# An End-to-End Framework for Video Multi-Person Pose Estimation

Zhihong Wei

University of Science and Technology of China

weizh588@mail.ustc.edu.cn

## Abstract

*Video-based human pose estimation models aim to address scenarios that cannot be effectively solved by static image models such as motion blur, out-of-focus and occlusion. Most existing approaches consist of two stages: detecting human instances in each image frame and then using a temporal model for single-person pose estimation. This approach separates the spatial and temporal dimensions and cannot capture the global spatio-temporal context between spatial instances for end-to-end optimization. In addition, it relies on separate detectors and complex post-processing such as RoI cropping and NMS, which reduces the inference efficiency of the video scene. To address the above problems, we propose VEPE (Video End-to-End Pose Estimation), a simple and flexible framework for end-to-end pose estimation in video. The framework utilizes three crucial spatio-temporal Transformer components: the Spatio-Temporal Pose Encoder (STPE), the Spatio-Temporal Deformable Memory Encoder (STDME), and the Spatio-Temporal Pose Decoder (STPD). These components are designed to effectively utilize temporal context for optimizing human body pose estimation. Furthermore, to reduce the mismatch problem during the cross-frame pose query matching process, we propose an instance consistency mechanism, which aims to enhance the consistency and discrepancy of the cross-frame instance query and realize the instance tracking function, which in turn accurately guides the pose query to perform cross-frame matching. Extensive experiments on the Posetrack dataset show that our approach outperforms most two-stage models and improves inference efficiency by 300%.*

## 1. Introduction

Human pose estimation is an important problem in computer vision, which aims to localize the position of keypoints (e.g., knees, ankles, etc.) or body parts. Currently, pose estimation finds extensive applications across various fields, such as action recognition [44], augmented reality, virtual humans, as well as surveillance and tracking



(a) Multi-stage framework of multi-person pose estimation in video



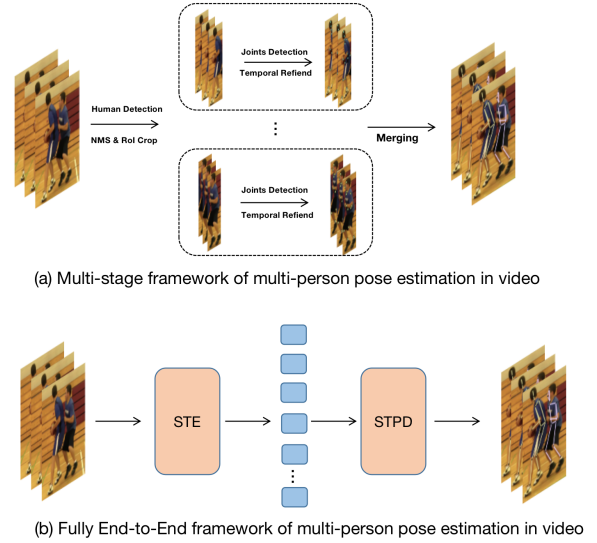(b) Fully End-to-End framework of multi-person pose estimation in video

Figure 1. Comparison of multi-stage and end-to-end frameworks. (a) Most of the existing methods are based on this framework: Multi-person pose estimation can be transformed into single-person pose estimation through the use of a human detector and a set of manual components and finally, all human instances are merged. (b) The end-to-end approach directly views the whole process as a sequence-to-sequence task, where the multi-person pose estimation task is accomplished by feature learning through an encoder and decoding through a decoder.

[23].Early approaches include techniques based on contouring, template matching, feature point detection, and statistical modeling[35, 37, 49], which typically rely on traditional computer graphics and image processing techniques and have limited adaptability to complex backgrounds and multi-pose situations. Recently, with the development of deep learning, a number of works[7, 10, 15, 18, 20, 24, 33, 36, 38, 39, 43, 46] have achieved excellent performance, advancing the development of human pose estimation.

Unfortunately, existing methods designed for static images fall short when it comes to dynamic video images. These methods lack consideration for cues between video frames and disregard the crucial temporal and geometric

consistency among frames. As a result, they fail to establish the dependency of human pose cues in temporal sequences, leading to significant performance degradation in scenarios involving motion blur, video out-of-focus, and pose occlusion. Consequently, effectively harnessing temporal sequence information in videos is essential for advancing human pose estimation algorithms in real-life scenarios.

