# INTERWEAVED GRAPH AND ATTENTION NETWORK FOR 3D HUMAN POSE ESTIMATION

*Ti Wang[1], Hong Liu[1], Runwei Ding[1*], Wenhao Li[1], Yingxuan You[1], Xia Li[2]*

Key Laboratory of Machine Perception, Peking University, Shenzhen Graduate School[1]
Department of Computer Science, ETH Zürich[2]

{tiwang, youyx}@stu.pku.edu.cn, {hongliu, dingrunwei, wenhaoli}@pku.edu.cn, xia.li@inf.ethz.ch

## ABSTRACT

Despite substantial progress in 3D human pose estimation from a single-view image, prior works rarely explore global and local correlations, leading to insufficient learning of human skeleton representations. To address this issue, we propose a novel Interweaved Graph and Attention Network (IGANet) that allows bidirectional communications between graph convolutional networks (GCNs) and attentions. Specifically, we introduce an IGA module, where attentions are provided with local information from GCNs and GCNs are injected with global information from attentions. Additionally, we design a simple yet effective U-shaped multi-layer perceptron (uMLP), which can capture multi-granularity information for body joints. Extensive experiments on two popular benchmark datasets (i.e. Human3.6M and MPI-INF-3DHP) are conducted to evaluate our proposed method. The results show that IGANet achieves state-of-the-art performance on both datasets. Code is available at https://github.com/xiu-cs/IGANet.
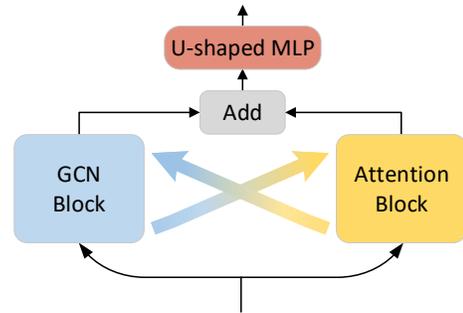
***Index Terms***— 3D Human Pose Estimation, Graph Convolutional Network, Attention

## 1. INTRODUCTION

Monocular 3D human pose estimation aims to recover the 3D positions of body joints from a single-view image. This task plays an important role in many applications, such as human-computer interaction, action recognition, and human mesh reconstruction. Typically, the pipeline can be divided into two parts: 1) estimating the locations of 2D keypoints from a monocular image, and 2) lifting the estimated 2D keypoints to 3D. In this paper, we focus on the problem of 2D-3D pose lifting, where the model input is a 2D pose detected from an image using off-the-shelf 2D pose detectors [1, 2].

With the development of graph convolutional networks (GCNs) in recent years, various GCN-based methods [3, 4] combined with human geometry priors have significantly improved the accuracy of predicting 3D human skeleton positions.
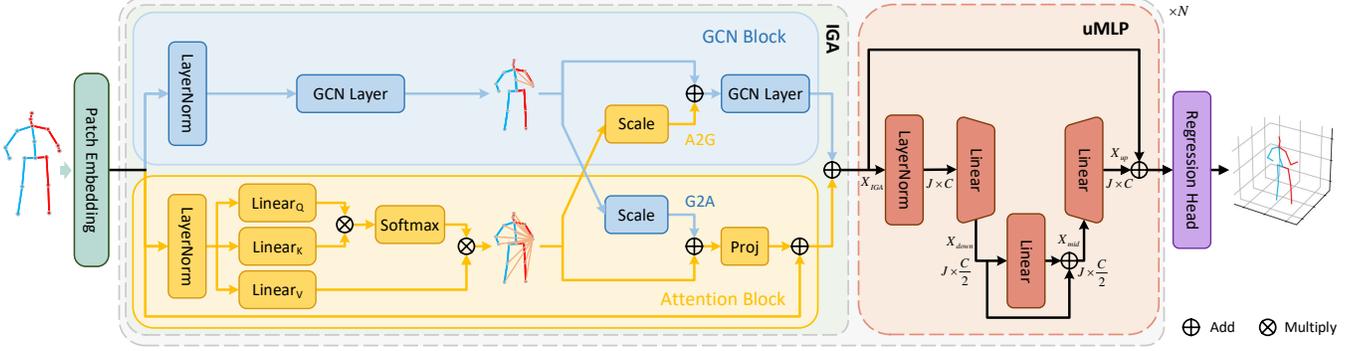
---

