41

42

43

44

45

46

47

48

49

50

51

52

53

54

58

1

# APP: Adaptive Pose Pooling for 3D Human Pose Estimation from Videos

Anonymous Authors

## ABSTRACT

Current advancements in 3D human pose estimation have attained notable success by converting 2D poses into their 3D counterparts. However, this approach is inherently influenced by the errors introduced by 2D pose detectors and overlooks the intrinsic spatial information embedded within RGB images. To address these challenges, we introduce a versatile module called Adaptive Pose Pooling (APP), compatible with many existing 2D-to-3D lifting models. The APP module includes three novel sub-modules: Pose-Aware Offsets Generation (PAOG), Pose-Aware Sampling (PAS), and Spatial Temporal Information Fusion (STIF). First, we extract latent features of the multi-frame lifting model. Then, a 2D pose detector is utilized to extract multi-level feature maps from the image. After that, PAOG generates offsets according to featuremaps. PAS uses offsets to sample featuremaps. Then, STIF can fuse PAS sampling features and latent features. This innovative design allows the APP module to simultaneously capture spatial and temporal information. We conduct comprehensive experiments on two widely used datasets: Human3.6M and MPI-INF-3DHP. Meanwhile, we employ various lifting models to demonstrate the efficacy of the APP module. Our results show that the proposed APP module consistently enhances the performance of lifting models, achieving state-of-the-art results. Significantly, our module achieves these performance boosts without necessitating alterations to the architecture of the lifting model. Code and checkpoints are available at: Anonymous Github.

# CCS CONCEPTS

• Computing methodologies  $\rightarrow$  Activity recognition and understanding.

### KEYWORDS

3D Human Pose Estimation, Adaptive Pose Pooling, Feature Fusion

# **1 INTRODUCTION**

3D human pose estimation (HPE) is a critical task in computer vision, and its goal is to estimate 3D human poses from images, videos, and 2D human poses. It attempts to estimate precise positions and orientations of 3D human keypoints. The applications of 3D HPE are pretty large, such as virtual/augmented reality [24, 29], motion analysis [2, 8, 43], human-computer interaction [1, 12], and autonomous driving [13] etc. Models for 3D HPE can be broadly

Unpublished working draft. Not for distribution.

on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM



(b) Our Method

Figure 1: Difference between our method and previous methods. (a) Previous methods freeze the 2D pose detector and train the lifting model. (b) Our proposed method utilizes the 2D pose features extracted by the lifting model and takes multi-level feature maps extracted by the 2D pose detector as inputs. Our method addresses the issue of underutilizing image features by leveraging multi-level feature maps extracted by a pretrained 2D pose detector. Furthermore, while single-frame models utilize image features, they lack temporal information. However, the APP module addresses this by extracting temporal and spatial information simultaneously.

categorized into two fundamental classes based on the number of viewpoints employed: multi-view and single-view. Multi-view approaches leverage information from multiple camera perspectives, as explored in studies such as [6, 14, 17, 34, 46]. They mitigate depth ambiguity and reach better results than monocular 3D HPE methods. However, their practical applicability is hindered by the necessity of employing multiple cameras.

In contrast, single-view models offer a more applicable alternative and can be subdivided into single-stage and two-stage methodologies. Single-stage methods [3, 18, 25, 31, 40, 50] directly infer 3D human pose from input images or video frames without intermediate steps. While two-stage like early methods [26, 32] adopt a sequential approach. As shown in Figure 1(a), off-the-shelf 2D pose detectors [5, 30, 39] extract 2D human poses from input images or video frames. After that, they can obtain corresponding 3D human poses based on these 2D estimations, so this stage is often called the lifting stage. Models in the lifting stage can further be categorized into single-frame or multi-frame variants. Single-frame models [7, 47] infer 3D human poses from individual images, while multi-frame models process sequences of 2D human poses. The latter one contains two categories: many-to-one and many-to-many

116

59

60

republish, to post on servers or to redistribute to lists, requires prior specific

AGNARA 2024 AS U

<sup>55</sup> ACM MM, 2024, Melbourne, Australia

<sup>56 © 2024</sup> Copyright held by the owner/author(s). Publication rights licensed to ACM.

<sup>57</sup> https://doi.org/10.1145/nnnnnnnnnnn

methods. Many-to-one models estimating the 3D human pose of
a central frame using multiple 2D human poses [21, 22, 48, 49],
and many-to-many, also called seq2seq models inferring every
frame of a 3D human pose from corresponding 2D human poses
[4, 27, 33, 37, 45, 51]. They try to address different problems, ranging from spatio-temporal information exploration to enhancing
model robustness.

In this paper, we focus on monocular 3D HPE. It uses only one 124 125 camera to obtain images of humans. Therefore, it introduces ambi-126 guity because of the lack of depth information, which is an ill-posed task. Besides, occlusion, including self-occlusion and occlusion by 127 other objects and persons, needs to be drawn to researchers' atten-128 tion. Hence, single-view models sometimes fail to get accurate and 129 robust 3D human poses owing to depth ambiguity and occlusion. 130 Many monocular 3D HPE models suffer from them because 131 the information of the input images fails to be fully used. 132 HEMlets Pose [50] and Context-Aware PoseFormer (CA-PF) [47] try 133 to handle this. HEMlets Pose explicitly lets the model learn relative 134 135 depth among different keypoints in 2D human poses. In comparison, CA-PF adopts Deformable Attention [52] to extract features 136 137 of feature maps obtained by the 2D pose detector. However, these 138 methods [47, 50] are single-frame models that cannot capture 139 temporal dynamics. Recently, multi-frame lifting models have reached significant improvement. Why not use both image and 140 temporal features? This work proposes a well-designed Adaptive 141 142 Pose Pooling (APP) module. Serving as a plug-and-play component that can be integrated with many existing lifting models. The APP 143 module facilitates extracting temporal and spatial information from 144 multi-frame 2D human pose data and feature maps. Figure 1 illus-145 trates the distinction between previous methods and our proposed 146 method. Previous methods freeze the 2D pose detector and only 147 148 take detected 2D poses as the model input to estimate correspond-149 ing 3D poses. Our method uses feature maps extracted by a 2D pose detector and takes multi-frame lifting models' calculated pose 150 151 features. Therefore, the APP module inherently can capture spatio-152 temporal features. Do we really need to use all the features of the feature maps? Features near 2D poses should be more important 153 than features at other positions. To enhance the model's flexibility, 154 155 the model should adjust the sampling points during training. In the meantime, getting spatial and temporal features and interacting 156 among features is also important for the model. 157

Our proposed APP module contains three novel submodules: 158 159 Pose-Aware Offsets Generation (PAOG), Pose-Aware Sampling (PAS), and Spatial Temporal Information Fusion (STIF). To address the 160 161 problem of insufficient utilization of spatial information in images, 162 we propose PAOG and PAS. PAOG leverages 2D human poses as indices to learn offsets during training. PAS adaptively extracts 163 features from feature maps generated by the 2D human pose de-164 tector. Though single-frame models use image features, they suffer 165 from lacking temporal features. We proposed STIF to solve this 166 problem by fusing image and temporal features. We conduct ex-167 tensive experiments on two popular datasets: Human3.6M and 168 MPI-INF-3DHP. Our proposed method achieves a new state-of-the-169 art performance, with 1.3mm (from 37.5mm to 36.2mm) and 1.1mm 170 (from 30.6mm to 29.5mm) improvement in terms of Mean Per-Joint 171 172 Position Error (MPJPE) and P-MPJPE (Procrustes-MPJPE) on Hu-173 man3.6M. There are 4.6 (from 85.9 to 89.5) and 4mm (from 16.7mm

174

to 12.7mm) improvements for Area Under Curve (AUC) and MPJPE on MPI-INF-3DHP.

Compared to current methods for monocular 3D HPE, our contributions are summarized as follows:

- We propose the APP module, enabling multi-frame lifting models to utilize information from the previous stage. This module is a plug-and-play solution compatible with most multi-frame lifting models and mitigates the problem of excessive reliance on 2D human pose outputs.
- We introduce Pose-Aware Offsets Generation (PAOG) and Pose-Aware Sampling (PAS). PAOG adjusts offsets relative to 2D human poses during training, allowing PAS to learn to adaptively capture spatial features from feature maps.
- We employ Spatial Temporal Information Fusion (STIF) to facilitate information interaction between the features sampled by PAS and the hidden feature of the lifting model. This enables the APP module to capture both temporal and spatial features simultaneously.
- We conduct experiments on two widely used datasets: Human3.6M and MPI-INF-3DHP. Our proposed approach achieves state-of-the-art performance on both datasets compared to previous methods. Importantly, it does not alter the lifting model's architecture.

