# IVT: An End-to-End Instance-guided Video Transformer for 3D Pose Estimation

Zhongwei Qiu*
University of Science and Technology Beijing
qiuzhongwei@xs.ustb.edu.cn

Qiansheng Yang*
Baidu Inc.
yangqiansheng@baidu.com

Jian Wang
Baidu Inc.
wangjian33@baidu.com

Dongmei Fu
University of Science and Technology Beijing
fdm_ustb@ustb.edu.cn

## ABSTRACT

Video 3D human pose estimation aims to localize the 3D coordinates of human joints from videos. Recent transformer-based approaches focus on capturing the spatiotemporal information from sequential 2D poses, which cannot model the contextual depth feature effectively since the visual depth features are lost in the step of 2D pose estimation. In this paper, we simplify the paradigm into an end-to-end framework, Instance-guided Video Transformer (IVT), which enables learning spatiotemporal contextual depth information from visual features effectively and predicts 3D poses directly from video frames. In particular, we firstly formulate video frames as a series of instance-guided tokens and each token is in charge of predicting the 3D pose of a human instance. These tokens contain body structure information since they are extracted by the guidance of joint offsets from the human center to the corresponding body joints. Then, these tokens are sent into IVT for learning spatiotemporal contextual depth. In addition, we propose a cross-scale instance-guided attention mechanism to handle the variational scales among multiple persons. Finally, the 3D poses of each person are decoded from instance-guided tokens by coordinate regression. Experiments on three widely-used 3D pose estimation benchmarks show that the proposed IVT achieves state-of-the-art performances.

## CCS CONCEPTS

• **Computing methodologies → Object recognition**.

## KEYWORDS

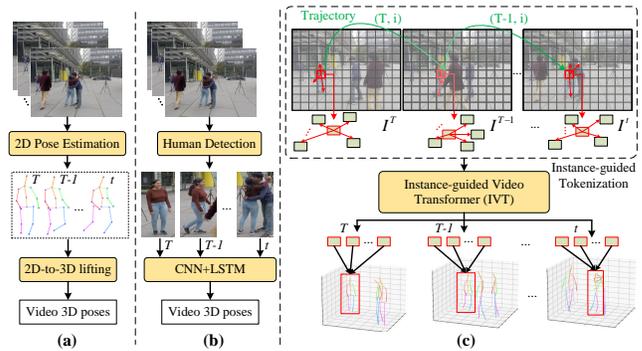Video Transformer, Human Pose Estimation

*Equal Contribution

Figure 1: Comparison of (a) The pipeline of 2D-to-3D video pose lifting [3, 51], (b) Recurrent structure based methods [7, 18] for video 3D pose estimation, (c) Our end-to-end Instance-guided Video Transformer (IVT). IVT is a single-stage framework while others are two-stage methods.

## 1 INTRODUCTION

3D human pose estimation aims to localize the 3D joints of person(s) from monocular images or videos. As a fundamental computer vision task, it has a lot of applications, including action recognition [27], human-robot interaction detection [21], and virtual reality [32], etc. Unfortunately, estimating 3D human poses from monocular 2D images or videos is very challenging because of the lack of depth information.

To tackle this problem, some image-based approaches learn depth information from image feature by depth map supervision [40] or 3D heatmap supervision [30, 37], while other image-based approaches [45, 49] firstly estimate 2D pose from image, and then lift 2D pose to 3D pose. However, the depth information implied in a single image is still limited. Compare to a single image, video can provide more motion cues which are quite helpful for the inference of depth. Thus, video-based approaches [1, 3, 6, 7, 38, 51] are developed rapidly in recent years. And we also focus on the video-based method in this paper.

By the benefit of the amazing performance of modern 2D pose estimators, many video-based approaches follow the 2D-to-3D lifting paradigm as Figure 1 (a), and exploit spatial and temporal modeling methods to improve the performance of 3D pose estimation. [26] and [13] apply convolution neural network (CNN) and recurrent neural network (RNN) respectively to model temporal dependency among the sequential 2D poses. But limited by the formulation

of CNN and RNN, they are not good at capturing the long-range dependency in both spatial and temporal dimensions. To alleviate this situation, some works [3, 6, 42] introduce graph neural network (GNN) to exploit spatial-temporal information between keypoints, which can capture both short-term and long-term dependency by setting an appropriate adjacent matrix. Besides GNN, other works [23, 51] use transformer to get more representative features from pose sequence, and also improve the performance of video 3D pose estimation significantly.

Despite the significant progress achieved by the above GNN or transformer-based methods, they still did not break out of the 2D-to-3D lifting paradigm. This paradigm only considers the structure of the 2D pose for depth estimation, while ignoring the contextual depth information contained in the semantic feature of video frames, since the semantic feature has been dropped at the 2D pose prediction stage. But we suppose that the context depth information embedded in the semantic feature is more effective than the 2D pose structure for 3D pose estimation.

As shown in Figure 1 (b), some video-based approaches [7, 38] with the recurrent neural network, firstly conduct human detection and then predict 3D pose directly from the cropped video patches, which can be regarded as exploiting contextual depth information from the semantic feature of video frames to some extent. But they apply the recurrent neural network to exchange features between different video frames for only temporal modeling, which is not effective compared with GNN or transformer-based spatial-temporal modeling methods mentioned above. Besides, the inputs of cropped patches bring a new problem of keypoints feature alignment.

To handle the above problems, we simplify the video 3D pose estimation into an end-to-end transformer-based framework as shown in Figure 1 (c), which aims to make full usage of the spatial-temporal depth feature and predicts 3D pose directly from video frames. In order to capture the effective contextual depth information and reduce the computational burden introduced by conducting self-attention on the dense semantic feature of video frames, we propose an Instance-guided Video Transformer to model the spatial-temporal depth information by the guidance of human instance.

Firstly, Instance-guided Video Transformer (IVT) introduces a series of instance-guided visual tokens and each token is capable for predicting 3D pose of an instance. These tokens are contructed by aggregating features from related spatial points by the guidance of a set of learned 2D offsets from human center to the corresponding human joints. This mechanism enables each token can capture the whole body information. For a query token, the attention is computed on both spatial and temporal dimensions to exchange the context depth information. Furthermore, we propose a cross-scale instance-guided attention mechanism to handle the variational scales among multiple persons. As a result, IVT enables effective spatial-temporal depth feature exchange and brings significant improvement to video 3D pose estimation.

In summary, IVT is a simple and unified framework that is suitable for both single-person and multi-person video 3D pose estimation tasks. And to the best of our knowledge, for multi-person video 3D pose estimation, IVT is the first end-to-end method that leverages transformer to directly capture multi-person depth information in the video. Our contributions can be summarized as follows:

- We propose a novel end-to-end transformer-based framework for both single-person and multi-person video 3D pose estimation, called instance-guided video transformer (IVT).
- We design a novel instance-guided attention mechanism to enable effective spatial-temporal depth information learning in videos.
- IVT achieves new state-of-the-art results on three widely-used video 3D pose estimation benchmarks, Human3.6M, 3DPW, and CMU Panoptic datasets.

## 2 RELATED WORK

### 2.1 Image-based 3D Pose Estimation

The image-based multi-person 3D pose estimation methods can be mainly divided into two kinds of paradigms: top-down [23, 30, 37, 45, 49] and bottom-up approaches [40, 48, 50].