To tackle this issue, various studies [2, 23, 37] suggest integrating shared sequence characteristics from adjacent frames (support frames). For instance, [23] utilizes a convolutional LSTM to capture both spatial and temporal features, ultimately predicting the pose sequence within videos.[37] proposed a 3DHRNet that utilizes 3D convolution to extract spatio-temporal features of video trajectories to estimate pose sequences. Additionally, some methods utilize optical flow or motion estimation to refine pose estimation for the current keyframe [21, 25, 30]. For instance, [25, 30] compute dense optical flow between frames and employ flow-based motion fields to temporally refine pose heatmaps. It is of concern that, optical flow estimation is computationally expensive and prone to fragility in the presence of severe image quality degradation. In another work[21], pose heatmaps of consecutive frames are aggregated, and motion residuals are modelled to enhance pose estimation for keyframes, leveraging implicit or explicit motion priors. Within these frameworks, attention is typically given to all visual motion cues, which often introduces unnecessary motion details, including irrelevant elements like surrounding people or the background. Moreover, most of the above methods use a two-stage detection strategy, where each human instance is first detected and then the instance region in each image frame is cropped, followed by feeding it into a temporal model for prediction. This approach separates the spatial and temporal dimensions and cannot capture the global spatio-temporal context between spatial instances for end-to-end optimization. Furthermore, the computational cost of the two-stage model increases with the number of instances in the image due to the need for additional human detectors and complex post-processing(such as RoI cropping and NMS), leading to less efficient inference, especially in dense scenes. In recent years, end-to-end models have attracted the attention of a wide range of researchers due to their efficient and concise architectures. Compared with traditional two-stage approaches, end-to-end models regress pixel-level poses directly from images without the need for predefined human detectors and complex artificial components, which significantly improves the inference speed and scalability of the models. A comparison of the multi-stage model and the full end-to-end framework structure is shown in Figure.1.

In this paper, we extend PETR[29] for video tasks and propose a simple video-based end-to-end pose estimation framework, which naturally treats pose instance objects in the video as a sequence-to-sequence task, and the whole process can be realized with full end-to-end training. The overall construction adopts the Encoder-Decoder structure, which does not require any hand-designed components and greatly improves inference efficiency. The Decoder decoded pose query contains only human pose instances, which naturally filters out background and other interfering factors, and the use of a Transformer to articulate the multi-frame pose query can be naturally generalized to the global context learning at a distance. Our approach incorporates three essential temporal components: the Spatio-Temporal Pose Encoder (STPE), the Spatio-Temporal Deformable Memory Encoder (STDME), and the Spatio-Temporal Pose Decoder (STPD). The Spatio-Temporal Pose Encoder focuses on establishing correlations between spatial pose queries across frames to learn the temporal context of pose instances. Meanwhile, the Spatio-Temporal Deformable Memory Encoder is responsible for extracting cross-frame spatio-temporal visual features for instances using the multi-scale spatial features generated by the Spatial Encoder. Additionally, the Spatio-Temporal Deformable Memory Encoder plays a crucial role in providing valuable human appearance and location cues to the Spatio-Temporal Decoder. To refine the pose results, the spatio-temporal information obtained from the STPE and STDME is decoded in a multi-level cascade using the Spatio-Temporal Pose Decoder (STPD). This decoder further enhances the accuracy and precision of the pose estimation. Furthermore, to reduce the mismatch problem during the cross-frame pose query matching process, we propose an instance consistency mechanism, which aims to enhance the consistency and discrepancy of cross-frame instance query, realize the instance tracking function, and to guide the pose query for cross-frame matching, which improves the accuracy of pose instance identification and tracking of Temporal Pose Encoder.

In summary, our contributions can be summarized as follows:

- This paper propose a simple and flexible end-to-end pose estimation framework based on video. The framework employs three key spatio-temporal Transformer components: the Spatio-Temporal Pose Encoder (STPE), Spatio-Temporal Deformable Memory Encoder (STDME), and Spatio-Temporal Pose Decoder (STPD), to efficiently utilize temporal information for optimizing human poses.
- To reduce the mismatch problem during the cross-frame pose query matching process, this paper proposes an instance consistency mechanism, which aims to enhance the consistency and discrepancy of the cross-frame instance query and realize the instance tracking function, which in turn guides the pose query to perform accurate cross-frame matching.

- Extensive experiments on the Posetrack dataset show that VEPE outperforms most two-stage models and improves inference efficiency by up to 300%.

