# GAITMIXER: SKELETON-BASED GAIT REPRESENTATION LEARNING VIA WIDE-SPECTRUM MULTI-AXIAL MIXER

*Ekkasit Pinyoanuntapong, Ayman Ali, Pu Wang, Minwoo Lee*      *Chen Chen*

University of North Carolina at Charlotte
{epinyoan, aali26, pu.wang, minwoo.lee}@uncc.edu

University of Central Florida
chen.chen@crcv.ucf.edu

## ABSTRACT

Most existing gait recognition methods are appearance-based, which rely on the silhouettes extracted from the video data of human walking activities. The less-investigated skeleton-based gait recognition methods directly learn the gait dynamics from 2D/3D human skeleton sequences, which are theoretically more robust solutions in the presence of appearance changes caused by clothes, hairstyles, and carrying objects. However, the performance of skeleton-based solutions is still largely behind the appearance-based ones. This paper aims to close such performance gap by proposing a novel network model, GaitMixer, to learn more discriminative gait representation from skeleton sequence data. In particular, GaitMixer follows a heterogeneous multi-axial mixer architecture, which exploits the spatial self-attention mixer followed by the temporal large-kernel convolution mixer to learn rich multi-frequency signals in the gait feature maps. Experiments on the widely used gait database, CASIA-B, demonstrate that GaitMixer outperforms the previous SOTA skeleton-based methods by a large margin while achieving a competitive performance compared with the representative appearance-based solutions. Code will be available at https://github.com/exitudio/gaitmixer

***Index Terms***— Gait Recognition, Self-Attention, Large-kernel Convolution, Multi-axial Mixer

## 1. INTRODUCTION

Unlike short-distance biometrics (e.g., fingerprints, facial, iris, palm, and finger vein patterns), gait can be recognized from a distance without the subject's cooperation or interference. Such long-distance biometrics has a huge potential to extend its applications to forensic identification, access control, and social security. The gait recognition methods are generally either appearance-based or skeleton-based. Appearance-based approaches [1][2][3][4] utilize background subtraction to obtain silhouettes from a video sequence, which are further analyzed using carefully-designed network models for gait representation learning. On the other hand, skeleton-based approaches [5][6][7] utilize the skeleton sequences extracted from 2D/3D pose estimators as the
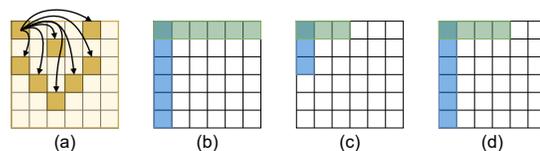


**Fig. 1**. (a) Global self-attention token mixing [8]. (b) Self-attention mixing along H and W axes [9]. (c) Convolution-mixing along H and W axes [6] (d) Heterogeneous multi-axial mixer (ours).

inputs to learn effective gait representations. Theoretically, skeleton-based methods are more robust to appearance variations caused by hairstyles, carrying objects, and clothes. However, the skeleton-based approaches, which still do not receive sufficient attention, yield a large performance gap compared with the appearance-based counterparts.

To close this gap, this paper tries to exploit more effective gait feature encoders by proposing the multi-axial mixer, which is a generic transformer-like architecture that mixes the feature patches (i.e., tokens) along each axis of the feature space, respectively, e.g., width-wise, height-wise, and channel-wise axes in image feature space. Many recent high-performance network backbones can be considered as the special cases of multi-axial mixer based on what types of mixing functions are applied, which mainly include convolution [6][10][11], self-attention [12][13], and multi-layer perceptrons (MLP) [14] [15] as shown in Fig. 1. Multi-axial mixers have been demonstrated to achieve SOTA performance in image classification and video recognition tasks, while significantly reducing computation complexities compared with other competitive network models, such as vision transformers. Despite their promising features, current multi-axial mixers generally exploit the homogeneous architecture design, where the same type of mixing functions (e.g., either convolution, self-attention, or MLP) is applied along each feature space axis. Such design, however, has limited capacity to learn multi-frequency features. In particular, it has been established that convolutions focus more on local information and therefore are good learners for high-frequency features [16]. Self-attentions, on the contrary, are designed to model
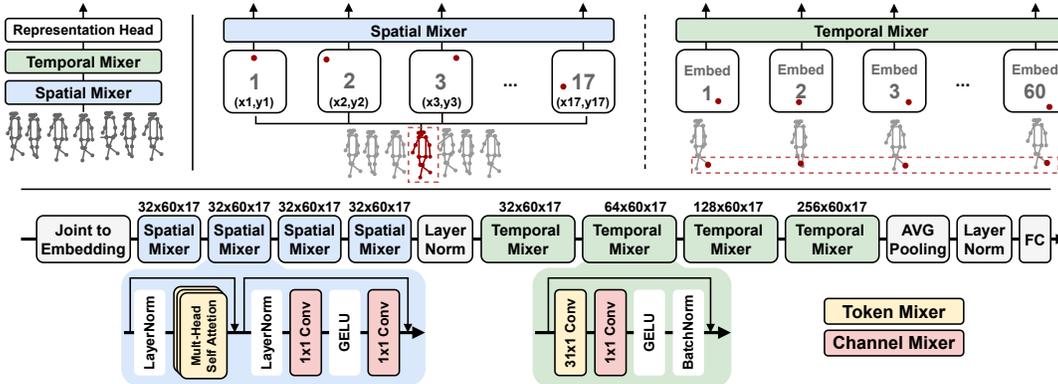
**Fig. 2**. (top) GaitMixer consists of a spatial mixer followed by a temporal mixer. (bottom) The detailed network architecture of the spatial self-attention mixer and large-kernel convolution mixer