The top-down paradigm follows a pipeline of conducting human detection firstly and performing single-person 3D pose estimation later. For the single-person cases, they predict 3D poses by learning 3D heatmaps [30], or estimating 2D poses by 2D pose estimator [34, 35] and lifting 2D poses to 3D poses [49]. Typically, PoseNet [30] predicts the root depths of each person at the stage of human detection, then estimates the 3D coordinates from 3D heatmaps. The bottom-up paradigm [40, 48, 50] follows a pipeline of firstly estimating the 3D coordinates for each human joint in an image and then assigning them to different human instances. For example, MubyNet [48] estimates keypoints and limb core at the same time and then integrates limb score to group keypoints into different persons. HMOR [40] propose hierarchical multi-person ordinal relations as an additional loss to help depth learning. However, the image-based approaches are not good at handling occlusion cases as the video-based approaches.

### 2.2 Video-based 3D Pose Estimation

Video-based multi-person 3D pose estimation aims to capture temporal information for 3D pose estimation. The ways of extracting temporal information can be divided into two categories: based on image visual features [7, 18, 38] and based on 2D poses [1, 3, 51].

The methods [7, 18, 38] based on image visual features usually crop the human features according to human bounding boxes, and then use 3D convolution or recurrent neural network to extract the temporal information from these cropped sequences features. Sun *et al.* [38] propose a skeleton-disentangling framework to separate 3D human pose and shape estimation into spatial and temporal dimensions. TCMR [7] uses ResNet to extract visual features from video frames, then captures temporal information on these deep features by the recurrent neural network. However, these methods essentially conduct single-person video 3D pose estimation, which brings a new problem of feature alignment due to crop images.

The methods [1, 3, 51] based on 2D coordinates usually estimate a sequence of 2D poses at first, then lift 2D coordinates sequence to 3D pose by a temporal lifting network. Typically, Cai *et al.* [3] exploit spatial-temporal relationships for 3D pose estimation via Graph Convolutional Networks (GCN). Zheng *et al.* [51] propose PoseFormer, a spatial-temporal transformer network to capture the spatial-temporal information among human joints. However, these methods cannot capture truly depth features from visual

images since the depth feature is lost in the stage of 2D human pose estimation. Besides, these methods disassemble the multi-person video task into a single-person video task. Thus, the depth information between different persons can not be captured. In this paper, we tackle the problems and build an end-to-end multi-person video 3D pose estimation framework.

## 2.3 Transformer in Human Pose Estimation

Recently, the transformer-based approaches [14, 20, 23, 29, 44, 51] have been proposed to improve the long-term modeling capabilities of sequence for human pose estimation. TransPose [44] and TF-Pose [29] formulate human joints as visual tokens and capture the relationship between human joints by self-attention. METRO [20] and PRTR [20] exploit the end-to-end transformer-based pose estimation network. PoseFormer [51] explores the spatial-temporal attention mechanism for 3D pose estimation. However, the Pose-Former didn't study the attention on real depth features from images since it lifts 3D poses from a sequence of 2D poses. Moreover, the existing transformer-based pose estimation methods are designed for single-person pose estimation, which limits their applications. In this paper, we study an end-to-end multi-person 3D pose estimation framework and explore to extract the relationship between multi-person joints in both spatial and temporal dimensions.

## 3 METHOD

In this section, we elaborate on the detail of the proposed Instance-guided Video Transformer (IVT). The framework of IVT is shown in Figure 2. Given a sequence of video frames $I = \{I_t \mid t \in [1, T]\}$, IVT firstly extracts deep features by a backbone network, which are used to learn instance 2D offsets (from body center to $J$ keypoints) and temporal feature motion (trajectory). Then, for each frame, the deep features are organized as visual tokens with the guidance of instance 2D offsets, called instance-guided tokens. And each instance-guided token is in charge of predicting the 3D pose of its corresponding instance by the assistant of the aggregated whole body information. This process is denoted as Instance-Guided Tokenization (IGT). After conducting IGT, instance-guided tokens are sent into a video transformer to capture context depth information in both spatial and temporal dimensions. Finally, the outputted visual tokens are further used to decode 3D poses for persons detected from a human center heatmap.

## 3.1 Formulation of Transformer Module

First of all, we review the formulation of the basic transformer module. Given query matrix $Q$, key matrix $\mathcal{K}$, value matrix $\mathcal{V}$, where the first dimension of them is sample dimension while the second is feature dimension, a typical attention formulation $A(Q, \mathcal{K}, \mathcal{V})$ can be expressed as:

$$A(Q, \mathcal{K}, \mathcal{V}) = softmax(\frac{Q \cdot \mathcal{K}^T}{\sqrt{d}}) \cdot \mathcal{V} \qquad (1)$$

where $d$ is the length of feature within these matrixs. If we split $Q$, $\mathcal{K}$, $\mathcal{V}$ into $h$ heads via feature axis and conduct attention for each head, denoted as $\{Q_1, ..., Q_h\}$, $\{\mathcal{K}_1, ..., \mathcal{K}_h\}$, $\{\mathcal{V}_1, ..., \mathcal{V}_h\}$, the

Multi-Head Attention $MHA(Q, \mathcal{K}, \mathcal{V})$ can be formulated as:

$$MHA(Q, \mathcal{K}, \mathcal{V}) = P(Concat_c(head_1, ..., head_h))$$
$$head_i = A(Q_i, \mathcal{K}_i, \mathcal{V}_i), i \in [1, h] \qquad (2)$$

where P is linear projection function, $Concat_c$ means concatenating matrixs along feature axis. Specially, when $Q$, $\mathcal{K}$, $\mathcal{V}$ are derived from a same input matrix $\mathcal{X}$, we can get Multi-Head Self-Attention $MHSA(\cdot)$ further, which is denoted as:

$$MHSA(\mathcal{X}) = MHA(Q, \mathcal{K}, \mathcal{V})$$
$$Q = P_q(\mathcal{X}), \mathcal{K} = P_k(\mathcal{X}), \mathcal{V} = P_v(\mathcal{X}) \qquad (3)$$

where $P_q$, $P_k$, $P_v$ are linear projection functions for generating $Q$, $\mathcal{K}$, $\mathcal{V}$, respectively.

Then, two basic transformer modules used in this paper can be formulated as:

$$\mathcal{X}_{out} = FFN(MHSA(\mathcal{X}_{in})) \qquad (4)$$
$$\mathcal{X}_{out} = FFN(MHA(Q, \mathcal{K}, \mathcal{V})) \qquad (5)$$

and they are served for self-attention and cross-attention respectively. While, in the above equations, $FFN(\cdot)$ represents feed forward network, which consists of two linear layers. For simple expression, the layer norm and shortcut path are ignored here.

## 3.2 Instance-guided Tokenization

In this section, we introduce the process of generating instance-guided tokens for each video frame $I_t$, named Instance-Guided Tokenization (IGT).

We firstly extract deep feature $F_t$ for $I_t$ by a backbone network $\phi(\cdot)$, and the shape of $F_t$ is $C \times H \times W$. Then we split $F_t$ into $N$ feature blocks, denote as $B_t$, which has a shape of $N \times C_b$, where $C_b = C \times K \times K$ and $K$ represents the block size. To extract visual tokens, traditional vision transformers [2, 10] take each block as a token and capture the spatial and temporal relationships among these block tokens. Then, $N$ tokens are generated from features $B_t$, denoted as $\tau_t$. The shape of $\tau_t$ is $N \times C_t$, where $C_t = C_b$ is feature dim of tokens. However, this tokenization is not fine-grained to capture the context information from the human body to predict depth for 3D human pose estimation. Here, we introduce the instance-guided tokenization (IGT) approach. Different from traditional tokenization, IGT considers the features from whole body as well as the relationship between human joints when extracting visual tokens, which enables each token to encode the body structure of its corresponding human instance.

