

Research Article

Optimizing Human Pose Estimation Through Focused Human and Joint Regions

Yingying Jiao^{1,2}, Zhigang Wang³, Zhenguang Liu^{4,5}, Shaojing Fan⁶, Sifan Wu^{1,2}, Zheqi Wu³, Zhuoyue Xu³

1. College of Computer Science and Technology, Jilin University, China; 2. Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education, Jilin University, China; 3. College of Computer Science and Technology, Zhejiang Gongshang University, China; 4. The State Key Laboratory of Blockchain and Data Security, Zhejiang University, China; 5. Hangzhou High-Tech Zone (Binjiang) Institute of Blockchain and Data Security, China; 6. School of Computing, National University of Singapore, Singapore

Human pose estimation has given rise to a broad spectrum of novel and compelling applications, including action recognition, sports analysis, as well as surveillance. However, accurate video pose estimation remains an open challenge. One aspect that has been overlooked so far is that existing methods learn motion clues from all pixels rather than focusing on the target human body, making them easily misled and disrupted by unimportant information such as background changes or movements of other people. Additionally, while the current Transformer-based pose estimation methods has demonstrated impressive performance with global modeling, they struggle with local context perception and precise positional identification.

In this paper, we try to tackle these challenges from three aspects: (1) We propose a bilayer Human-Keypoint Mask module that performs coarse-to-fine visual token refinement, which gradually zooms in on the target human body and keypoints while masking out unimportant figure regions. (2) We further introduce a novel deformable cross attention mechanism and a bidirectional separation strategy to adaptively aggregate spatial and temporal motion clues from constrained surrounding contexts. (3) We mathematically formulate the deformable cross attention, constraining that the model focuses solely on the regions centered at the target person body. Empirically, our method achieves state-of-the-art performance on three large-scale benchmark datasets. A remarkable highlight is that our method achieves an 84.8 mean Average Precision (mAP) on the challenging wrist joint, which significantly outperforms the 81.5 mAP achieved by the current state-of-the-art method on the PoseTrack2017 dataset.

Introduction

Human pose estimation, as a fundamental problem in the realm of computer vision and artificial intelligence^{[1][2]}, involves accurately identifying the anatomical keypoints of human bodies. Precise pose estimation is the key for the success of a machine as it paves the way for machines to accurately interpret human movements and behaviors. Accordingly, human pose estimation spans a wide range of applications from action recognition, movement tracking, to augmented reality^{[3][4][5][6][7]}.

A plethora of research has been dedicated to the field of pose estimation on still images, evolving from early methods employing tree-based and random forest models^{[8][9]} to current methodologies utilizing convolutional neural networks^[10] and Transformers^[11]. Despite their excellent performance on still images, applying these methods directly to video pose estimation leads to significant performance degradation due to the exclusive characteristics in videos, such as rapid movement and video defocus, which are frequently encountered in videos but absent in static images^[12].

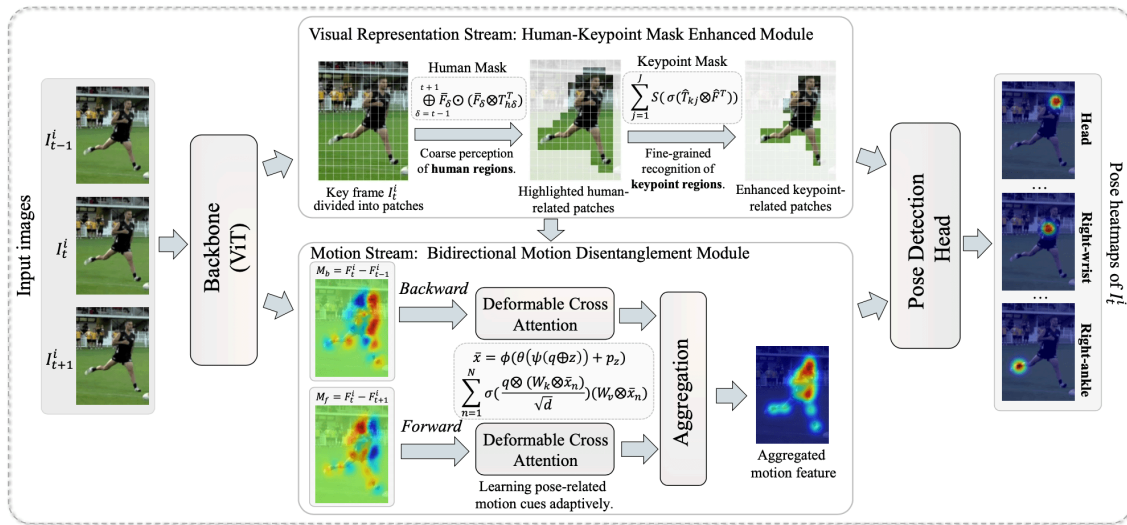


Figure 1. A high-level overview of our proposed VREMD, which utilizes a dual-stream architecture to collaboratively process and integrate complementary visual and motion features. The visual representation stream executes progressive enhancement of human keypoint-related features to achieve precise location recognition. The motion stream performs adaptive pose-related motion disentanglement through the novel deformable cross attention. $\{\mathbf{F}_{t-1}^i, \mathbf{F}_t^i, \mathbf{F}_{t+1}^i\}$ denote the visual features of three input frames $\{I_{t-1}^i, I_t^i, I_{t+1}^i\}$ output by backbone network.

To address this issue, substantial studies have emerged that leverage temporal continuity to extract rich semantic visual contexts for human pose estimation in videos. Current methods can be roughly categorized into two main branches. One line of research^{[13][14]} aggregates temporal information from neighboring frames for video pose estimation, employing CNN-based architectures and pose calibration. Fueled by the development of Transformers^{[15][16]}, another line of studies^{[17][18]} strive to integrate attention mechanisms into model construction, yielding impressive results and showcasing their immense potential. However, a limitation inherent in existing Transformer-based methods^[17] lies in their inability to effectively manage local dependencies. This limitation poses a notable challenge for visual perception tasks such as pose estimation, which require precise local positioning.

Following thorough experimentation and empirical investigation, we uncover two insights: (1) Existing methods^{[19][20][21]} struggle to handle subtle pose changes, particularly in challenging scenarios with occlusions or motion blur. This may stem from the fact that current methods tend to capture temporal dynamics pixel-by-pixel rather than focusing solely on target human regions, leading to them being distracted by unuseful cues such as background changes or pixels far from the target person. (2) Additionally, previous studies^{[14][19]} adopting multiple sets of fixed deformable convolutions with varying dilation rates, which neglect the importance of adaptive scale selection.