long-range interactions and are more capable to capture low-frequency signals (global information) in feature map [17].

In this paper, we propose GaitMixer, a novel *heterogeneous* spatial-temporal axial mixer, which can effectively learn the discriminative gait representation by capturing both high-frequency and low-frequency features. In particular, GaitMixer consists of a spatial self-attention mixer and a temporal large-kernel axial mixer (Fig. 2). The spatial axial mixer learns interactions among the joints within each skeleton frame. The temporal axial mixer models the interactions among the temporal tokens of each single joint at different time indices. Experiments on the widely used gait database, CASIA-B, demonstrate that GaitMixer outperforms the previous SOTA skeleton-based [6] methods by 12% on average, while achieving a competitive performance compared with the representative appearance-based solutions.

## 2. RELATED WORK

Appearance-based approaches extract binary images of a human silhouette from the source images by subtracting static background [18]. GaitNet [1] integrates silhouette extraction into the model as an end-to-end network for gait recognition. GaitSet [2] decouples the temporal continuous sequence by learning identity information from the set of independent frames to be immune to permutation of frames and be able to integrate frames from different videos. While the majority of methods [1][2][4] take the entire shape as input, more recent approaches GaitPart [3] focus on each part of the body individually assuming that each part of human body needs its own spatial-temporal learning by separating silhouette into several parts horizontally. Skeleton-based approaches (i.e., model-based approaches) use skeleton data as the model inputs. In the early work, pose-based temporal-spatial network (PTSN) [5] utilizes a long-short term memory (LSTM) to capture the dynamic information and CNN to learn static information of a gait sequence in parallel. PoseGait [19] utilizes 3D pose es-

timated from images in order to be invariant to view changes, along with hand-crafted features including joint angle, limb length, and joint motion. The most recent methods, Gait-Graph [6] and GaitGraph2 [7], adopt graph convolution neural networks (GCNs) for gait recognition, inspired by the successes of GCNs in action recognition tasks.

## 3. GAITMIXER: HETEROGENEOUS WIDE-SPECTRUM SPATIAL-TEMPORAL MIXER

To effectively learn both high-frequency and low-frequency gait features, we introduce the GaitMixer, a heterogeneous spatial-temporal axial mixer architecture. As shown in Fig. 2, GaitMixer consists of a spatial self-attention mixer and a temporal large-kernel axial mixer (Fig. 2). The spatial axial mixer only learns interactions among the joints within each skeleton frame. The $d_y$-dimensional spatial representation $y^t \in \mathbb{R}^{|J| \times d_y}$ for each skeleton frame $t$ with $|J|$ joints is learned after $B_S$ self-attention blocks. Then, the representations of $T$ skeleton frames within a gait sequence are concatenated into $z \in \mathbb{R}^{|J| \times T \times d_y}$, which is then forwarded to a temporal axial mixer to capture the interactions among the tokens of each single joint at different temporal indices. The temporal axial mixer consists of $B_T$ one-dimensional large-kernel convolution blocks. Moreover, to simplify our Gait-Mixer architecture, both spatial and temporal mixers adopt the isotropic design, which does not perform feature down-sampling and maintains the same feature resolutions at all layers.