For the $i$th block in $B_t$, denoted as $B_t[i]$, the process of instance-guided tokenization includes two steps. Firstly, it gathers features from $J$ corresponding blocks of $B_t[i]$ with the guidance of instance 2D offsets. As shown in Figure 2, instance 2D offsets are the joint offsets from body center to $J$ joints, and these offsets are preserved in an offset map $M_o^{2D}$, which is predicted by several convolutional layers based on deep feature $F_t$. $M_o^{2D}$ contains the whole body information of each instance and indicates the feature locations of relative keypoints. The gathered blocks are concatenated into one feature vector, denoted as $\widetilde{B}_t[i]$, and its length is $J \times C_b$. Secondly, a multi-head self-attention module is used to encode $\widetilde{B}_t[i]$ and generate instance-guided token $\tau_t[i]$. This enables feature exchange between different joints of a human instance and makes the token feature be aware of body structure, which improves the
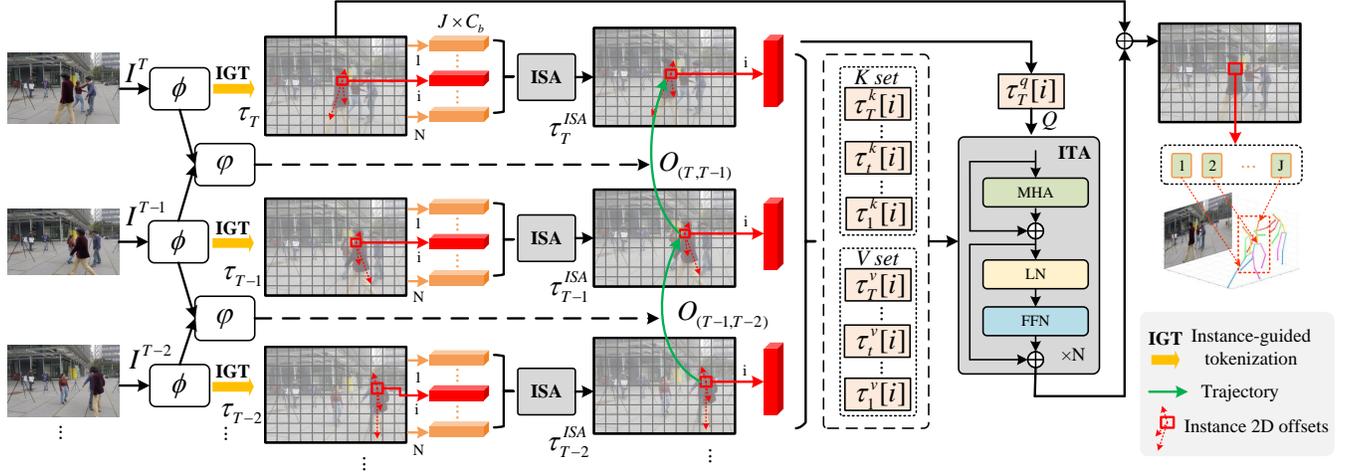
**Figure 2: The overview of Instance-guided Video Transformer (IVT), which includes instance-guided tokenization (IGT), instance-guided spatial attention (ISA), and instance-guided temporal attention (ITA). Given video frames of height $H$ and width $W$, deep features are extracted by embedding network $\phi(\cdot)$, and are further used to estimate trajectory motions $O$ by network $\varphi(\cdot)$. The keypoints features are extracted from deep features according to the trajectory and instance 2D offsets and to further learn instance-aware tokens of shape $J \times C_b$ by IGT. For a query token $\tau_T^q[i]$ at the $i^{th}$ block in $T^{th}$ frame, ISA is computed in each frame and ITA is computed among temporal frames to capture depth information. Each token in the final layer of IVT outputs the 3D coordinates of a person. MHA, LN, and FFN denote multi-head attention, layer norm, and feed-forward network, respectively. $\oplus$ means element-wise addition. $J$ represents the joints number. N represents token numbers in each frame.**

robustness of instance-guided token for predicting whole body 3D pose. Concretely, this self-attention can be formulated as:

$$f = reshape(\widetilde{B}_t[i], (J, C_b))$$
$$f = FFN(MHSA(f)) \tag{6}$$
$$\tau_t[i] = reshape(f, J \times C_b)$$

where $MHSA$ and $FFN$ are multi-head self-attention module and feed forward module respectively. $reshape(f, shape)$ means reshaping the input data $f$ into target $shape$. As a result, each generated instance-guided token $\tau_t[i]$ can encode the global context from a human instance.

Once the token feature $\tau_t$ for each video frame $I_t$ is extracted, we pass all the token features $\tau = \{\tau_t | t \in [1, T]\}$ of $T$ video frames to the instance-guided video transformer for 3D pose estimation, which will be elaborated in the next section.

### 3.3 Instance-guided Video Transformer

Instance-guided video transformer (IVT) takes the instance-guided tokens $\tau_t$ as inputs and conducts spatial-temporal attention to capture context depth information in both spatial and temporal dimensions. It can be divided into two sequential attention stages, Instance-guided Spatial Attention(ISA) and Instance-guided Temporal Attention(ITA). ISA computes the correlation between all tokens within one frame, which can gather the context depth features from other human instances or objects. Based on the output of ISA, ITA calculates attention among a group of corresponding tokens in the temporal dimension, which can aggregate depth information of the same instance from different video frames. Since the inputs are whole images, the human instances in images suffer variational

scales. To tackle this problem, we propose a cross-scale attention mechanism for IVT. It enables IVT to be more robust to handle the different scales of human instances.

*3.3.1 Instance-guided Spatial Attention.* Instance-guided Spatial Attention (ISA) conducts spatial self-attention within one frame. Here, for the $t$th frame, it is tokenized as instance-guided tokens $\tau_t$. The tokens are sent into ISA and output $\tau_t^{ISA}$, which can be formulated as:

$$\tau_t^{ISA} = A^{ISA}(\tau_t) = FFN(MHSA(\tau_t)) \tag{7}$$

where $MHSA$ and $FFN$ are multi-head self-attention and feed-forward network, respectively. $A^{ISA}(\cdot)$ means instance-guided spatial attention. Due to the instance-guided tokens, ISA can capture more fine-grained keypoints relationships between the same person and different human instances at the same time. We compute ISA on each video frame to obtain a sequence of token maps denoted as $\tau^{ISA} = \{\tau_t^{ISA} | t \in [1, T]\}$.

*3.3.2 Instance-guided Temporal Attention.* Temporal information is important for 3D human pose estimation, especially in handling occlusion problems. To capture global depth information from temporal features, we introduce instance-guided temporal attention (ITA) here, which computes the cross-attention on instance-guided tokens from different video frames.