**Fig. 1**. A schematic diagram of IGANet. The local relations extracted by the GCN block are interweaved with the global information obtained by the attention block.

For a given graph node, GCN-based methods can exploit semantic information by aggregating the features of its neighbor nodes. However, some distant joints can also provide useful information, such as the relations between the joints of hands and feet when running, which are difficult to capture by the GCN-based methods. Recently, transformers have been widely applied to 3D human pose estimation [5, 6, 7]. Thanks to the attention mechanism of transformers [8], global relations between body joints can be well captured. However, without the topological structure information of the human skeleton, attentions can not fully explore the geometrical connections between body joints. There are few existing works that combine GCNs and attentions in a serial or parallel manner [9, 10, 11]. Nevertheless, these simple combination manners cannot sufficiently utilize the complementarity of GCNs and attentions for effective global and local modeling.

To overcome this issue, a novel Interweaved Graph and Attention Network (IGANet) is proposed to learn bidirectional communications between GCNs and attentions, as shown in Fig. 1. Specifically, we introduce an IGA module based on two guidance strategies: a) *G2A*: the topology information of the human skeleton extracted by the GCN block is injected into the attention block, so that the attention can be guided by the GCN to better learn the structure information of the human body. b) *A2G*: the global information of body joints captured by the attention block is passed to the GCN block, so that the GCN can take the global connections into account while paying attention to its neighbor nodes. Moreover, inspired

**Fig. 2**. Structures of Interweaved Graph and Attention Network (IGANet). The input 2D body joints are embedded into high-dimensional features through the patch embedding. In our IGA module, the captured global and local human skeleton features are interweaved to provide complementary clues. After that, our uMLP module further captures multi-granularity information of body joints. Finally, the predicted results are obtained by the regression head.

by the U-shaped structure [12, 13] that captures features at multiple scale, we design a U-shaped multi-layer perception (uMLP) module with a bottleneck structure along the channel dimension, which can learn multi-granularity features via down-projection and up-projection layers. We find that such a simple design outperforms the original MLP in transformer for learning skeletal representations.

Overall, our contributions can be summarized as follows: (1) We propose a new network named IGANet for single-frame 3D human pose estimation, which allows GCNs and attentions to complement each other. (2) A simple yet effective U-shaped MLP is designed to capture multi-granularity information. (3) The proposed IGANet outperforms existing state-of-the-art methods on two popular benchmark datasets, *i.e.*, Human3.6M [14] and MPI-INF-3DHP [15].

## 2. METHOD

As shown in Fig. 2, IGANet contains two main components: the Interweaved Graph and Attention module (IGA) and the U-shaped MLP module (uMLP). Through patch embedding, the 2D skeleton is mapped into $X \in \mathbb{R}^{J \times C}$, and then sent to IGA and uMLP to extract the features of body joints. Finally, regressions are performed to predict the 3D pose $Y \in \mathbb{R}^{J \times 3}$. IGANet can learn bidirectional communications between GCNs and attentions to mutually improve each other, thus enhancing the capability of modeling human skeleton relationships.