## 2. Related Work

### 2.1. End-to-end human pose estimation model

Current end-to-end multi-person pose estimation[19, 29, 45] frameworks are built by following the designs of DETR[6] and its variants [50]. PETR[29]views pose estimation as a hierarchical set prediction problem and propose the first fully end-to-end pose estimation framework with the advent of DETR. The process begins by predicting a series of human poses through a pose decoder, after which each pose's keypoints are further refined via a joint (keypoint) decoder. QueryPose[41] adopts a method similar to Sparse R-CNN[32], developing two parallel decoders based on RoIAlign, each dedicated to human detection and pose estimation. ED-POSE[45] reinterprets the multi-person pose estimation task as two distinct processes of explicit box detection.It unifies the global person and local keypoints into the same box representation, and they can be optimized by the consistent regression loss in a fully end-to-end manner. GroupPose[19] introduces a straightforward approach for multi-person pose estimation in an end-to-end manner. By treating each keypoint as an object, it leverages N × K keypoint queries to predict the N × K keypoint positions. Additionally, the method utilizes N instance queries, with each query representing a pose consisting of K keypoints, to provide scoring for the predicted K-keypoint pose. Subsequently, the traditional self-attention is replaced with two consecutive grouped self-attention for performing intra-instance and cross-instance interactions, respectively.

### 2.2. Video-based human pose estimation model

Current image-based methods struggle to adapt effectively to video streams due to inherent challenges in grasping the temporal dynamics across sequential frames.A direct method for cross-frame modelling and utilizing temporal, context is to use the convolutional LSTM proposed in [2, 23]. One primary limitation of this model is its tendency to misalign features across varying frameworks, which consequently diminishes the efficacy of those supporting frameworks. Several studies [25, 30] utilize optical flow to incorporate motion priors. Generally, these methods calculate dense optical flow between frames and use these motion cues to enhance the accuracy of the predicted pose heatmaps. Nevertheless, estimating optical flow is resource-intensive and can be sensitive to significant declines in image quality. Another line of research [5, 21, 37] explores implicit motion compensation by applying deformable convolutions or 3D CNNs. For instance, [5, 21, 22] suggest modeling joint movements at multiple granularities based on heatmap residuals and using techniques like pose resampling or pose warping through deformable convolutions. According to [21], motion offsets are calculated between the keyframe and additional frames, which serve as a foundation for resampling pose heatmaps across successive frames. In both scenarios, the precision of pose estimation is significantly influenced by the effectiveness of optical flow or motion offset calculations. All of the above methods are based on a two-stage approach that requires pre-training of the human detector. The prediction quality of human instances depends on the human detector. The model cannot be trained end-to-end, in which case the lack of effective supervision of the intermediate feature layer may lead to inaccurate pose estimation or even error propagation. In addition, two-stage-based models require the use of artificial components such as NMS, and ROI components, which will greatly increase the burden of inference efficiency.

## 3. Method

### 3.1. Review of PETR

Inspired by the paradigm emerging from object detection [50], Shi et al.[29] have proposed for the first time a fully end-to-end Transformer-based multi-person pose estimation framework, called PETR. The proposed method reduces the problem of pose estimation to a hierarchical ensemble prediction problem by transforming pose estimation into a hierarchical ensemble prediction problem that combines the human body instances and the keypoints in a unified manner. Firstly, a multi-scale feature map is extracted through a mainstream backbone network, and a visual feature encoder takes the spread image features as input and refines them to obtain a refined multi-scale visual feature memory. Next, it operates through multiple randomly initialized trainable pose queries, and the pose decoder is responsible for learning and inferring interactions between objects while evaluating the pose of the instances in the overall image environment. Moreover, the query-based framework is trained utilizing a bipartite matching technique, thereby eliminating heuristic label assignments and removing the necessity for NMS-based post-processing.

### 3.2. VEPE Framework

The overall flow of VEPE is shown in Fig.2. The goal of the framework is to temporally learn static pose results and output fine-tuned, high-quality pose results.