The query, key, and value for ITA are generated from the token maps $\tau^{ISA}$ outputted from ISA. For a query token $\tau_T^{ISA}[i]$ at the $i$th block in $T$th frame, ITA is computed on the same block places among different frames, and we denote the query, key and value for updating this token as $Q_T[i]$, $\mathcal{K}_T[i]$ and $\mathcal{V}_T[i]$ respectively, which
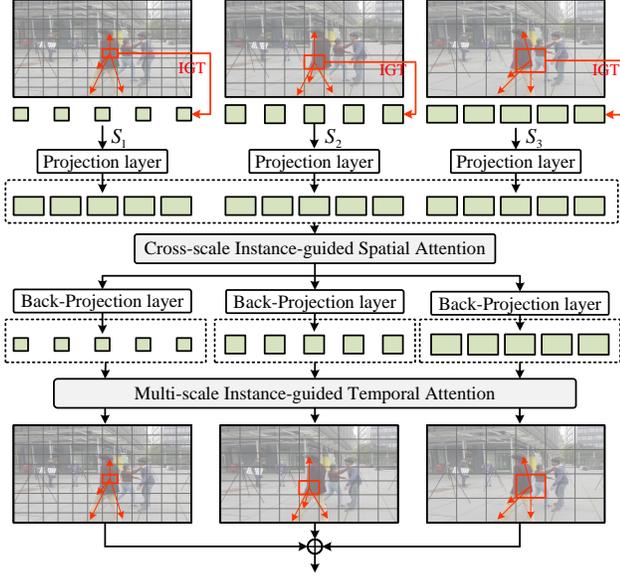
**Figure 3: The illustration of the cross-scale instance-guided attention for IVT. $S_1$, $S_2$, and $S_3$ represents three different scales for cross-scale attention. IGT means instance-guided tokenization. The projection layer is a Linear layer to project visual tokens from different scales into the same size. Back-projection is the inverse operation of the projection layer.**

are formulated as:

$$Q_T[i] = P_q(\tau_T^{ISA}[i]),$$
$$\mathcal{K}_T[i] = P_k(Concat_n(\tau_1^{ISA}[i], ..., \tau_{T-1}^{ISA}[i], \tau_T^{ISA}[i])), \quad (8)$$
$$\mathcal{V}_T[i] = P_v(Concat_n(\tau_1^{ISA}[i], ..., \tau_{T-1}^{ISA}[i], \tau_T^{ISA}[i]))$$

where $P_q$, $P_k$, and $P_v$ are linear projection layers for generating the query, key, and value, respectively. $Concat_n$ means concatenating matrixes along sample dimension. $i \in [1, N]$ represents the block index. It is worth to note that, before conducting temporal attention, all the token maps from the video sequence are aligned with optical flow, which is calculated between each pair of adjacent frames in advance. Therefore, the ITA can be computed on the tokens from different frames with the same block index $i$. As shown in Figure 2, the green line represents the corresponding relationship between adjacent frames. Then, for the $i$th token at $T$th frame, the instance-guided temporal attention(ITA) is computed as:

$$\tau_T^{ITA}[i] = A^{ITA}(\tau_T^{ISA}[i]) = FFN(MHA(Q_T[i], \mathcal{K}_T[i], \mathcal{V}_T[i]))$$
$$(9)$$

where $A^{ITA}(\cdot)$ means instance-guided temporal attention and it is repeated for all $t \in [1, T]$ and $i \in [1, N]$ for outputting token maps $\tau^{ITA} = \{\tau_t^{ITA} | t \in [1, T]\}$.

Combined with ISA and ITA, the final output of IVT can be formulated as:

$$\tau_t^{IVT} = A^{ITA}(A^{ISA}(\tau_t)) + \tau_t, \quad t \in [1, T] \quad (10)$$

### 3.3.3 Cross-scale Attention.
Scale information also matters in depth estimation because the relative depth of different subjects is correlated with their scales. Therefore, to improve the accuracy of depth estimation, we design a cross-scale attention mechanism for the instance-guided video transformer. This cross-scale mechanism encourages the video transformer to calculate attention between blocks with different scales, which can make the relative depth information from other persons be aggregated better.

Concretely, we split the deep feature map $F_t$ with different block sizes at the instance-guided tokenization stage. In this paper, we apply three block sizes $s \in \{2, 4, 8\}$. Then, we pass all the tokens from three scales into cross-scale instance-guided spatial attention (CISA), as shown in Figure 3. Due to the feature dimension of tokens from different scales being different, we add projection layers before the cross-scale attention for aligning the feature dimension. Meanwhile, back-projection layers are added after cross-scale attention for restoring the feature dimension of those tokens to their original state.

After cross-scale instance-guided spatial attention for each frame, three sequences of token maps with different scales are outputted. Then these three sequences are passed into ITA to perform temporal attention individually. Finally, token maps from different scales are added up into one token map frame by frame. For clearity, we name the above operation as multi-scale instance-guided temporal attention (MITA).

In a nutshell, combined with CISA and MITA, the whole process of IVT in Equation 10 can be recapped as

$$\tau_t^{IVT} = A^{MITA}(A^{CISA}(\tau_t)) + \tau_t, \quad t \in [1, T] \quad (11)$$

where $A^{CISA}(\cdot)$ denotes cross-scale instance-guided spatial attention and $A^{MITA}(\cdot)$ denotes multi-scale instance-guided temporal attention.

## 3.4 Loss Function
For $t$th frame in video, the outputted token map $\tau_t^{IVT}$ are used to learn root keypoints heatmaps $M_h$ of size $H \times W$ and 3D offset $M_o$ of size $J \times H \times W$ by a convolutional layer. For each point $p$ in $M_h$ with high confidence, the corresponding 3D offsets at the $p$ place in $M_o$ are extracted to decode a whole 3D pose of size $J \times 3$ for a person, $J$ represents joints number. Finally, NMS is used to remove the superfluous 3D poses. The decoding process is same as previous work [31, 52].

During training, we use $L_1$ loss for 3D offsets regression and $L_2$ loss for heatmap learning. Meanwhile, the instance 2D offsets $M_o^{2D}$ are learned with the supervision of ground-truth instance offsets $\hat{M}_o^{2D}$. The loss function $\mathcal{L}$ is

$$\mathcal{L} = L_1(M_o, \hat{M}_o) + L_1(M_o^{2D}, \hat{M}_o^{2D}) + \alpha L_2(M_h, \hat{M}_h), \quad (12)$$

where $\hat{M}_o^j$ and $\hat{M}_h$ means the ground-truth of $M_o$ and $M_h$, respectively. $\alpha$ represents a loss weight.

## 4 EXPERIMENTS
In this section, we elaborate the experiment results of IVT. We firstly introduce the implemental details of IVT, and then report results and compare with SOTA methods on three widely-used datasets: Human3.6M, 3DPW, and CMU Panoptic. All ablation studies are

**Table 1: Quantitative comparison with state-of-the-art methods on Human3.6M under Protocol 1 (MPJPE) and Protocol 2 (PA-MPJPE). $f$ denotes the number of input frames used in each method, and $*$ represents a Transformer-based model. Bold indicates the best and underline indicates the second best.**