Inspired by these, we propose a dual-stream framework, which executes Visual Representation Enhancement and Motion Disentanglement (VREMD) for human pose estimation in videos. Technically, we embrace three novel designs to tackle the challenge. (1) We propose a two-step human-keypoint mask module for coarse-to-fine visual enhancement, which progressively refines extracted representations from the human body and keypoints perspectives. (2) We further introduce a bidirectional decoupled module tailored for adaptively disentangling motion cues of the target person from unnecessary visual elements. (3) Furthermore, we mathematically formulate a deformable cross attention mechanism that constrains the model to focus exclusively on regions circumscribing the target human body.

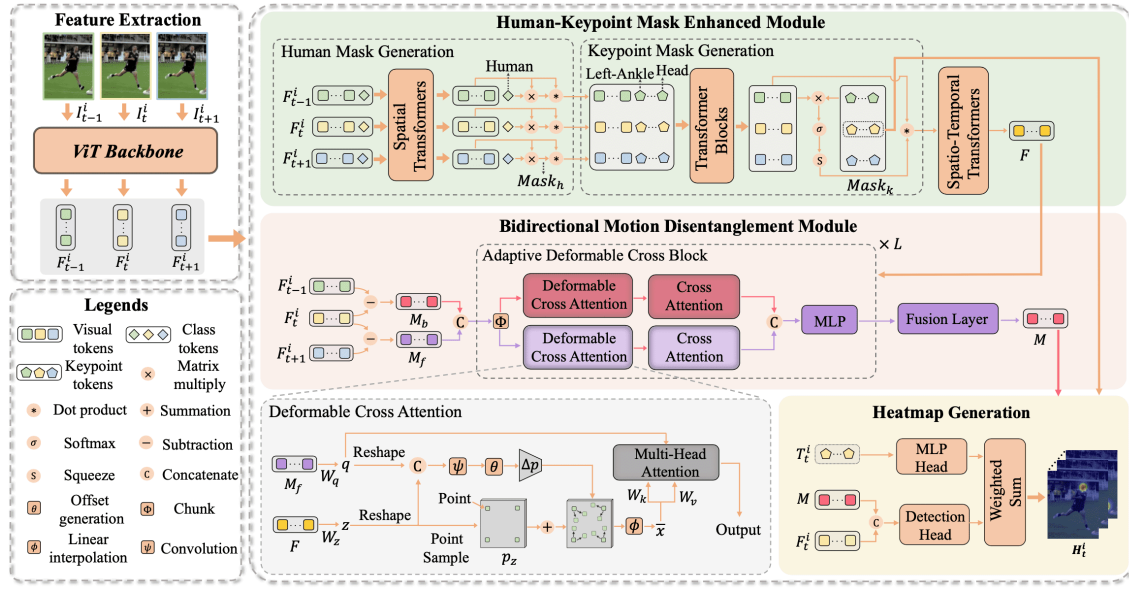


Figure 2. The overall pipeline of our VREMD framework. Given an input sequence $\{I_{t-1}^i, I_t^i, I_{t+1}^i\}$, our goal is to estimate the human pose of the key frame I_t^i . We initially extract the visual features via a ViT backbone, and then feed them into the Human-Keypoint Enhanced module and the Bidirectional Motion Disentanglement module to obtain T_t^i and M . Finally, the outputs derived from different heads are combined through a weighted sum to arrive at the final predicted pose heatmap H_t^i .

Our framework exemplifies the collaborative advantage between local spatial focus and adaptive temporal clues extraction, opening up possibilities for rethinking the pose estimation task from emphasizing on the target human body and masking out the irrelevant spatio-temporal contexts. To evaluate the efficacy of our method, we conduct extensive experiments on three public benchmarks, achieving state-of-the-art performance. The key contributions of our method are summarized as follows:

- We present a dual-stream framework that integrates visual enhancement and motion disentanglement to highlight target human areas and filter other non-essential regions for human pose estimation.
- We creatively introduce a deformable cross attention to disentangle pose-related motion cues, harnessing bidirectional temporal dynamics and enabling the model to robustly handle complex pose variations of the target human.
- Empirically, our method achieves state-of-the-art performance on three large-scale benchmarks, and overall provides insights into integrating Transformer-based methods with region-specific enhancement strategies to boost their local localization capabilities.

Our Method

Preliminaries. Our method follows the top-down paradigm, which first extracts each individual person from an image and then estimates their poses. Specifically, we first utilize an object detector to extract the bounding box for person i in a video frame I_t that is to be detected. Subsequently, we expand the bounding box by 25% and crop the same person in the adjacent frames (i.e., I_{t-1} and I_{t+1}). As a result, we obtain a sequence of consecutive frames for person i : $\mathcal{I}_t^i = \{I_{t-1}^i, I_t^i, I_{t+1}^i\}$. Given a sequence of video frames \mathcal{I}_t^i that includes the key frame I_t^i and the auxiliary frames I_{t-1}^i and I_{t+1}^i , our target is to detect the human pose within I_t^i . We aim to strengthen the utilization of supplementary temporal information in auxiliary frames by employing incremental visual representation enhancement and adaptively disentangling useful motion information, thus tackling the common issue of existing methods being interfered with by irrelevant information regarding the target human.

Method overview. The overview pipeline of our proposed VREMD is depicted in Figure 2. VREMD constructs a dual-stream architecture with inter-module communication that enhances both visual features and captures meaningful motion cues. Specifically, VREMD incorporates two distinct modules: a Human-Keypoint Mask Enhanced module (HKME) and a Bidirectional Motion Disentanglement module (BMD). First, we utilize a Vision Transformer backbone to extract visual features $\{\mathbf{F}_{t-1}^i, \mathbf{F}_t^i, \mathbf{F}_{t+1}^i\}$ from the input frame sequence \mathcal{I}_t^i , which are then simultaneously fed into both the HKME and BMD modules. The HKME generates dual masks for a coarse-to-fine representation refinement, resulting in enhanced feature \mathbf{F} and key frame keypoint tokens \mathbf{T}_t^i . The BMD computes the motion features and, utilizing \mathbf{F} as a constraint, dynamically derives joint-related motion contexts to produce the filtered \mathbf{M} . Finally, the keypoint heatmaps \mathbf{H}_k from key frame tokens \mathbf{T}_t^i via an MLP and the heatmaps \mathbf{H}_m decoded from \mathbf{M} and key frame features \mathbf{F}_t^i are weighted, summed, and combined to produce the final pose estimation \mathbf{H}_t^i . The following sections will elaborate on the two key components in detail.

Human-Keypoint Mask Enhanced Module

Despite the Transformers architecture achieving remarkable success in various fields^{[15][16]}, its application in video pose estimation has been limited. Given the significant potential demonstrated by this architecture in other visual perception tasks^{[22][23]}, we seek to design a novel Transformer-based framework specially tailored for video pose detection. A naive approach to aggregate unique temporal cues from a video would be to concatenate features across multiple frames for full-token computation.