### 3.1. Spatial Mixer with Axial Self-attention

The spatial mixer module aims to learn a high dimensional representation embedding from each skeleton frame. Although self-attention tends to capture low-frequency (or global) features, our experiments demonstrate that self-attention is sufficient to learn both high-frequency and low-

frequency signals in the feature map along spatial axis. This also indicates that self-attention can effectively model both short-range and long-range inter-joint dependencies. Given a 2D skeleton with joints $J$, we consider each joint (*i.e.*, $x$ and $y$ coordinates) as a spatial token (with 2 channels) and perform the feature extraction among all $|J|$ spatial tokens by following the isotropic transformer pipeline. Specifically, the spatial taken $x_i \in \mathbb{R}^2$ is passed through a trainable linear projection, which maps each token to a high dimension embedding $\mathbf{x}_i \in \mathbb{R}^{d_x}$. Then, the spatial token embeddings of each skeleton frame $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{|J|})$ are mixed by inter-token dot product attentions to generate an output sequence $\mathbf{y} = (\mathbf{y}_1, \mathbf{y}_2, \ldots, \mathbf{y}_{|j|})$ where $\mathbf{y}_i \in \mathbb{R}^{d_y}$. Running $h$ self-attentions in parallel leads to the multihead self-attention with $h$ heads, where the outputs of the attention heads are concatenated and projected into the expected dimensions.

### 3.2. Temporal Mixer with Large-kernel Convolution

The essence of a walking sequence is composed of multiple short repeated cycles. In the temporal axis, self-attention may not be able to capture wide-band multi-frequency features, considering that the global receptive field of self-attention is much easier to capture low-frequency features. It demands a large amount of data for self-attention to establish the desirable locality inductive bias that is the key to learn high-frequency features. To learn both high-frequency and low-frequency temporal data, we utilize large kernel depth-wise separable convolution in the temporal mixer as illustrated in Fig. 2. In general, convolution neural networks tend to capture high-frequency (local) features, however, the large kernel allows the model to also learn low-frequency features. 60 frames are used in the temporal model which covers around 4 cycles of walking. A large one-dimension kernel with size of $31 \times 1$ is used to capture mid-range information (around 2 walking cycles). A reverse padding with size of 30 is applied to keep the temporal dimension the same. The temporal mixer only communicates with all frames along the temporal axis of the same joints. A temporal mixer is composed of two types of axial mixers. First, a token mixer is implemented by a depth-wise convolution that learns all embeddings only in the same channel. Next, a channel mixer is a $1 \times 1$ convolution that learns only specific embedding along all channels.

### 3.3. Representation head and Loss Function

We apply average pooling along spatial and temporal dimensions to reduce the output from $\mathbf{x}_{temp} \in \mathbb{R}^{F \cdot (J \times c)}$ to $\mathbf{x}_{hidden} \in \mathbb{R}^c$, where $F$ and $J$ are number of frames and joints respectively. The number of channels $c$ is set to 256. Finally, layer norm, fully connected layer, and $l^2$-norm are applied respectively and return the feature embedding in 128 dimensions. To learn the discriminative gait representation, we apply the triplet loss with multi-similarity miner.
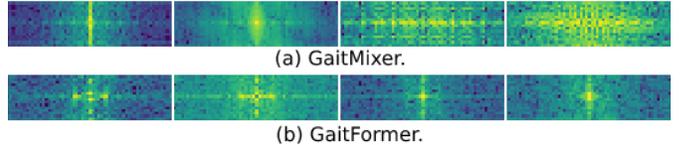


**Fig. 3**. 2D FFT of GaitMixer and GaitFormer feature maps of 4 channels. Higher temperature indicates larger magnitude. Pixels closer to the center represent lower frequencies

## 4. EXPERIMENTS

### 4.1. Dataset

**CASIA-B** [20] has been widely adopted as a multi-view, RGB, and silhouette gait dataset. The data acquisition is performed by 124 individuals from 11 viewing angles ranging from 0 to 180 with 18 angle differences. To mimic typical daily walking conditions, each subject performs six sequences of normal walking (NM), two sequences of walking with a coat (CL), and two sequences of walking with carrying a bag (BG). For each individual, ten sequences are captured from each view angle. This paper follows a widely-used test protocol [1][2][3][6][7][19], which uses the data of the first 74 subjects' sequences for the training and the remaining 50 subjects' sequences for testing. Furthermore, the test dataset is divided into gallery and probe sets. The gallery set includes the first four sequences of the normal walking condition. The probe set consists of the last two sequences of normal walking, two walking with a coat on, and walking with carrying a bag. Finally, the results are reported for all viewing angles.