| Protocol 1 | | Dir. | Disc. | Eat. | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Somke | Wait | WalkD. | Walk | WalkT. | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Dabral et al. [8] (f=243) | ECCV'18 | 44.8 | 50.4 | 44.7 | 49.0 | 52.9 | 61.4 | 43.5 | 45.5 | 63.1 | 87.3 | 51.7 | 48.5 | 52.2 | 37.6 | 41.9 | 52.1 |
| Cai et al. [3] ($f=7$) | 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. [33] ($f=243$) | 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. [22] ($f=50$) | BMVC'19 | 42.5 | 44.8 | 42.6 | 44.2 | 48.5 | 57.1 | 52.6 | 41.4 | 56.5 | 64.5 | 47.4 | 43.0 | 48.1 | 33.0 | 35.1 | 46.6 |
| Yeh et al. [46] ($f=243$) | NeurIPS'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 |
| Liu et al. [26] ($f=243$) | CVPR'20 | 41.8 | 44.8 | 41.1 | 44.9 | 47.4 | 54.1 | 43.4 | 42.2 | 56.2 | 63.6 | <u>45.3</u> | 43.5 | 45.3 | 31.3 | 32.2 | 45.1 |
| Zeng et al. [49] ($f=243$) | ECCV'20 | 46.6 | 47.1 | 43.9 | <u>41.6</u> | 45.8 | <u>49.6</u> | 46.5 | **40.0** | 53.4 | 61.1 | 46.1 | 42.6 | <u>43.1</u> | 31.5 | 32.6 | 44.8 |
| Wang et al. [42] ($f=96$) | ECCV'20 | <u>41.3</u> | 43.9 | 44.0 | 42.2 | 48.0 | 57.1 | 42.2 | 43.2 | 57.3 | 61.3 | 47.0 | 43.5 | 47.0 | 32.6 | 31.8 | 45.6 |
| Chen et al. [4] ($f=81$) | TCSVT'21 | 42.1 | <u>43.8</u> | 41.0 | 43.8 | <u>46.1</u> | 53.5 | 42.4 | 43.1 | 53.9 | 60.5 | 45.7 | <u>42.1</u> | 46.2 | 32.2 | 33.8 | 44.6 |
| Lin et al. [23] ($f=1$)* | CVPR'21 | - | - | - | - | - | - | - | - | - | - | - | - | - | - | - | 54.0 |
| Liu et al. [26] ($f=243$)* | ICRA'21 | 43.3 | 46.1 | 40.9 | 44.6 | 46.6 | 54.0 | 44.1 | 42.9 | 55.3 | **57.9** | 45.8 | 43.4 | 47.3 | 30.4 | **30.3** | 44.9 |
| Zheng et al. [51] ($f=81$)* | ICCV'21 | 41.5 | 44.8 | <u>39.8</u> | 42.5 | 46.5 | 51.6 | <u>42.1</u> | 42.0 | <u>53.3</u> | 60.7 | 45.5 | 43.3 | 46.1 | **31.8** | <u>32.2</u> | <u>44.3</u> |
| **Ours (IVT) ($f=5$)*** | | **36.5** | **40.1** | **38.4** | **40.7** | **42.6** | **42.8** | **30.1** | 43.4 | **46.1** | <u>58.0</u> | **40.2** | **37.1** | **40.8** | <u>32.1</u> | 33.5 | **40.2** |
| **Protocol 2** | | Dir. | Disc. | Eat. | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Somke | Wait | WalkD. | Walk | WalkT. | Avg. |
| Hossain et al. [13] ($f=243$) | 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 |
| Cai et al. [3] ($f=7$) | 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 | 32.3 | 39.0 |
| Lin et al. [22] ($f=50$) | 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. [33] ($f=243$) | 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 |
| Liu et al. [26] ($f=243$) | CVPR'20 | <u>32.3</u> | 35.2 | 33.3 | 35.8 | 35.9 | 41.5 | 33.2 | 32.7 | 44.6 | 50.9 | 37.0 | <u>32.5</u> | 37.0 | <u>25.2</u> | 27.2 | 35.6 |
| Wang et al. [42] ($f=96$) | ECCV'20 | 32.9 | 35.2 | 35.6 | <u>34.4</u> | 36.4 | 42.7 | <u>31.2</u> | 32.5 | 45.6 | 50.2 | 37.3 | 32.8 | 36.3 | 26.0 | **23.9** | 35.5 |
| Chen et al. [4] ($f=81$) | TCSVT'21 | 33.1 | 35.3 | 33.4 | 35.9 | 36.1 | 41.7 | 32.8 | 33.3 | <u>42.6</u> | 49.4 | 37.0 | 32.7 | 36.5 | 25.5 | 27.9 | 35.6 |
| Liu et al. [26] ($f=243$)* | ICRA'21 | 32.7 | 36.2 | 33.4 | 36.5 | 36.0 | 41.5 | 33.6 | 33.1 | 44.1 | 46.8 | 36.7 | 33.1 | 35.8 | 24.2 | 24.8 | 35.2 |
| Zheng et al. [51] ($f=81$)* | ICCV'21 | 32.5 | <u>34.8</u> | <u>32.6</u> | 34.6 | <u>35.3</u> | <u>39.5</u> | 32.1 | <u>32.0</u> | 42.8 | <u>48.5</u> | <u>34.8</u> | **32.4** | <u>35.3</u> | **24.5** | 26.0 | <u>34.6</u> |
| **Ours (IVT) ($f=5$)*** | | **27.0** | **24.8** | **32.2** | **30.1** | **27.8** | **32.1** | **22.3** | **28.7** | **30.7** | **24.4** | **32.7** | 37.8 | **21.9** | 31.1 | <u>24.7</u> | **28.5** |

based on CMU Panoptic dataset. Meanwhile, some visualization results on Human3.6M are given for presenting the superiority of IVT in an intuitionistic way.

### 4.1 Implemental Details

We use HRNet-32 [41] pre-trained on 2D pose estimation dataset COCO [24] as the backbone network of IVT and SPyNet [36] as the motion estimation network between different frames. In our experiments, IVT is stacked with 3 layers, and it is trained on 8 V100 GPUs with a batch size of 4 sequences/GPU, while the sequence length is 5 frames and the input size is $512 \times 512$. The total training epochs is 50. Adam optimizer is adopted and the initial learning rate is 5e-4, which decreases 10× at 30 and 40 epochs. The loss weight $\alpha$ equals 10 during training.

### 4.2 Datasets and Metrics

*4.2.1 Human3.6M dataset.* Human3.6 [15] is the largest indoor benchmark for single-person 3D pose estimation, which includes 7 subjects that performing 15 actions. Following the previous works [23, 30, 49, 51], we use two protocols for evaluation. For **Protocol 1**, IVT is trained on the subjects S1, S5, S6, S7, and S8, and tested on the subjects S9 and S11 by using Mean-Per-Joint-Position-Error (MPJPE), which measures the Euclidean distances between the ground truth joints and the predicted joints. For **Protocol 2**, subjects S1, S5, S7, S8, and S9 are used for training, and S11 is used for testing by PA-MPJPE. It calculates the Euclidean distance between predicted and ground-truth 3D joint coordinates after root joint alignment and further rigid alignment by Procrustes analysis [12].

*4.2.2 3DPW dataset.* 3DPW [39] is a multi-person outdoor 3D pose estimation dataset, which contains 22K images for training and 35K images for testing. Following the previous works [7, 18, 23], we train IVT on the training set and evaluate IVT on the testing set in PA-MPJPE.

*4.2.3 CMU Panoptic dataset.* CMU Panoptic [16] is a larger-scale multi-person dataset, captured by multiple cameras. Following the settings of previous works [40, 50], we use 160K images from different videos as the training set and the videos from two cameras (16, 30) as the testing set. For comparison, MPJPE is used for evaluation.

### 4.3 Comparison with SOTA Methods

*4.3.1 Results on Human3.6M Dataset.* The comparisons with state-of-the-art methods on the Human3.6M dataset are shown in Table 1. Our IVT with $f = 5$ achieves new state-of-the-art results with an MPJPE of 41.3mm and a PA-MPJPE of 28.5mm in Protocol 1 and Protocol 2, respectively. The relative gains are 6.8% and 17.6%, respectively. The results demonstrate the effectiveness of the proposed IVT. Compared with other transformer-based methods [23, 25, 51], IVT outperforms them. Even the PoseFormer [51] is based on a frame number 81, IVT with $f = 5$ obtains better results since the PoseFormer loses the visual depth feature in the process of temporal modeling.

