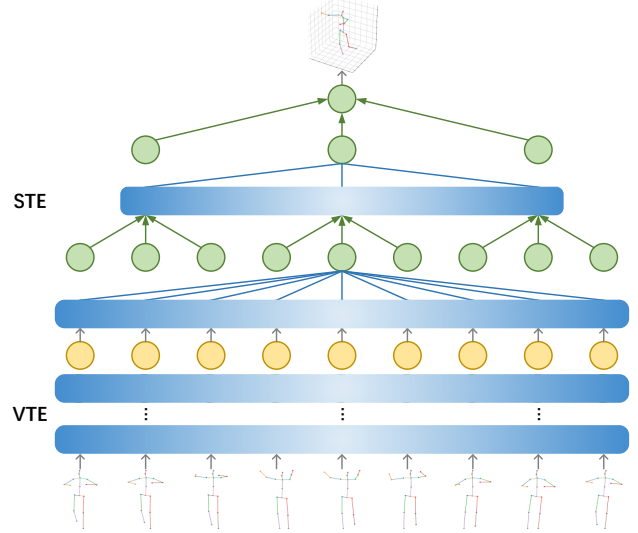*Index Terms*—3D human pose estimation, Transformer, Strided convolution.



Fig. 1: Our Strided Transformer Encoder (STE) takes the outputs of Vanilla Transformer Encoder (VTE) as input (yellow) and generates a 3D pose for the target frame as output (top). The self-attention mechanism (blue) concentrates on global context and the strided convolution (green) aggregates information from local contexts.

## I. INTRODUCTION

**3**D human pose estimation is a classic computer vision task that aims to estimate 3D joint locations of a human body from images or videos. This task has drawn tremendous attention in the past decades [1]–[4] since it plays a significant role in wide applications, such as clinic [5], computer animation [6], action recognition [7]–[16], and human-robot interaction [17], [18]. Many state-of-the-art approaches adopt a two-stage pipeline (*i.e.*, 2D-to-3D lifting method) [19]–[21], which first estimates 2D keypoints and then lifts them to 3D space. Although the 2D-to-3D lifting method benefits from the reliable performance of 2D pose detectors, it is still a highly ill-posed problem due to the inherent ambiguity in depth, since multiple 3D interpretations can be projected to the same 2D pose in the image space.

To alleviate this problem, temporal context information has been investigated by many researchers. Some methods [22]–[24] leverage past and future data in the sequence to predict the 3D pose of the target frame. For instance, Cai *et al.* [24] presented a local-to-global graph convolutional network to exploit spatio-temporal relations to estimate 3D keypoints from a 2D pose sequence. However, these approaches have small temporal receptive fields and limited temporal correlation windows, thus suffering from modeling long-range dependencies.

Vanilla Transformer [25] is developed for exploiting long-range dependencies and achieves tremendous success in natural language processing [26], [27] and computer vision [28]–[32]. It consists of a self-attention module and a position-wise feed-forward network (FFN). The self-attention module computes pairwise dot-product among all input elements to
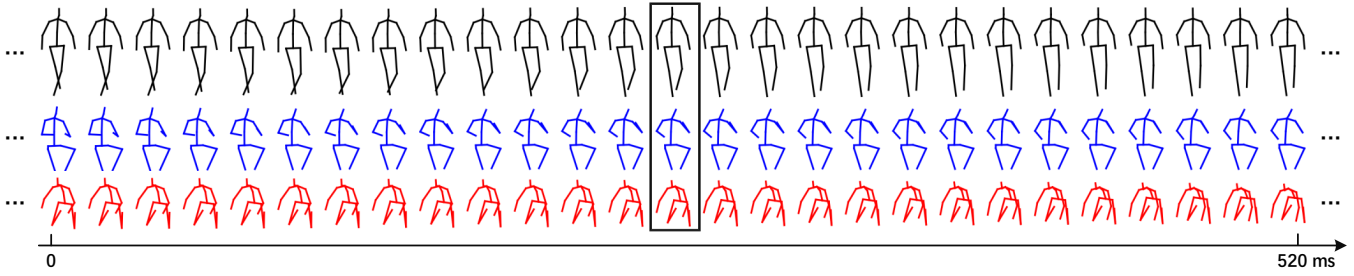
Fig. 2: Example of 2D pose sequences of 27 consecutive frames (520 ms) on Human3.6M dataset (captured from 50 Hz cameras). It contains huge redundant information as nearby poses are same. The rectangle denotes the center frame.

capture global-context information, and the FFN acts as pattern detectors over the input across all layers [33]. Such a design looks like a good choice for the 2D-to-3D pose lifting method to capture long-range dependencies. However, there are several shortcomings in the Vanilla Transformer Encoder (VTE) [25]: (i) The full-length sequence in the forward pass across all layers contains significant redundancy for video-based pose estimation as nearby poses are quite similar, as illustrated in Fig. 2. (ii) The time and memory complexity of the attention operation grows quadratically with the input length, making it very expensive to process long sequences. Thus, the receptive field may be forced to decrease in real-time applications, whereas a large receptive field is important to enhance the estimation consistency [34]. (iii) The VTE architecture is less capable to extract fine-grained local feature patterns, which is well-known to be crucial for computer vision tasks. To mitigate these issues, we propose to gradually merge nearby poses to shrink the sequence length until one representation of the target pose is acquired. An alternative is to perform pooling operation after the FFN [27]. However, lots of valuable information will be lost using pooling operation, and the local information can not be well exploited. Motivated by the previous methods [20], [34] that are able to elegantly handle variable-length sequences via temporal convolutions, we propose to replace fully-connected layers in FFN with strided convolutions to progressively reduce the sequence length. The modified Transformer is dubbed Strided Transformer Encoder (STE), as shown in Fig. 1. With the proposed STE, we can model both global and local information in a hierarchical architecture, and the computation in FFN can be traded off for constructing a deeper model to boost the model capacity.

Although the STE can aggregate long-range information to a single-pose representation, it remains a question whether this single representation is enough to represent a long sequence and how to make this representation work in improving the performance. We observe that directly supervising the model at a single target frame scale always breaks temporal smoothness among video frames, while only supervising at a full sequence scale cannot explicitly learn a specific representation for the target frame. These observations encourage us to develop a method that can effectively embed both scales into a learnable framework. Therefore, based on the outputs of VTE and STE, a full-to-single supervision scheme is designed at both full and single scales, which can impose extra temporal smoothness constraints at the full sequence scale and refine the estimation

at the single target frame scale. This scheme brings great benefits in producing smoother and more accurate 3D poses.

The proposed architecture is called Strided Transformer, as shown in Fig. 3. Extensive experiments are conducted on two standard 3D human pose estimation datasets, *i.e.*, Human3.6M [35] and HumanEva-I [36]. Experimental results show that the proposed method achieves state-of-the-art performance.

Our contributions are summarized as follows:

- We propose a new Transformer-based architecture for 3D human pose estimation called Strided Transformer, which can simply and effectively lift a long 2D pose sequence to a single 3D pose.
- To reduce the sequence redundancy and computation cost, Strided Transformer Encoder (STE) is introduced to gradually reduce the temporal dimensionality and aggregate long-range information into a single-vector representation of pose sequences in a hierarchical global and local fashion.
- A full-to-single supervision scheme is designed to impose extra temporal smoothness constraints during training at the full sequence scale and further refine the estimation at the single target frame scale.
- State-of-the-art results are achieved with fewer parameters on two commonly used benchmark datasets, making our method a strong baseline for Transformer-based 3D pose estimation.

## II. RELATED WORK

At the early stage of applying deep neural networks on 3D pose estimation task, many methods [37]–[40] learned the direct mapping from RGB images to 3D poses (*i.e.*, one-stage pose estimation). However, these methods require sophisticated architectures with high computation costs, which are impractical in realistic applications.

**Two-stage pose estimation.** Two-stage methods formulate the problem of 3D human pose estimation as 2D keypoint detection followed by 2D-to-3D lifting estimation [19], [41], [42]. Recent works show that 3D locations of body joints can be efficiently and effectively recovered using detected 2D poses from state-of-the-art 2D pose detectors, and this 2D-to-3D pose lifting method outperforms one-stage approaches. For example, Martinez *et al.* [19] lifted 2D joint locations to 3D space via a fully-connected residual network. Fang *et al.* [41] proposed a pose grammar model to encode the human
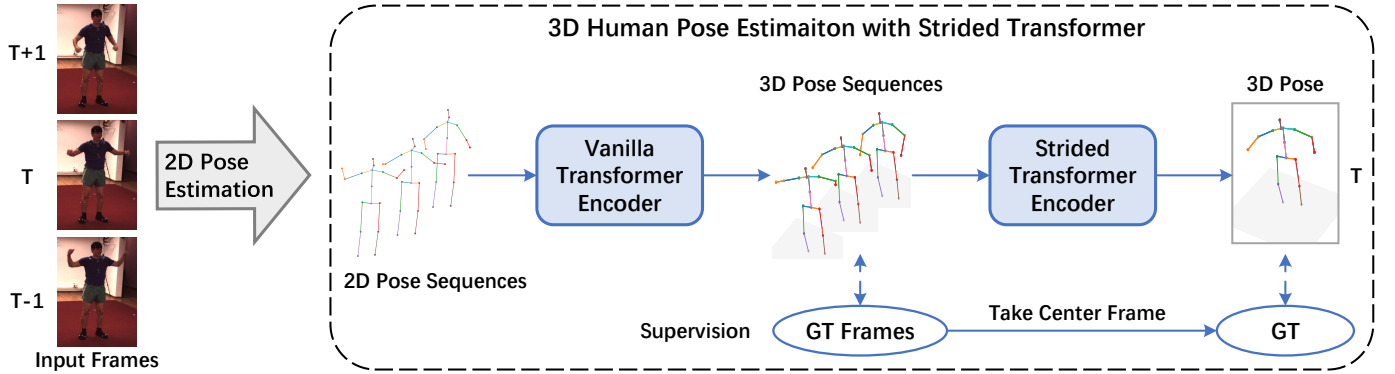
Fig. 3: Overview of our proposed Strided Transformer for predicting the 3D joint locations of the target frame (center frame) from the estimated 2D pose sequences. It mainly consists of a Vanilla Transformer Encoder (VTE) and a Strided Transformer Encoder (STE). The network first models long-range information via VTE and then aggregates the information into one target pose representation from the proposed STE. The model is trained end-to-end at both full sequence and single target frame scales.

body configuration of human poses from 2D space to 3D space. To improve the generalization of the trained 2D-to-3D pose estimator, Gong [43] introduced a pose augmentation framework that is differentiable. We also follow this two-stage pipeline because it is widely adopted among the state-of-the-art methods in this domain.