### 4.2. Implementation Details

**Training Details** HRNet [21] is used as a 2D human pose estimator. We follow data augmentations from GaitGraph[6] and add normalization of the joint position in $(x, y)$-coordinates by dividing 320 which is the width of the original videos to input data while keeping the aspect ratio. Adam optimizer is used with $6e-3$ learning rate with 1-cycle learning rate and $1e-5$ weight decay. We are using a balanced batch sampler to sample the number of walking data per person equally. The batch size is (74, 4), denoting 74 people and 4 walking samples per person. **Testing**. Each gait testing sample contains 60 frames selected from the middle of the sequence data. The test set is separated into probe and gallery. Both are fed into the model to obtain the feature representations. The ID of the gallery representation that has the smallest cosine distance from the probe will be the predicted ID of the probe.

### 4.3. Comparison with the SOTA Methods

To demonstrate the superior performance of GaitMixer as a heterogeneous multi-axial mixer model, we also build GaitFormer, which is a homogeneous multi-axial mixer that

**Table 1**. Averaged Rank-1 accuracies on CASIA-B per probe angle excluding identical-view cases.

| Gallery NM#1-4 | | 0°-180° | | | | | | | | | | | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Probe | 0° | 18° | 36° | 54° | 72° | 90° | 108° | 126° | 144° | 162° | 180° | |
| NM#5-6 | PoseGait [19] | 55.3 | 69.6 | 73.9 | 75.0 | 68.0 | 68.2 | 71.1 | 72.9 | 76.1 | 70.4 | 55.4 | 68.7 |
| | GaitGraph [6] | 85.3 | 88.5 | 91.0 | 92.5 | 87.2 | 86.5 | 88.4 | 89.2 | 87.9 | 85.9 | 81.9 | 87.7 |
| | GaitGraph2 [7] | 78.5 | 82.9 | 85.8 | 85.6 | 83.1 | 81.5 | 84.3 | 83.2 | 84.2 | 81.6 | 71.8 | 82.0 |
| | GaitFormer (ours) | 90.9 | 91.2 | 93.7 | 91.9 | 91.9 | 92.7 | 93.3 | 91.8 | 92.5 | 90.5 | 85.5 | 91.5 |
| | **GaitMixer(ours)** | **94.4** | **94.9** | **94.6** | **96.3** | **95.3** | **96.3** | **95.3** | **94.7** | **95.3** | **94.7** | **92.2** | **94.9** |
| BG#1-2 | PoseGait [19] | 35.3 | 47.2 | 52.4 | 46.9 | 45.5 | 43.9 | 46.1 | 48.1 | 49.4 | 43.6 | 31.1 | 44.5 |
| | GaitGraph [6] | 75.8 | 76.7 | 75.9 | 76.1 | 71.4 | 73.9 | 78.0 | 74.7 | 75.4 | 75.4 | 69.2 | 74.8 |
| | GaitGraph2 [7] | 69.9 | 75.9 | 78.1 | 79.3 | 71.4 | 71.7 | 74.3 | 76.2 | 73.2 | 73.4 | 61.7 | 73.2 |
| | GaitFormer (ours) | 82.5 | 83.2 | 85.7 | 85.7 | 84.2 | 80.2 | 78.9 | 82.6 | 82.2 | 78.6 | 71.3 | 81.4 |
| | **GaitMixer(ours)** | **83.5** | **85.6** | **88.1** | **89.7** | **85.2** | **87.4** | **84.0** | **84.7** | **84.6** | **87.0** | **81.4** | **85.6** |
| CL#1-2 | PoseGait [19] | 24.3 | 29.7 | 41.3 | 38.8 | 38.2 | 38.5 | 41.6 | 44.9 | 42.2 | 33.4 | 22.5 | 36.0 |
| | GaitGraph [6] | 69.6 | 66.1 | 68.8 | 67.2 | 64.5 | 62.0 | 69.5 | 65.6 | 65.7 | 66.1 | 64.3 | 66.3 |
| | GaitGraph2 [7] | 57.1 | 61.1 | 68.9 | 66 | 67.8 | 65.4 | 68.1 | 67.2 | 63.7 | 63.6 | 50.4 | 63.6 |
| | GaitFormer (ours) | 76.1 | 80.3 | 81.0 | 78.2 | 77.7 | 76.6 | 77.4 | 75.8 | 76.5 | 75.7 | 77.2 | 77.2 |
| | **GaitMixer(ours)** | **81.2** | **83.6** | **82.3** | **83.5** | **84.5** | **84.8** | **86.9** | **88.9** | **87.0** | **85.7** | **81.6** | **84.5** |