We also give a fine-grained analysis of videos from the Human3.6M dataset in Figure 4. As shown in Figure 4 (a), given the video inputs, GAST-Net [25] and PoseFormer [51] fail on these

**Table 2: Quantitative comparison with state-of-the-art methods on multi-person 3D human pose estimation dataset (3DPW) in PA-MPJPE. Frame denotes the number of input frames used in each method. * denotes transformer-based methods. Lower is better.**

| Methods | | Frames | PA-MPJPE ↓ |
|---|---|---|---|
| Doersch et al. [9] | NeurIPS'19 | 31 | 74.4 |
| Kanazawa et al. [17] | CVPR'19 | 10 | 72.6 |
| Cheng et al. [6] | AAAI'20 | >90 | 71.8 |
| Sun et al. [38] | ICCV'19 | 45 | 69.5 |
| Kolotouros et al. [19] | ICCV'19 | 1 | 59.2 |
| Kocabas et al. [18] | CVPR'20 | 16 | 57.6 |
| Luo et al. [28] | ACCV'20 | 90 | 54.7 |
| Cheng et al. [5] | CVPR'21 | >90 | 62.9 |
| Choi et al. [7] | CVPR'21 | 16 | 52.7 |
| **Ours (IVT)*** | | 5 | **46.0** |

hard cases with complex postures or occlusions, but our IVT performs well on these cases since the captured temporal depth context. The red circle denotes the wrong pose predicted by GAST-Net and PoseFormer.

As shown in Figure 4 (b) and (c), the MPJPE of IVT is lower than GAST-Net and PoseFormer, and the depth error of IVT is lower than GAST-Net and PoseFormer. It shows that the improvements of IVT mainly benefit from better depth prediction. Combined with Figure 4 (a), (b), and (c), we can found that GAST-Net, PoseFormer, and IVT have similar results on frame 235. The MPJPE and depth error show the similar results at frame 235 since the human poses in this period are clear.

*4.3.2 Results on 3DPW Dataset.* The comparisons with state-of-the-art methods on the 3DPW dataset are shown in Table 2. In the multi-person 3DPW dataset, IVT outperforms previous video-based methods and achieves 46.0mm in PA-MPJPE. Compared with METRO [23], the relative gain is 6% in PA-MPJPE. These results verify the effectiveness and generalization ability of the proposed IVT since 3DPW is an in-the-wild dataset.

*4.3.3 Results on CMU Panoptic Dataset.* The comparisons with SOTA methods on CMU Panoptic dataset are shown in Table 3. It shows that, IVT obtains 48.4mm in MPJPE and achieves a new state-of-the-art result with a relative gain of 10% compared with the the DAS [43]. The results on CMU Panoptic dataset further demonstrate the effectiveness of the proposed IVT framework.

## 4.4 Ablation Study

In this section, we verify the effectiveness of the proposed ISA, ITA, and cross-scale attention mechanism in instance-guided video transformer (IVT). Based on IVT, we also study the influence of frame numbers on video transformer. Then, we compare the parameters and computational costs of different types of IVT.

*4.4.1 Effectiveness of proposed attention mechanisms.* We conduct the ablation study on CMU Panoptic dataset to verify the effectiveness of proposed modules in the instance-guided video transformer.

First of all, to verify the different types of attention mechanisms, we build an end-to-end multi-person 3D pose estimation baseline

**Table 3: Comparison with SOTA methods on multi-person 3D human pose estimation dataset (CMU Panoptic) in MPJPE. † means an extra refining network is used. Lower is better.**

| Methods | | MPJPE (mm) ↓ |
|---|---|---|
| SFB [47] | CVPR'18 | 153.4 |
| PoseNet [30] | ICCV'19 | 87.6 |
| MubyNet [48] | NeurIPS'18 | 78.1 |
| SMAP [50] | ECCV'20 | 73.1 |
| LoCO [11] | CVPR'20 | 69.0 |
| SMAP† [50] | ECCV'20 | 61.8 |
| DAS [43] | CVPR'22 | 53.8 |
| **Ours(IVT)** | | **48.4** |

**Table 4: Ablation study of IVT on CMU Panoptic. SA means traditional spatial attention. ISA means using instance-guided spatial attention in IVT. ITA means using instance-guided temporal attention in IVT. CISA represents cross-scale ISA. MITA means multi-scale ITA. Flops are computed on two input images with a size of $512 \times 512$.**

| Methods | Feature | Params (M) | Flops (T) | MPJPE (mm)↓ | Δ |
|---|---|---|---|---|---|
| Baseline | + SA | 32.96 | 0.127 | 52.0 | - |
| IVT | + ISA | 34.88 | 0.123 | 50.8 | ↓ 2.3% |
| IVT | + ISA + ITA | 35.81 | 0.126 | 49.5 | ↓ 2.6% |
| IVT | + CISA + MITA | 40.85 | 0.134 | **48.4** | ↓ 2.3% |

**Table 5: Ablation study of IVT (ISA-ITA) on frame numbers on CMU Panoptic dataset. Note that each video is sampled with a sampling rate of 5 frames. Thus, the temporal receptive field is $r = f \times 5$, where $f$ means the used frame number.**

| Frames | 1 | 3 | 5 | 7 | 9 |
|---|---|---|---|---|---|
| MPJPE (mm) | 51.8 | 50.3 | 49.5 | 49.5 | 50.0 |

with traditional simple spatial attention as [2, 10], noted as SA. As shown in Table 4, SA achieves 52.0mm in MPJPE. Combined with the instance-guided attention, IVT with ISA obtains 50.8mm in MPJPE and achieves a relative gain of 2.3%. Compared with only using ISA, the ITA obtains 49.5mm in MPJPE and achieves a relative gain of 2.6%. Besides, the cross-scale attention mechanism brings a relative gain of 2.3% and achieves 48.4mm in MPJPE. These results show that the proposed ISA, ITA, and cross-scale attention mechanism are useful for 3D human pose estimation.

*4.4.2 The ablation study on frame number.* To explore the influence of frame numbers for IVT, we conduct the ablation study of frame numbers based on IVT with ITA. As shown in Table 5, the frames number means the used frames, but the truth receptive field on time is $f \times 5$ since the sampling rate of the video is 5 frames. For example, $f = 5$ means we used 5 frames but the interval is 25 from the first frame to the last frame. As shown in Table 5, given more frames from 1 to 5, the performance of IVT improves to 49.5mm from 50.8mm. But with the frame number increasing to 9, the performance of IVT decreases to 50.00mm. The result shows that long-term frames could damage the performance of IVT since the long-term motion is hard to estimate.