### 2 RELATED WORK

Transformers[42] initially dominated the domain of natural language processing. Owing to its expressive performance, it has gained significant attention in computer vision as well [10]. Therefore, Transformer-based models have been increasingly applied across various domains, including 3D human pose estimation. Pose-Former [49] stands out as the pioneering model to integrate Transformers into 3D human pose estimation. Building upon this foundation, PoseFormerV2 [48] represents an advancement, introducing a frequency domain approach to enhance model robustness against fast motion changes. Furthermore, MHFormer [22] addresses the challenge of depth ambiguity by learning the representations of multiple hypothesis poses.

In the quest for improved efficiency and accuracy, novel architectural modifications have emerged. StridedTrans [21] utilizes fully connected layers in the Transformer encoder with stride convolutions. It reduces computational complexity and progressively shortens sequence lengths, thus enhancing pose estimation accuracy for the central frame. Meanwhile, MixSTE [45] focuses on modeling strong intra-frame correspondences between individual joints, thereby facilitating the learning of spatio-temporal correlations. P-STMO [37] adopts a strategy of occlusion masking to enable the model to learn more resilient pose representations. Additionally, it leverages a large volume of unlabeled data for self-supervised pretraining, thereby augmenting generalization capabilities. HoT [23] achieves this by implementing feature compression and restoration techniques, rendering it compatible with various lifting models while reducing computational overhead.

Advancements in attention mechanisms have also been pivotal. HDFormer [4], for instance, integrates self-attention and higherorder attention mechanisms to construct a hierarchical attention module aimed at mitigating complexities and handling heavily 227

228

229

230

231



Figure 2: Our proposed APP module. (a) APP module contains three novel sub-modules: Pose-Aware Offsets Generation (PAOG), Pose-Aware Sampling (PAS), and Spatial Temporal Information Fusion (STIF). STIF comprises Feature Updating and Multi-Head Cross Attention, which is not shown in this figure. Feature Updating is used to obtain  $\bar{F}^l$  and  $X^l$ . MHCA fuses the hidden features  $X^l$  with pose features *P*. (b) PAOG learns sampling offsets and utilizes these offsets to guide feature extraction. (c) PAS leverages learned offsets of PAOG and 2D human poses to sample multi-level feature maps.

occluded scenarios. STCFormer [41] decomposes correlation learning into spatial and temporal dimensions, thereby reducing model complexity. Moreover, MotionBERT [51] emerges as a dual-stream spatio-temporal Transformer capable of concurrently capturing spatio-temporal information, showcasing state-of-the-art performance across multiple downstream tasks. Similarly, MotionAG-Former [27] introduces a hybrid architecture composed of multiple dual-stream modules, effectively capturing both spatial and temporal features of keypoints in tandem. KTPFormer[33] tries to overcome the weakness in existing transformer-based methods by proposing two plug-and-play attention modules, namely Kinematics Prior Attention (KPA) and Trajectory Prior Attention (TPA). These modules take advantage of the structure of the human body and motion trajectory information.

For the 2D-to-3D lifting stage, given a sequence of 2D human poses  $\boldsymbol{p} \in \mathbb{R}^{T \times J \times C_{in}}$ , the goal is to estimate the corresponding 3D human poses  $\hat{\boldsymbol{p}} \in \mathbb{R}^{T \times J \times C_{out}}$ , where *T* is the sequence length of 2D pose frames, *J* is the number of keypoints of the 2D human pose, *C<sub>in</sub>* and *C<sub>out</sub>* are the dimension of the 2D and 3D human pose keypoints, respectively.

These advancements collectively underscore the profound impact and ongoing evolution of Transformer-based approaches in advancing the field of 3D human pose estimation. Unlike our proposed method, which integrates image features into estimating 3D human poses, none of the existing transformer-based multi-frame methods take advantage of such information. This presents a unique opportunity for our approach to bridge this gap and potentially unlock new avenues for improving the accuracy and robustness of 3D human pose estimation methods.

### 3 APPROACH

### 3.1 Framework Overview

Figure 2 illustrates the overview of our proposed method. Given input images  $I \in \mathbb{R}^{t \times 3 \times H \times W}$ , where *t* is the number of image frames and *H*, *W* are the height and width of the image, note that  $t \ll T$ . We first use a 2D pose detector to extract the multi-level feature maps denoted by  $\{H_i\}_{s=1}^S$ , where *S* is the number of scales. *H*<sub>1</sub> is a heatmap to get the 2D pose *p*. Then pose features  $P \in \mathbb{R}^{T \times J \times dP}$  is calculated in the multi-frame lifting models without regression head, where *d*<sub>*P*</sub> is the dimension of pose features. This process is formulated as:

 $MLM(\cdot)$  denotes multi-frame lifting models without regression head; the backbone is a 2D pose detector. In Figure 2, the values of averaged image features  $\bar{F} \in \mathbb{R}^{T \times (k^2 \times J) \times d}$  and hidden features  $X \in \mathbb{R}^{T \times J \times d}$  are updated by each layer of the APP module, where k is the window size.

To get the initialized offsets  $\Delta^0 \in \mathbb{R}^{T \times (k^2 \times J) \times C_{in}}$ , we repeat the values by taking each point in the 2D human poses p as a center and generating a  $k \times k$  window around it. Then the p is resized to  $\mathbb{R}^{T \times (k^2 \times J) \times C_{in}}$ . Now, the 2D poses p data is a 2D vector, which means the x and y coordinates of the human keypoints.

We denote  $\bar{F}^{l-1}$  and  $X^{l-1}$  as the last layer outputs of the APP module. As for the first layer of APP, the offsets  $\Delta^0$  are used to extract the initialized  $\bar{F}^0$  and  $X^0$ :

$$\{H_{s}\}_{s=1}^{S} = \theta_{f}^{0}(\{H_{s}\}_{s=1}^{S}),$$
  

$$\{F_{s}\}_{s=1}^{S} = GridSample(\{H_{s}\}_{s=1}^{S}, \Delta^{0}),$$
  

$$\bar{F}^{0} = \frac{1}{S}\sum_{s=1}^{S} F_{s},$$
(2)

 $S_{s=1}^{o}$  $X^{0} = AveragePool(\bar{F}^{0})$ 

Where  $\theta_f$  is a group of 2D convolution operators to project multilevel feature maps into the same dimension d, the implementation of *GridSample* is consistent with **nn.functional.grid\_sample** in PyTorch, *AveragePool* is to get the mean value of each  $k \times k$  window.

Our proposed APP module consists of *L* stacked layers. In detail, the calculation process of each layer APP module is defined as follows:

$$\bar{F}^{l}, X^{l} = APP^{l}(\bar{F}^{l-1}, X^{l-1}, \boldsymbol{p}, \{H_{i}\}_{s=1}^{S}, \boldsymbol{P})$$
(3)

In the last layer of the APP module, a simple regression head consisting of two fully connected layers estimates the final 3D human pose given by  $X^{l}$ .

### 3.2 Adaptive Pose Pooling

3.2.1 Pose-Aware Offsets Generation. Inspired by deformable convolution and deformable attention [9, 52], we propose the PAOG module as shown in Figure 2 (b). PAOG adaptively learns offset points during training. It takes the previous layer features  $\bar{F}^{l-1}$  and  $X^{l-1}$ , and the 2D poses  $P_p^{2D}$  as input to generate learned offsets  $\Delta^l$ , which is defined as follows:

$$\Delta_p^l = PAOG^l(\bar{F}^{l-1}, X^{l-1}, p)$$
  
=  $\theta_s^l(\bar{F}^{l-1}) \otimes \theta_w^l(X^{l-1}) + p$  (4)

Where  $\otimes$  denotes element-wise multiplication,  $\theta_s^l$  and  $\theta_w^l$  generate offset points based on the image center and their corresponding weights. Note that they are both fully connected layers.

3.2.2 Pose-Aware Sampling. Using the offset points  $\Delta_p^l$  generated by PAOG and the previous averaged image features  $\bar{F}^{l-1}$ , PAS also utilizes the multi-level feature maps  $\{H\}_{s=1}^S \in \mathbb{R}^{t \times C_s \times H_s \times W_s}$ extracted from the 2D pose detector. We follow Deformable-DETR [52] before sampling features from the multi-level feature maps F. They are projected into the same dimension d as in Eq. (2). The sampling process is then defined as:

$$\{H_{s}\}_{s=1}^{S} = \theta_{f}^{l}(\{H_{s}\}_{s=1}^{S}),$$
  

$$\{F_{s}\}_{s=1}^{S} = GridSample(\{H_{s}\}_{s=1}^{S}, \Delta^{l}),$$
  

$$\bar{F}^{l} = \frac{1}{S}\sum_{s=1}^{S} F_{s}$$
(5)

*3.2.3 Feature Updating.* At the *l* layer of the APP module, we update the averaged image features  $\bar{F}^l$  and hidden features  $X^l$  as

follows:

$$=\theta_k^l(\bar{F}^l),\tag{6}$$

$$\bar{F}^{l} = \alpha \bar{F}^{l-1} + (1-\alpha) \bar{F}^{l}$$
(6) 409
410
410
410

Where  $\theta_k^l(\cdot)$  represents a convolution operation with a kernel size of  $k^2$ , whose size is the same as the previously mentioned  $k \times k$  window, this operation allows the PAS to learn the weights of  $k^2$  features within the window to calculate  $X^l$  given by  $\bar{F}^l$ .  $\alpha$  is a parameter controlling the magnitude of the update for  $\bar{F}^l$ .