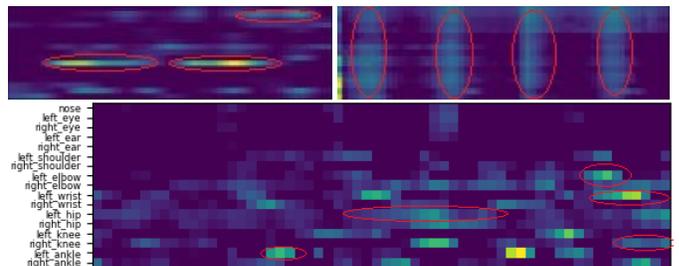**Table 2**. Averaged Rank-1 accuracies on CASIA-B comparison with both appearance-based and skeleton-based methods

| | Method | Probe | | | |
|---|---|---|---|---|---|
| | | NM | BG | CL | Mean |
| Appearance based | GaitNet [1] | 91.6 | 85.7 | 58.9 | 78.7 |
| | GaitSet [2] | 95.0 | 87.2 | 70.4 | 84.2 |
| | GaitPart [3] | **96.2** | **91.5** | 78.7 | **88.8** |
| Skeleton based | PoseGait | 68.7 | 44.5 | 36.0 | 49.7 |
| | GaitGraph | 87.7 | 74.8 | 66.3 | 76.3 |
| | GaitGraph2 | 82.0 | 73.2 | 63.6 | 72.9 |
| | **GaitFormer (ours)** | 91.5 | 81.4 | 77.2 | 83.4 |
| | **GaitMixer (ours)** | 94.9 | 85.6 | **84.5** | 88.3 |



**Fig. 4**. **Grad-CAM** [22] visualizations. **Top-left**: GaitGraph [6]. **Top-right**: GaitFormer. **Bottom**: GaitMixer. X-axis represents frames 1 to 60 and Y-axis represents 17 joints. Features with higher contributions have higher heat temperatures.

adopts self-attention for both spatial and temporal axes. In Fig. 3, we visualize the frequency magnitude of the output feature maps from GaitMixer and GaitFormer, respectively. It can be observed that GaitMixer concentrates on both high-frequency and low-frequency components along both temporal and spatial axes in feature maps. This confirms the superior capacity of GaitMixer to capture features in wide-spectrum bands. GaitFormer, however, cannot effectively model the high-frequency feature components. The performance comparisons between our approaches and the SOTA skeleton-based methods are shown in Table 1. It is shown that our multi-axial mixer models outperform the existing solutions by a large margin in both cross-view and cross-walking-condition cases. Moreover, GaitMixer achieves better recognition accuracy than GaitFormer because GaitMixer can jointly exploit the heterogeneous mixing at different feature space dimensions. Table 2 shows a competitive performance of our skeleton-based methods, compared with the representative appearance-based methods. GaitMixer achieves much higher accuracy than all appearance-based approaches in wearing coat condition. It is due to the inherent robustness of skeleton data against large appearance changes.

### 4.4. Visualization

We use class activation map (Grad-CAM) [22] to show which parts of the input gait sequence contribute most to the final recognition result. As shown in Fig. 4 (bottom), GaitMixer focuses on continuous joint sequences with a variety of differ-

ent temporal windows, thus capturing short-, mid, and- long-range temporal feature interactions. Moreover, GaitMixer also pays attentions to a diverse set of joints except ears, eyes, and nose, which is also as expected because the landmarks on face are not relevant to the gait dynamics. As shown in Fig. 4 (top-left), GaitGraph tends to focus on some specific joints over a large temporal window and it also exploits the features from face landmarks for gait recognition. Both limitations could degrade the performance of GaitGraph. GaitFormer (Fig. 4 (top-right)) pays more attention to certain skeleton frameworks without capturing rich spatial-temporal feature interactions. This can be the key contributing factor that affects its performance.

## 5. CONCLUSION

In this paper, we present GaitMixer model, a novel heterogeneous multi-axial architecture combining a spatial self-attention mixer and a large kernel temporal convolution mixer to capture both high-frequency and low-frequency dynamics of gait data. Our approach achieves the best accuracy on the well-known CASIA-B gait dataset for all conditions when compared to previous skeleton-based methods and is superior to appearance-based approaches with coats conditions.