The framework includes the following two parts: the spatial component and the temporal component. The spatial component uses PETR as a spatial feature extractor to extract spatial visual feature memories as well as spatial pose queries through its internal Spatial Encoder and Decoder. The spatio-temporal components include the Spatio-Temporal Pose Encoder (STPE) for aligning and refin-
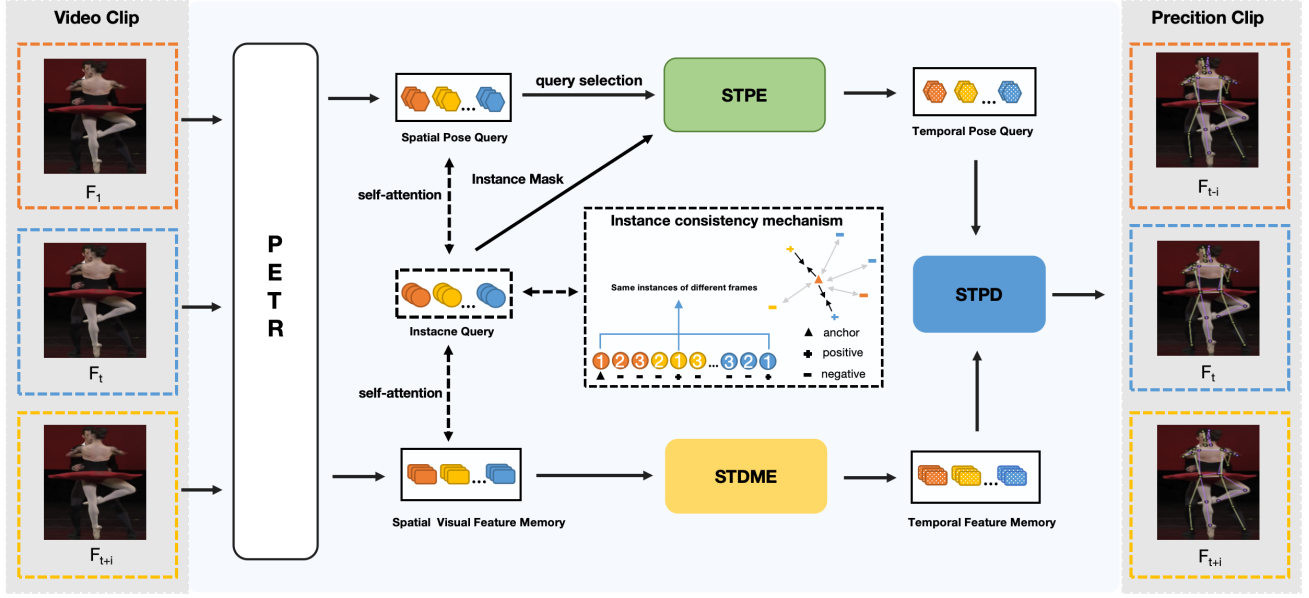
Figure 2. VEPE model pipeline. The VEPE model pipeline consists of two phases: spatial and temporal. In the spatial phase: spatial visual feature memories and spatial pose queries are obtained by the spatial feature encoder and spatial pose decoder of the PETR model, respectively. In the temporal phase: Initiated with spatial pose queries from PETR first applies query selection to obtain high quality spatial pose queries. Subsequently, high-quality spatio-temporal pose queries are generated through a spatio-temporal learning process with the help of Spatio-Temporal Pose Encoder (STPE) and instance tracking of instance query. Concurrently, the Spatio-Temporal Deformable Memory Encoder (STDME) processes the spatial visual feature memories, temporally encoding them to produce temporal feature memories. Finally, the Spatio-Temporal Pose Decoder (STPD) outputs the refined pose results by combining the spatio-temporal cues from STPE and STDME. In addition, an instance consistency mechanism is introduced, which aims to optimize the consistency and differentiation of feature representations between instance queries across frames.

ing spatial pose queries, the Spatio-Temporal Deformable Memory Encoder (STDME) for fusing spatial visual feature memories, and the Spatio-Temporal Pose Decoder (STPD) for extracting spatio-temporal pose results. In addition, the instance consistency mechanism aims to optimize the consistency and uniqueness of cross-frame instance query feature representations for instance tracking, which ultimately leads to accurate cross-frame matching of pose queries during temporal encoding.

### 3.3. Spatio-Temporal Pose Encoder

The learnable pose query represents a high-level semantic keypoint for capturing human instances. By training this query, we can filter out distractions such as backgrounds, resulting in a query that focuses only on specific instance information. This learnable query facilitates direct temporal interaction with instance objects during cross-frame learning. To facilitate this process, we propose a simple but powerful encoder to facilitate the interaction between cross-frame pose instances. This encoder aligns their feature spaces to convey and interact with information about the same instances in different frames. its idea is to es-

tablish correlations between pose queries across the frame space so as to understand the temporal context of the pose instances. We refer to this module as the Spatio-Temporal Pose Encoder (STPE).

The design of the Spatio-Temporal Pose Encoder (STPE) module is shown in Figure.3 (a), where the STPE includes a self-attention layer, a cross-attention layer, and a Feed-forward Neural Network (FFN). To simplify the explanation of the process in the figure, we designate one frame as the keyframe and the remaining frames as reference frames. Specifically, firstly, the pose query of the keyframe learns the relationship between the instances in the frame through the self-attention layer, and then interacts with the spatial pose object query extracted from the reference frame through the cross-attention layer to compute the common attention between the reference query and the features of the key query, and gradually extracts the useful information to be aggregated to the pose query of the keyframe. In the cross-attention mechanism, it's noteworthy that we use an instance mask to isolate interference from unrelated instances to ensure accuracy. This instance mask is derived from similarity calculations among instance queries,

effectively filtering out the most matched targets for each instance in every frame. This process ensures that our algorithm is not impacted by irrelevant instances, which helps reduce unnecessary noise and enhances the overall accuracy and reliability of the model, particularly in complex and highly dynamic scenarios.