 $X^l$ 

3.2.4 *Multi-Head Cross Attention.* MHCA enables the APP module to harness the temporal features extracted by the multi-frame lifting models while interacting with the spatial features extracted by PAS. We first combine the pose feature P extracted with positional embeddings *PE*, initialized to zero, and updated during training. We take  $X^l$  as Q, P as K, and V, then we perform as follows:

$$Q = X^{l}W_{Q}, K = PW_{K}, V = PW_{V},$$

$$X^{l} = MHCA^{l}(Q, K, V)$$

$$= Carrect(head = head)$$
(7)

$$= Concat(head_1, ..., head_h), \tag{7}$$

$$ad_j = Softmax\left(\frac{Q_j K_j}{\sqrt{d_P}}\right) V_j$$

Where  $W_Q \in \mathbb{R}^{d \times d}$ ,  $W_K \in \mathbb{R}^{d_P \times d}$ , and  $W_V \in \mathbb{R}^{d_P \times d}$ , j(j = 1, ..., h) represents the *j*-th attention head, and *h* is the number of attention heads.

3.2.5Feed Forward Network. Like many monocular 3D HPE transformer.434based methods, we apply an FFN at the end to finish computing<br/>the current layer. We use GeLU [15] as the activation function, as435outlined in Eq. (8).437

$$H^{l} = FFN^{l} = FC(GeLU(FC(X^{l})))$$
(8)

3.2.6 Loss Function. The loss function of our proposed APP module is the same as [27, 51]. In addition to the widely used MPJPE loss, we also incorporate the Mean Per-Joint Velocity Error (MPJVE) loss and the normalized MPJPE [36] loss. Let  $\mathcal{L}_p$ ,  $\mathcal{L}_v$ , and  $\mathcal{L}_s$  denote these three losses.  $\mathcal{L}_p$  calculates the  $L_2$  distance between estimated 3D human poses  $\hat{p}$  and ground truth g, while  $\mathcal{L}_v$  minimizes the difference between  $\Delta \hat{p}$  and  $\Delta g$ .  $\mathcal{L}_p$  and  $\mathcal{L}_v$  are thus given by:

$$\mathcal{L}_{p} = \sum_{t=1}^{T} \sum_{j=1}^{J} \|\hat{p}_{tj} - g_{tj}\|_{2},$$
(9)

$$\mathcal{L}_{v} = \sum_{t=2}^{I} \sum_{j=1}^{J} ||\Delta \hat{\boldsymbol{p}}_{tj} - \Delta \boldsymbol{g}_{tj}||_{2}$$

Where  $\Delta \hat{p}_t = \hat{p}_t - \hat{p}_{t-1}$ ,  $\Delta g_t = g_t - g_{t-1}$ . For  $\mathcal{L}_s$ , estimated 3D poses  $\hat{p}$  is scaled according to ground truth g. Then, we calculate MPJPE loss between scaled  $\hat{p}$  and g. The total loss is formulated as follows:

$$\mathcal{L} = \gamma_p \mathcal{L}_p + \gamma_v \mathcal{L}_v + \gamma_s \mathcal{L}_s \tag{10}$$

where  $\gamma_p$ ,  $\gamma_v$ , and  $\gamma_s$  are weights for  $\mathcal{L}_p$ ,  $\mathcal{L}_v$ , and  $\mathcal{L}_s$ , respectively.

he

### Table 1: Qualitative comparisons of 3D human pose estimation per action on Human3.6M. The best and second-best results are bolded and blue, respectively. T: Number of the input frames. Seq2Seq: Estimate 3D human pose for every frame.