**Video pose estimation.** Recently, many approaches tried to exploit temporal information [20], [23], [24], [44] to improve the accuracy and the smoothness of the estimated 3D pose sequence. To predict temporally consistent 3D poses, Hossain *et al.* [23] designed a sequence-to-sequence network with LSTM. Pavllo *et al.* [20] introduced a fully convolutional model based on dilated temporal convolutions. Cai *et al.* [24] directly chose the 3D pose of the target frame from the outputs of the proposed graph-based method and then fed it to a refinement model. To produce smoother 3D sequences, Wang *et al.* [44] designed an U-shaped graph convolutional network and involved motion modeling into learning. However, the temporal connectivity of these architectures is inherently limited and is mainly constrained to simple sequential correlations. Different from most existing works that employed LSTM-based [23], graph-based [24], [44], or temporal convolutional networks [20], [34], [45] to exploit temporal information, we propose a Transformer-based architecture to capture long-range dependencies from input 2D pose sequences. Furthermore, compared with previous methods [24], [44] that either utilize a refinement model or use a motion loss to improve estimations, we design a full-to-single supervision scheme that refines the intermediate predictions to produce smoother and more accurate estimations.

**Visual Transformers.** Transformer models first proposed in [25] are commonly used in various language tasks. Recently, Transformers have shown promising performance in many computer vision tasks, such as object detection [46], [47] and image classification [48], [49]. DETR [46] presented a new Transformer-based design for object detection systems. ViT [48] proposed to apply a standard Transformer architecture directly to sequential image patches for image classification. METRO [50] introduced a Transformer frame-

work to reconstruct 3D human pose and mesh from a single image. However, METRO focused on the one-stage pose estimation and ignores the temporal information across frames. Unlike DETR [46], ViT [48], or METRO [50] that directly apply Transformer to images, we utilize a Transformer-based architecture to effectively map 2D keypoints to 3D poses. Additionally, efficient strided convolutions are incorporated into Transformer models to address the redundancy problem for the video-based 3D pose estimation task.

## III. METHOD

In this section, we first present an overview of the proposed Strided Transformer for 3D human pose estimation from a 2D video stream, and then show how our Transformer-based architecture learns a representative single-pose representation from redundant sequences resulting in an enhanced estimation. Finally, the complexity analysis of our network is presented.

### A. Overview

The overall framework of our proposed method is illustrated in Fig. 3. Given a sequence of the estimated 2D poses $P = \{p_1, \ldots, p_T\}$ from videos, we aim at reconstructing 3D joint locations $X \in \mathbb{R}^{J \times 3}$ for a target frame (center frame), where $p_t \in \mathbb{R}^{J \times 2}$ denotes the 2D joint locations at frame $t$, $T$ is the number of video frames, and $J$ is the number of joints. The network contains a Vanilla Transformer Encoder (VTE) followed by a Strided Transformer Encoder (STE), which is trained in a full-to-single prediction scheme at both full sequence and single target frame scales. Specifically, VTE is first used to model long-range information and is supervised by the full sequence scale to enforce temporal smoothness. Then, the proposed STE aggregates the information to generate one target pose representation and is supervised by the single target frame scale to produce more accurate estimations.

### B. Strided Transformer Encoder

Despite the substantial performance gains achieved by Transformers [25] in many computer vision tasks, the full-length token representation makes it unsuitable for many

video-based vision tasks that only require a single-vector representation of a sequence. To this end, STE is proposed to gradually compress the sequence of hidden states and model both global and local information in a hierarchical architecture. Each layer of the proposed STE consists of a multi-head self-attention (MSA) and a convolutional feed-forward network (CFFN).

*1) Multi-head self-attention:* The core mechanism of the Transformer model is MSA [25]. Suppose there are a set of queries ($Q$), keys ($K$), and values ($V$) of dimension $d_m$. Then the MSA can be computed as:

$$head_i = \text{Self-Attn}\left(QW_i^Q, KW_i^K, VW_i^V\right), \quad (1)$$

$$\text{MSA}(Q, K, V) = \text{Concat}\left(head_1, \ldots, head_h\right)W^O, \quad (2)$$

where $\text{Self-Attn}(Q, K, V) = \text{softmax}\left(QK^T/\sqrt{d_k}\right)V$ and $W_i^Q \in \mathbb{R}^{d_m \times d_k}, W_i^K \in \mathbb{R}^{d_m \times d_k}, W_i^V \in \mathbb{R}^{d_m \times d_v}$, and $W^O \in \mathbb{R}^{hd_v \times d_m}$ are parameter matrices. The hyperparameter $h$ is the number of multi-attention heads, $d_m$ is the dimension of the model, and $d_k = d_v = d_m/h$ in our implementation.

*2) Convolutional feed-forward network:* In the existing fully-connected (FC) layers in the FFN of VTE (Eq. (3)), it always maintains a full-length sequence of hidden representations across all layers with a high computation cost. It contains significant redundancy for video-based pose estimation, as nearby poses are quite similar. However, to reconstruct more accurate 3D body joints of the target frame, crucial information should be extracted from the entire pose sequences. Therefore, it requires selectively aggregating useful information.

To tackle this issue, inspired by the previous works [20], [34] that employ temporal convolutions to effectively shrink the sequence length, we make modifications to the generic FFN. Given the input feature vector $Z \in \mathbb{R}^{T \times D_{in}}$ with $T$ sequences and $D_{in}$ channels to generate an output of $(\tilde{T}, D_{out})$ features, the operation performed by FC in FFN can be formulated as:

$$\text{FC}_{t,d_{out}}(z) = \sum_{i}^{D_{in}} w_{d_{out},i} * z_{t,i}. \quad (3)$$

If 1D convolution is considered with kernel size $K$ and strided factor $S$, a strided convolution in CFFN can be computed as:

$$\text{Conv}_{S(t),c_{out}}(z) = \sum_{i}^{D_{in}} \sum_{k}^{K} w_{d_{out},i,k} * z_{S(t-\frac{K-1}{2}+k),i}. \quad (4)$$

In this way, fully-connected layers in FFN of VTE are replaced with strided convolutions. The modified VTE is termed as Strided Transformer Encoder (STE), which can be represented as:

$$\hat{Z}^{n-1} = Z^{n-1} + \text{MSA}(\text{LN}(Z^{n-1})), \quad (5)$$

$$Z^n = \text{MaxPool}(\hat{Z}^{n-1}) + \text{CFFN}(\text{LN}(\hat{Z}^{n-1})), \quad (6)$$

where $\text{LN}(\cdot)$ denotes the layer normalization, $\text{MaxPool}(\cdot)$ denotes the max pooling operation, and $n \in [1, \ldots, N]$ is the index of STE layers.

The STE is a hierarchical global and local architecture, where the self-attention mechanism models global context
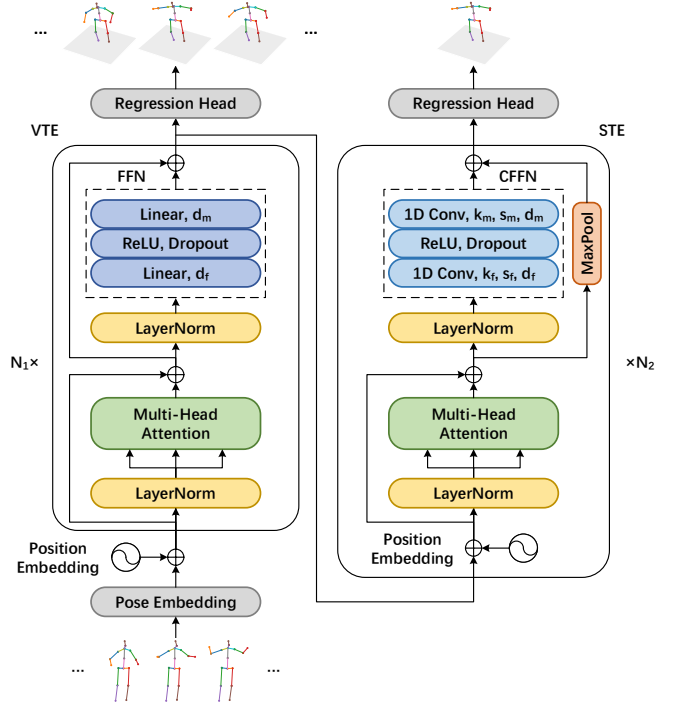


Fig. 4: The network architecture of our proposed Strided Transformer. The left is the VTE and the right is the STE. Here, $N_1$ and $N_2$ denote the number of layers of the two modules, respectively. The hyperparameters $k$, $s$, $d_m$ and $d_f$ are the kernel size, the strided factor, the dimension, and the number of hidden units. The max pooling operation is applied to the residuals to match the temporal dimensions.

and the strided convolution helps capture local contexts, as presented in Fig. 4 (right). It gradually merges the nearby poses to a short sequence length representation, as illustrated in Fig. 5. Importantly, through such a hierarchical design, the redundancy of the sequence and the computation cost can be reduced.