### 2.1. Preliminary

**Graph Convolutional Network (GCN).** An undirected graph can be represented as $G = \{V, E\}$, where $V$ is the set of nodes, and $E$ is the set of edges. The edges can be encoded in an adjacency matrix $A \in \{0, 1\}^{N \times N}$. For the $l^{th}$ layer feature $X_l$, the vanilla graph convolution aggregates neighboring node features, which can be formulated as:

$$X_l = \sigma(W_l X_{l-1} \tilde{A}), \qquad (1)$$

where $W_l \in \mathbb{R}^{C' \times C}$ is the layer-specific trainable weight matrix, $\tilde{A} = A + I_N$ is the adjacency matrix of the graph with added self-connections, and $I_N$ is the identity matrix. The stacking of multiple graph convolutional layers aggregates neighboring nodes to obtain enhanced feature representations.

**Attention.** The input tokens $X \in \mathbb{R}^{J \times C}$ are first projected to queries $Q \in \mathbb{R}^{J \times d}$, keys $K \in \mathbb{R}^{J \times d}$, and values $V \in \mathbb{R}^{J \times d}$, and then $Q, K, V$ are fed to a scaled dot-product attention [8]:

$$\text{Attention}(Q, K, V) = \text{Softmax}(QK^T / \sqrt{d})V, \qquad (2)$$

where $d$ is the dimension of $Q, K, V$. Multi-head self-attention (MSA) [8] splits $Q, K, V$ into multiple heads, each of which applies scaled dot-product attention in parallel. This enables the model to efficiently utilizes information from various representation subspaces with different locations.

### 2.2. Interweaved Graph and Attention

The attention block can capture long-range relationships between nodes in the graph, but it is difficult to model the inherent patterns of the human body structure, such as the left-right symmetry of the human body. The GCN focuses on the local information of human body joints but is limited in capturing global dependencies. To relieve these limitations, we propose an IGA module, which contains a GCN block and an attention block with two guidance strategies (*i.e.*, Graph2Attention (G2A) and Attention2Graph (A2G)) to interchange complementary information with each other.

**Graph2Attention (G2A).** To guide the attention to learn the topology priors of the human skeleton, we inject the skeletal information $f_{graph}$ captured by the first layer of the GCN block into the attention block. This can be formulated as:

$$X_{G2A} = \text{Softmax}(QK^T / \sqrt{d})V + s_{G2A} \cdot f_{graph}, \qquad (3)$$

where $s_{G2A}$ is the scale factor for $f_{graph}$. With the guidance of skeleton information from the GCN block, attention can fully capture the relations between human body joints.

**Table 1**. Quantitative comparisons on Human3.6M under MPJPE. The 2D pose detected by cascaded pyramid network (CPN) [1] is used as input. We use § to highlight methods that use refinement module [4, 13]. The top two best results for each action are highlighted in bold and underlined, respectively.