In dense scenes, where human bodies often share similar static features, pose queries tend to exhibit significant overlap and resemblance. To tackle this redundancy and leverage the spatio-temporal pose encoder (STPE) for effective learning, we've introduced the Pose Query Selection (PQS) technique. In PQS, we employ a trained instance pose evaluator during the spatial phase to streamline the process. By aggregating spatial pose queries from all frames, the evaluator can then select high-quality pose queries based on specific thresholds. This approach ensures that only the most informative and non-redundant queries are retained, reducing the impact of low-quality pose instances and lowering computational complexity.
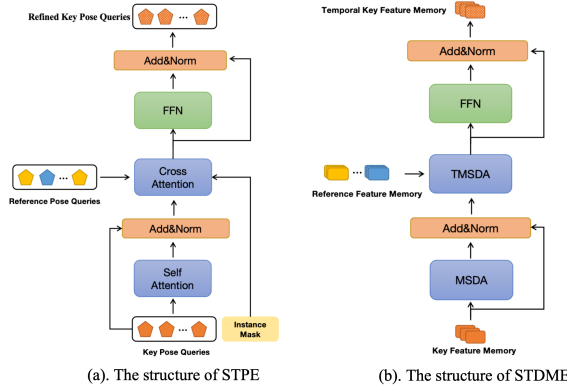


(a). The structure of STPE      (b). The structure of STDME

Figure 3. The structure of (a). Spatio-Temporal Pose Encoder (STPE) and (b).Spatio-Temporal Deformable Memory Encoder (STDME). **STPE**: It takes a keyframe pose query and a reference frame pose query as input and outputs an updated keyframe pose query by aggregating temporal features. **STDME**: It takes a multi-scale visual feature memory for keyframes and a multi-scale visual feature memory for reference frames as input and outputs an updated multi-scale visual feature memory through temporal feature aggregation.

## 3.4. Spatio-Temporal Deformable Memory Encoder

The purpose of the Spatio-Temporal Deformable Memory Encoder (STDME) is to encode spatio-temporal features from the multi-scale spatial feature memory of the Spatial Encoder (SE) and provide human appearance and position cues to the Spatio-Temporal Pose Decoder. The ultimate goal is to temporally aggregate visual cues of human pose across frames. Using a plain Transformer encoder[34] directly can lead to excessive computations due to the sim-

ilarity of neighbouring features containing similar appearance and background information. Additionally, irrelevant information like the human body background can disrupt the capture of spatio-temporal information. In contrast, deformable attention[50] samples only a portion of the information based on the learned offset field. The key idea is to focus on a small portion of the surrounding key sampling points and compute similarity for interaction, reducing the computational load of attentional similarity computation and mitigating interference from redundant information.

Multi-scale visual features play an important role in performing human pose estimation tasks. Using multi-scale visual features, pose information can be captured at different scales, including whole-body poses and localized poses, which can provide rich contextual information and help to understand the relationships among human body parts. Therefore, by utilizing extended temporal multi-scale deformable attention (TMSDA), we can link these multi-scale spatial features in the temporal dimension through this manipulation to provide more cross-frame multi-scale visual information for the visual features of the keyframes, which in turn provides additional visual cues for situations such as motion blur and keypoint occlusion.

To simplify the explanation of the process in the figure, we designate one frame as the keyframe and the remaining frames as reference frames. STDME receives as input the multi-scale feature memories of the reference and keyframes and outputs the multi-scale feature memories of the keyframes of the temporal sequence. The specific flow is represented in Figure.3 (b). The Multi-Scale Deformable Attention (MSDA) captures the contextual relationships among multi-scale visual feature tokens within frames. The Temporal Multi-Scale Deformable Attention (TMSDA) aligns the multi-scale visual feature memory of keyframes and reference frame pose queries, aggregating the most salient temporal information to the keyframe using a partial sampling strategy. Where Temporal multi-scale deformable attention(TMSDA) is realized as shown in the following equation:

$$\text{TMSDA}(z_q, \hat{p}_q, \mathbf{X}) = \sum_{m=1}^{M} W_m \Big[ \sum_{t=1}^{T} \sum_{l=1}^{L} \sum_{k=1}^{K} A_{mtlqk} \cdot W'_m x^{lt}(\phi_l(\hat{p}_q) + \Delta p_{mtlqk}) \Big]. \quad (1)$$

where m indexes the attention head, $t$ indexes the frame sampled from the same video clip, and $l$ indexes multi-scale feature maps from the backbone, and $k$ indexes the sampling points, and $\Delta p_{mtlqk}$ and $A_{mtlqk}$ indicate the sampling offset and attention weights of the $k^{th}$ sampling point in the $l^{th}$ feature map of $t^{th}$ frame and the $m^{th}$ attention head, respectively. Each reference point, represented by normal-

ized coordinates $\hat{p}_q \in [0,1]^2$, is rescaled using $\phi l$ to enable sampling across feature maps $l$ that vary in resolution. The scalar attention weight, denoted by $A_{mtlqk}$, is normalized such that $\sum_{t=1}^{T} \sum_{l=1}^{L} \sum_{k=1}^{K} A_{mtlqk} = 1$. The temporal multi-scale deformable attention samples $TLK$ points from $TL$ feature maps instead of $K$ points from single-frame feature maps.