| MPIPF                                | Sea2Sea      | Ιт  | Dir  | Disc   | Fat  | Greet  | Phone    | Photo  | Pose         | Purch        | Sit          | SitD   | Smoke   | Wait         | WalkD  | Walk             | WalkT        | Avo   |
|--------------------------------------|--------------|-----|------|--------|------|--------|----------|--------|--------------|--------------|--------------|--------|---------|--------------|--------|------------------|--------------|-------|
|                                      | 1 3042304    | 1   |      | 10.4   | Lat  | 01001  | 1 110110 | 111010 | 1030         | 1 41 0       | 50.5         | 511.D. | 51110KC | wait         | waikD. | waik             | waik i.      | 1108  |
| MHFormer [22] CVPR 22                |              | 351 | 39.2 | 43.1   | 40.1 | 40.9   | 44.9     | 51.2   | 40.6         | 41.3         | 53.5         | 60.3   | 43.7    | 41.1         | 43.8   | 29.8             | 30.6         | 43.0  |
| Strided Irans [21] IMM 22            |              | 351 | 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.   |
| MixSTE [45] CVPR 22                  | V            | 243 | 37.6 | 40.9   | 37.3 | 39.7   | 42.3     | 49.9   | 40.1         | 39.8         | 51.7         | 55.0   | 42.1    | 39.8         | 41.0   | 27.9             | 27.9         | 40.   |
| MixSTE [45] CVPR 22                  | ✓            | 243 | 36.7 | 39.0   | 36.5 | 39.4   | 40.2     | 44.9   | 39.8         | 36.9         | 47.9         | 54.8   | 39.6    | 37.8         | 39.3   | 29.7             | 30.6         | 39.   |
| P-STMO [37] ECCV 22                  |              | 243 | 38.9 | 42.7   | 40.4 | 41.1   | 45.6     | 49.7   | 40.9         | 39.9         | 55.5         | 59.4   | 44.9    | 42.2         | 42.7   | 29.4             | 29.4         | 42.   |
| CA-PF-HRNet-48 [47] NeurIPS 23       |              | 1   | -    | -      | -    | -      | -        | -      | -            | -            | -            | -      | -       | -            | -      | -                |              | 39.   |
| HDFormer [4] IJCAI 23                | <b>√</b>     | 96  | 38.1 | 43.1   | 39.3 | 39.4   | 44.3     | 49.1   | 41.3         | 40.8         | 53.1         | 62.1   | 43.3    | 41.8         | 43.1   | 31.0             | 29.7         | 42.   |
| HDFormer [4] IJCAI 23                | ✓            | 96  | 34.7 | 41.7   | 36.0 | 38.4   | 41.1     | 45.3   | 39.6         | 37.4         | 49.0         | 63.1   | 39.8    | 38.9         | 40.2   | 29.3             | 29.1         | 40.   |
| PoseFormerV2 [48] (f=27) CVPR'23     |              | 243 | -    | -      | -    | -      | -        | -      | -            | -            | -            | -      | -       | -            | -      | -                | -            | 45.   |
| STCFormer [41] CVPR'23               | √            | 243 | 39.6 | 41.6   | 37.4 | 38.8   | 43.1     | 51.1   | 39.1         | 39.7         | 51.4         | 57.4   | 41.8    | 38.5         | 40.7   | 27.1             | 28.6         | 41.   |
| STCFormer-L [41] CVPR'23             | √            | 243 | 38.4 | 41.2   | 36.8 | 38.0   | 42.7     | 50.5   | 38.7         | 38.2         | 52.5         | 56.8   | 41.8    | 38.4         | 40.2   | 26.2             | 27.7         | 40.   |
| UPS [11] CVPR'23                     |              | 243 | 37.5 | 39.2   | 36.9 | 40.6   | 39.3     | 46.8   | 39.0         | 41.7         | 50.6         | 63.5   | 40.4    | 37.8         | 44.2   | 26.7             | 29.1         | 40.   |
| D3DP [38] (H=20,K=10, J-Agg) ICCV'23 | $\checkmark$ | 243 | 37.3 | 39.5   | 35.6 | 37.8   | 41.3     | 48.2   | 39.1         | 37.6         | 49.9         | 52.8   | 41.2    | 39.2         | 39.4   | 27.2             | 27.1         | 39.   |
| GLA-GCN [44] ICCV'23                 |              | 243 | 41.3 | 44.3   | 40.8 | 41.8   | 45.9     | 54.1   | 42.1         | 41.5         | 57.8         | 62.9   | 45.0    | 42.8         | 45.9   | 29.4             | 29.9         | 44.   |
| MotionBERT [51] (scratch) ICCV'23    | $\checkmark$ | 243 | 36.3 | 38.7   | 38.6 | 33.6   | 42.1     | 50.1   | 36.2         | 35.7         | 50.1         | 56.6   | 41.3    | 37.4         | 37.7   | 25.6             | 26.5         | 39    |
| MotionBERT [51] (finetune) ICCV'23   | √            | 243 | 36.1 | 37.5   | 35.8 | 32.1   | 40.3     | 46.3   | 36.1         | 35.3         | 46.9         | 53.9   | 39.5    | 36.3         | 35.8   | 25.1             | 25.3         | 37    |
| HoT [23] CVPR'24                     |              | 243 | -    | -      | -    | -      | -        | -      | -            | -            | -            | -      | -       | -            | -      | -                | -            | 39    |
| KTPFormer [33] CVPR'24               | $\checkmark$ | 243 | 37.3 | 39.2   | 35.9 | 37.6   | 42.5     | 48.2   | 38.6         | 39.0         | 51.4         | 55.9   | 41.6    | 39.0         | 40.0   | 27.0             | 27.4         | 40    |
| MotionAGFormer-B [27] WACV'24        | $\checkmark$ | 243 | 36.4 | 38.4   | 36.8 | 32.9   | 40.9     | 48.5   | 36.6         | 34.6         | 51.7         | 52.8   | 41.0    | 36.4         | 36.5   | 26.7             | 27.0         | 38    |
| APP w. MotionBERT (finetune) (Ours)  | √            | 243 | 34.6 | 36.7   | 37.4 | 31.5   | 38.9     | 45.7   | 35.9         | 33.3         | 48.4         | 53.4   | 39.1    | 36.3         | 34.4   | 24.2             | 24.8         | 37.   |
| APP w. MotionAGFormer-B (Ours)       | $\checkmark$ | 243 | 33.7 | 37.2   | 35.0 | 32.2   | 37.2     | 40.7   | 34.6         | 34.3         | 45.1         | 51.4   | 37.0    | 34.6         | 36.8   | 27.0             | 26.8         | 36.   |
| P-MPJPE                              |              | T   | Dir. | Disc   | Eat  | Greet  | Phone    | Photo  | Pose         | Purch.       | Sit          | SitD.  | Smoke   | Wait         | WalkD. | Walk             | WalkT.       | Av    |
| MHFormer [22] CVPR'22                |              | 351 | 31.5 | 34.9   | 32.8 | 33.6   | 35.3     | 39.6   | 32.0         | 32.2         | 43.5         | 48.7   | 36.4    | 32.6         | 34.3   | 23.9             | 25.1         | 34.   |
| StridedTrans [21] TMM'22             |              | 351 | -    | -      | -    | -      | -        | -      | -            | -            | -            | -      | -       | -            | -      | -                | -            | -     |
| MixSTE [45] CVPR'22                  | $\checkmark$ | 243 | 30.8 | 33.1   | 30.3 | 31.8   | 33.1     | 39.1   | 31.1         | 30.5         | 42.5         | 44.5   | 34.0    | 30.8         | 32.7   | 22.1             | 22.9         | 32    |
| MixSTE [45] CVPR'22                  | $\checkmark$ | 243 | 28.0 | 30.9   | 28.6 | 30.7   | 30.4     | 34.6   | 28.6         | 28.1         | 37.1         | 47.3   | 30.5    | 29.7         | 30.5   | 21.6             | 20.0         | 30.   |
| P-STMO [37] ECCV'22                  |              | 243 | 31.3 | 35.2   | 32.9 | 33.9   | 35.4     | 39.3   | 32.5         | 31.5         | 44.6         | 48.2   | 36.3    | 32.9         | 34.4   | 23.8             | 23.9         | 34    |
| CA-PF-HRNet-48 [47] NeurIPS'23       |              | 1   | -    | -      | -    | -      | -        | -      | -            | -            | -            | -      | -       | -            | -      | -                | -            | 32.   |
| HDFormer [4] IICAI'23                | 1            | 96  | 29.6 | 33.8   | 31.7 | 31.3   | 33.7     | 37.7   | 30.6         | 31.0         | 41.4         | 47.6   | 35.0    | 30.9         | 33.7   | 25.3             | 23.6         | 33    |
| HDFormer [4] IICAI'23                | 1            | 96  | 27.9 | 32.8   | 29.7 | 30.6   | 32.5     | 35.0   | 28.9         | 29.2         | 38.3         | 50.0   | 32.9    | 30.1         | 31.8   | 23.6             | 22.8         | 31    |
| PoseFormerV2 [48] (f=27) CVPR'23     |              | 243 | -    | -      | -    | -      | -        | -      | -            | -            | -            | -      | -       | -            | -      | -                | -            | -     |
| STCFormer [41] CVPR'23               | 1            | 243 | 29.5 | 33.2   | 30.6 | 31.0   | 33.0     | 38.0   | 30.4         | 29.4         | 41.8         | 45.2   | 33.6    | 29.5         | 31.6   | 21.3             | 22.6         | 32    |
| STCFormer-L [41] CVPR'23             | 1            | 243 | 29.3 | 33.0   | 30.7 | 30.6   | 32.7     | 38.2   | 29.7         | 28.8         | 42.2         | 45.0   | 33.3    | 29.4         | 31.5   | 20.9             | 22.3         | 31    |
| UPS [11] CVPR'23                     |              | 243 | 30.3 | 32.2   | 30.8 | 33.1   | 31.1     | 35.2   | 30.3         | 32.1         | 39.4         | 49.6   | 32.9    | 29.2         | 33.9   | 21.6             | 24.5         | 32    |
| D3DP [38] (H=20 K=10 I-Agg) ICCV'23  |              | 243 | 30.6 | 32.4   | 29.2 | 30.9   | 31.9     | 37.4   | 30.2         | 293          | 40.4         | 43.2   | 33.2    | 30.4         | 31.3   | 21.5             | 22.3         | 31    |
| GLA-GCN [44] ICCV'23                 | ľ            | 243 | 32.4 | 35.3   | 32.6 | 34.2   | 35.0     | 42.1   | 32.1         | 31.9         | 45.5         | 49.5   | 36.1    | 32.4         | 35.6   | 23.5             | 24 7         | 34    |
| MotionBERT [51] (scratch) ICCV'23+   | 1            | 243 | 30.8 | 32.8   | 32.0 | 28.7   | 34.3     | 38.9   | 30.1         | 30.0         | 42.5         | 49.7   | 36.0    | 30.8         | 31.7   | 22.0             | 23.0         | 32    |
| MotionBERT [51] (finetune) ICCV'23   |              | 243 |      | - 52.0 | -    | - 20.7 | -        | -      |              | -            |              | -      | -       | -            | -      |                  | - 23.0       | - 52. |
| HoT [23] CVPR'24                     |              | 243 |      | _      |      | _      | _        | _      | _            | _            | _            | _      | _       | _            | _      | _                | _            |       |
| KTPFormer [33] CVPP'24               |              | 243 | 30.1 | 322    | 20.6 | 30.8   | 323      | 373    | 30.0         | 30.2         | 41.0         | 45.3   | 33.6    | 20.0         | 31.4   | 21.5             | 22.6         | 31    |
| MotionAGFormer-B [27] WACV'24        | 1            | 243 | 30.6 | 32.5   | 32.2 | 28.2   | 33.8     | 38.6   | 30.5         | 29.9         | 43.3         | 47.0   | 35.2    | 29.8         | 31.4   | 21.5             | 22.0         | 32    |
| ADD w MotionBEDT (finature) (Orrec)  |              | 242 | 20.0 | 21.2   | 21.0 | 26.0   | 22.9     | 27.0   | 20.8         | 20.2         | 40.7         | 47.0   | 24.1    | 20.0         | 20.0   | 21.0             | 21.7         | 21    |
| APP w Motion AGEormer-B (Ours)       |              | 243 | 28.2 | 30.2   | 285  | 26.5   | 30.1     | 33.2   | 29.6<br>27.2 | 29.2<br>28.1 | 40.7<br>36.4 | 47.9   | 30.8    | 30.0<br>27.7 | 29.7   | 21.0<br>21.4     | 21./<br>21.9 | 20    |
| ALL W. MOUDIAGIOI MEI D (OUIS)       | 1 Y          | 243 | 40.2 | 30.2   | 20.5 | 20.5   | 30.1     | 55.4   | 41.4         | 20.1         | 30.4         | 44.0   | 50.0    | 41.1         | 49.1   | 41. <del>4</del> | 41.7         | 49.   |

#### EXPERIMENTS

#### **Datasets and Evaluation Metrics** 4.1

4.1.1 Human3.6M. Human3.6M [16] is a widely used large-scale 3D human pose estimation dataset. The dataset consists of images captured by four high-resolution cameras placed at the corners of a room to capture the human body from multiple angles. It includes 11 subjects performing 17 action scenarios (e.g., talking on the phone, taking photos), making it a benchmark dataset for 3D human pose estimation. We follow previous studies [27, 51] by using subjects 1, 5, 6, 7, and 8 for model training and subjects 9 and 11 for testing.