| Method | Dire. | Disc. | Eat | Greet | Phone | Photo | Pose | Purch. | Sit | SitD. | Smoke | Wait | WalkD. | Walk | WalkT. | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Martinez *et al.* [16] | 51.8 | 56.2 | 58.1 | 59.0 | 69.5 | 78.4 | 55.2 | 58.1 | 74.0 | 94.6 | 62.3 | 59.1 | 65.1 | 49.5 | 52.4 | 62.9 |
| Zhao *et al.* [3] | 47.3 | 60.7 | 51.4 | 60.5 | 61.1 | **49.9** | 47.3 | 68.1 | 86.2 | **55.0** | 67.8 | 61.0 | **42.1** | 60.6 | 45.3 | 57.6 |
| Ci *et al.* [17] | 46.8 | 52.3 | <u>44.7</u> | 50.4 | 52.9 | 68.9 | 49.6 | 46.4 | 60.2 | 78.9 | 51.2 | 50.0 | 54.8 | 40.4 | 43.3 | 52.7 |
| Xu *et al.* [18] | 45.2 | 49.9 | 47.5 | 50.9 | 54.9 | 66.1 | 48.5 | 46.3 | 59.7 | 71.5 | 51.4 | 48.6 | 53.9 | 39.9 | 44.1 | 51.9 |
| Zhao *et al.* [19] | 45.2 | 50.8 | 48.0 | 50.0 | 54.9 | 65.0 | 48.2 | 47.1 | 60.2 | 70.0 | 51.6 | 48.7 | 54.1 | 39.7 | 43.1 | 51.8 |
| Cai *et al.* [13]§ | 46.5 | 48.8 | 47.6 | 50.9 | 52.9 | 61.3 | 48.3 | 45.8 | 59.2 | 64.4 | 51.2 | 48.4 | 53.5 | 39.2 | 41.2 | 50.6 |
| Zeng *et al.* [20] | 44.5 | 48.2 | 47.1 | 47.8 | 51.2 | 56.8 | 50.1 | 45.6 | 59.9 | 66.4 | 52.1 | <u>45.3</u> | 54.2 | 39.1 | 40.3 | 49.9 |
| Zou *et al.* [4]§ | 45.4 | 49.2 | 45.7 | 49.4 | 50.4 | 58.2 | 47.9 | 46.0 | 57.5 | 63.0 | 49.7 | 46.6 | 52.2 | 38.9 | 40.8 | 49.4 |
| IGANet (Ours) | <u>42.9</u> | <u>47.9</u> | **44.9** | <u>47.7</u> | <u>49.8</u> | 58.4 | <u>46.0</u> | <u>44.7</u> | <u>56.8</u> | 61.4 | <u>49.3</u> | 46.1 | 52.0 | <u>37.6</u> | <u>39.8</u> | <u>48.3</u> |
| IGANet (Ours)§ | **42.6** | **47.3** | 45.4 | **47.6** | **49.5** | <u>56.2</u> | **45.9** | **44.1** | **56.3** | <u>59.2</u> | **48.6** | **45.1** | <u>50.3</u> | **37.3** | **39.6** | **47.7** |

**Table 2**. Quantitative comparisons on Human3.6M under MPJPE. Ground truth 2D pose is used as input.

| Method | MPJPE ↓ |
|---|---|
| Martinez *et al.* [16] | 45.5 |
| Zhao *et al.* [3] | 43.8 |
| Cai *et al.* [13] | 38.1 |
| Gong *et al.* [21] | 38.2 |
| Zou *et al.* [4] | 37.4 |
| Zeng *et al.* [20] | 36.4 |
| Xu *et al.* [18] | 35.8 |
| IGANet (Ours) | **32.7 (↑ 8.7%)** |

**Table 3**. Quantitative comparisons with state-of-the-art methods on MPI-INF-3DHP.

| Method | Outdoor | All (PCK) ↑ | All (AUC) ↑ |
|---|---|---|---|
| Martinez *et al.* [16] | 31.2 | 42.5 | 17.0 |
| Ci *et al.* [17] | 77.3 | 74.0 | 36.7 |
| Li *et al.* [22] | 66.6 | 67.9 | - |
| Zeng *et al.* [20] | 80.3 | 77.6 | 43.8 |
| Xu *et al.* [18] | 75.2 | 80.1 | 45.8 |
| Liu *et al.* [23] | 80.1 | 79.3 | 47.6 |
| Zou *et al.* [4] | 85.7 | 86.1 | 53.7 |
| IGANet (Ours) | **86.4** | **86.1** | **54.2** |

**Attention2Graph (A2G).** Similarly, we feed back the global relations of human skeleton $f_{global}$ captured by the attention block into the GCN block, which enables the GCN to have a better awareness of the global relations of human body joints. This operation can be defined as:

$$X_{A2G} = G_1 + s_{A2G} \cdot f_{global},\qquad(4)$$

where $G_1$ is the output of the first GCN layer in GCN block, $s_{A2G}$ denotes the scale factor for $f_{global}$. In this way, the global information can be perceived by the GCN block.

With the interweaved complementary information, the modeling abilities of the GCN block and the attention block are both enhanced. Finally, the outputs from these two blocks are added together, which can be expressed as:

$$X_{IGA} = G_2(X_{A2G}) + \text{Proj}(X_{G2A}),\qquad(5)$$

where $\text{Proj}(\cdot)$ is the projection head containing a linear layer.