# 6. REFERENCES

[1] Chunfeng Song, Yongzhen Huang, Yan Huang, Ning Jia, and Liang Wang, "Gaitnet: An end-to-end network for gait based human identification," *Pattern Recognition*, vol. 96, pp. 106988, 2019.

[2] Hanqing Chao, Kun Wang, Yiwei He, Junping Zhang, and Jianfeng Feng, "Gaitset: Cross-view gait recognition through utilizing gait as a deep set," 2021.

[3] C. Fan, Y. Peng, C. Cao, X. Liu, S. Hou, J. Chi, Y. Huang, Q. Li, and Z. He, "Gaitpart: Temporal part-based model for gait recognition," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 14213–14221.

[4] Z. Wu, Y. Huang, L. Wang, X. Wang, and T. Tan, "A comprehensive study on cross-view gait based human identification with deep cnns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 2, pp. 209–226, 2017.

[5] Rijun Liao, Chunshui Cao, Edel B. Garcia, Shiqi Yu, and Yongzhen Huang, "Pose-based temporal-spatial network (ptsn) for gait recognition with carrying and clothing variations," in *Biometric Recognition*. 2017, pp. 474–483, Springer International Publishing.

[6] Torben Teepe, Ali Khan, Johannes Gilg, Fabian Herzog, Stefan Hörmann, and Gerhard Rigoll, "Gaitgraph: Graph convolutional network for skeleton-based gait recognition," in *2021 IEEE International Conference on Image Processing (ICIP)*, 2021, pp. 2314–2318.

[7] Torben Teepe, Johannes Gilg, Fabian Herzog, Stefan Hörmann, and Gerhard Rigoll, "Towards a deeper understanding of skeleton-based gait recognition," 2022.

[8] 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, "An image is worth 16x16 words: Transformers for image recognition at scale," 2020.

[9] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid, "Vivit: A video vision transformer," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 6816–6826.

[10] Sijie Yan, Yuanjun Xiong, and Dahua Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," 2018.

[11] Yi-Fan Song, Zhang Zhang, Caifeng Shan, and Liang Wang, "Stronger, faster and more explainable: A graph convolutional baseline for skeleton-based action recognition," in *Proceedings of the 28th ACM International Conference on Multimedia*. oct 2020, ACM.

[12] Chiara Plizzari, Marco Cannici, and Matteo Matteucci, "Spatial temporal transformer network for skeleton-based action recognition," in *Pattern Recognition. ICPR International Workshops and Challenges*, pp. 694–701. Springer International Publishing, 2021.

[13] C. Zheng, S. Zhu, M. Mendieta, T. Yang, C. Chen, and Z. Ding, "3d human pose estimation with spatial and temporal transformers," 2021.

[14] I. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit, M. Lucic, and A. Dosovitskiy, "Mlp-mixer: An all-mlp architecture for vision," 2021.

[15] D. J. Zhang, K. L., Y. Chen, Y. Wang, S. Chandra, Y. Qiao, L. Liu, and M. Z. Shou, "Morphmlp: A self-attention free, mlp-like backbone for image and video," *arXiv preprint arXiv:2111.12527*, 2021.

[16] X. Ding, X. Zhang, Y. Zhou, J. Han, G. Ding, and J. Sun, "Scaling up your kernels to 31x31: Revisiting large kernel design in cnns," 2022.

[17] Zizheng Pan, Jianfei Cai, and Bohan Zhuang, "Fast vision transformers with hilo attention," 2022.

[18] Liang Wang, Tieniu Tan, Huazhong Ning, and Weiming Hu, "Silhouette analysis-based gait recognition for human identification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 25, no. 12, pp. 1505–1518, 2003.

[19] Rijun Liao, Shiqi Yu, Weizhi An, and Yongzhen Huang, "A model-based gait recognition method with body pose and human prior knowledge," *Pattern Recognition*, vol. 98, pp. 107069, 2020.

[20] Shiqi Yu, Daoliang Tan, and Tieniu Tan, "A framework for evaluating the effect of view angle, clothing and carrying condition on gait recognition," in *18th International Conference on Pattern Recognition (ICPR'06)*, 2006, vol. 4, pp. 441–444.

[21] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang, "Deep high-resolution representation learning for human pose estimation," 2019.

[22] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," *International Journal of Computer Vision*, vol. 128, no. 2, pp. 336–359, oct 2019.