### C. Network Architecture

In this section, we describe how to use the proposed Transformer-based network architecture to estimate 3D human poses from a sequence of 2D poses. As shown in Fig. 5, the proposed network is composed of four components: a pose embedding, a Vanilla Transformer Encoder (VTE), a Strided Transformer Encoder (STE), and a regression head.

*1) Pose embedding:* Given a sequence of the estimated 2D poses $P \in \mathbb{R}^{T \times J \times 2}$, the pose embedding first concatenates $(x, y)$ coordinates of the $J$ joints for each frame to tokens $P' \in \mathbb{R}^{T \times (J \cdot 2)}$, and then embeds each token to a high dimensional feature $Z_0 \in \mathbb{R}^{T \times d_m}$ using a 1D convolutional layer with $d_m$ channels, followed by batch normalization, dropout, and a ReLU activation.

*2) Vanilla Transformer Encoder:* Suppose that the VTE consists of $N_1$ layers, the learnable position embedding $E_1 \in \mathbb{R}^{T \times d_m}$ is used before the first layer of VTE, which can be formulated as follows:
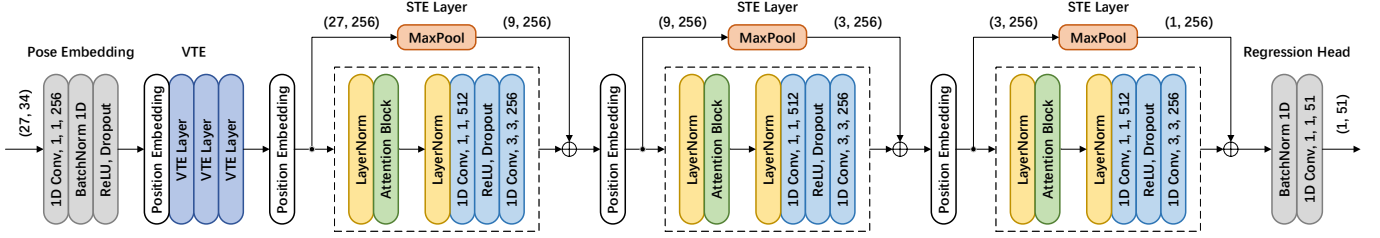
$$Z_1^0 = Z_0 + E_1. \quad (7)$$

Fig. 5: An instantiation of the proposed Strided Transformer network. It reconstructs the target 3D body joints by progressively reducing the sequence length. The input consists of 2D keypoints for a receptive field of 27 frames with $J = 17$ joints. Convolutional feed-forward networks are in blue where $(3, 3, 256)$ denotes kernels of size 3 with strided factor 3 and 256 output channels. The tensor sizes are shown in parentheses, *e.g.*, $(27, 34)$ denotes 27 frames and 34 channels. Due to strided convolutions, the max pooling operation is applied to the residuals to match the shape of subsequent tensors.

Then, given the embedded feature $Z_1^0$, the VTE layers can be represented as:

$$\hat{Z}_1^{n-1} = Z_1^{n-1} + \text{MSA}(\text{LN}(Z_1^{n-1})), \tag{8}$$
$$Z_1^n = \hat{Z}_1^{n-1} + \text{FFN}(\text{LN}(\hat{Z}_1^{n-1})), \tag{9}$$

where $n \in [1, \ldots, N_1]$ is the index of VTE layers. It can be expressed by using a function of a VTE layer $\text{VTE}(\cdot)$:

$$Z_1^n = \text{VTE}(Z_1^{n-1}). \tag{10}$$

*3) Strided Transformer Encoder:* For the STE, it is built upon the outputs of VTE and takes the $Z_1^{N_1} \in \mathbb{R}^{T \times d_m}$ as input. The learnable position embeddings $E_2 \in \mathbb{R}^{S(t) \times d_m}$ with strided factor $S$ are used for every layer of STE due to the different sequence lengths. Then, the STE layers can be represented as follows:

$$Z_2^n = \text{STE}(Z_2^{n-1} + E_2^n), \tag{11}$$

where $n \in [1, \ldots, N_2]$ is the index of STE layers, $Z_2^0 = Z_1^{N_1}$, and $\text{STE}(\cdot)$ denotes the function of an STE layer whose details can be found in Eq. (5) and Eq. (6).

*4) Regression head:* In order to perform the regression, a batch normalization and a 1D convolutional layer are applied to the outputs of VTE and STE, $Z_1^{N_1} \in \mathbb{R}^{T \times d_m}$ and $Z_2^{N_2} \in \mathbb{R}^{1 \times d_m}$, respectively. Finally, the outputs of 3D pose prediction are $\tilde{X}$ and $X$, where $\tilde{X} \in \mathbb{R}^{T \times J \times 3}$ and $X \in \mathbb{R}^{J \times 3}$ are predictions of the 3D pose sequence and the 3D joint locations of the target frame, respectively.

### D. Full-to-Single Prediction

The iterative refinement scheme, aimed at producing predictions in multiple processing stages, is effective for 3D pose estimation [24], [37]. Motivated by the success of such iterative processing, we also consider a refinement scheme. A full-to-single scheme is proposed to incorporate both full sequence and single target frame scales constraints into the framework. This scheme further refines the intermediate predictions to produce more accurate estimations rather than using a single component with a single output. More precisely, the full sequence scale can enforce temporal smoothness and the single target frame scale helps learn a specific representation for the target frame.

*1) Full sequence scale:* The first step is to supervise at full sequence scale by imposing extra temporal smoothness constraints during training from the output of VTE followed by a regression head. A sequence loss $\mathcal{L}_f$ is used to improve upon single frame predictions for temporal consistency over a sequence. This loss ensures that the estimated 3D pose sequences $\tilde{X} \in \mathbb{R}^{T \times J \times 3}$ coincide with the ground truth 3D joint sequences $Y \in \mathbb{R}^{T \times J \times 3}$:

$$\mathcal{L}_f = \sum_{t=1}^{T} \sum_{i=1}^{J} \left\| Y_i^t - \tilde{X}_i^t \right\|_2, \tag{12}$$

where $\tilde{X}_i^t$ and $Y_i^t$ represent the sequence of estimated 3D poses and ground truth 3D joint locations of joint $i$ at frame $t$, respectively.

*2) Single target frame scale:* In the second step, the supervision is adopted on the output of STE followed by a regression head. A single-frame loss $\mathcal{L}_s$ is used to refine the estimation at the single target frame scale. It minimizes the distance between the estimated 3D pose $X \in \mathbb{R}^{J \times 3}$ and the target ground truth 3D joint annotation $Y \in \mathbb{R}^{J \times 3}$:

$$\mathcal{L}_s = \sum_{i=1}^{J} \left\| Y_i - X_i \right\|_2, \tag{13}$$

where $X_i$ and $Y_i$ represent the target frame's estimated 3D pose and ground truth 3D joint locations of joint $i$, respectively.

*3) Loss function:* In our implementation, the model is supervised at both full sequence scale and single target frame scale. We train the entire network in an end-to-end manner with the total loss:

$$\mathcal{L} = \lambda_f \mathcal{L}_f + \lambda_s \mathcal{L}_s, \tag{14}$$

where $\lambda_f$ and $\lambda_s$ are weighting factors.

### E. Complexity Analysis

In this section, we use floating-point operations (FLOPs) to measure the computational cost and analyze the compression ratio of our proposed Strided Transformer network. Given the sequence length $t$, dimension $d_m = d_f/2 = d$, strided factor

$s$, and kernel size $k$, the FLOPs of a VTE layer $\mathcal{F}_{VTE}^n$ and an STE layer $\mathcal{F}_{STE}^n$ can be computed by:

$$\begin{aligned} \mathcal{F}_{VTE}^n(t,d) &= \mathcal{F}_{MSA}^n(t,d) + \mathcal{F}_{FFN}^n(t,d) \\ &= 8td^2 + 2t^2d, \end{aligned} \quad (15)$$

$$\begin{aligned} \mathcal{F}_{STE}^n(t,d,s) &= \mathcal{F}_{MSA}^n(t,d,s) + \mathcal{F}_{CFFN}^n(t,d,s) \\ &= (6 + 2s^{-1}k)td^2 + 2t^2d, \end{aligned} \quad (16)$$

where $\mathcal{F}_{MSA}^n$, $\mathcal{F}_{FFN}^n$, and $\mathcal{F}_{CFFN}^n$ are the FLOPs of the MSA, FFN, and CFFN, respectively.

Then if we consider $N$ layers of VTE and STE with input sequence length $T$, dimension $D$, strided factor $S$, and kernel size $K$, the encoder-wise FLOPs of VTE $\mathcal{F}_{VTE}$ can be formulated as:

$$\mathcal{F}_{VTE} = N\mathcal{F}_{VTE}^n = N(8TD^2 + 2T^2D), \quad (17)$$

the encoder-wise FLOPs of STE $\mathcal{F}_{STE}$ can be formulated as:

$$\begin{aligned} \mathcal{F}_{STE} &= \sum_{n=1}^{N} \mathcal{F}_{STE}^n \\ &= \sum_{n=1}^{N} \left[ \left( \frac{6 + 2KS^{-1}}{S^{n-1}} \right) TD^2 + \frac{2}{S^{2(n-1)}} T^2D \right]. \end{aligned} \quad (18)$$

For our 27-frame Strided Transformer, which contains $N_1$ VTE layers and $N_2$ STE layers with $N_1 = N_2 = N = 3$, $S = 3$, and $K = 3$. In this case, the compression ratio $\alpha$ can be computed by:

$$\alpha = \frac{2\mathcal{F}_{VTE}}{\mathcal{F}_{VTE} + \mathcal{F}_{STE}} = \frac{2}{1 + \beta}, \quad (19)$$

where

$$\beta = \frac{\mathcal{F}_{STE}}{\mathcal{F}_{VTE}} = \frac{468D + 91T}{972D + 243T}. \quad (20)$$

We have $\lim_{D \to \infty} \alpha = 1.35$ with a fixed $T$. Thus, the compression ratio $\alpha$ of our 27-frame Strided Transformer is 1.35.