## 2.3. U-shaped MLP

Unlike the original MLP of transformer that uses an inverted bottleneck structure, we design a simple yet effective U-shaped MLP (uMLP) that adopts a bottleneck structure along the channel dimension. The input $X_{IGA}$ is first fed to a down-projection layer, followed by a middle layer that keeps the same dimension, and finally to an up-projection layer. To preserve the original information, shortcuts are used between layers of the same dimension. This above is formulated as:

$$\begin{aligned} X_{down} &= \text{MLP}_{down}(\text{LN}(X_{IGA})), \\ X_{mid} &= \text{MLP}_{mid}(X_{down}) + X_{down}, \\ X_{up} &= \text{MLP}_{up}(X_{mid}) + X_{IGA}, \end{aligned}\qquad(6)$$

where $\text{MLP}(\cdot)$ consists of a linear layer and a GELU activation function, $\text{LN}(\cdot)$ denotes the Layer Normalization.

## 3. EXPERIMENTS

### 3.1. Datasets and Evaluation Metrics

**Human3.6M.** Human3.6M [14] is the most widely used indoor dataset for 3D human pose estimation. Following previous works [4, 16, 18], our model is trained on 5 subjects (S1, S5, S6, S7, S8) and tested on 2 subjects (S9 and S11). We adopt the most commonly used evaluation metric MPJPE [24], which calculates the mean Euclidean distance between the estimated joints and the ground truth in millimeters.

**MPI-INF-3DHP.** MPI-INF-3DHP [15] is a more challenging dataset that contains both indoor and complex outdoor scenes. This dataset records 8 actors performing 8 activity sets each from 14 camera views, covering more diverse poses than Human3.6M. Following previous works [4, 20], the percentage of correct keypoints (PCK) under 150 mm radius and the area under the curve (AUC) are used as evaluation metrics.

**Table 4**. Ablation study on different designs of our model.

| Attention | GCN | G2A | A2G | uMLP | MPJPE ↓ |
|:---:|:---:|:---:|:---:|:---:|:---:|
| ✓ | | | | | 51.7 |
| ✓ | ✓ | | | | 49.6 |
| ✓ | ✓ | ✓ | | | 49.2 |
| ✓ | ✓ | | ✓ | | 49.4 |
| ✓ | ✓ | ✓ | | ✓ | 48.8 |
| ✓ | ✓ | | ✓ | ✓ | 48.8 |
| ✓ | ✓ | ✓ | ✓ | ✓ | **48.3** |

## 3.2. Implementation Details

Our IGANet containing IGA and uMLP modules loops $N$=3 times. The scale factor of $f_{global}$ and $f_{graph}$ are set as 0.8 and 0.5. Following [4], we adopt 2D joints detected by CPN [1] on Human3.6M and ground truth 2D joints on MPI-INF-3DHP. Horizontal flipping is used for data augmentation during training and testing following [4, 13]. Our model is trained for 20 epochs on a single RTX 1080 Ti GPU and the batch size is set to 128. The Adam optimizer is adopted with an initial learning rate of 0.001 and a decay factor of 0.95 per epoch.

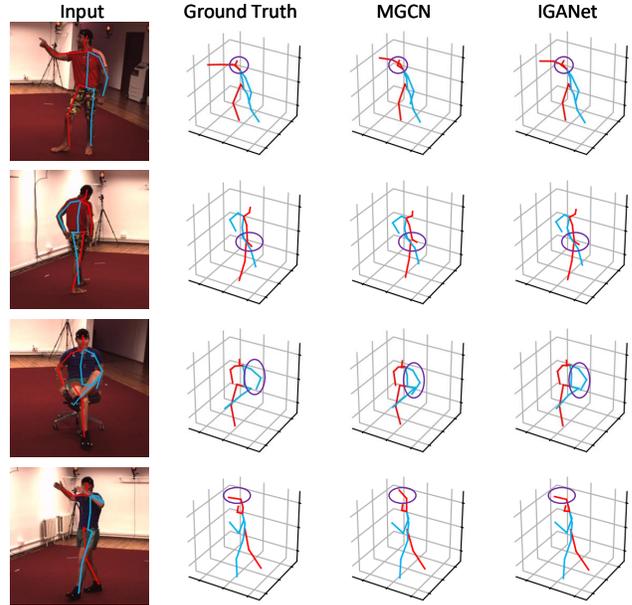## 3.3. Comparison with State-of-the-Art Methods