4.1.2 MPI-INF-3DHP. MPI-INF-3DHP [28] dataset comprises indoor and outdoor scenes. It contains over 1.3 million frames captured from 14 camera viewpoints, recording eight actions performed by eight subjects. The test set includes seven action classes and three different scenarios: green screen, non-green screen, and outdoor scenes.

4.1.3 Evaluation Metrics. For the Human3.6M dataset, we report the MPJPE and P-MPJPE. Before computing the error, the former

aligns the estimated 3D pose with the ground truth root node (hip). At the same time, the latter requires rigid alignment based on the ground truth, including translation and rotation, before error calculation. For the MPI-INF-3DHP dataset, consistent with previous works [4, 27, 38, 47, 48], we report MPJPE, Percentage of Correct Keypoints (PCK) with a threshold of 150mm, and Area Under Curve (AUC).

## 4.2 Implementation Details

We adopt Transformer [42], MotionAGFormer-B [27], MixSTE [45], MotionBERT [51], and HoT [23] as the foundation models. These models' dimensions dp are 128, 256, 512, and 128, respectively. For a fair comparison, we use pretrained lifting models provided by the authors.

YOLOv3 [35] is used to extract bounding boxes of humans in the images, and then we crop the images to  $192 \times 256$  according to the bounding boxes. However, in [47], the 3D ground truth human pose is directly used to extract bounding boxes, which is not feasible in real-world scenarios.

Table 2: Qualitative comparisons of 3D human pose estimation per action using 2D Ground Truth (GT) human poses on Human3.6M. The best and second-best results are bolded and blue, respectively. *T*: Number of the input frames. Seq2Seq: Estimate 3D human pose for every frame.

| MPJPE (GT)                           | Seq2Seq      | Т   | Dir. | Disc | Eat  | Greet | Phone | Photo | Pose | Purch. | Sit  | SitD.       | Smoke | Wait | WalkD. | Walk | WalkT. | Av  |
|--------------------------------------|--------------|-----|------|------|------|-------|-------|-------|------|--------|------|-------------|-------|------|--------|------|--------|-----|
| MHFormer [22] CVPR'22                |              | 351 | 27.7 | 32.1 | 29.1 | 28.9  | 30.0  | 33.9  | 33.0 | 31.2   | 37.0 | 39.3        | 30.0  | 31.0 | 29.4   | 22.2 | 23.0   | 30. |
| StridedTrans [21] TMM'22             |              | 351 | 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. |
| MixSTE [45] CVPR'22                  | $\checkmark$ | 81  | 25.6 | 27.8 | 24.5 | 25.7  | 24.9  | 29.9  | 28.6 | 27.4   | 29.9 | 29.0        | 26.1  | 25.0 | 25.2   | 18.7 | 19.9   | 25  |
| MixSTE [45] CVPR'22                  | $\checkmark$ | 243 | 21.6 | 22.0 | 20.4 | 21.0  | 20.8  | 24.3  | 24.7 | 21.9   | 26.9 | 24.9        | 21.2  | 21.5 | 20.8   | 14.7 | 15.6   | 21  |
| P-STMO [37] ECCV'22                  |              | 243 | 28.5 | 30.1 | 28.6 | 27.9  | 29.8  | 33.2  | 31.3 | 27.8   | 36.0 | 37.4        | 29.7  | 29.5 | 28.1   | 21.0 | 21.0   | 29  |
| D3DP [38] (H=20,K=10, J-Agg) ICCV'23 | $\checkmark$ | 243 | 19.9 | 19.4 | 19.4 | 19.0  | 19.8  | 22.0  | 21.4 | 19.1   | 24.8 | 23.2        | 19.6  | 18.7 | 18.6   | 14.0 | 14.5   | 19  |
| GLA-GCN [44] ICCV'23                 |              | 243 | 20.1 | 21.2 | 20.0 | 19.6  | 21.5  | 26.7  | 23.3 | 19.8   | 27.0 | 29.4        | 20.8  | 20.1 | 19.2   | 12.8 | 13.8   | 21  |
| STCFormer [41] CVPR'23               | $\checkmark$ | 81  | 25.9 | 25.9 | 22.7 | 24.0  | 24.6  | 27.5  | 27.6 | 23.1   | 30.1 | 31.5        | 25.1  | 24.7 | 23.8   | 18.4 | 19.6   | 25  |
| STCFormer-L [41] CVPR'23             | $\checkmark$ | 243 | 20.8 | 21.8 | 20.0 | 20.6  | 23.4  | 25.0  | 23.6 | 19.3   | 27.8 | 26.1        | 21.6  | 20.6 | 19.5   | 14.3 | 15.1   | 21  |
| HDFormer [4] IJCAI'23                | $\checkmark$ | 96  | -    | -    | -    | -     | -     | -     | -    | -      | -    | -           | -     | -    | -      | -    | -      | 21  |
| MotionBERT [51] (scratch) ICCV'23    | $\checkmark$ | 243 | 16.7 | 19.9 | 17.1 | 16.5  | 17.4  | 18.8  | 19.3 | 20.5   | 24.0 | 22.1        | 18.6  | 16.8 | 16.7   | 10.8 | 11.5   | 17  |
| MotionBERT [51] (finetune) ICCV'23   | $\checkmark$ | 243 | 15.9 | 17.3 | 16.9 | 14.6  | 16.8  | 18.6  | 18.6 | 18.4   | 22.0 | 21.8        | 17.3  | 16.9 | 16.1   | 10.5 | 11.4   | 16  |
| KTPFormer [33] CVPR'24               | $\checkmark$ | 243 | 18.8 | 17.4 | 18.1 | 17.7  | 18.3  | 20.6  | 19.6 | 17.7   | 23.3 | <b>22.0</b> | 18.7  | 17.0 | 16.8   | 12.4 | 13.5   | 18  |
| MotionAGFormer-B [27] WACV'24        | $\checkmark$ | 243 | -    | -    | -    | -     | -     | -     | -    | -      | -    | -           | -     | -    | -      | -    | -      | 19  |
| APP w. MotionAGFormer-B (Ours)       | $\checkmark$ | 243 | 18.2 | 20.6 | 18.4 | 17.9  | 19.5  | 21.3  | 20.7 | 20.6   | 25.2 | 25.7        | 19.3  | 18.2 | 17.4   | 11.3 | 12.1   | 19  |

The lengths of 2D human pose sequences *T* on Human3.6M and MPI-INF-3DHP are set to 243 and 81, respectively, and images are uniformly sampled into sequence lengths  $t = \frac{T}{9}$  for a tradeoff between performance and efficiency. Additionally, *t* must be divisible by *T*. We extract feature maps from images using HRNetw32 [39] with four resolutions. It is chosen as the pose detector,  $C_s = 32^{s+1} \in \{32, 64, 128, 256\}, H_s = \frac{H}{4^{s+1}} \in \{\frac{H}{4}, \frac{H}{8}, \frac{H}{16}, \frac{H}{32}\}, W_s = \frac{W}{4^{s+1}} \in \{\frac{W}{4}, \frac{W}{8}, \frac{W}{16}, \frac{W}{32}\}.$ 

The parameters of our proposed APP module are mainly determined by the number of layers *L*, the number of attention heads *h*, the kernel size *k*, the dimension *d* of the APP module, and the dimension of hidden feature  $F_p$  used in the selected lifting model. Also, we employ DropPath [20] with the probability *p*. For Human3.6M and MPI-INF-3DHP, *L*, *h*, *k*,  $\alpha$ , *p*, and *d* are set to 6, 8, 3, 0.9, 0.2, and 128, respectively. We train the APP module using the Adam optimizer [19]. 2D ground truth human pose is sent to the model for MPI-INF-3DHP. In Eq. (9), the weights of each part of the loss are set to 1.0, 20.0, and 0.5, respectively. All experiments are conducted on one NVIDIA RTX 3090 GPU.

### 4.3 Qualitative Comparison of Human3.6M

We conducted comparative evaluations with recent state-of-the-art (SOTA) models on the Human3.6M dataset. The results, summarized in Table 1, underscore the superiority of our proposed APP module when integrated with MotionAGFormer-B, attaining SOTA performance. Specifically, in comparison with lifting models MixSTE [45] and HDFormer [4], which also utilize HRNet-detected 2D human pose as input, our approach achieves notable improvements of 3.6mm (from 39.8mm to 36.2mm) and 4.1mm (from 40.3mm to 36.2mm) in MPJPE, and 1.1mm (from 30.6mm to 29.5mm) and 2.2mm (from 31.7mm to 29.5mm) in terms of P-MPJPE, respectively. What's more, our proposed method surpasses CA-PF-HRNet-48 [47] 3.6mm (from 39.8mm to 36.2mm) in MPJPE, and 3.2mm (from 32.7mm to 29.5mm).