## IV. EXPERIMENTS

### A. Datasets and Evaluation

The proposed method is evaluated on two challenging benchmark datasets, *i.e.*, Human3.6M [35] and HumanEva-I [36]. Human3.6M dataset is the largest publicly available dataset for 3D human pose estimation, which consists of 3.6 million images captured from 4 synchronized cameras with 50 Hz. There are 7 professional subjects performing 15 daily activities such as "Waiting", "Smoking", and "Posing". Following the standard protocol in prior works [20], [56], [57], 5 subjects (S1, S5, S6, S7, S8) are used for training and 2 subjects (S9 and S11) are used for evaluation. The frames from all views are trained by a single model for all actions. HumanEva-I is a much smaller dataset with fewer subjects and actions compared to Human3.6M. Following [20], [22], our model is trained for all subjects (S1, S2, S3) and all actions (Walk, Jog, Box).

Three standard evaluation protocols are used in the experiments. The mean per joint position error (MPJPE) is the average Euclidean distance between the ground truth and predicted positions of the joints, which is referred to as protocol #1 in many works [41], [58]. Procrustes analysis MPJPE (P-MPJPE) is adopted, where the estimated 3D pose is aligned to the ground truth in translation, rotation, and scale. This protocol is referred to as protocol #2 [19], [23]. Following [20], [44], [45], we also report the mean per joint velocity error (MPJVE) corresponding to the MPJPE of the first derivative of the 3D pose sequences. This metric measures the smoothness of predictions over time and is vital for video-based 3D pose estimation.

### B. Implementation Details

In our experiments, the proposed Strided Transformer contains $N_1 = N_2 = 3$ encoder layers, $h = 8$ attention heads, $d_m = 256$ dimensions, and $d_f = 512$ hidden units for both VTE and STE. The kernel sizes $k_f$ and $k_m$ are set to 1 and 3 in all STE layers, respectively. The strided factor $s_f$ is set to 1, and $s_m$ is set to $\{3, 3, 3\}$ for the receptive field of 27 frames, $\{9, 3, 3\}$ for 81, $\{3, 9, 9\}$ for 243, and $\{3, 9, 13\}$ for 351. The weighting factors $\lambda_f$ and $\lambda_s$ are set to 1.

All experiments are conducted on the PyTorch framework with one GeForce GTX 3090 GPU. The network is trained using Amsgrad optimizer. An initial learning rate of 0.001 is used with a shrink factor of 0.95 applied after each epoch. The same refine module as [24], [44] is adopted. We only apply horizontal flip augmentation during training/test stages. The 2D poses can be obtained by performing any classic 2D pose detections or directly using the 2D ground truth. Following [20], [59], the cascaded pyramid network (CPN) [60] is used for Human3.6M and Mask R-CNN [61] is adopted for HumanEva-I to obtain 2D poses for a fair comparison.

### C. Comparison with State-of-the-art Results

Our method is compared with previous state-of-the-art approaches on Human3.6M dataset. The performance of our 351-frame model with CPN input is reported in Table I. Our method outperforms the state-of-the-art methods on Human3.6M under all metrics (43.7 mm on protocol #1 and 35.2 mm on protocol #2).

Table II compares the computational complexity, MPJPE, and frame per second (FPS) with several state-of-the-art methods in different receptive fields on Human3.6M. Our model is lightweight and the number of parameters hardly increases with the increased receptive fields, which is practical for real-time applications. Compared with temporal convolutional networks [20], [45], our proposed Transformer-based network requires fewer total parameters with competitive performance for 3D pose estimation in videos. Besides, even though the inference speed of the proposed model is lower than [20], [45], it still has an acceptable FPS for real-time inference. Fig. 7 shows some qualitative comparisons with state-of-the-art methods [20], [34], which indicates that our methods can produce more accurate 3D predictions.

To further explore the upper bound of our method, the results from 2D ground truth inputs are reported in Table III. It can be seen that our method achieves the best result (28.5 mm in MPJPE), outperforming all other methods. This

TABLE I: Quantitative comparisons on Human3.6M under protocol #1 and protocol #2, where † indicates the temporal information used in each method. Best in bold, second-best underlined.

| **Protocol #1** | Dir. | Disc | Eat | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Smoke | Wait | WalkD. | Walk | WalkT. | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Martinez et al. [19] 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 | 49.5 | 52.4 | 62.9 |
| Fang et al. [41] AAAI'18 | 50.1 | 54.3 | 57.0 | 57.1 | 66.6 | 73.3 | 53.4 | 55.7 | 72.8 | 88.6 | 60.3 | 57.7 | 62.7 | 47.5 | 50.6 | 60.4 |
| Lee et al. [22] 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 |
| Xu et al. [42] CVPR'21 | 45.2 | 49.9 | 47.5 | 50.9 | 54.9 | 66.1 | 48.5 | 46.3 | 59.7 | 71.5 | 51.4 | 48.6 | 53.9 | 39.9 | 44.1 | 51.9 |
| Gong et al. [43] CVPR'21 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 50.2 |
| Cai et al. [24] 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 |
| Pavllo et al. [20] 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 |
| Lin et al. [51] 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 |
| Xu et al. [52] 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 |
| Liu et al. [34] CVPR'20 † | 41.8 | 44.8 | 41.1 | 44.9 | 47.4 | 54.1 | 43.4 | 42.2 | 56.2 | 63.6 | 45.3 | 43.5 | 45.3 | 31.3 | 32.2 | 45.1 |
| Zeng et al. [53] ECCV'20 † | 46.6 | 47.1 | 43.9 | 41.6 | 45.8 | 49.6 | 46.5 | 40.0 | 53.4 | 61.1 | 46.1 | 42.6 | 43.1 | 31.5 | 32.6 | 44.8 |
| Wang et al. [44] ECCV'20 † | 40.2 | 42.5 | 42.6 | 41.1 | 46.7 | 56.7 | 41.4 | 42.3 | 56.2 | 60.4 | 46.3 | 42.2 | 46.2 | 31.7 | 31.0 | 44.5 |
| Chen et al. [45] TCSVT'21 † | 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 |
| Ours † | 40.3 | 43.3 | 40.2 | 42.3 | **45.6** | 52.3 | 41.8 | 40.5 | 55.9 | 60.6 | **44.2** | 43.0 | 44.2 | **30.0** | 30.2 | 43.7 |
| **Protocol #2** | Dir. | Disc | Eat | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Smoke | Wait | WalkD. | Walk | WalkT. | Avg. |
| Martinez et al. [19] 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 |
| Pavlakos et al. [54] 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 |
| Liu et al. [55] ECCV'20 | 35.9 | 40.0 | 38.0 | 41.5 | 42.5 | 51.4 | 37.8 | 36.0 | 48.6 | 56.6 | 41.8 | 38.3 | 42.7 | 31.7 | 36.2 | 41.2 |
| Gong et al. [43] CVPR'21 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 39.1 |
| Cai et al. [24] 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 |
| Lin et al. [51] 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 | 28.9 | 36.8 |
| Pavllo et al. [20] 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 |
| Xu et al. [52] 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 |
| Liu et al. [34] CVPR'20 † | 32.3 | 35.2 | 33.3 | 35.8 | 35.9 | 41.5 | 33.2 | 32.7 | 44.6 | 50.9 | 37.0 | 32.4 | 37.0 | 25.2 | 27.2 | 35.6 |
| Wang et al. [44] ECCV'20 † | 32.9 | 35.2 | 35.6 | 34.4 | 36.4 | 42.7 | **31.2** | 32.5 | 45.6 | 50.2 | 37.3 | 32.8 | 36.3 | 26.0 | **23.9** | 35.5 |
| Ours † | 32.7 | 35.5 | 32.5 | 35.4 | 35.9 | 41.6 | 33.0 | 31.9 | 45.1 | 50.1 | 36.3 | 33.5 | 35.1 | 23.9 | 25.0 | 35.2 |

TABLE II: Quantitative comparisons with state-of-the-art methods in different receptive fields on Human3.6M. The computational complexity, MPJPE, and frame per second (FPS) are reported. FPS is computed on a single GeForce GTX 2080 Ti GPU.

| Model | $T$ | Param (M) | FLOPs (G) | MPJPE (mm) | FPS |
|---|---|---|---|---|---|
| Pavllo et al. [20] | 27 | 8.56 | 0.017 | 48.8 | 1492 |
| Pavllo et al. [20] | 81 | 12.75 | 0.025 | 47.7 | 1121 |
| Pavllo et al. [20] | 243 | 16.95 | 0.033 | 46.8 | 863 |
| Chen et al. [45] | 27 | 31.88 | 0.061 | 45.3 | 410 |
| Chen et al. [45] | 81 | 45.53 | 0.088 | 44.6 | 315 |
| Chen et al. [45] | 243 | 59.18 | 0.116 | 44.1 | 264 |
| Ours (27 frames) | 27 | 4.01 | 0.128 | 46.9 | 118 |
| Ours (81 frames) | 81 | 4.06 | 0.392 | 45.4 | 112 |
| Ours (243 frames) | 243 | 4.23 | 1.372 | 44.0 | 108 |
| Ours (351 frames) | 351 | 4.34 | 2.142 | **43.7** | 105 |

demonstrates if a more robust 2D pose detection is available, our Strided Transformer can produce more accurate 3D poses.