Yet, such a straightforward treatment strategy faces two issues: excessive capture of redundant information between adjacent frames, and a lack of focus on task-relevant tokens.

Inspired by previous work^{[24][11]}, we propose a Human-Keypoint Mask Enhanced module with a progressive refinement architecture, addressing the aforementioned issues through three steps: (1) We generate a human mask to coarsely enhance the perception of the target human. (2) We produce a keypoint mask to achieve finer filtering of keypoint-related features. (3) We utilize spatio-temporal networks to aggregate the highlighted spatio-temporal cues of these visual features. This step-by-step optimization strategy can discern articular visual tokens, simulating the capability of localized identification, which promotes precise pose estimation.

Human mask. Given a visual feature sequence $\{\mathbf{F}_{t-1}^i, \mathbf{F}_t^i, \mathbf{F}_{t+1}^i\} \in \mathbb{R}^{3 \times N \times D}$ output by the ViT backbone, we concatenate a learnable class token $\mathbf{T}_h \in \mathbb{R}^{3 \times 1 \times D}$ with a category of human to each feature. These features then individually pass through cascaded Transformer blocks for intra-frame spatial similarity computation. We separate the result into human token \mathbf{T}_h and visual features $\bar{\mathbf{F}} \in \mathbb{R}^{3 \times N \times D}$. After transposing the human token, we perform matrix multiplication to obtain the human mask $\mathbf{Mask}_h \in \mathbb{R}^{3 \times N \times 1}$. Finally, we secure a coarsely selected feature $\mathbf{F}_c \in \mathbb{R}^{3 \times N \times D}$ by executing element-wise dot product between the human mask \mathbf{Mask}_h and the visual feature $\bar{\mathbf{F}}$, utilizing broadcasting. The above operations can be formulated as:

$$\mathbf{F}_c = \bigoplus_{\delta=t-1}^{t+1} \bar{\mathbf{F}}_\delta \odot \underbrace{(\bar{\mathbf{F}}_\delta \otimes \mathbf{T}_{h\delta}^T)}_{\mathbf{Mask}_{h\delta}}, \quad (1)$$

where \bigoplus , δ , \odot , \otimes , and \mathbf{T}^T denote concatenation, temporal index of frames, dot product, matrix multiplication, the transpose of \mathbf{T} , respectively.

Keypoint mask. In pursuit of more precise keypoint-related feature enhancement, we employ additional auxiliary tokens to accurately localize spatial positions by integrating multi-frame representations in the spatio-temporal domain. We concatenate the learnable keypoint tokens $\mathbf{T}_k \in \mathbb{R}^{3 \times J \times D}$ (Note that J is the number of keypoints) to the coarsely selected feature \mathbf{F}_c and separate the multi-frame features, which are then linked along the token dimension and fed into Transformer blocks for spatio-temporal learning. Subsequently, we split the visual features and keypoint tokens from the output and gather them over multiple frames, resulting in multi-frame features $\hat{\mathbf{F}} \in \mathbb{R}^{(3 \cdot N) \times D}$ and multi-frame keypoint tokens $\hat{\mathbf{T}}_k \in \mathbb{R}^{(3 \cdot J) \times D}$. After transposing the multi-frame features, we perform matrix multiplication with the multi-frame keypoint tokens to produce the keypoint confidence map $\mathbf{Map} \in \mathbb{R}^{(3 \cdot J) \times (3 \cdot N)}$. We apply the

softmax function to compute element-wise weights for the map **Map**, and summing along the second-to-last dimension followed by transposition yields the keypoint mask $\mathbf{Mask}_k \in \mathbb{R}^{(3 \cdot N) \times 1}$:

$$\mathbf{Mask}_k = \sum_{j=1}^J S(\sigma(\hat{\mathbf{T}}_{kj} \otimes \hat{\mathbf{F}}^T)), \quad (2)$$

where $j \in \{1, \dots, J\}$, $S(\cdot)$, $\sigma(\cdot)$, and \otimes denote the keypoint index, squeeze operation, softmax function, and matrix multiplication, respectively. The keypoint mask is element-wise multiplied with the multi-frame features $\hat{\mathbf{F}}$ to create the refined filtered features $\mathbf{F}_f \in \mathbb{R}^{(3 \cdot N) \times D}$.

Spatio-temporal aggregation. To fully leverage the refined representation information, we perform decoupled spatio-temporal feature aggregation through the spatio-temporal Transformers. Specifically, we first separate the refined filtered features \mathbf{F}_f and undertake frame-level spatial modulation. Then, each token is concatenated with its corresponding token in the temporal domain to undergo temporal modulation, resulting in $\bar{\mathbf{F}}_f \in \mathbb{R}^{(3 \cdot N) \times D}$. Finally, we adopt an MLP to execute token dimensionality reduction on $\bar{\mathbf{F}}_f$ to attain spatio-temporal aggregation of multi-frame features, leading to the enhanced feature $\mathbf{F} \in \mathbb{R}^{N \times D}$.

Bidirectional Motion Disentanglement Module

To extract useful complementary information from auxiliary frames, prior methods^{[14][20]} implicitly model feature residuals to capture motion evidence. The common practice among these paradigms is to directly concatenate the computed multiple motion features for convolution after their calculation, which considers temporal continuity but overlooks insights from the temporal direction. We observe that, from the perspective centered around the key frame, the essential temporal details that need to be focused on actually originate from two different directions, namely forward and backward. Considering this intrinsic factor, we design a bidirectional separation strategy to decouple the continuous motion into parallel forward and backward motion trajectories. Furthermore, existing methods do not differentiate motion clues in the spatial dimension, which can lead to learning pose-irrelevant information (*e.g.*, background, other people, etc.) that can disrupt detection. Moreover, existing methods heavily rely on deformable convolutions for local motion calibration, potentially leading to models that are overly tailored and limiting their compatibility with Transformer-based architectures. To tackle these challenges, we introduce deformable cross attention (DCA) for the first time and create the Adaptive Deformable Cross block by employing it, which adaptively captures pose-related motion dynamics.