**(a) Visualization results on video frames**



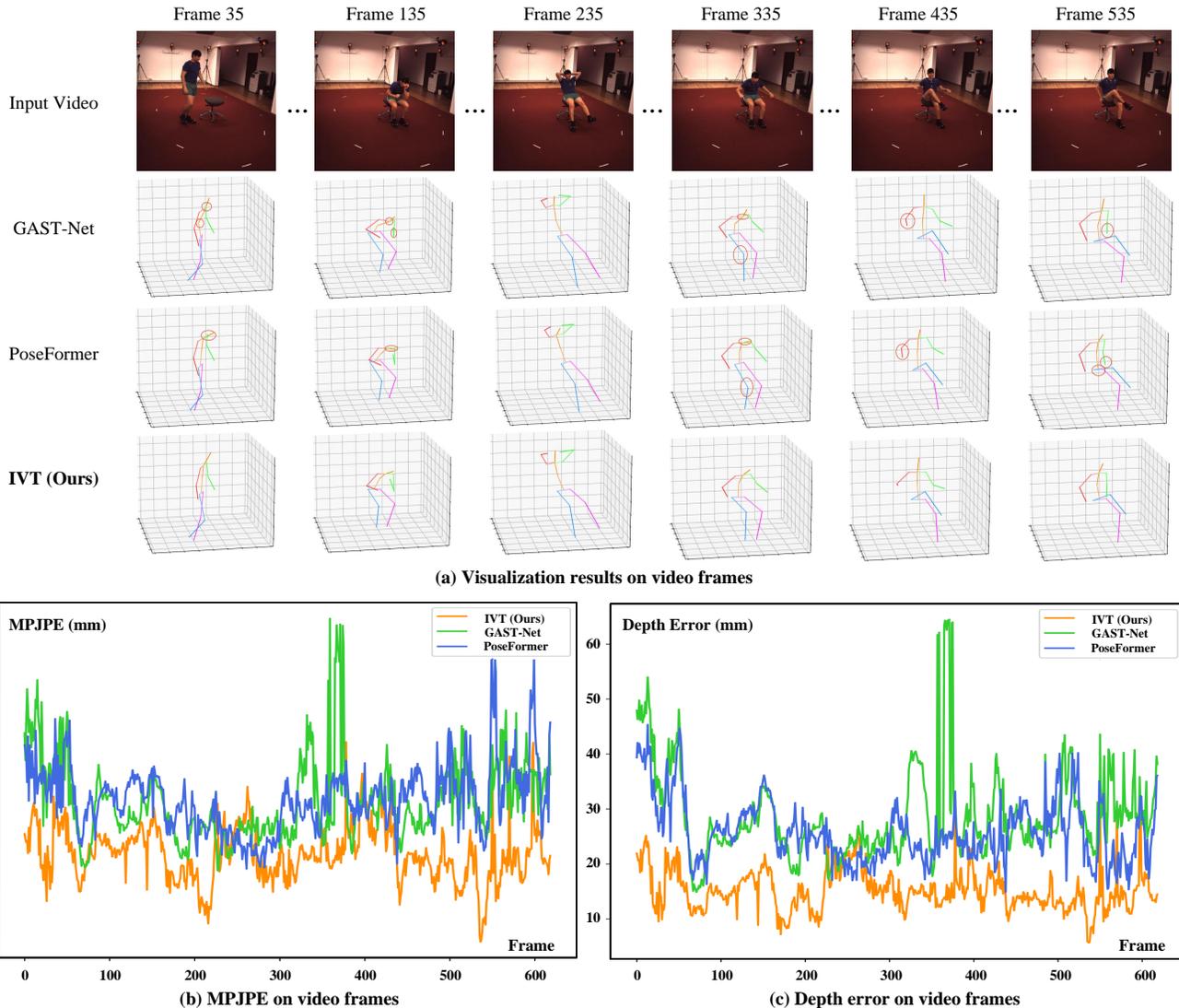**(b) MPJPE on video frames**                 **(c) Depth error on video frames**

**Figure 4: Qualitative comparison between IVT and the SOTA video methods (GAST-Net [25] and PoseFormer [51]) on video from Human3.6M dataset. (a) visualization results (Red circle denotes the wrong prediction). (b) The curve of MPJPE-Frame shows the MPJPE on each frame. (c) The curve of Depth Error-Frame shows the MPJPE in depth dimension on each frame. The depth prediction by IVT is better than GAST-Net and PoseFormer. Best viewed in color.**

*4.4.3 Parameters and computational costs.* The comparisons of different attention mechanisms on parameters and computational costs are shown in Table 4. Compared with the baseline with spatial attention, instance-guided attention (ISA) obtains a relative gain of 2.3% by adding 1.92MB parameters, while the flops of ISA decrease since the instance-aware tokens are based on image blocks. Even for the IVT with CISA and ITA, the increasing parameters and the computational costs are acceptable.

## 5 CONCLUSIONS

In this paper, we propose a novel end-to-end instance-guided video transformer (IVT) for video 3D human pose estimation to capture

global depth context. To capture the depth context in both spatial and temporal dimensions, we propose instance-guided spatial attention (ISA) and instance-guided temporal attention (ITA) mechanisms. To further tackle the variational human scales in video, we propose cross-scale attention for IVT. Combined with ISA, ITA, and cross-scale attention, IVT outperforms state-of-the-art methods on three widely-used 3D human pose estimation datasets.

## 6 ACKNOWLEDGEMENT

# REFERENCES

[1] Anurag Arnab, Carl Doersch, and Andrew Zisserman. 2019. Exploiting temporal context for 3D human pose estimation in the wild. In *CVPR*. 3395–3404.

[2] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. 2021. Is Space-Time Attention All You Need for Video Understanding?. In *ICML*. PMLR, 813–824.

[3] Yujun Cai, Liuhao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann. 2019. Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks. In *ICCV*. 2272–2281.

[4] Tianlang Chen, Chen Fang, Xiaohui Shen, Yiheng Zhu, Zhili Chen, and Jiebo Luo. 2021. Anatomy-aware 3d human pose estimation with bone-based pose decomposition. *TCSVT* 32, 1 (2021), 198–209.

[5] Yu Cheng, Bo Wang, Bo Yang, and Robby T Tan. 2021. Monocular 3D multi-person pose estimation by integrating top-down and bottom-up networks. In *CVPR*. 7649–7659.

[6] Yu Cheng, Bo Yang, Bo Wang, and Robby T Tan. 2020. 3d human pose estimation using spatio-temporal networks with explicit occlusion training. In *AAAI*, Vol. 34. 10631–10638.

[7] Hongsuk Choi, Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. 2021. Beyond static features for temporally consistent 3d human pose and shape from a video. In *CVPR*. 1964–1973.

[8] Rishabh Dabral, Anurag Mundhada, Uday Kusupati, Safeer Afaque, Abhishek Sharma, and Arjun Jain. 2018. Learning 3d human pose from structure and motion. In *ECCV*. 668–683.

[9] Carl Doersch and Andrew Zisserman. 2019. Sim2real transfer learning for 3d human pose estimation: motion to the rescue. *NeurIPS* 32 (2019).

[10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *ICLR*.

[11] Matteo Fabbri, Fabio Lanzi, Simone Calderara, Stefano Alletto, and Rita Cucchiara. 2020. Compressed volumetric heatmaps for multi-person 3d pose estimation. In *CVPR*. 7204–7213.

[12] John C Gower. 1975. Generalized procrustes analysis. *Psychometrika* 40, 1 (1975), 33–51.

[13] Mir Rayat Imtiaz Hossain and James J Little. 2018. Exploiting temporal information for 3d human pose estimation. In *ECCV*. 68–84.

[14] Yupan Huang, Hongwei Xue, Bei Liu, and Yutong Lu. 2021. Unifying multimodal transformer for bi-directional image and text generation. In *ACM MM22*. 1138–1147.

[15] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu. 2013. Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. *TPAMI* 36, 7 (2013), 1325–1339.

[16] Hanbyul Joo, Tomas Simon, Xulong Li, Hao Liu, Lei Tan, Lin Gui, Sean Banerjee, Timothy Godisart, Bart Nabbe, Iain Matthews, et al. 2017. Panoptic studio: A massively multiview system for social interaction capture. *TPAMI* 41, 1 (2017), 190–204.