As shown in Table IV, with the supervision of full sequence scale, our method reduces the MPJVE by 15.4% (from 2.6 mm to 2.2 mm), achieving smoother predictions with lower MPJVE than other models. It indicates that the full-to-single supervision scheme can enhance temporal smoothness and produce vastly smoother poses.

To evaluate the generalizability of our model to smaller datasets, experiments are conducted on HumanEva-I based on Mask R-CNN 2D detections and 2D ground truth. The results in Table V demonstrate that our method achieves promising results on all kinds of actions.

## D. Ablation Studies

**Input sequence length.** The MPJPE results of our model with different sequence lengths (between 1 and 351) on Human3.6M are shown in Fig. 6 (a). It can be seen that our proposed method obtains larger gains under both 2D pose inputs (CPN and GT) with more input frames used for predictions, but the error saturates past a certain point. This is expected since directly lifting 3D poses from disjointed 2D poses leads to temporally incoherent outputs [62]. It is worth mentioning that our method gets a better result with $T = 351$ (43.7 mm) than $T = 243$ (44.0 mm), while the performance decreases with longer inputs ($T > 243$) in [34]. This indicates that our method equipped with the global self-attention mechanism is powerful in modeling long-range dependencies. Meanwhile, with the help of STE, our method can learn the representative representation from long sequences. Next, we choose $T = 27$ on Human3.6M in the following ablation experiments as a compromise between the accuracy and computational complexity.

**2D detections.** For the 2D-to-3D pose lifting task, the accuracy of the 2D detections directly influences the results of 3D pose estimation [19]. To show the effectiveness of our method on different 2D pose detectors, we carry out experiments with the detections from Stack Hourglass (SH) [63], Detectron [20], and CPN [60]. Moreover, to test the tolerance of our method to different levels of noise, we also train our network by 2D ground truth (GT) with various levels of additive Gaussian noises. The results are shown in Fig. 6 (b). It can be observed that the MPJPE of 3D poses increases linearly with the two-norm errors of 2D poses. Besides, our method performs well on different 2D inputs, indicating the effectiveness and robustness of our method.

**Model hyperparameters.** As shown in Table VI, we first

TABLE III: Quantitative comparisons of MPJPE in millimeter on Human3.6M under protocol #1, using ground truth 2D joint locations as input. † means the method utilizing temporal information. Best in bold.

| Protocol #1 | Dir. | Disc | Eat | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Smoke | Wait | WalkD. | Walk | WalkT. | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Martinez et al. [19] ICCV'17 | 37.7 | 44.4 | 40.3 | 42.1 | 48.2 | 54.9 | 44.4 | 42.1 | 54.6 | 58.0 | 45.1 | 46.4 | 47.6 | 36.4 | 40.4 | 45.5 |
| Lee et al. [22] ECCV'18 † | 32.1 | 36.6 | 34.3 | 37.8 | 44.5 | 49.9 | 40.9 | 36.2 | 44.1 | 45.6 | 35.3 | 35.9 | 30.3 | 37.6 | 35.5 | 38.4 |
| Pavllo et al. [20] CVPR'19 † | 35.2 | 40.2 | 32.7 | 35.7 | 38.2 | 45.5 | 40.6 | 36.1 | 48.8 | 47.3 | 37.8 | 39.7 | 38.7 | 27.8 | 29.5 | 37.8 |
| Cai et al. [24] ICCV'19 † | 32.9 | 38.7 | 32.9 | 37.0 | 37.3 | 44.8 | 38.7 | 36.1 | 41.0 | 45.6 | 36.8 | 37.7 | 37.7 | 29.5 | 31.6 | 37.2 |
| Xu et al. [42] CVPR'21 | 35.8 | 38.1 | 31.0 | 35.3 | 35.8 | 43.2 | 37.3 | 31.7 | 38.4 | 45.5 | 35.4 | 36.7 | 36.8 | 27.9 | 30.7 | 35.8 |
| Liu et al. [34] CVPR'20 † | 34.5 | 37.1 | 33.6 | 34.2 | 32.9 | 37.1 | 39.6 | 35.8 | 40.7 | 41.4 | 33.0 | 33.8 | 33.0 | 26.6 | 26.9 | 34.7 |
| Chen et al. [45] TCSVT'21 † | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 32.3 |
| Zeng et al. [53] ECCV'20 † | 34.8 | 32.1 | 28.5 | 30.7 | 31.4 | 36.9 | 35.6 | 30.5 | 38.9 | 40.5 | 32.5 | 31.0 | 29.9 | 22.5 | 24.5 | 32.0 |
| Ours | **27.1** | **29.4** | **26.5** | **27.1** | **28.6** | **33.0** | **30.7** | **26.8** | **38.2** | **34.7** | **29.1** | **29.8** | **26.8** | **19.1** | **19.8** | **28.5** |

TABLE IV: Results show the velocity error (MPJPV) of our methods and other state-of-the-arts on Human3.6M. Here, * denotes our result without the supervision of full sequence scale. Best in bold.

| MPJPV | Dir. | Disc | Eat | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Smoke | Wait | WalkD. | Walk | WalkT. | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pavllo et al. [20] CVPR'19 | 3.0 | 3.1 | 2.2 | 3.4 | 2.3 | 2.7 | 2.7 | 3.1 | 2.1 | 2.9 | 2.3 | 2.4 | 3.7 | 3.1 | 2.8 | 2.8 |
| Lin et al. [51] BMVC'19 | 2.7 | 2.8 | 2.1 | 3.1 | 2.0 | 2.5 | 2.5 | 2.9 | 1.8 | 2.6 | 2.1 | 2.3 | 3.7 | 2.7 | 3.1 | 2.7 |
| Chen et al. [45] TCSVT'21 | 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 |
| Wang et al. [44] ECCV'20 | **2.3** | **2.5** | 2.0 | **2.7** | 2.0 | 2.3 | **2.2** | **2.5** | 1.8 | 2.7 | 1.9 | 2.0 | **3.1** | **2.2** | 2.5 | 2.3 |
| Ours * | 2.8 | 2.8 | 2.1 | 3.2 | 2.2 | 2.5 | 2.6 | 2.8 | 1.8 | 2.4 | 2.1 | 2.3 | 3.5 | 3.0 | 2.6 | 2.6 |
| Ours | 2.4 | **2.5** | **1.8** | 2.8 | **1.8** | **2.2** | **2.2** | **2.5** | **1.5** | **2.0** | **1.8** | **1.9** | 3.2 | 2.5 | **2.1** | **2.2** |

TABLE V: Quantitative results on HumanEva-I dataset under protocol #2. Best in bold, second-best underlined. (MRCNN) - Mask-RCNN; (GT) - 2D ground truth.

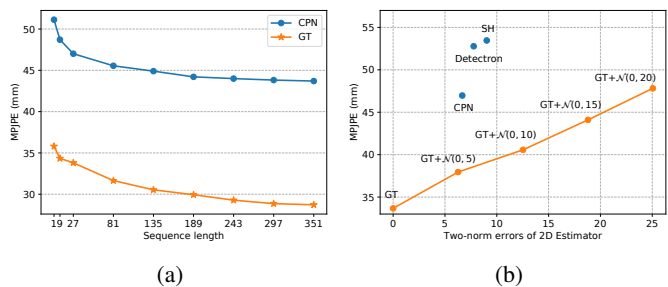| | Walk | | | Jog | | | Box | | | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|
| | S1 | S2 | S3 | S1 | S2 | S3 | S1 | S2 | S3 | |
| Martinez et al. [19] | 19.7 | 17.4 | 46.8 | 26.9 | 18.2 | 18.6 | - | - | - | - |
| Pavlakos et al. [37] | 22.3 | 19.5 | **29.7** | 28.9 | 21.9 | 23.8 | - | - | - | - |
| Lee et al. [22] | 18.6 | 19.9 | _30.5_ | 25.7 | 16.8 | 17.7 | 42.8 | 48.1 | 53.4 | 30.3 |
| Pavllo et al. [20] | **13.9** | _10.2_ | 46.6 | _20.9_ | **13.1** | 13.8 | _23.8_ | _33.7_ | _32.0_ | _23.1_ |
| Ours (T = 27 MRCNN) | _14.0_ | **10.0** | 32.8 | **19.5** | _13.6_ | _14.2_ | **22.4** | **21.6** | **22.5** | **18.9** |
| Ours (T = 27 GT) | 9.7 | 7.6 | 15.8 | 12.3 | 9.4 | 11.2 | 14.8 | 12.9 | 16.5 | 12.2 |

Fig. 6: (a) Ablation studies on different sequence lengths of our method on Human3.6M with the MPJPE metric. (b) The impact of 2D detections on Human3.6M. Here, $\mathcal{N}(0, \sigma^2)$ represents the Gaussian noise with mean zero and $\sigma$ is the standard deviation. (CPN) - Cascaded Pyramid Network; (SH) Stack Hourglass; (GT) - 2D ground truth.

analyze the effect of the number of VTE layers. Empirically, it can be found that the performance cannot be improved when naively stacking multiple standard Transformer encoder layers. Notably, our model equipped with STE is more accurate at the same number of Transformer encoder layers and comparable model parameters. For example, our method ($N_1 = 3$ and $N_2 = 3$) has better performance and fewer FLOPs than the model of $N_1 = 6$ at the same $d_m = 256$ and $d_f = 512$ (46.9 mm vs. 47.9 mm, 0.128G vs. 0.174G). In addition, our STE ($N_2 = 3$, 0.041G) also has fewer FLOPs than standard Transformer encoder ($N_1 = 3$, 0.087G) with similar parameters, which achieves $2.1\times$ less computation. It verifies the effectiveness of our proposed STE in reducing computation cost and boosting performance. Then, we investigate the influence of various hyperparameters combinations to find the optimal network architecture. It can be observed that using 3 encoder layers of both VTE and STE modules, 256 dimensions, and 512 hidden units achieves the best performance.