Method	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	Mean
PoseTracker ^[25]	67.5	70.2	62.0	51.7	60.7	58.7	49.8	60.6
PoseFlow ^[26]	66.7	73.3	68.3	61.1	67.5	67.0	61.3	66.5
JointFlow ^[27]	-	-	-	-	-	-	-	69.3
FastPose ^[28]	80.0	80.3	69.5	59.1	71.4	67.5	59.4	70.3
TML++ ^[29]	-	-	-	-	-	-	-	71.5
Simple (R-50) ^[30]	79.1	80.5	75.5	66.0	70.8	70.0	61.7	72.4
Simple (R-152) ^[30]	81.7	83.4	80.0	72.4	75.3	74.8	67.1	76.7
STEmbedding ^[31]	83.8	81.6	77.1	70.0	77.4	74.5	70.8	77.0
HRNet ^[10]	82.1	83.6	80.4	73.3	75.5	75.3	68.5	77.3
MDPN ^[32]	85.2	88.5	83.9	77.5	79.0	77.0	71.4	80.7
CorrTrack ^[33]	86.1	87.0	83.4	76.4	77.3	79.2	73.3	80.8
Dynamic-GNN ^[34]	88.4	88.4	82.0	74.5	79.1	78.3	73.1	81.1
PoseWarper ^[13]	81.4	88.3	83.9	78.0	82.4	80.5	73.6	81.2
DCPose ^[14]	88.0	88.7	84.1	78.4	83.0	81.4	74.2	82.8
DetTrack ^[35]	89.4	89.7	85.5	79.5	82.4	80.8	76.4	83.8
SLT-Pose ^[36]	88.9	89.7	85.6	79.5	84.2	83.1	75.8	84.2
HANet ^[37]	90.0	90.0	85.0	78.8	83.1	82.1	77.1	84.2
KPM ^[38]	89.5	90.0	87.6	81.8	81.1	82.6	76.1	84.6
M-HANet ^[39]	90.3	90.7	85.3	79.2	83.4	82.6	77.8	84.8
FAMI-Pose ^[19]	89.6	90.1	86.3	80.0	84.6	83.4	77.0	84.8
DSTA ^[18]	89.3	90.6	87.3	82.6	84.5	85.1	77.8	85.6
TDMI-ST ^[20]	90.6	91.0	87.2	81.5	85.2	84.5	78.7	85.9
VREMD (Ours)	89.9	91.4	88.8	84.8	88.5	87.8	81.0	87.6

Table 1. Comparisons with the state-of-the-art methods for video pose estimation on the validation sets of the PoseTrack2017^[40] dataset. Note that we aggregate temporal information from neighboring frames (i.e., one frame to the left and one to the right).

Adaptive Deformable Cross block. Given the features $\{\mathbf{F}_{t-1}^i, \mathbf{F}_t^i, \mathbf{F}_{t+1}^i\}$ from the backbone, we subtract \mathbf{F}_t^i from both \mathbf{F}_{t+1}^i and \mathbf{F}_{t-1}^i to obtain $\{\mathbf{M}_f, \mathbf{M}_b\}$. Adaptive Deformable Cross blocks (ADC) take the concatenation of \mathbf{M}_f and \mathbf{M}_b , along with the enhanced feature \mathbf{F} from HKME. After entering the ADC block, \mathbf{M}_f and \mathbf{M}_b are first split, and then pass through a dual-branch structure that includes a deformable cross attention (DCA) and a cross attention. The results from the dual branches are concatenated and sent into an MLP for nonlinear transformation. After the final block, a fusion layer is applied to integrate the bidirectional motion features to obtain an aggregated motion representation \mathbf{M} .

Deformable cross attention. Our deformable cross attention (DCA) predicts multiple offsets at a single point, rather than predicting offsets at each point of the kernel as in the case of deformable convolution. This endows it with a stronger ability to characterize the relationships between elements and to flexibly handle different scales. The concept of our cross mechanism is realized by incorporating the enhanced feature \mathbf{F} as a constraint to control the generation of offsets in the spatial domain, ensuring that only a subset of motion features are selected as keys and values for attention computation. Specifically, the DCA can be represented by the following formulas:

$$\begin{aligned} q &= W_q \otimes x, \quad z = W_z \otimes \mathbf{F}, \\ \triangle p &= \theta(\psi(q \oplus z)), \quad \bar{x} = \phi(\triangle p + p_z), \end{aligned} \quad (3)$$

$$\text{DCA}(x, z, p_z) = \sum_{n=1}^N \sigma \left(\frac{q \otimes (W_k \otimes \bar{x}_n)^T}{\sqrt{d}} \right) (W_v \otimes \bar{x}_n), \quad (4)$$

where $x, q, \triangle p, p_z, \bar{x}, N$, and d are motion features \mathbf{M}_f or \mathbf{M}_b , query, point offset, reference points from z , sample features, number of sampling points, and embedding dimension, respectively. $\otimes, \oplus, \psi(\cdot), \theta(\cdot), \phi(\cdot)$, and $\sigma(\cdot)$ denote the operations of matrix multiplication, concatenation, convolution, offset generation, bilinear interpolation, softmax, respectively. W_q, W_k, W_v , and W_z are all learnable mapping matrices. The offset $\triangle p$ generated under the constraint of \mathbf{F} , ensures the filtering of spatial regions related to the human joints within the global domain, thereby facilitating adaptive motion cue extraction from motion features.

Heatmap generation. We first split the key frame keypoint tokens \mathbf{T}_t^i from $\hat{\mathbf{T}}_k$ and then transform them into \mathbf{H}_k through an MLP and reshaping. By aggregating \mathbf{M} and \mathbf{F}_t^i and up-sampling, we obtain \mathbf{H}_m . The final pose heatmaps \mathbf{H}_t^i are derived by adding \mathbf{H}_k and \mathbf{H}_m with equal weights.

Loss function. We adopt the established pose heatmap loss \mathcal{L}_H to supervise the final predicted pose heatmaps \mathbf{H}_t^i to converge to the ground truth pose heatmaps \mathbf{G}_t^i :

$$\mathcal{L}_H = \|\mathbf{H}_t^i - \mathbf{G}_t^i\|_2^2. \quad (5)$$

Experiments

Experimental Settings

Datasets. PoseTrack has become a crucial dataset in video-based human pose estimation benchmarks. **PoseTrack2017**^[40] introduces 250 training videos and 50 validation videos, with 80,144 pose annotations across 15 key points. **PoseTrack2018**^[41] expands to 593 training and 170 validation videos, totaling 153,615 annotations. **PoseTrack2021**^[42] further enriches the dataset, particularly improving the representation of smaller figures and crowded scenes, reaching 177,164 pose annotations, with recalibrated joint visibility flags to better address occlusions.

Evaluation metric. To evaluate the efficacy of our proposed model in pose estimation, we calculate the average precision (AP) for each joint and then aggregate these values to obtain the mean average precision (mAP).

Implementation details. Our VREMD framework is realized utilizing PyTorch. For feature extraction on single frames, we adopt the most primitive Vision Transformer (ViT-L) architecture^{[15][43]}, pre-trained on the COCO dataset^[44], as our backbone. The input image size is fixed at 256×192 . We integrate a series of data augmentation techniques, consistent with methodologies employed in previous works^{[13][14]}, comprising random rotation $[-45^\circ, 45^\circ]$, random scale $[0.65, 1.35]$, truncation (half body), and flipping during training. The number of input frames is set to 3, consisting of one key frame accompanied by two auxiliary frames sourced from preceding and succeeding neighbors, respectively. This configuration mirrors that of DCPose^[14], rather than employing the five frame input as seen in TDMI^[20] and FAMI-Pose^[19]. Our model is trained on a single RTX 4090 GPU for 20 epochs with the backbone frozen. We utilize the AdamW optimizer with an initial learning rate of $2e-3$, which is then reduced by a factor of ten at the 16th epoch.