**Human3.6M.** Table 1 shows the comparisons of our IGANet with previous single-frame methods, using 2D keypoints detected by CPN [1]. IGANet significantly improves the performance from 49.9 mm to 48.3 mm in MPJPE. We note that some works [4, 13] adopt refinement module to further boost the performance. Following them, IGANet achieves 47.7 mm MPJPE surpassing MGCN [4] by a large margin (relative 3.4% improvement). To explore the upper bound of our method, comparison results using ground truth 2D keypoints are shown in Table 2. Our method achieves an obvious improvement (from 35.8 mm to 32.7 mm, relative 8.7% improvements), which further demonstrate its effectiveness.

**MPI-INF-3DHP.** We evaluate the generalization ability of IGANet by testing its performance on the test set of MPI-INF-3DHP after pre-training on Human3.6M. Comparison results in Table 3 show our model achieves the best performance over all metrics, which demonstrates the superior generalization ability of our model in unseen scenarios.

## 3.4. Ablation Study

To investigate the design choices of IGANet, experiments are conducted on Human3.6M for analysis and verification. The 2D pose extracted by CPN [1] is used as input. The results are reported in Table 4.

**IGA.** The combination of attention block and original MLP in transformer is adopted as the baseline. After introducing the GCN block, the performance is improved. Further, we try to make connections between GCN and attention blocks. We find that either injecting the skeleton information from GCN block into attention block (G2A) or introducing global



**Fig. 3**. Qualitative comparisons with the previous state-of-the-art method (MGCN [4]) on Human3.6M dataset.

information from attention block to GCN block (A2G) can boost the performance. Finally, the information from GCN and attention blocks are interweaved (both G2A and A2G), and the performance achieves the best (48.3 mm).

**uMLP.** For some model designs in Table 4, we replace the original MLP in transformer with our designed uMLP module. Experimental results show that uMLP module can better learn sleletal representations, which further proves its effectiveness.

## 3.5. Qualitative Results

Fig. 3 shows visualization results on Human3.6M dataset. Our IGANet is able to predict 3D poses that are close to ground truth under various actions performed by different people. Compared with the MGCN [4], IGANet can predict more reasonable and accurate 3D poses. For example, in the third row, the 3D pose of the sitting action predicted by our method is close to the ground truth. As a comparison, the left arm predicted by MGCN is at a lower position.

## 4. CONCLUSION

In this paper, we present a novel Interweaved Graph and Attention Network (IGANet) for 3D human pose estimation. The key idea is to interweave GCNs and attentions to better capture both global and local relations of human body joints. To interchange complementary information with each other, we introduce an IGA module with a bidirectional guidance mechanism. Moreover, a uMLP module is designed to gather multi-granularity information. Experimental results demonstrate the effectiveness of our IGANet for estimating 3D human pose.

# 5. REFERENCES

[1] Yilun Chen, Zhicheng Wang, Yuxiang Peng, Zhiqiang Zhang, Gang Yu, and Jian Sun, "Cascaded pyramid network for multi-person pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7103–7112.

[2] Alejandro Newell, Kaiyu Yang, and Jia Deng, "Stacked hour-glass networks for human pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016, pp. 483–499.

[3] Long Zhao, Xi Peng, Yu Tian, Mubbasir Kapadia, and Dimitris N Metaxas, "Semantic graph convolutional networks for 3D human pose regression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3425–3435.