### 3.5. Spatio-Temporal Pose Decoder

In contrast to the spatial pose decoder, the goal of the Spatio-Temporal Pose Decoder (STPD) is to take the temporal feature memory of Spatio-Temporal Deformable Memory Encoder and the temporal pose query of Spatio-Temporal Pose Encoder as inputs, and then output the decoded results of the pose query of frames that are updated by the temporal cues.

The implementation is shown in Figure.4. To simplify the explanation of the process in the figure, we designate one frame as the keyframe. First, temporal key pose queries from STPE are fed into Self Pose-Pose Attention to capture the relationship between N instances in the frame, which can provide context information between these objects. After Self Pose-Pose Attention, the output query can be represented as follows:

$$\hat{\mathcal{Q}}_{pose}^t = MHA(\bar{\mathcal{Q}}_{pose}^t, \bar{\mathcal{Q}}_{pose}^t, \bar{\mathcal{Q}}_{pose}^t) + \bar{\mathcal{Q}}_{pose}^t \qquad (2)$$

where HMA stands for the multi-head attention mechanism and $t$ stands for the keyframe.

Then the output $\hat{\mathcal{Q}}_{pose}^t$ as Query interacts with the temporal key multi-scale feature memory $\bar{\mathcal{F}}^t$ from STDME with the cross-attention mechanism to extract the temporal image features and aggregate them on $\tilde{\mathcal{Q}}_{pose}^t$, and outputs $\tilde{\mathcal{Q}}_{pose}^t$. The process can be expressed as:

$$\tilde{\mathcal{Q}}_{pose}^t = MSDA(\hat{\mathcal{Q}}_{pose}^t, \bar{\mathcal{F}}^t, \bar{\mathcal{F}}^t) + \hat{\mathcal{Q}}_{pose}^t \qquad (3)$$

where MSDA is multi-scale deformable attention, and $\bar{\mathcal{F}}^t$ is temporal multi-scale feature memroy.

It is worth noting that the cross-attention uses a multi-scale deformable attention module, which we name Deformable Feature-to-Pose Attention. After $\tilde{\mathcal{Q}}_{pose}^t$ passes through the cross-attention mechanism, it passes through the FFN network, which performs nonlinear transformations and pattern extractions, thus allowing the model to perform deeper representation learning. Finally, the classification head predicts the pose confidence score for each query by a linear projection layer.

In the STPD approach, three decoder layers are utilized in a step-by-step manner. Unlike the conventional method of using only the final decoder layer to determine pose coordinates, our model employs all decoder layers for gradual refinement of these coordinates. Specifically, each successive layer reinterprets the pose according to the predictions

generated by the preceding layer. Formally, if $Q_{d-1}$ represents the normalized pose predicted by the $(d-1)^{th}$ decoder layer, then the $d^{th}$ decoder layer further refines this pose as follows:

$$\mathcal{Q}_d = \sigma(\sigma^{-1}(\mathcal{Q}_{d-1}) + \Delta\mathcal{Q}_p) \qquad (4)$$

In the $d^{th}$ layer, the offsets are predicted as $\Delta\mathcal{Q}_p$, where $\sigma$ represents the sigmoid function, and $\sigma^{-1}$ denotes its inverse, the inverse sigmoid function. This progressive output prediction with coarseness to fineness can effectively resolve regressions of fine coordinates to prevent misalignment.
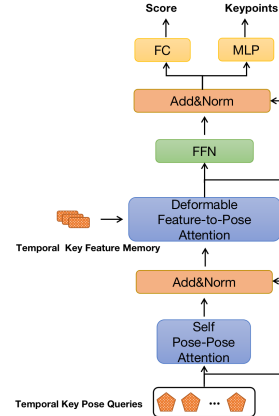


Figure 4. Detailed Architecture of STPD. STPD employs a self-attention module to capture the relationship between spatial objects and cross-attention to capture the interaction between temporal pose query and temporal feature memory. A deformable cross-attention module is designed to focus on visual features closely associated with the target keypoints, and, finally, two separate heads decode the confidence scores and keypoint coordinates, respectively.