Method	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	Mean
AlphaPose ^[45]	63.9	78.7	77.4	71.0	73.7	73.0	69.7	71.9
TML++ ^[29]	-	-	-	-	-	-	-	74.6
MDPN ^[32]	75.4	81.2	79.0	74.1	72.4	73.0	69.9	75.0
PGPT ^[46]	-	-	-	72.3	-	-	72.2	76.8
Dynamic-GNN ^[34]	80.6	84.5	80.6	74.4	75.0	76.7	71.8	77.9
PoseWarper ^[13]	79.9	86.3	82.4	77.5	79.8	78.8	73.2	79.7
PT-CPN++ ^[47]	82.4	88.8	86.2	79.4	72.0	80.6	76.2	80.9
DCPose ^[14]	84.0	86.6	82.7	78.0	80.4	79.3	73.8	80.9
DetTrack ^[35]	84.9	87.4	84.8	79.2	77.6	79.7	75.3	81.5
FAMI-Pose ^[19]	85.5	87.7	84.2	79.2	81.4	81.1	74.9	82.2
HANet ^[37]	86.1	88.5	84.1	78.7	79.0	80.3	77.4	82.3
M-HANet ^[37]	86.7	88.9	84.6	79.2	79.7	81.3	78.7	82.7
KPM ^[38]	85.1	88.9	86.4	80.7	80.9	81.5	77.0	83.1
DSTA ^[18]	85.9	88.8	85.0	81.1	81.5	83.0	77.4	83.4
TDMI-ST ^[20]	86.7	88.9	85.4	80.6	82.4	82.1	77.6	83.6
VREMD (Ours)	86.7	89.3	85.6	82.1	85.0	83.9	79.3	84.6

Table 2. Comparisons with the state-of-the-art methods for video pose estimation on the validation sets of the PoseTrack2018^[41] dataset.

Method	Head	Shoulder	Elbow	Wrist	Hip	Knee	Ankle	Mean
Tracktor++ w. poses ^[48]	-	-	-	-	-	-	-	71.4
CorrTrack ^[33]	-	-	-	-	-	-	-	72.3
Tracktor++ w. corr. ^[48]	-	-	-	-	-	-	-	73.6
DCPose ^[14]	83.2	84.7	82.3	78.1	80.3	79.2	73.5	80.5
FAMI-Pose ^[19]	83.3	85.4	82.9	78.6	81.3	80.5	75.3	81.2
DSTA ^[18]	87.5	87.0	84.2	81.4	82.3	82.5	77.7	83.5
TDMI-ST ^[20]	86.8	87.4	85.1	81.4	83.8	82.7	78.0	83.8
VREMD (Ours)	87.2	89.1	85.2	82.4	85.1	83.4	79.2	84.5

Table 3. Comparisons with the state-of-the-art methods for video pose estimation on the validation sets of the PoseTrack2021^[42] dataset.

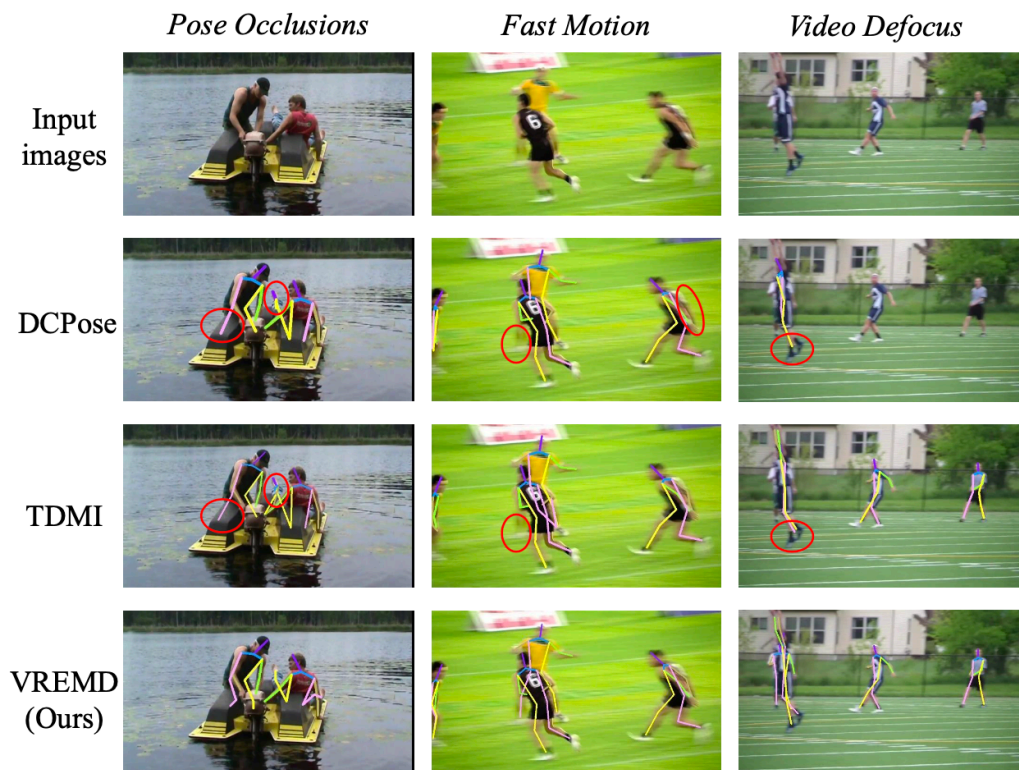


Figure 3. Qualitative comparison of our VREMD, DCPose^[14], and TDMI^[20] on the PoseTrack2017 dataset, featuring challenges such as pose occlusions, fast motion, and video defocus. Red solid circles denote the inaccurate pose predictions.-