[4] Zhiming Zou and Wei Tang, "Modulated graph convolutional network for 3D human pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 11477–11487.

[5] Wenhao Li, Hong Liu, Runwei Ding, Mengyuan Liu, Pichao Wang, and Wenming Yang, "Exploiting temporal contexts with strided transformer for 3D human pose estimation," *IEEE Transactions on Multimedia (TMM)*, 2022.

[6] Ce Zheng, Sijie Zhu, Matias Mendieta, Taojiannan Yang, Chen Chen, and Zhengming Ding, "3D human pose estimation with spatial and temporal transformers," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 11656–11665.

[7] Wenhao Li, Hong Liu, Hao Tang, Pichao Wang, and Luc Van Gool, "MHFormer: Multi-hypothesis transformer for 3D human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 13147–13156.

[8] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998–6008.

[9] Kevin Lin, Lijuan Wang, and Zicheng Liu, "Mesh graphormer," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 12939–12948.

[10] Ce Zheng, Matias Mendieta, Pu Wang, Aidong Lu, and Chen Chen, "A lightweight graph transformer network for human mesh reconstruction from 2D human pose," in *Proceedings of the 30th ACM International Conference on Multimedia (ACM MM)*, 2022, pp. 5496–5507.

[11] Wenhao Li, Hong Liu, Tianyu Guo, Hao Tang, and Runwei Ding, "GraphMLP: A graph MLP-like architecture for 3D human pose estimation," *arXiv preprint arXiv:2206.06420*, 2022.

[12] Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015, pp. 234–241.

[13] Yujun Cai, Liuhao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann, "Exploiting spatial-temporal relationships for 3D pose estimation via graph convolutional networks," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 2272–2281.

[14] Catalin Ionescu, Dragos Papava, Vlad Olaru, and Cristian Sminchisescu, "Human3.6M: Large scale datasets and predictive methods for 3D human sensing in natural environments," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 36, no. 7, pp. 1325–1339, 2013.

[15] Dushyant Mehta, Helge Rhodin, Dan Casas, Pascal Fua, Oleksandr Sotnychenko, Weipeng Xu, and Christian Theobalt, "Monocular 3D human pose estimation in the wild using improved cnn supervision," in *International Conference on 3D Vision (3DV)*, 2017, pp. 506–516.

[16] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little, "A simple yet effective baseline for 3D human pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2640–2649.

[17] Hai Ci, Chunyu Wang, Xiaoxuan Ma, and Yizhou Wang, "Optimizing network structure for 3D human pose estimation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 2262–2271.

[18] Tianhan Xu and Wataru Takano, "Graph stacked hourglass networks for 3D human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 16105–16114.

[19] Weixi Zhao, Weiqiang Wang, and Yunjie Tian, "GraFormer: Graph-oriented transformer for 3D pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 20438–20447.

[20] Ailing Zeng, Xiao Sun, Fuyang Huang, Minhao Liu, Qiang Xu, and Stephen Lin, "SRNet: Improving generalization in 3D human pose estimation with a split-and-recombine approach," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 507–523.

[21] Kehong Gong, Jianfeng Zhang, and Jiashi Feng, "PoseAug: A differentiable pose augmentation framework for 3D human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 8575–8584.

[22] Chen Li and Gim Hee Lee, "Generating multiple hypotheses for 3D human pose estimation with mixture density network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 9887–9895.

[23] Kenkun Liu, Rongqi Ding, Zhiming Zou, Le Wang, and Wei Tang, "A comprehensive study of weight sharing in graph networks for 3D human pose estimation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 318–334.

[24] Qi Dang, Jianqin Yin, Bin Wang, and Wenqing Zheng, "Deep learning based 2D human pose estimation: A survey," *Tsinghua Science and Technology*, vol. 24, no. 6, pp. 663–676, 2019.