**Strided factor.** We observe that the strided factor of STE used in our Strided Transformer has an impact on the estimation performance. Here, we study the influence of using different design choices of strided factor of STE. The experimental results are depicted in Table VII. It shows that using a strided factor $s_m = \{3, 3, 3\}$ has the best performance. This

demonstrates the benefit of gradually reducing the temporal dimensionality with a small strided factor.

**Prediction scheme.** We further examine the proposed prediction scheme of full sequence scale and single target frame scale by using five different designs: (i) Full: the STE of our proposed method is replaced with VTE, and the new architecture is only supervised by the full sequence scale (the sequence loss). (ii) Single: the proposed method is only supervised by the single target frame scale (single-frame loss). (iii) Full-to-full: the architecture consists of six VTE layers, whose first three layers and final three layers are both supervised by the sequence loss. (iv) Single-to-single: VTE and STE of the proposed method are both supervised by the single-frame loss. (v) Full-to-single: our proposed method. In Table VIII, it can be observed that the schemes of considering only one prediction manner (i, ii, iii, iv) decay performance, and our full-to-single prediction scheme (v) is the best. The empirical results indicate that our proposed full-to-single mechanism is
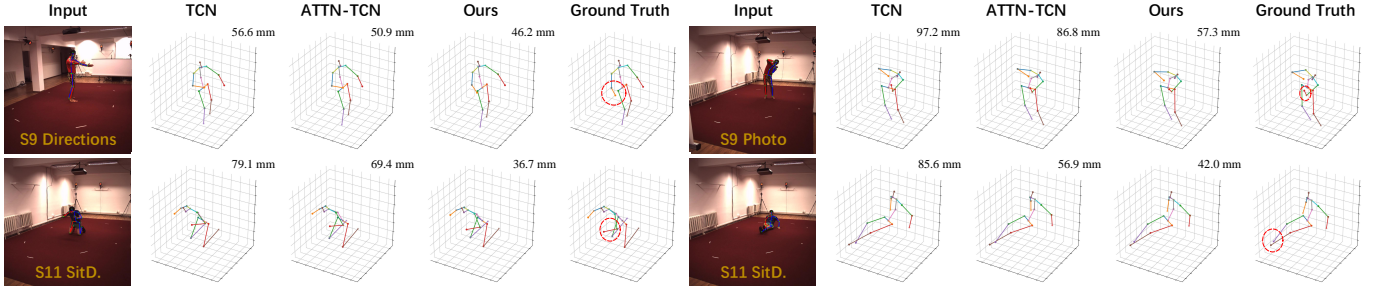
Fig. 7: Qualitative comparisons with the previous state-of-the-art methods, TCN [20] and ATTN-TCN [34] on Human3.6M dataset. Wrong estimations are highlighted by red circles.

TABLE VI: Ablation study on the hyperparameters of our model on Human3.6M under protocol #1. $N_1$ and $N_2$ are the number of VTE and STE layers, respectively. $d_m$ and $d_f$ are the dimensions and the number of hidden units.

| $N_1$ | $N_2$ | $d_m$ | $d_f$ | Param (M) | FLOPs (G) | MPJPE (mm) |
|---|---|---|---|---|---|---|
| 2 | - | 512 | 2048 | 6.36 | 0.342 | 47.9 |
| 3 | - | 512 | 2048 | 9.51 | 0.514 | 47.8 |
| 4 | - | 512 | 2048 | 12.66 | 0.685 | 48.0 |
| 5 | - | 512 | 2048 | 15.82 | 0.856 | 48.4 |
| 6 | - | 512 | 2048 | 18.97 | 1.028 | 49.3 |
| 2 | - | 256 | 512 | 1.08 | 0.058 | 47.8 |
| 3 | - | 256 | 512 | 1.61 | 0.087 | 47.6 |
| 4 | - | 256 | 512 | 2.13 | 0.116 | 47.8 |
| 5 | - | 256 | 512 | 2.66 | 0.145 | 47.7 |
| 6 | - | 256 | 512 | 3.19 | 0.174 | 47.9 |
| - | 3 | 256 | 512 | 2.42 | 0.041 | 48.0 |
| 2 | 3 | 256 | 512 | 3.48 | 0.099 | 47.4 |
| 3 | 3 | 256 | 512 | 4.01 | 0.128 | **46.9** |
| 2 | 3 | 512 | 2048 | 22.18 | 0.589 | 47.4 |
| 3 | 3 | 512 | 2048 | 25.33 | 0.761 | 47.3 |

TABLE VII: Ablation study on the strided factor of STE with the receptive field $T = 3 \times 3 \times 3 = 27$. The evaluation is performed on Human3.6M under protocol #1.

| Layers | Strided factor | MPJPE (mm) |
|---|---|---|
| 3 | 3, 3, 3 | **46.9** |
| 3 | 3, 9, 1 | 47.5 |
| 3 | 9, 3, 1 | 47.3 |
| 2 | 3, 9 | 47.2 |
| 2 | 9, 3 | 47.1 |
| 1 | 27 | 47.7 |

TABLE VIII: Ablation study on different prediction schemes. The evaluation is performed on Human3.6M under protocol #1. $\Delta$ represents the performance gap between the methods and ours.

| Prediction scheme | MPJPE (mm) | $\Delta$ |
|---|---|---|
| Full | 47.9 | 1.0 |
| Single | 48.3 | 1.4 |
| Full-to-full | 47.4 | 0.5 |
| Single-to-single | 48.5 | 1.6 |
| Full-to-single | **46.9** | - |

TABLE IX: Ablation study on each component of our network architecture on Human3.6M under protocol #1.

| Method | MPJPE (mm) |
|---|---|
| Ours, proposed | **46.9** |
| Ours, intermediate predictions | 48.1 |
| Ours, Pooling Transformer | 47.3 |
| w/o VTE | 48.0 |
| w/o STE | 47.6 |

local contexts to aggregate information. Removing VTE (only trained with single-frame loss) leads to a 1.1 mm increase in MPJPE error. Besides, removing STE (only trained with sequence loss) increases the MPJPE to 47.6 mm. These results validate the importance of both VTE and STE modules in our Strided Transformer, where VTE mainly models long-range information and STE focuses on aggregating information in a hierarchical global and local fashion.

crucial for performance improvement.

**Model components.** As shown in Table IX, an ablation study is performed to assess the effectiveness of different components of our method. We select the center frame of intermediate predictions from VTE as final results, which increases the MPJPE by 1.2 mm (from 46.9 mm to 48.1 mm). It proves that the scheme of intermediate supervision can further improve estimation accuracy. Next, we perform pooling operation after FFN of VTE following [27] and then replace STE of our proposed method with it. The new architecture is termed as Pooling Transformer, and its error increases by 0.4 mm, which highlights that our STE can preserve more valuable information than Pooling Transformer by exploiting

### E. Qualitative Results

**Attention visualization.** Our method is easily interpretable through visualizing the attention score across frames to explain what the target frame relies on. Visualization results of the multi-head attention maps of the first attention layers from VTE and STE (243-frame model) are shown in Fig. 8. The left map shows strong attention close to the input frames [64], [65], while the right map mainly pays strong attention to the center frame across all the sequences. This is expected since the proposed full-to-single strategy enables the VTE and STE modules to learn different representations: (i) VTE selectively identifies important sequences that are close to the input frames and enforces temporal consistency across frames. (ii) STE learns a specific representation from the input sequences
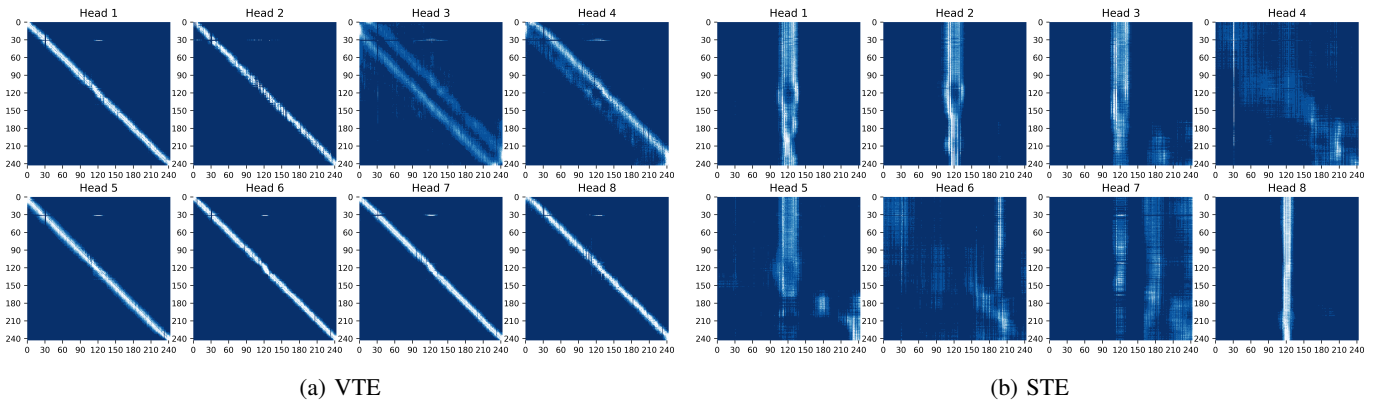
(a) VTE