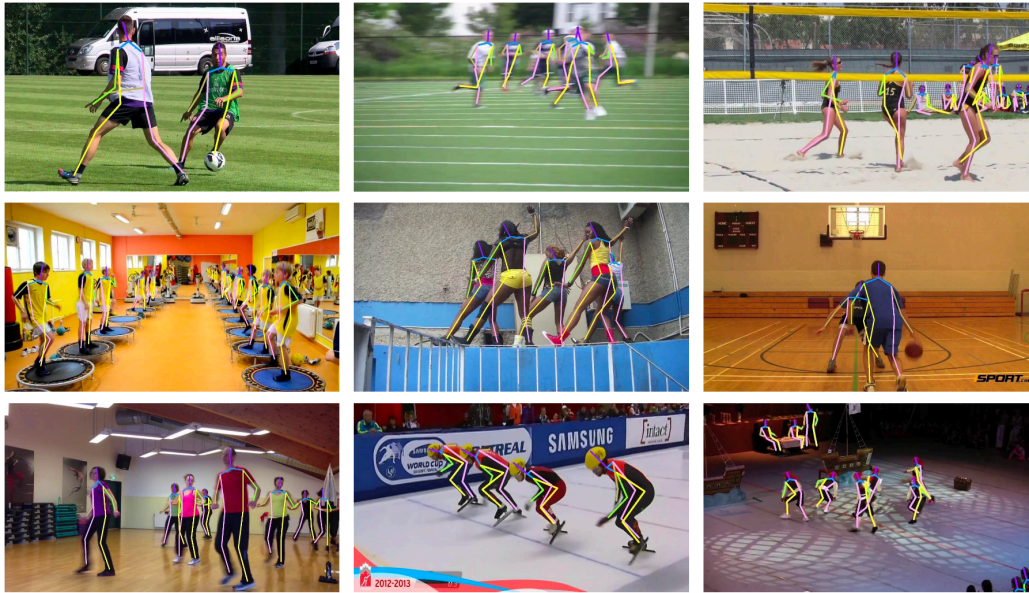


Figure 4. Visual results of our VREMD on PoseTrack datasets. Challenging Scenarios such as fast motion or pose occlusion are involved.

Method	HKME	BMD	Mean
Baseline			80.2
(a)	✓		85.3
(b)		✓	85.6
(c)	✓	✓	87.6

Table 4. Ablation of different components in our VREMD.

Comparison with State-of-the-art Approaches

Results on the PoseTrack2017 Dataset. We first benchmark our method on the PoseTrack2017^[40] dataset. A total of 22 methods are compared and their performances on the PoseTrack2017 validation set are summarized in Table 1. Our proposed VREMD consistently outperforms existing state-of-the-art methods, reaching an mAP of 87.6. Compared to the latest top-performing method TDMI-ST^[20], our VREMD obtain a 1.7 mAP gain. The performance boost for challenging joints

(i.e., wrist, ankle) is also promising: we attain an mAP of 84.8 (\uparrow 3.3) for wrists and an mAP of 81.0 (\uparrow 2.3) for ankles. It is noteworthy that our VREMD operates effectively with fewer input video frames than the most recent works^{[19][20]}, requiring just three frames as opposed to five. These consistent and substantial improvements in effectiveness indicate the importance of reinforcing the positional attributes of visual representations and integrating joint-related motion dynamics. In addition, we present the visualized results, which include a comparison with existing methods, for scenarios involving complex spatio-temporal interactions (*e.g.*, pose occlusion, blur) in Fig 3, demonstrating our method’s robustness. More visualization results are shown in Figure 4.

Results on the PoseTrack2018 Dataset. We further evaluate our VREMD on the PoseTrack2018 dataset, and the detailed validation set results are showcased in Table 2. Once again, as illustrated in this table, our VREMD surpasses all prior state-of-the-art methods, achieving the most exceptional outcomes. We obtain the final performance of 84.6 mAP. The precision for wrists and ankles also shows a noticeable improvement compared to TDMI-ST, scoring 82.1 (\uparrow 1.5) and 79.3 (\uparrow 1.7) respectively.

Results on the PoseTrack2021 Dataset. Performance comparisons of our model and previous state-of-the-art methods on the PoseTrack21 dataset are provided in Table 3. When evaluated on the PoseTrack2021 validation dataset, the results highlight the outstanding performance of our model. Achieving new state-of-the-art results, our model records an overall mAP of 84.5, outperforming TDMI-ST by a margin of 0.7 mAP. Encouragingly, our method yields a 1.0 mAP improvement over the previous best, attaining 82.4 at the wrist, and shows a 1.2 mAP advance, achieving 79.2 at the ankle, which are recognized as difficult joints to accurately predict. These results, once again, underscore the robustness and superiority of our method in this domain.

Method	Human mask	Keypoint mask	Mean
(a)			85.9
(b)	✓		86.5
(c)		✓	86.8
(d)	✓	✓	87.6

Table 5. Ablation of various designs in the HKME module.

Method	DC	DA	DCA (Ours)	BS (Ours)	Mean
(a)	✓				84.7
(b)		✓			85.8
(c)			✓		87.1
(d)			✓	✓	87.6

Table 6. Ablation of various designs in the BMD module.

Ablation Study

We carry out extensive ablation studies centered on assessing the impact of individual components within our VREMD architecture, encompassing the Human-Keypoint Mask Enhancement module (HKME) and the Bidirectional Motion Disentanglement module (BMD). We additionally probe into the efficacy of diverse micro-designs incorporated in each module. All experiments are performed on the PoseTrack2017 validation set.

Study on components of VREMD. We experimentally evaluate the effectiveness of each component in our VREMD framework, detailing the quantitative results in Table 4. Firstly, we establish a baseline for this experiment by coupling a Vision Transformer (ViT) Backbone with a pose detection head. (a) Integrating the Human-Keypoint Mask Enhanced module (HKME) into the baseline yields a substantial gain of 5.1 mAP. This substantial progress indicates that the dual-mask mechanism, offering a coarse-to-fine representation refinement, facilitates improvements in human pose estimation. (b) In the next setup, we exclusively incorporate the Bidirectional Motion Disentanglement module (BMD) into the baseline system. Notably, the Adaptive Deformable Cross (ADC) block, which originally utilized enhanced features from the HKME, now receives backbone output features instead. The outcome achieves an mAP of 85.6, marking an increase of 5.4 mAP. Such a significant boost in performance unequivocally validates the BMD module’s proficiency in adaptively excavating bidirectional temporal information, guiding accurate pose estimation. (c) Finally, we incorporate both the HKME and BMD modules into our framework, attaining a culminating performance of 87.6 mAP, which indicates that the synergy of these two components can lead to further enhancements.

Study on Human-Keypoint Mask Enhanced module. We then investigate the impact of the two mask generation techniques in HKME on overall performance. We conduct four experiments and presented them in Table 5. (a) Generating visual representations using only the spatio-temporal Transformers network. (b) Producing a human mask for coarse filtering of human-related tokens. (c) Calculating a keypoint mask for basic joint token screening. (d) Utilizing dual masks, derived from methods (b) and (c), for the progressive refinement and enhancement of visual tokens, transitioning from coarse to fine detail. This table illustrates that method (a), which does not generate any masks, offers a slight improvement of 0.3 mAP over the setting that removes HKME. Subsequently, applying the human mask alone (b) and the keypoint mask alone (c) achieves respective performances of 86.5 mAP and 86.8 mAP. Although utilizing these masks individually can yield certain accuracy gains, simultaneously employing both for coarse-to-fine representation refinement (d) leads to the optimal results. This promising outcome attests to the superiority of our dual-mask paradigm, which provides a prompt of human joints to the framework, enabling more accurate keypoint localization.