Although our method exhibits a slightly higher MPJPE of 1.3mm
 compared to MotionBERT [51] when ground truth 2D human pose
 is employed as input, it surpasses MotionAGFormer-B by 0.3mm.

These findings collectively attest to the efficacy and competitiveness of our proposed approach in advancing the state-of-the-art in monocular 3D human pose estimation.

Table 3: Qualitative comparisons of 3D human pose estimation per action on MPI-INF-3DHP. We report AUC, MPJPE, and PCK. The best and second-best results are bolded and blue, respectively. *T*: Number of the input frames.

| Method                               | Т   | AUC↑ | MPJPE↓ | PCK↑        |
|--------------------------------------|-----|------|--------|-------------|
| D3DP [38] (H=20,K=10, J-Agg) ICCV'23 | 243 | 78.2 | 29.7   | 97.7        |
| HDFormer [4] IJCAI'23                | 32  | 64.0 | 51.5   | 96.8        |
| HDFormer [4] IJCAI'23                | 96  | 72.9 | 37.2   | 98.7        |
| PoseFormerV2 [48] CVPR'23            | 81  | 78.8 | 27.8   | 97.9        |
| CA-PF-HRNet-32 [47] NeurIPS'23       | 1   | 75.4 | 32.7   | 98.0        |
| CA-PF-HRNet-48 [47] NeurIPS'23       | 1   | 76.3 | 31.4   | 98.2        |
| MotionAGFormer-B [27] WACV'24        | 81  | 84.2 | 18.2   | <b>98.3</b> |
| KTPFormer [33] CVPR'24               | 27  | 84.4 | 19.2   | 98.9        |
| KTPFormer [33] CVPR'24               | 81  | 85.9 | 16.7   | 98.9        |
| APP w. MotionAGFormer-B (Ours)       | 81  | 89.5 | 12.7   | 98.9        |

# 4.4 Qualitative Comparison of MPI-INF-3DHP

In evaluating the MPI-INF-3DHP dataset, we aligned the sequence length *T* to 81 frames to maintain consistency with the experimental setup of MotionAGFormer-B [27]. As detailed in Table 2, our method demonstrates incremental enhancements across various evaluation metrics. Specifically, we observe marginal improvements in the AUC and MPJPE upon KTPFormer [33], with gains of 3.6 (from 85.9 to 89.5) and 4mm (from 16.7mm to 12.7mm), respectively. Moreover, the PCK of our method is the same as KTPFormer [33] because the threshold of PCK is set to 150mm. When comparing our method with the baseline model MotionAGFormer-B[27], there are huge improvements for AUC (from 84.2 to 89.5) and MPJPE (from 18.2mm to 12.7mm). These findings underscore the efficacy and robustness of our proposed method in enhancing monocular 3D human pose estimation performance on the MPI-INF-3DHP dataset.

# 4.5 Ablation Study

We conduct extensive experiments to substantiate the effectiveness of our proposed APP module. As shown in Table 4, our experimental results manifest a consistent enhancement in the performance of the lifting models across diverse configurations of 2D pose detectors and multiple-frame lifting models, affirming the efficacy of the APP module. Transformer [42] is a simple spatial-temporal vanilla transformer with 12 layers.

Table 4: In comparison with baselines on Human3.6M, we choose different lifting models and 2D poses.

| 709 | Method                              | Т   | 2D Pose | Param  | MACs    | MACs/frame | MPJPE↓ | P-MPJPE↓ |
|-----|-------------------------------------|-----|---------|--------|---------|------------|--------|----------|
| 710 | Transformer [42] NeurIPS'17         | 243 | SH      | 9.62M  | 46.16G  | 189.95M    | 39.2   | 32.8     |
|     | APP w. Transformer (Ours)           | 243 | SH      | 13.20M | 70.71G  | 290.97M    | 39.0   | 32.7     |
| 711 | MixSTE [45] CVPR'22                 | 243 | CPN     | 33.78M | 147.60G | 607.38M    | 40.9   | 32.6     |
| 712 | APP w. MixSTE (Ours)                | 243 | CPN     | 39.54M | 174.85G | 719.54M    | 39.9   | 32.1     |
| /12 | MotionBERT [51] (finetune) ICCV'23  | 243 | SH      | 42.47M | 185.60G | 763.77M    | 37.5   | -        |
| 713 | APP w. MotionBERT (finetune) (Ours) | 243 | SH      | 48.23M | 211.77G | 871.48M    | 37.0   | 31.6     |
| 714 | MotionAGFormer-B [27] WACV'24       | 243 | CPN     | 11.72M | 64.78G  | 266.60M    | 39.2   | 31.5     |
| /14 | APP w. MotionAGFormer-B (Ours)      | 243 | CPN     | 15.30M | 90.14G  | 370.97M    | 38.9   | 31.3     |
| 715 | MotionAGFormer-B [27] WACV'24       | 243 | SH      | 11.72M | 64.78G  | 266.60M    | 38.4   | 32.6     |
|     | APP w. MotionAGFormer-B (Ours)      | 243 | SH      | 15.30M | 90.14G  | 370.97M    | 38.4   | 32.3     |
| 716 | MotionAGFormer-B [27] WACV'24       | 243 | HRNet   | 11.72M | 64.78G  | 266.60M    | 36.6   | 29.5     |
|     | APP w. MotionAGFormer-B (Ours)      | 243 | HRNet   | 15.30M | 90.14G  | 370.97M    | 36.2   | 29.5     |
| 717 | MotionAGFormer-B [27] WACV'24       | 243 | GT      | 11.72M | 64.78G  | 266.60M    | 19.4   | -        |
|     | APP w. MotionAGFormer (Ours)        | 243 | GT      | 15.30M | 90.14G  | 370.97M    | 19.1   | 18.4     |
| 718 | HoT [23] CVPR'24                    | 243 | SH      | 16.35M | 33.48G  | 137.79M    | 39.8   | 33.5     |
| 719 | APP w. HoT (Ours)                   | 243 | SH      | 20.66M | 50.80G  | 209.06M    | 39.7   | 33.4     |
|     |                                     |     |         |        |         |            |        |          |

Notably, the APP module contributes to performance improvements while maintaining small parameters, accounting for approximately 37% (from 9.62M to 13.2M), 17% (from 33.78M to 39.54M), 14% (from 42.47M to 48.23M), 30% (from 11.72M to 15.3M), 26% (from 16.35M to 20.66M) parameters of Transformer [42], MixSTE [45], MotionBERT [51], MotionAGFormer-B [27], HoT [23], respectively. When it comes to MACs, the gains are 34% (from 46.16G to 70.71G), 16% (from 147.6G to 174.85G), 12% (from 185.6G to 211.77G), 28% (from 64.78G to 90.14G), 34% (from 33.48G to 50.8G).

All ablation studies are conducted on the Human3.6M dataset, with MotionAGFormer-B [27] chosen as the baseline model and 2D human body pose detected using HRNet [39]. Initially, we conducted experiments on parameter selections for the APP module. The results of Table 5 presents indicate that our proposed method achieves the best performance when L, k,  $\alpha$ , p, and d are selected as 6, 3, 0.9, 0.2, and 128, respectively.

Table 5: Ablation study of the parameters on Human3.6M. L: The number of layers of our proposed APP module. p: Probability of the DropPath. d: The Number of the dimensions. We report MPJPE and P-MPJPE. The best and second-best results are bolded and blue, respectively.

| L | k | α   | p   | d   | Params | MPJPE↓ | P-MPJPE↓ |
|---|---|-----|-----|-----|--------|--------|----------|
| 6 | 3 | 0.9 | 0.2 | 128 | 3.58M  | 36.2   | 29.5     |
| 4 | 3 | 0.9 | 0.2 | 256 | 7.50M  | 36.5   | 29.7     |
| 4 | 3 | 0.9 | 0.4 | 128 | 2.76M  | 36.3   | 29.6     |
| 6 | 5 | 0.9 | 0.2 | 128 | 5.15M  | 36.5   | 29.7     |
| 6 | 3 | 0.6 | 0.2 | 128 | 3.58M  | 36.6   | 30.0     |
| 6 | 3 | 0.9 | 0   | 128 | 3.58M  | 36.5   | 29.9     |

Then, we conducted ablation experiments on the submodules of our proposed APP module. In Table 6, the first row presents partial

components of three submodules: MHCA, PAOG, and PAS. The table shows that MHCA has the most significant impact on the APP model's performance. However, the other three components also enhance the model's performance to varying degrees.