(b) STE

Fig. 8: Multi-head attention maps ($h = 8$) from VTE and STE of our 243-frame model. It illustrates that the self-attention mechanism systematically assigns a weight distribution to frames, all of which might benefit the inference. Brighter color indicates higher attention score.
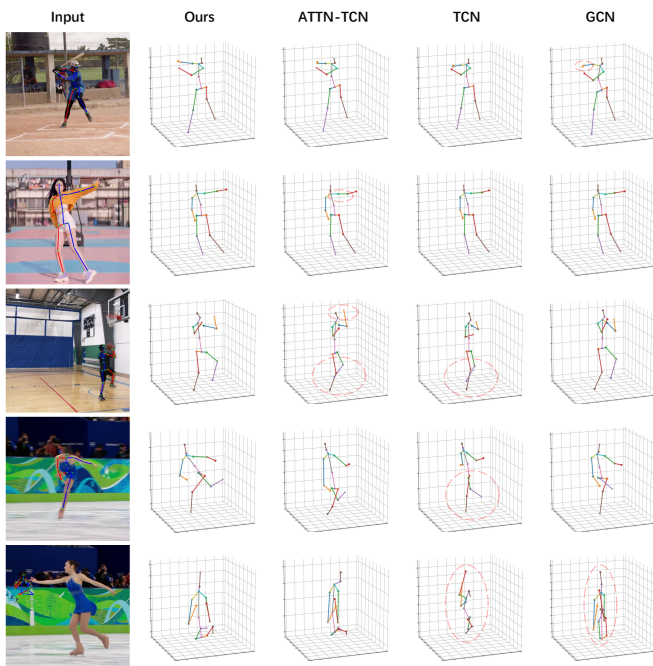


Fig. 9: Qualitative comparisons on challenging in-the-wild videos with previous state-of-the-art methods, ATTN-TCN [34], TCN [20], and GCN [24]. The last row shows the failure case, where the 2D detector has failed badly.

using both past and future data, improving the representation ability of features to reach an optimal inference for the target frame. Note that a few attention head maps are sparse due to the different temporal patterns or semantics.

**3D reconstruction visualization.** We further evaluate our method on challenging in-the-wild videos from YouTube. Fig. 9 shows the qualitative comparisons with the previous state-of-the-art methods [20], [24], [34]. We use the same 2D detector (cascaded pyramid network [60]) to obtain 2D poses and then feed them to the models for a fair comparison. Despite the challenging samples with complex actions and fast movements, the proposed method can produce realistic and

structurally plausible 3D predictions outperforming previous works. This demonstrates our method is robust to partial occlusions and tolerant to depth ambiguity. The last row shows the failure case caused by a big 2D detection error.

## V. CONCLUSION

In this work, we investigate the suitableness of applying a Transformer-based network to the task of video-based 3D human pose estimation. From the proposed Strided Transformer with Strided Transformer Encoder (STE) and full-to-single supervision scheme, we show how the representative single-pose representation can be learned from redundant sequences. The key is to reasonably use strided convolutions in the Transformer architecture to aggregate long-range information into a single-vector pose in a hierarchical global and local fashion. Meanwhile, the computation cost can be reduced significantly. Moreover, our full-to-single supervision scheme enhances temporal smoothness and further refines the representation for the target frame. Comprehensive experiments on two benchmark datasets demonstrate that our method achieves superior performance compared with state-of-the-art methods.

Although our method can reduce the computation cost of Transformers, the computational complexity and runtime cost of our method are still larger than temporal convolutional networks [20], [45], indicated in Table II. It is well acknowledged that the strong performance of Transformers comes at high computational costs. Note that the scope of this paper only targets improving FFN in the Transformer model. Future works may include designing a more efficient self-attention mechanism and extending our Strided Transformer to solve multi-view 3D human pose estimation. In addition, we hope that our approach would bring inspiration to the field of skeleton-based representation learning, *e.g.*, action recognition, motion prediction, pose tracking, and so on.

## REFERENCES

[1] I. Radwan, A. Dhall, and R. Goecke, "Monocular image 3d human pose estimation under self-occlusion," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2013, pp. 1888–1895.

[2] S. Li and A. B. Chan, "3d human pose estimation from monocular images with deep convolutional neural network," in *Asian Conference on Computer Vision (ACCV)*, 2014, pp. 332–347.

[3] T. Zhao, S. Li, K. N. Ngan, and F. Wu, "3-d reconstruction of human body shape from a single commodity depth camera," *IEEE Transactions on Multimedia*, vol. 21, no. 1, pp. 114–123, 2018.

[4] P. Hu, E. S.-l. Ho, and A. Munteanu, "3dbodynet: Fast reconstruction of 3d animatable human body shape from a single commodity depth camera," *IEEE Transactions on Multimedia*, 2021.

[5] A. Kadkhodamohammadi and N. Padoy, "A generalizable approach for multi-view 3d human pose regression," *Machine Vision and Applications*, vol. 32, no. 1, pp. 1–14, 2021.

[6] K. Pullen and C. Bregler, "Motion capture assisted animation: Texturing and synthesis," in *Proceedings of the 29th annual conference on Computer graphics and interactive techniques*, 2002, pp. 501–508.

[7] P. Wang, W. Li, Z. Gao, C. Tang, and P. O. Ogunbona, "Depth pooling based large-scale 3-d action recognition with convolutional neural networks," *IEEE Transactions on Multimedia*, vol. 20, no. 5, pp. 1051–1061, 2018.

[8] M. Liu, H. Liu, and C. Chen, "Robust 3d action recognition through sampling local appearances and global distributions," *IEEE Transactions on Multimedia*, vol. 20, no. 8, pp. 1932–1947, 2017.

[9] M. Liu and J. Yuan, "Recognizing human actions as the evolution of pose estimation maps," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1159–1168.

[10] P. Wei, H. Sun, and N. Zheng, "Learning composite latent structures for 3d human action representation and recognition," *IEEE Transactions on Multimedia*, vol. 21, no. 9, pp. 2195–2208, 2019.

[11] Y.-F. Song, Z. Zhang, C. Shan, and L. Wang, "Constructing stronger and faster baselines for skeleton-based action recognition," *arXiv preprint arXiv:2106.15125*, 2021.

[12] T. Chen, D. Zhou, J. Wang, S. Wang, Y. Guan, X. He, and E. Ding, "Learning multi-granular spatio-temporal graph network for skeleton-based action recognition," in *Proceedings of the 29th ACM International Conference on Multimedia (ACMMM)*, 2021, pp. 4334–4342.

[13] C. Li, C. Xie, B. Zhang, J. Han, X. Zhen, and J. Chen, "Memory attention networks for skeleton-based action recognition," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

[14] D. Yang, Y. Wang, A. Dantcheva, L. Garattoni, G. Francesca, and F. Bremond, "Unik: A unified framework for real-world skeleton-based action recognition," *arXiv preprint arXiv:2107.08580*, 2021.

[15] B. Zhang, Y. Yang, C. Chen, L. Yang, J. Han, and L. Shao, "Action recognition using 3d histograms of texture and a multi-class boosting classifier," *IEEE Transactions on Image processing*, vol. 26, no. 10, pp. 4648–4660, 2017.

[16] C. Chen, M. Liu, H. Liu, B. Zhang, J. Han, and N. Kehtarnavaz, "Multi-temporal depth motion maps-based local binary patterns for 3-d human action recognition," *IEEE Access*, vol. 5, pp. 22 590–22 604, 2017.

[17] M. Garcia-Salguero, J. Gonzalez-Jimenez, and F.-A. Moreno, "Human 3d pose estimation with a tilting camera for social mobile robot interaction," *Sensors*, vol. 19, no. 22, p. 4943, 2019.

[18] L. Gui, K. Zhang, Y. Wang, X. Liang, J. M. Moura, and M. Veloso, "Teaching robots to predict human motion," in *Proceedings of the IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 562–567.

[19] 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 (ICCV)*, 2017, pp. 2640–2649.

[20] D. Pavllo, 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 (CVPR)*, 2019, pp. 7753–7762.

[21] G. Hua, W. Li, Q. Zhang, R. Ding, and H. Liu, "Weakly-supervised cross-view 3d human pose estimation," *arXiv preprint arXiv:2105.10882*, 2021.

[22] 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.

[23] 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.

[24] 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 (ICCV)*, 2019, pp. 2272–2281.

[25] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems (NIPS)*, 2017, pp. 5998–6008.

[26] Y. Tay, M. Dehghani, D. Bahri, and D. Metzler, "Efficient transformers: A survey," *arXiv preprint arXiv:2009.06732*, 2020.

[27] D. Zihang, L. Guokun, Y. Yiming, and Q. L. V., "Funnel-transformer: Filtering out sequential redundancy for efficient language processing," in *Advances in Neural Information Processing Systems (NIPS)*, 2020.

[28] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu *et al.*, "A survey on visual transformer," *arXiv preprint arXiv:2012.12556*, 2020.

[29] S. He, H. Luo, P. Wang, F. Wang, H. Li, and W. Jiang, "Transreid: Transformer-based object re-identification," *arXiv preprint arXiv:2102.04378*, 2021.

[30] X. Li, Y. Hou, P. Wang, Z. Gao, M. Xu, and W. Li, "Trear: Transformer-based rgb-d egocentric action recognition," *arXiv preprint arXiv:2101.03904*.

[31] L. Han, P. Wang, Z. Yin, F. Wang, and H. Li, "Exploiting better feature aggregation for video object detection," in *Proceedings of the 28th ACM International Conference on Multimedia (ACMMM)*, 2020, pp. 1469–1477.

[32] X. Li, Y. Hou, P. Wang, Z. Gao, M. Xu, and W. Li, "Transformer guided geometry model for flow-based unsupervised visual odometry," *Neural Computing and Applications*, pp. 1–12, 2021.

[33] M. Geva, R. Schuster, J. Berant, and O. Levy, "Transformer feed-forward layers are key-value memories," *arXiv preprint arXiv:2012.14913*, 2020.

[34] R. Liu, J. Shen, H. Wang, C. Chen, S.-c. Cheung, and V. Asari, "Attention mechanism exploits temporal contexts: Real-time 3d human pose reconstruction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 5064–5073.

[35] 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.