Study on Bidirectional Motion Disentanglement module. Additionally, we explore the influence of our deformable cross attention (DCA) and bidirectional separation strategy. Four experiments are performed and displayed in Table 6. (a) We first replace our Adaptive Deformable Cross (ADC) block with the deformable conv (DC)^[49], as adopted in previous works^{[14][19][20]}. We observe a slight performance decline, that is, a 0.6 mAP decrease. We speculate that the reason might be the feature map obtained through the attention mechanism is more spatially dispersed and structurally diverse, which is incompatible with the local adaptive variation characteristics of deformable conv. (b) We further employ plain deformable attention (DA)^[50] and achieve an 85.8 mAP, which proves that deformable attention is more suitable for our frameworks based on attention mechanisms. (c) We propose a novel deformable cross attention (DCA), which integrates the advantages of adaptive receptive field of deformable attention and selective feature highlighting of cross attention, achieving an 87.1 mAP. (d) Finally, we apply a bidirectional separation (BS) strategy to independently capture bidirectional motion dynamics, resulting in a 0.5 mAP improvement, unlike previous methods that concatenate and jointly process bidirectional motion features. These results strongly demonstrate that our method can more effectively capture task-relevant motion cues to facilitate pose estimation.

Related Work

Image-based human pose estimation. Recent progress in deep learning architectures, as chronicled in^{[15][51]}, coupled with the proliferation of extensive datasets referenced in^{[44][52]}, has catalyzed the development of a multitude of deep learning methodologies. These methodologies, delineated in^[10]^[53] and proposed for the purpose of image-based human pose estimation, predominantly align with two distinct paradigms: bottom-up and top-down. Bottom-up approaches^[54] initiate with the detection of individual body parts in an image and subsequently attempt to aggregate these parts into a comprehensive human pose. The top-down paradigm^{[10][43]} start by detecting the bounding box around the human body and then localize the target human's keypoints within that area. However, these image-based methods struggle when applied to video streams, since they fail to effectively incorporate the temporal changes between frames. our research builds upon previous image-based approaches, extending them with temporal dynamics capture specifically tailored for video pose estimation.

Video-based human pose estimation. In the early stages, substantial approaches involve utilizing optical flow to establish motion-based assumptions^[55]. These approaches commonly generate dense optical flow across frames to improve pose heatmap predictions, yet the technique is computationally demanding and prone to errors when faced with marked deterioration in image quality. Recent methods^[18] have shifted towards attempting to implicitly capture motion evidence from temporal information by employing deformable convolutions. DCPose^[14] and PoseWarper^[13] model and process pose temporal residuals and re-refine keypoint detection via multi-scale deformable convolutions for accurate pose estimation. TDM^[20] introduces a multi-stage framework that encodes temporal differences for dynamic context modeling, leveraging mutual information to uncover useful temporal clues. Contrary to prior approaches that directly execute feature difference learning in the global space, we strive to enhance visual representations through the aggregation of joint positions, and to dissect representative joint-associated motion dynamics for more robust human pose estimation.

Conclusion and Future Work

Conclusion. In this paper, we investigate the video-based human pose estimation task from the perspective of local spatial perception and temporal cues disentanglement. A dual-stream architecture is designed to effectively capture spatio-temporal dependencies by collaboratively executing gradual human joint focus and adaptive motion decoupling. Specifically, we present a Human-Keypoint Mask

Enhanced module that performs a coarse-to-fine selective representation enhancement to assist the framework in exploring human and joint regions. Additionally, we create a Bidirectional Motion Disentanglement module to adaptively mine pose-related motion evidence. Our method significantly and consistently outperforms state-of-the-art performances on three benchmark datasets: PoseTrack2017, PoseTrack2018, and PoseTrack2021.

Limitations and future works. We identified two limitations in our model: (1) The accuracy of our head joint localization is suboptimal. We believe this is due to good spatial separation of joints but imperfect recognition of their relationships, causing interference from nearby joints, such as the shoulder. We plan to address this by incorporating Graph Neural Networks (GNNs) to better capture these interrelationships. (2) When the target person is severely occluded by others, our method may mistakenly incorporate temporal cues from non-target individuals, reducing pose estimation accuracy. We plan to optimize our visual and motion features using clustering techniques to address this issue.

Acknowledgments

This work is supported by the National Natural Science Foundation of China (No. 62372402), and the Key R&D Program of Zhejiang Province (No. 2023C01217).

References

1. [△]Wang D, Zhang S (2022). "Contextual instance decoupling for robust multi-person pose estimation." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11060–11068.
2. [△]Geng Z, Wang C, Wei Y, Liu Z, Li H, Hu H (2023). "Human pose as compositional tokens". *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 660–671.
3. [△]Yang Y, Chen H, Liu Z, Lyu Y, Zhang B, Wu S, Wang Z, Ren K (2023). "Action recognition with multi-stream motion modeling and mutual information maximization". *arXiv preprint arXiv:2306.07576*.
4. [△]Tse THE, De Martini D, Marchegiani L (2019). "No need to scream: Robust sound-based speaker localisation in challenging scenarios." In: *Social Robotics: 11th International Conference, ICSR 2019, Madrid, Spain, November 26--29, 2019, Proceedings 11*. Springer. pp. 176--185.
5. [△]Su P, Liu Z, Wu S, Zhu L, Yin Y, Shen X (2021). "Motion prediction via joint dependency modeling in phase space." In: *Proceedings of the 29th ACM international conference on multimedia*, 713–721.