Table 6: Ablation study of the components on Human3.6M. MHCA: Use MHCA. Weights: Apply weights to the learned offsets. PoseAware: Choose the detected 2D pose as the pivot; otherwise, the center of the image will be selected. Conv: Convolution is performed on the sampled features; otherwise, the mean operation is carried out directly. We report MPJPE and P-MPJPE. The best and second-best results are bolded and blue, respectively.

| MHCA         | Weights      | PoseAware    | Conv         | Params | MPJPE↓ | P-MPJPE↓ |
|--------------|--------------|--------------|--------------|--------|--------|----------|
| $\checkmark$ |              |              |              | 2.69M  | 36.6   | 29.8     |
|              | $\checkmark$ |              |              | 2.30M  | 62.0   | 43.1     |
|              |              | $\checkmark$ |              | 2.30M  | 68.0   | 44.7     |
|              |              |              | $\checkmark$ | 3.18M  | 67.2   | 45.7     |
| $\checkmark$ | $\checkmark$ |              |              | 2.69M  | 36.2   | 29.6     |
| $\checkmark$ |              | $\checkmark$ |              | 2.69M  | 36.7   | 29.9     |
| $\checkmark$ |              |              | $\checkmark$ | 3.58M  | 36.4   | 29.6     |
|              | $\checkmark$ | $\checkmark$ |              | 2.30M  | 63.6   | 43.2     |
|              | $\checkmark$ |              | $\checkmark$ | 3.18M  | 59.1   | 40.0     |
|              |              | $\checkmark$ | $\checkmark$ | 3.18M  | 71.0   | 45.1     |
| $\checkmark$ | $\checkmark$ | $\checkmark$ |              | 2.69M  | 36.5   | 30.0     |
| $\checkmark$ | $\checkmark$ |              | $\checkmark$ | 3.58M  | 36.3   | 29.9     |
| $\checkmark$ |              | $\checkmark$ | $\checkmark$ | 3.58M  | 36.4   | 29.6     |
|              | $\checkmark$ | $\checkmark$ | $\checkmark$ | 3.18M  | 57.0   | 39.1     |
| $\checkmark$ | $\checkmark$ | $\checkmark$ | $\checkmark$ | 3.58M  | 36.2   | 29.5     |

Finally, we conducted ablaton experiments on image frames t setting. Although the MACs (24.83G) of the APP module are not affected by the number of frames, yet image frames t need more computation in training stage because we need multi-level featuremaps extracted by 2D pose detector. Furthermore, when  $t \ge 81$ , the training time becomes very long. Therefore,  $t = 27 = \frac{T}{9}$ is tradeoff between accuracy and efficiency. In inference stage, this problem no longer appears, owing to that 2D pose detector is used to estimate 2D pose for every frame.

Table 7: Ablation study of the image frames t on Human3.6M. We report MPJPE and P-MPJPE. The best and second-best results are bolded and blue, respectively.

| t  | Params | MPJPE↓ | P-MPJPE↓ |
|----|--------|--------|----------|
| 1  | 3.58M  | 36.5   | 30.0     |
| 3  | 3.58M  | 36.5   | 29.8     |
| 9  | 3.58M  | 36.3   | 29.6     |
| 27 | 3.58M  | 36.2   | 29.6     |

### 4.6 Visualization

To validate the generalization ability of the model, we selected the MotionBERT [51] (finetune version) and MotionAGFormer-B [27] as foundation models, which were trained on 2D poses extracted by SH [30] on Human3.6M dataset, for comparison with our proposed APP module. We chose two video segments from the internet to

visually analyze the APP module. For the visual analysis in Figures 3 and 4, MotionAGFormer-B was selected as the foundation model.

4.6.1 Analysis of MHCA Module. We visualized the attention maps of the MHCA submodule at each layer of our proposed APP module. The values of each attention map were normalized to [0, 1]. Each attention map extracted at every layer is a  $17 \times 17$  matrix. Figure 3 shows that although each layer has different focuses, the APP module can capture the relationships between key points.



Figure 3: Visualization of the attention maps of each layer of MHCA in our proposed APP module.



Figure 4: Visualization of the initialized/learned sampling points learned by each layer, where the initialized/learned sampling points are marked in orange and magenta, respectively. Each sampling point is represented by a cross (x).

*4.6.2 Analysis of PAOG Module.* Figure 4 illustrates the sampling points learned by our proposed PAOG submodule. The first row in the figure is appended for observation, while the second to fifth

rows represent feature maps extracted at different resolutions. Each layer of the PAOG submodule can adaptively learn sampling points based on the distribution of feature maps at different resolutions.

*4.6.3* Analysis of In-the-wild Videos. We selected six frames from two videos to compare the predicted 3D human poses by two lifting models, MotionBERT [51] and MotionAGFormer-B [27], with those predicted by our proposed model. The differences have been circled in orange.



Figure 5: Qualitative comparison with SOTA methods on in-the-wild images. We select MotionBERT[51] and MotionAGFormer-B [27] as foundation models.

# 5 CONCLUSION

Presenting APP, short for Adaptive Pose Pooling, a flexible module crafted to mesh with various multi-frame lifting models. This versatile module utilizes the hidden features extracted by lifting models with the feature maps obtained from the 2D pose detector. APP excels in simultaneously extracting spatial and temporal features through this collaborative approach module, and our visualization demonstrates this. Extensive experimentation highlights the reliability of our module, showcasing its ability to enhance the performance metrics of lifting models while consistently achieving state-of-the-art results across diverse scenarios. Notably, APP module presents remarkable adaptability, maintaining its effectiveness even when paired with different 2D poses. A noteworthy feature is its capacity to enhance performance without requiring structural changes to the existing lifting model architecture. APP: Adaptive Pose Pooling for 3D Human Pose Estimation from Videos

ACM MM, 2024, Melbourne, Australia

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

1043

1044

### 929 **REFERENCES**

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

- Peter Bauer, Arij Bouazizi, Ulrich Kressel, and Fabian B Flohr. 2023. Weakly supervised multi-modal 3d human body pose estimation for autonomous driving. In 2023 IEEE Intelligent Vehicles Symposium (IV). IEEE, 1–7.
- [2] Tobias Baumgartner and Stefanie Klatt. 2023. Monocular 3D human pose estimation for sports broadcasts using partial sports field registration. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. 5108–5117.
- [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 Proceedings of the IEEE/CVF International Conference on Computer Vision. 2272–2281.
- [4] Hanyuan Chen, Jun-Yan He, Wangmeng Xiang, Zhi-Qi Cheng, Wei Liu, Hanbing Liu, Bin Luo, Yifeng Geng, and Xuansong Xie. 2023. HDFormer: high-order directed transformer for 3D human pose estimation. In Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence. 581–589.
- [5] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun. 2018. Cascaded pyramid network for multi-person pose estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 7103–7112.
- [6] Sungho Chun, Sungbum Park, and Ju Yong Chang. 2023. Learnable human mesh triangulation for 3D human pose and shape estimation. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. 2850–2859.
- [7] Hai Ci, Chunyu Wang, Xiaoxuan Ma, and Yizhou Wang. 2019. Optimizing network structure for 3d human pose estimation. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 2262–2271.
- [8] Huseyin Coskun, David Joseph Tan, Sailesh Conjeti, Nassir Navab, and Federico Tombari. 2018. Human motion analysis with deep metric learning. In Proceedings of the European conference on computer vision (ECCV). 667–683.
- [9] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. 2017. Deformable convolutional networks. In Proceedings of the IEEE international conference on computer vision. 764–773.
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *ICLR*.
- [11] Lin Geng Foo, Tianjiao Li, Hossein Rahmani, Qiuhong Ke, and Jun Liu. 2023. Unified pose sequence modeling. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 13019–13030.
- [12] Jia Gong, Zhipeng Fan, Qiuhong Ke, Hossein Rahmani, and Jun Liu. 2022. Meta agent teaming active learning for pose estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 11079–11089.
- [13] Renshu Gu, Gaoang Wang, and Jenq-Neng Hwang. 2019. Efficient multi-person hierarchical 3d pose estimation for autonomous driving. In 2019 IEEE conference on multimedia information processing and retrieval (MIPR). IEEE, 163–168.
- [14] Yihui He, Rui Yan, Katerina Fragkiadaki, and Shoou-I Yu. 2020. Epipolar transformers. In Proceedings of the ieee/cvf conference on computer vision and pattern recognition. 7779–7788.
- [15] Dan Hendrycks and Kevin Gimpel. 2016. Bridging Nonlinearities and Stochastic Regularizers with Gaussian Error Linear Units. (2016).
- [16] 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. *IEEE TPAMI* 36, 7 (2013), 1325–1339.
- [17] Karim Iskakov, Egor Burkov, Victor Lempitsky, and Yury Malkov. 2019. Learnable triangulation of human pose. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 7718–7727.
- [18] Angjoo Kanazawa, Michael J Black, David W Jacobs, and Jitendra Malik. 2018. Endto-end recovery of human shape and pose. In Proceedings of the IEEE conference on computer vision and pattern recognition. 7122–7131.
- [19] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014).
- [20] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. 2016. FractalNet: Ultra-Deep Neural Networks without Residuals. In International Conference on Learning Representations.
- [21] Wenhao Li, Hong Liu, Runwei Ding, Mengyuan Liu, Pichao Wang, and Wenming Yang. 2022. Exploiting Temporal Contexts with Strided Transformer for 3D Human Pose Estimation. *IEEE TMM* 25 (2022), 1282–1293.
- [22] Wenhao Li, Hong Liu, Hao Tang, Pichao Wang, and Luc Van Gool. 2022. MH-Former: Multi-Hypothesis Transformer for 3D Human Pose Estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 13147–13156.
- [23] Wenhao Li, Mengyuan Liu, Hong Liu, Pichao Wang, Jialun Cai, and Nicu Sebe. 2024. Hourglass Tokenizer for Efficient Transformer-Based 3D Human Pose Estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- [24] Huei-Yung Lin and Ting-Wen Chen. 2010. Augmented reality with human body interaction based on monocular 3d pose estimation. In *International Conference*