[36] L. Sigal, A. O. Balan, and M. J. Black, "Humaneva: Synchronized video and motion capture dataset and baseline algorithm for evaluation of articulated human motion," *International Journal of Computer Vision*, vol. 87, no. 12, pp. 4–27, 2010.

[37] 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 (CVPR)*, 2017, pp. 7025–7034.

[38] 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.

[39] L. Zhao, X. Peng, Y. Tian, M. Kapadia, and D. N. Metaxas, "Semantic graph convolutional networks for 3d human pose regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3425–3435.

[40] J. Liu, H. Ding, A. Shahroudy, L.-Y. Duan, X. Jiang, G. Wang, and A. C. Kot, "Feature boosting network for 3d pose estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 2, pp. 494–501, 2019.

[41] H.-S. Fang, Y. Xu, W. Wang, X. Liu, and S.-C. Zhu, "Learning pose grammar to encode human body configuration for 3d pose estimation," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[42] T. Xu and W. Takano, "Graph stacked hourglass networks for 3d human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 16 105–16 114.

[43] K. Gong, J. Zhang, and J. Feng, "Poseaug: A differentiable pose augmentation framework for 3d human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8575–8584.

[44] J. Wang, S. Yan, Y. Xiong, and D. Lin, "Motion guided 3d pose estimation from videos," *arXiv preprint arXiv:2004.13985*, 2020.

[45] T. Chen, C. Fang, X. Shen, Y. Zhu, Z. Chen, and J. Luo, "Anatomy-aware 3d human pose estimation with bone-based pose decomposition," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.

[46] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European Conference on Computer Vision (ECCV)*, 2020, pp. 213–229.

[47] X. Zhu, W. Su, L. Lu, B. Li, X. Wang, and J. Dai, "Deformable detr: Deformable transformers for end-to-end object detection," *arXiv preprint arXiv:2010.04159*, 2020.

[48] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*,

"An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[49] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, F. E. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *arXiv preprint arXiv:2101.11986*, 2021.

[50] K. Lin, L. Wang, and Z. Liu, "End-to-end human pose and mesh reconstruction with transformers," *arXiv preprint arXiv:2012.09760*, 2020.

[51] J. Lin and G. H. Lee, "Trajectory space factorization for deep video-based 3d human pose estimation," *arXiv preprint arXiv:1908.08289*, 2019.

[52] 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 Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 899–908.

[53] A. Zeng, X. Sun, F. Huang, M. Liu, Q. Xu, and S. Lin, "Srnet: Improving generalization in 3d human pose estimation with a split-and-recombine approach," in *European Conference on Computer Vision (ECCV)*, 2020, pp. 507–523.

[54] 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 (CVPR)*, 2018, pp. 7307–7316.

[55] K. Liu, R. Ding, Z. Zou, L. Wang, and W. Tang, "A comprehensive study of weight sharing in graph networks for 3d human pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 318–334.

[56] 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 (CVPR)*, 2019, pp. 10 895–10 904.

[57] D. Tome, M. Toso, L. Agapito, and C. Russell, "Rethinking pose in 3d: Multi-stage refinement and recovery for markerless motion capture," in *2018 International Conference on 3D Vision (3DV)*, 2018, pp. 474–483.

[58] 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 (CVPR)*, 2019, pp. 1077–1086.

[59] 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 (ICCV)*, 2019, pp. 723–732.

[60] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun, "Cascaded pyramid network for multi-person pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7103–7112.

[61] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision (ICCV)*, 2017, pp. 2961–2969.

[62] R. Dabral, A. Mundhada, U. Kusupati, S. Afaque, 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.

[63] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 483–499.

[64] Z. Wu, Z. Liu, J. Lin, Y. Lin, and S. Han, "Lite transformer with long-short range attention," in *International Conference on Learning Representations (ICLR)*, 2020.

[65] Z. Jiang, W. Yu, D. Zhou, Y. Chen, J. Feng, and S. Yan, "Convbert: Improving bert with span-based dynamic convolution," in *Advances in Neural Information Processing Systems (NIPS)*, 2020.
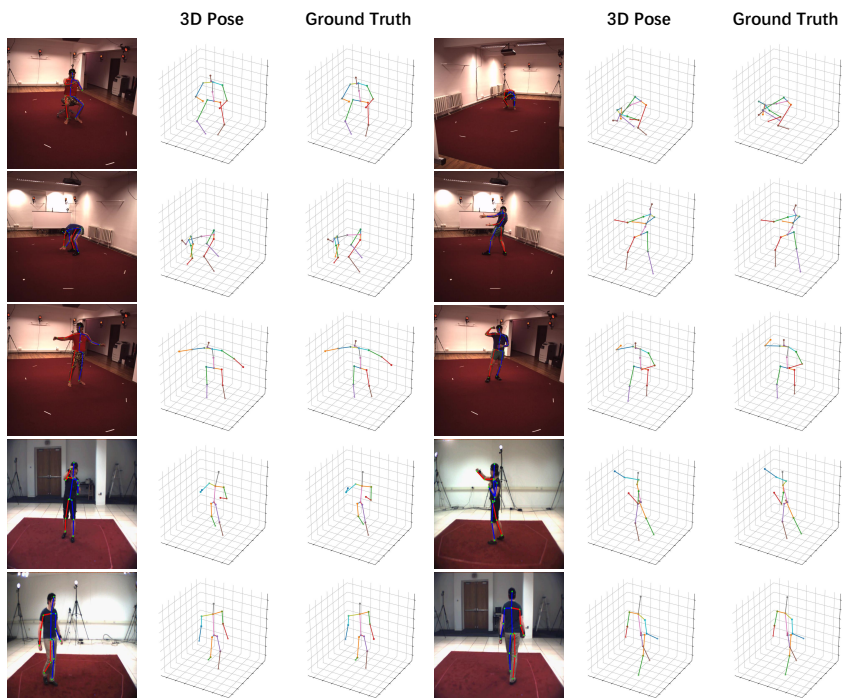
## VI. APPENDIX

Fig. 10: Visual results of our proposed method on Human3.6M dataset (first 3 rows) and HumanEva-I dataset (last 2 rows).
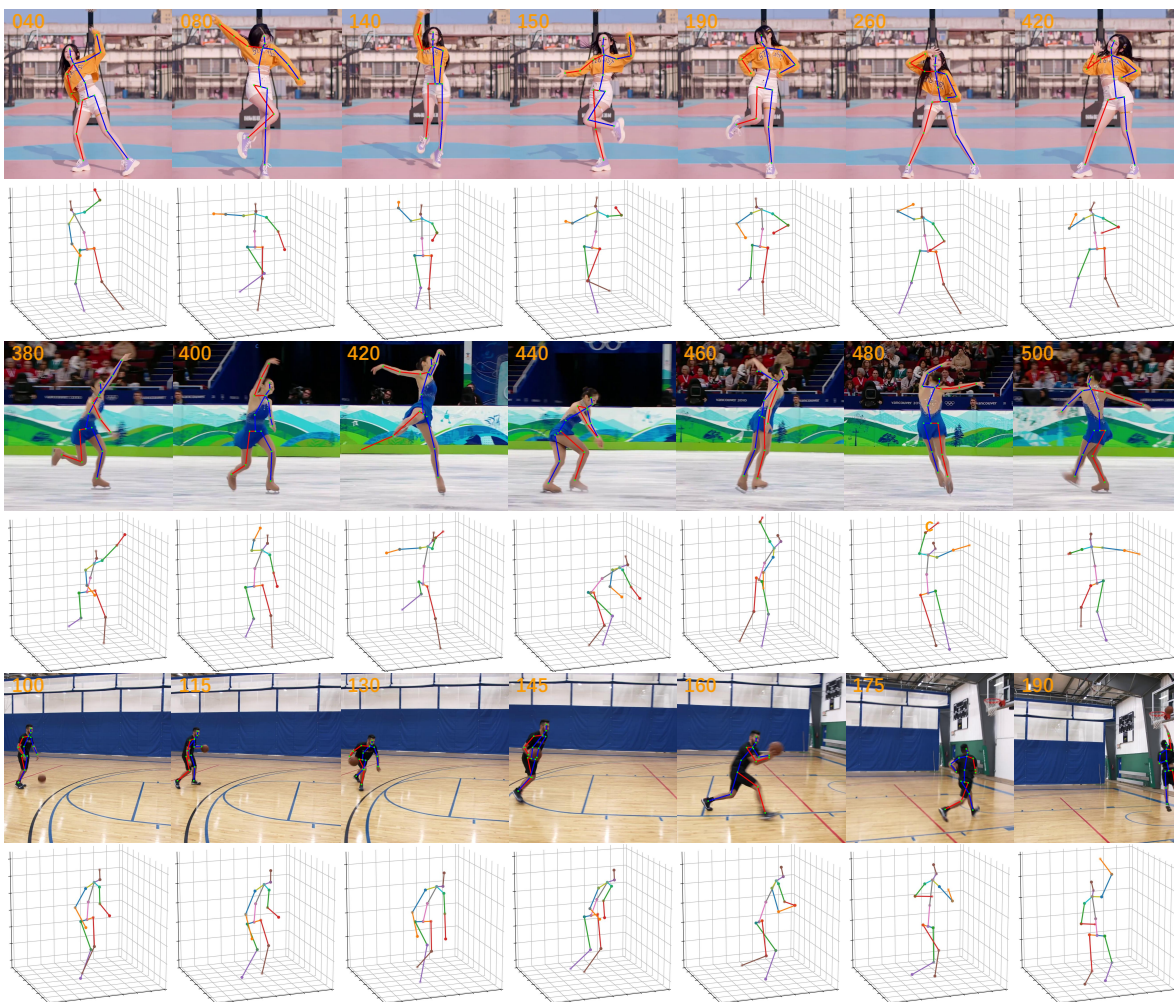


Fig. 11: Qualitative results on challenging wild videos. The number is the frame index of input videos.