6. [△]Wu S, Liu Z, Zhang B, Zimmermann R, Ba Z, Zhang X, Ren K (2024). "Do as I Do: Pose Guided Human Motion Copy". *IEEE Transactions on Dependable and Secure Computing*.
7. [△]Liu Z, Wu S, Xu C, Wang X, Zhu L, Wu S, Feng F (2022). "Copy motion from one to another: Fake motion video generation". *arXiv preprint arXiv:2205.01373*. Available from: [arXiv:2205.01373](https://arxiv.org/abs/2205.01373).
8. [△]Wang Y, Mori G (2008). "Multiple tree models for occlusion and spatial constraints in human pose estimation." In: *Computer Vision--ECCV 2008: 10th European Conference on Computer Vision, Marseille, France, October 12-18, 2008, Proceedings, Part III* 10. Springer. pp. 710–724.
9. [△]Sapp B, Toshev A, Taskar B (2010). "Cascaded models for articulated pose estimation." In: *Computer Vision--ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part II* 11. Springer. pp. 406–420.
10. [△]_a, [△]_b, [△]_c, [△]_d Sun K, Xiao B, Liu D, Wang J (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.
11. [△]_a, [△]_b Li Y, Zhang S, Wang Z, Yang S, Yang W, Xia S-T, Zhou E (2021). "Tokenpose: Learning keypoint tokens for human pose estimation". *Proceedings of the IEEE/CVF International conference on computer vision*. 11313–11322.
12. [△]Zhao Y, Xiong Y, Lin D (2018). "Recognize actions by disentangling components of dynamics". *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6566–6575.
13. [△]_a, [△]_b, [△]_c, [△]_d Bertasius G, Feichtenhofer C, Tran D, Shi J, Torresani L (2019). "Learning temporal pose estimation from sparsely-labeled videos". *Advances in neural information processing systems*. 32.
14. [△]_a, [△]_b, [△]_c, [△]_d, [△]_e, [△]_f, [△]_g, [△]_h, [△]_i, [△]_j, [△]_k Liu Z, Chen H, Feng R, Wu S, Ji S, Yang B, Wang X (2021). "Deep dual consecutive network for human pose estimation". *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 525–534.
15. [△]_a, [△]_b, [△]_c, [△]_d Dosovitskiy A, Beyer L, Kolesnikov A, Weissenborn D, Zhai X, Unterthiner T, Dehghani M, Minderer M, Heigold G, Gelly S, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*. 2020.
16. [△]_a, [△]_b Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser Ł, Polosukhin I (2017). "Attention is all you need". *Advances in neural information processing systems*. 30.
17. [△]_a, [△]_b Jin KM, Lee GH, Lee SW. OTPose: occlusion-aware transformer for pose estimation in sparsely-labeled videos. In: *2022 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*. IEEE; 2022. p. 3255–3260.

3.

32. ^aGuo H, Tang T, Luo G, Chen R, Lu Y, Wen L (2018). "Multi-domain pose network for multi-person pose estimation and tracking". In: *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*. 0–0.
33. ^aRafi U, Doering A, Leibe B, Gall J (2020). "Self-supervised keypoint correspondences for multi-person pose estimation and tracking in videos." In: *Computer Vision--ECCV 2020: 16th European Conference, Glasgow, UK, August 23--28, 2020, Proceedings, Part XX 16*. Springer; 2020. p. 36--52.
34. ^aYang Y, Ren Z, Li H, Zhou C, Wang X, Hua G (2021). "Learning dynamics via graph neural networks for human pose estimation and tracking". In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8074--8084.
35. ^aWang M, Tighe J, Modolo D (2020). "Combining detection and tracking for human pose estimation in videos." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11088--11096.
36. ^ΔGai D, Feng R, Min W, Yang X, Su P, Wang Q, Han Q (2023). "Spatiotemporal learning transformer for video-based human pose estimation". *IEEE Transactions on Circuits and Systems for Video Technology*. 33(9): 4564--4576.
37. ^a^ΔJin KM, Lim BS, Lee GH, Kang TK, Lee SW (2023). "Kinematic-aware hierarchical attention network for human pose estimation in videos." In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 5725--5734.
38. ^a^ΔFu Z, Zuo W, Hu Z, Liu Q, Wang Y (2023). "Improving Multi-Person Pose Tracking with A Confidence Network". *IEEE Transactions on Multimedia*.
39. ^ΔJin KM, Lee GH, Nam WJ, Kang TK, Kim HW, Lee SW (2024). "Masked Kinematic Continuity-aware Hierarchical Attention Network for pose estimation in videos". *Neural Networks*. 169: 282--292.
40. ^a^ΔIqbal U, Milan A, Gall J (2017). "Posetrack: Joint multi-person pose estimation and tracking." In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2011--2020.
41. ^a^ΔAndriluka M, Iqbal U, Insafutdinov E, Pishchulin L, Milan A, Gall J, Schiele B (2018). "Posetrack: A benchmark for human pose estimation and tracking." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 5167--5176.
42. ^a^ΔDoering A, Chen D, Zhang S, Schiele B, Gall J (2022). "Posetrack21: A dataset for person search, multi-object tracking and multi-person pose tracking". *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20963--20972.

43. ^a Xu Y, Zhang J, Zhang Q, Tao D (2022). "Vitpose: Simple vision transformer baselines for human pose estimation". *Advances in Neural Information Processing Systems*. 35: 38571–38584.
44. ^a Lin T-Y, Maire M, Belongie S, Hays J, Perona P, Ramanan D, Dollár P, Zitnick CL. Microsoft coco: Common objects in context. In: *Computer Vision--ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V* 13. Springer; 2014. p. 740-755.
45. ^Δ Fang HS, Xie S, Tai YW, Lu C (2017). "Rmpe: Regional multi-person pose estimation." In: *Proceedings of the IEEE international conference on computer vision*. pp. 2334–2343.
46. ^Δ Bao Q, Liu W, Cheng Y, Zhou B, Mei T (2020). "Pose-guided tracking-by-detection: Robust multi-person pose tracking". *IEEE Transactions on Multimedia*. 23: 161–175.
47. ^Δ Yu D, Su K, Sun J, Wang C (2018). "Multi-person pose estimation for pose tracking with enhanced cascaded pyramid network". In: *Proceedings of the european conference on computer vision (ECCV) Workshops*. 0–0.
48. ^a Bergmann P, Meinhardt T, Leal-Taixe L (2019). "Tracking without bells and whistles". *Proceedings of the IEEE/CVF international conference on computer vision*. 941–951.
49. ^Δ Dai J, Qi H, Xiong Y, Li Y, Zhang G, Hu H, Wei Y (2017). "Deformable convolutional networks." *Proceedings of the IEEE international conference on computer vision*. 764–773.
50. ^Δ Zhu X, Su W, Lu L, Li B, Wang X, Dai J (2020). "Deformable detr: Deformable transformers for end-to-end object detection". *arXiv preprint arXiv:2010.04159*.
51. ^Δ He K, Zhang X, Ren S, Sun J (2016). "Deep residual learning for image recognition." In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778.
52. ^Δ Shuai C, Zhong J, Wu S, Lin F, Wang Z, Ba Z, Liu Z, Cavallaro L, Ren K (2023). "Locate and verify: A two-stream network for improved deepfake detection." In: *Proceedings of the 31st ACM International Conference on Multimedia*. pp. 7131–7142.
53. ^Δ Wei SE, Ramakrishna V, Kanade T, Sheikh Y (2016). "Convolutional pose machines". *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*. 4724–4732.
54. ^Δ Cao Z, Simon T, Wei S-E, Sheikh Y (2017). "Realtime multi-person 2d pose estimation using part affinity fields". *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7291–7299.
55. ^Δ Pfister T, Charles J, Zisserman A (2015). "Flowing convnets for human pose estimation in videos." In: *Proceedings of the IEEE international conference on computer vision*. pp. 1913–1921.

Declarations

Funding: This work is supported by the National Natural Science Foundation of China (No. 62372402), and the Key R&D Program of Zhejiang Province (No. 2023C01217).

Potential competing interests: No potential competing interests to declare.