- on Advanced Concepts for Intelligent Vision Systems. Springer, 321–331.
- [25] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision* and Pattern Recognition. 5064–5073.
- [26] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little. 2017. A simple yet effective baseline for 3D human pose estimation. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 2640–2649.
- [27] Soroush Mehraban, Vida Adeli, and Babak Taati. 2024. MotionAGFormer: Enhancing 3D Human Pose Estimation with a Transformer-GCNFormer Network. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. 6920–6930.
- [28] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt. 2017. Monocular 3D human pose estimation in the wild using improved CNN supervision. In 3DV. 506–516.
- [29] Dushyant Mehta, Srinath Sridhar, Oleksandr Sotnychenko, Helge Rhodin, Mohammad Shafiei, Hans-Peter Seidel, Weipeng Xu, Dan Casas, and Christian Theobalt. 2017. Vnect: Real-time 3d human pose estimation with a single rgb camera. Acm transactions on graphics (tog) 36, 4 (2017), 1–14.
- [30] Alejandro Newell, Kaiyu Yang, and Jia Deng. 2016. Stacked hourglass networks for human pose estimation. In *Proceedings of the European conference on computer* vision (ECCV). 483–499.
- [31] Georgios Pavlakos, Xiaowei Zhou, Konstantinos G Derpanis, and Kostas Daniilidis. 2017. Coarse-to-fine volumetric prediction for single-image 3D human pose. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 7025–7034.
- [32] Dario Pavllo, Christoph Feichtenhofer, David Grangier, and Michael Auli. 2019. 3D human pose estimation in video with temporal convolutions and semi-supervised training. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 7753–7762.
- [33] Jihua Peng, Yanghong Zhou, and PY Mok. 2024. KTPFormer: Kinematics and Trajectory Prior Knowledge-Enhanced Transformer for 3D Human Pose Estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.
- [34] N Dinesh Reddy, Laurent Guigues, Leonid Pishchulin, Jayan Eledath, and Srinivasa G Narasimhan. 2021. Tessetrack: End-to-end learnable multi-person articulated 3d pose tracking. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 15190–15200.
- [35] Joseph Redmon and Ali Farhadi. 2018. YOLOv3: An Incremental Improvement. arXiv (2018).
- [36] Helge Rhodin, Mathieu Salzmann, and Pascal Fua. 2018. Unsupervised geometryaware representation for 3d human pose estimation. In Proceedings of the European conference on computer vision (ECCV). 750–767.
- [37] Wenkang Shan, Zhenhua Liu, Xinfeng Zhang, Shanshe Wang, Siwei Ma, and Wen Gao. 2022. P-STMO: Pre-Trained Spatial Temporal Many-to-One Model for 3D Human Pose Estimation. In Proceedings of the European conference on computer vision (ECCV).
- [38] Wenkang Shan, Zhenhua Liu, Xinfeng Zhang, Zhao Wang, Kai Han, Shanshe Wang, Siwei Ma, and Wen Gao. 2023. Diffusion-based 3d human pose estimation with multi-hypothesis aggregation. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 14761–14771.
- [39] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. 2019. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE/CVF* conference on computer vision and pattern recognition. 5693–5703.
- [40] Xiao Sun, Bin Xiao, Fangyin Wei, Shuang Liang, and Yichen Wei. 2018. Integral human pose regression. In Proceedings of the European conference on computer vision (ECCV). 529–545.
- [41] Zhenhua Tang, Zhaofan Qiu, Yanbin Hao, Richang Hong, and Ting Yao. 2023. 3D Human Pose Estimation With Spatio-Temporal Criss-Cross Attention. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 4790–4799.
- [42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NeurIPS*. 5998–6008.
- [43] Jingwei Xu, Zhenbo Yu, Bingbing Ni, Jiancheng Yang, Xiaokang Yang, and Wenjun Zhang. 2020. Deep kinematics analysis for monocular 3d human pose estimation. In Proceedings of the IEEE/CVF Conference on computer vision and Pattern recognition. 899–908.
- [44] Bruce XB Yu, Zhi Zhang, Yongxu Liu, Sheng-hua Zhong, Yan Liu, and Chang Wen Chen. 2023. Gla-gcn: Global-local adaptive graph convolutional network for 3d human pose estimation from monocular video. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 8818–8829.
- [45] Jinlu Zhang, Zhigang Tu, Jianyu Yang, Yujin Chen, and Junsong Yuan. 2022. MixSTE: Seq2seq Mixed Spatio-Temporal Encoder for 3D Human Pose Estimation in Video. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 13232–13242.
- [46] Zhe Zhang, Chunyu Wang, Weichao Qiu, Wenhu Qin, and Wenjun Zeng. 2021. Adafuse: Adaptive multiview fusion for accurate human pose estimation in the

|                          | [50] | Kun Zhou, Xiaoguang Han, Nianjuan Jiang, Kui Jia, and Jiangbo Lu. 2019. Hem-  | 1103 |
|--------------------------|------|---|------|
| le 2D Pose               |      | lets pose: Learning part-centric heatmap triplets for accurate 3d human pose  | 1104 |
| ny-sevenin               |      | Vision. 2344–2353.  | 1105 |
| Chen. 2023.              | [51] | Wentao Zhu, Xiaoxuan Ma, Zhaoyang Liu, Libin Liu, Wayne Wu, and Yizhou  | 1106 |
| 3D Human<br>outer Vision |      | Wang, 2023. MotionBER I: A Unified Perspective on Learning Human Motion Rep-<br>resentations. In Proceedings of the IEEE/CVF International Conference on Computer | 1107 |
|                          |      | Vision. 15085–15099.  | 1108 |
| and Zheng-               | [52] | Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2020.   | 1109 |
| Computer                 |      | In International Conference on Learning Representations.  | 1110 |
|                          |      |   | 1111 |
|                          |      |   | 1112 |
|                          |      |   | 1113 |
|                          |      |   | 1114 |
|                          |      |   | 1115 |
|                          |      |   | 1116 |
|                          |      |   | 1117 |
|                          |      |   | 1118 |
|                          |      |   | 1119 |
|                          |      |   | 1120 |
|                          |      |   | 1121 |
|                          |      |   | 1122 |
|                          |      |   | 1123 |
|                          |      |   | 1124 |
|                          |      |   | 1125 |
|                          |      |   | 1126 |
|                          |      |   | 1127 |
|                          |      |   | 1128 |
|                          |      |   | 1129 |
|                          |      |   | 1130 |
|                          |      |   | 1131 |
|                          |      |   | 1132 |
|                          |      |   | 1133 |
|                          |      |   | 1134 |
|                          |      |   | 1135 |
|                          |      |   | 1136 |
|                          |      |   | 1137 |
|                          |      |   | 1138 |
|                          |      |   | 1139 |
|                          |      |   | 1140 |

- wild. International Journal of Computer Vision 129 (2021), 703-718.
- [47] Qitao Zhao, Ce Zheng, Mengyuan Liu, and Chen Chen. 2023. A Single 2D Pose with Context is Worth Hundreds for 3D Human Pose Estimation. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- [48] Qitao Zhao, Ce Zheng, Mengyuan Liu, Pichao Wang, and Chen Chen. 2023
   PoseFormerV2: Exploring Frequency Domain for Efficient and Robust 3D Human Pose Estimation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 8877–8886.
- [49] Ce Zheng, Sijie Zhu, Matias Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding. 2021. 3D Human Pose Estimation with Spatial and Temporal Transformers. In Proceedings of the IEEE/CVF International Conference on Computer Vision. 11656–11665.