[17] Angjoo Kanazawa, Jason Y Zhang, Panna Felsen, and Jitendra Malik. 2019. Learning 3d human dynamics from video. In *CVPR*. 5614–5623.

[18] Muhammed Kocabas, Nikos Athanasiou, and Michael J Black. 2020. Vibe: Video inference for human body pose and shape estimation. In *CVPR*. 5253–5263.

[19] Nikos Kolotouros, Georgios Pavlakos, Michael J Black, and Kostas Daniilidis. 2019. Learning to reconstruct 3D human pose and shape via model-fitting in the loop. In *ICCV*. 2252–2261.

[20] Ke Li, Shijie Wang, Xiang Zhang, Yifan Xu, Weijian Xu, and Zhuowen Tu. 2021. Pose recognition with cascade transformers. In *CVPR*. 1944–1953.

[21] Yong-Lu Li, Xinpeng Liu, Han Lu, Shiyi Wang, Junqi Liu, Jiefeng Li, and Cewu Lu. 2020. Detailed 2d-3d joint representation for human-object interaction. In *CVPR*. 10166–10175.

[22] Jiahao Lin and Gim Hee Lee. 2019. Trajectory space factorization for deep video-based 3d human pose estimation. In *BMVC*.

[23] Kevin Lin, Lijuan Wang, and Zicheng Liu. 2021. End-to-end human pose and mesh reconstruction with transformers. In *CVPR*. 1954–1963.

[24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*. Springer, 740–755.

[25] Junfa Liu, Juan Rojas, Yihui Li, Zhijun Liang, Yisheng Guan, Ning Xi, and Haifei Zhu. 2021. A graph attention spatio-temporal convolutional network for 3D human pose estimation in video. In *ICRA*. IEEE, 3374–3380.

[26] Ruixu Liu, Ju Shen, He Wang, Chen Chen, Sen-ching Cheung, and Vijayan Asari. 2020. Attention mechanism exploits temporal contexts: Real-time 3d human pose reconstruction. In *CVPR*. 5064–5073.

[27] Ziyu Liu, Hongwen Zhang, Zhenghao Chen, Zhiyong Wang, and Wanli Ouyang. 2020. Disentangling and unifying graph convolutions for skeleton-based action recognition. In *CVPR*. 143–152.

[28] Zhengyi Luo, S Alireza Golestaneh, and Kris M Kitani. 2020. 3d human motion estimation via motion compression and refinement. In *ACCV*.

[29] Weian Mao, Yongtao Ge, Chunhua Shen, Zhi Tian, Xinlong Wang, and Zhibin Wang. 2021. Tfpose: Direct human pose estimation with transformers. *arXiv preprint arXiv:2103.15320* (2021).

[30] Gyeongsik Moon, Ju Yong Chang, and Kyoung Mu Lee. 2019. Camera distance-aware top-down approach for 3d multi-person pose estimation from a single rgb image. In *ICCV*. 10133–10142.

[31] Xuecheng Nie, Jiashi Feng, Jianfeng Zhang, and Shuicheng Yan. 2019. Single-stage multi-person pose machines. In *ICCV*. 6951–6960.

[32] Mathias Parger, Chengcheng Tang, Yuanlu Xu, Christopher David Twigg, Lingling Tao, Yijing Li, Robert Wang, and Markus Steinberger. 2021. UNOC: Understanding occlusion for embodied presence in virtual reality. *TVCG* (2021).

[33] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 2019. 3d human pose estimation in video with temporal convolutions and semi-supervised training. In *CVPR*. 7753–7762.

[34] Zhongwei Qiu, Kai Qiu, Jianlong Fu, and Dongmei Fu. 2019. Learning recurrent structure-guided attention network for multi-person pose estimation. In *ICME*. IEEE, 418–423.

[35] Zhongwei Qiu, Kai Qiu, Jianlong Fu, and Dongmei Fu. 2020. Dgcn: Dynamic graph convolutional network for efficient multi-person pose estimation. In *AAAI*, Vol. 34. 11924–11931.

[36] Anurag Ranjan and Michael J Black. 2017. Optical flow estimation using a spatial pyramid network. In *CVPR*. 4161–4170.

[37] Yu Sun, Qian Bao, Wu Liu, Yili Fu, Michael J Black, and Tao Mei. 2021. Monocular, one-stage, regression of multiple 3d people. In *ICCV*. 11179–11188.

[38] Yu Sun, Yun Ye, Wu Liu, Wenpeng Gao, Yili Fu, and Tao Mei. 2019. Human mesh recovery from monocular images via a skeleton-disentangled representation. In *ICCV*. 5349–5358.

[39] Timo von Marcard, Roberto Henschel, Michael J Black, Bodo Rosenhahn, and Gerard Pons-Moll. 2018. Recovering accurate 3d human pose in the wild using imus and a moving camera. In *ECCV*. 601–617.

[40] Can Wang, Jiefeng Li, Wentao Liu, Chen Qian, and Cewu Lu. 2020. Hmor: Hierarchical multi-person ordinal relations for monocular multi-person 3d pose estimation. In *ECCV*. Springer, 242–259.

[41] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. 2020. Deep high-resolution representation learning for visual recognition. *TPAMI* 43, 10 (2020), 3349–3364.

[42] Jingbo Wang, Sijie Yan, Yuanjun Xiong, and Dahua Lin. 2020. Motion guided 3d pose estimation from videos. In *ECCV*. Springer, 764–780.

[43] Zitian Wang, Xuecheng Nie, Xiaochao Qu, Yunpeng Chen, and Si Liu. 2022. Distribution-Aware Single-Stage Models for Multi-Person 3D Pose Estimation. In *CVPR*.

[44] Sen Yang, Zhibin Quan, Mu Nie, and Wankou Yang. 2021. Transpose: Keypoint localization via transformer. In *ICCV*. 11802–11812.

[45] Wei Yang, Wanli Ouyang, Xiaolong Wang, Jimmy Ren, Hongsheng Li, and Xiaogang Wang. 2018. 3d human pose estimation in the wild by adversarial learning. In *CVPR*. 5255–5264.

[46] Raymond Yeh, Yuan-Ting Hu, and Alexander Schwing. 2019. Chirality nets for human pose regression. *NeurIPS* 32 (2019).

[47] Andrei Zanfir, Elisabeta Marinoiu, and Cristian Sminchisescu. 2018. Monocular 3d pose and shape estimation of multiple people in natural scenes-the importance of multiple scene constraints. In *CVPR*. 2148–2157.

[48] Andrei Zanfir, Elisabeta Marinoiu, Mihai Zanfir, Alin-Ionut Popa, and Cristian Sminchisescu. 2018. Deep network for the integrated 3d sensing of multiple people in natural images. *NeurIPS* 31 (2018).

[49] Ailing Zeng, Xiao Sun, Fuyang Huang, Minhao Liu, Qiang Xu, and Stephen Lin. 2020. Srnet: Improving generalization in 3d human pose estimation with a split-and-recombine approach. In *ECCV*. Springer, 507–523.

[50] Jianan Zhen, Qi Fang, Jiaming Sun, Wentao Liu, Wei Jiang, Hujun Bao, and Xiaowei Zhou. 2020. Smap: Single-shot multi-person absolute 3d pose estimation. In *ECCV*. Springer, 550–566.

[51] Ce Zheng, Sijie Zhu, Matias Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding. 2021. 3d human pose estimation with spatial and temporal transformers. In *ICCV*. 11656–11665.

[52] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. 2019. Objects as points. *arXiv preprint arXiv:1904.07850* (2019).