### 3.6. Instance Consistency Mechanism

In scenes where subtle human differences in appearance exist, matching pose instances across frames can be challenging due to the highly similar features among different individuals. This similarity can cause errors in matching, which hinders the learning efficiency of the spatio-temporal pose encoder.

To address this issue, we introduced a new learnable embedding (instance query) to guide pose instances in complex environments for cross-frame matching, thereby improving the performance of the Spatio-Temporal Pose Encoder. Specifically, we created learnable embeddings with the same quantity and feature dimension as the pose queries, where each instance query is linked one-to-one with its corresponding pose query. Using a self-attention mechanism, the instance query interacts with pose queries and

visual features to extract and capture each instance's most distinguishing visual characteristics, achieving cross-frame instance tracking through similarity matching. Ultimately, the instance query generates an instance mask by calculating similarity, effectively filtering out irrelevant instances to identify the target that best matches across each frame. This process enhances cross-frame pose matching accuracy and significantly improves the learning efficacy of the Spatio-Temporal Pose Encoder.

At the same time, we introduce the Instance Consistency Mechanism (ICM), designed to mitigate mismatches among instances through the implementation of Instance Consistency Loss (ICL). The essence of this mechanism is to ensure that instances of the same individual across different frames maintain consistent feature representations. Simultaneously, it aims to increase the spatial feature distance between instances of different individuals, thereby enhancing their distinguishability. As depicted in Fig.2, the Instance Consistency Mechanism is crucial for differentiating between instances. Our method begins with the Hungarian algorithm[17] to establish a one-to-one correspondence between the predicted results and the ground truth (GT) data, thus forming a set of candidate results. We use trackId information to identify instance query of the same individual across multiple frames as positive samples, and instance query of other individuals as negative samples, using a designated instance as the anchor. Following this, the Instance Consistency Loss (ICL) is introduced to supervise the learning within this candidate set. The formula for instance consistency loss, denoted as

$$L_{ic} = \sum_{i=1}^{N} \left[ \max\left(0, d(a_i, p_i) - d(a_i, n_i) + margin\right) \right] \quad (5)$$

Where $N$ is the number of instances in the candidate set, $a_i$ denotes the anchor, i.e., the $i^{th}$ instance, $p_i$ denotes the positive sample, i.e., the same individual's pose instance object in other frames, and $n_i$ denotes the negative sample, i.e., the different individual's pose instance object. $d(x, y)$ denotes the distance between $x$ and $y$, and cosine similarity is used as the metric distance here. The parameter Margin plays a crucial role, as it regulates the model's ability to effectively discriminate between positive and negative sample pairs. In this way, the instance consistency loss aims to optimize the model performance by ensuring more accurate identification and tracking of the same instance in consecutive frames and reducing the number of mismatch events.

Finally, with the above technique, we successfully realize instance query with high differentiation and one-to-one correspondence with pose query.

# 4. Experiments

## 4.1. Datasets

The PoseTrack dataset represents an extensive collection for human pose tracking and estimation across video sequences. It captures demanding scenarios, incorporating intricate motions and densely populated scenes. In the 2017 version of PoseTrack[14], there are 514 video segments, containing a total of 16,219 pose labels, split into 250 clips for training, 50 for validation, and 214 for testing. Later, PoseTrack2018[1] augmented this dataset, amassing 1,138 video segments and 153,615 pose annotations, allocated into 593 training clips, 170 validation clips, and 375 testing clips. Both datasets annotate 15 distinct joints, along with additional tags indicating joint visibility. PoseTrack21[9] extends the annotations, particularly for small individuals and crowded scenes, with 177,164 pose annotations. In our model evaluation, we focus on visible joints and use average precision (AP).

## 4.2. Experimental setup

### 4.2.1. Training details.

Based on prior research, we adopt HRNet-W48[31] as the backbone network in our approach. We first performed pre-training on the MS COCO dataset, and the second phase performed 20 epochs of temporal phase training on the Posetrack dataset. Similar to PETR[29], we employ the AdamW optimizer with an initial learning rate of $2 \times 10^{-4}$ for Transformers and weight decay is set to $10^{-4}$. We initialize the number of pose queries as 100. Throughout training, we use a batch size of 8, and to ensure comparability with prior models, we maintain a fixed number of 3 frames in all experiments. Our training and testing processes are conducted on 8 Tesla V100 GPUs.

### 4.2.2. Testing details.

The input images are adjusted to ensure their shortest edges measure 800, while their longest edges do not exceed 1333. The evaluation time is recorded with a single NVIDIA Tesla V100 GPU.

### 4.2.3. details on inference speed

The data is presented in Table.1 illustrates the average inference durations across the PoseTrack2017 dataset for various models, calculated per frame. On average, each frame of the PoseTrack2017 dataset contains six human instances. The inference times for all evaluated methods were tested using a V100 GPU. Notably, the top-down model employs the YOLOv3[28] model as its human detector.

## 4.3. Comparison with State-of-the-art Approaches

### 4.3.1. Results on the PoseTrack2017 Dataset.

We initially test our model on the PoseTrack2017 validation dataset. In total, 15 different methods are com-

Table 1. Quantitative outcomes on the PoseTrack2017 validation set. Note that all of the above except VEPE are multi-stage models. The time indicated represents the model's inference duration across the entire PoseTrack dataset, averaged per frame. The inference times for all methods were tested on a V100, with the top-down models utilizing the YOLOv3[28] model as the human detector.

| Method | Shoulder | Head | Elbow | Wrist | Hip | Ankle | Knee | **Mean** | Time [ms] |
|---|---|---|---|---|---|---|---|---|---|
| PoseTracker [11] | 70.2 | 67.5 | 62.0 | 51.7 | 60.7 | 49.8 | 58.7 | 60.6 | - |
| PoseFlow [42] | 73.3 | 66.7 | 68.3 | 61.1 | 67.5 | 61.3 | 67.0 | 66.5 | - |
| JointFlow [8] | - | - | - | - | - | - | - | 69.3 | - |
| FastPose [48] | 80.3 | 80.0 | 69.5 | 59.1 | 71.4 | 59.4 | 67.5 | 70.3 | - |
| TML++ [13] | - | - | - | - | - | - | - | 71.5 | - |
| Simple (ResNet-50) [40] | 80.5 | 79.1 | 75.5 | 66.0 | 70.8 | 61.7 | 70.0 | 72.4 | - |
| Simple (ResNet-152) [40] | 83.4 | 81.7 | 80.0 | 72.4 | 75.3 | 67.1 | 74.8 | 76.7 | 560 |
| STEmbedding [16] | 81.6 | 83.8 | 77.1 | 70.0 | 77.4 | 70.8 | 74.5 | 77.0 | - |
| HRNet [31] | 83.6 | 82.1 | 80.4 | 73.3 | 75.5 | 68.5 | 75.3 | 77.3 | 580 |
| MDPN [12] | 88.5 | 85.2 | 83.9 | 77.5 | 79.0 | 71.4 | 77.0 | 80.7 | - |
| Dynamic-GNN [47] | 88.4 | 88.4 | 82.0 | 74.5 | 79.1 | 73.1 | 78.3 | 81.1 | - |
| PoseWarper [5] | 88.3 | 81.4 | 83.9 | 78.0 | 82.4 | 73.6 | 80.5 | 81.2 | 1293 |
| DCPose [21] | 88.7 | 88.0 | 84.1 | 78.4 | 83.0 | 74.2 | 81.4 | 82.8 | 1390 |
| FAMI-Pose [22] | 90.1 | 89.6 | 86.3 | 80.0 | 84.6 | 77.0 | 83.4 | 84.8 | 1577 |
| **VEPE (Ours)** | **88.5** | **87.6** | **84.5** | **78.8** | **83.3** | **74.7** | **81.7** | **83.0** | 334 |

pared, including PoseTracker [11], PoseFlow [42], Joint-Flow [8], FastPose [48], TML++ [13], and SimpleBaseline (both ResNet-50 and ResNet-152). Additionally, comparisons are made with STEmbedding [16], HRNet [31], MDPN [12], Dynamic-GNN [47], PoseWarper [5], DCPose [21], FAMI-Pose[22], as well as our VEPE model. Table 1 provides the results of these methods on the PoseTrack2017 validation set.

From the table, we can observe that our method outperforms most of the two-stage models and slightly outperforms the DCPose model, and our inference speed substantially outperforms the DCPose and FAMIPose models. In addition, we obtain excellent accuracy at challenging joints such as ankles, wrists, knees, etc., which indicates that our model is robust to fast-motion and body-obscuring scenarios and has excellent inference efficiency.

### 4.3.2. Results on the PoseTrack2018 Dataset.

We conducted an evaluation of our model using the PoseTrack2018 validation set, making comparisons across a total of 10 different methods. The methods included STAF[26], PGPT[3], AlphaPose[10], TML++ [13],Dynamic-GNN[47], MDPN [12], PoseWarper[5], DCPose[21], FAMIPose[22], as well as our proposed VEPE. From the table **??**, VEPE has a simple and flexible end-to-end structure that outperforms most of the two-stage models and improves 0.2 mAP over DCPose.



Figure 5. Visualization of Posetrack21 validation results. The first row shows the scene in fast motion, the row column shows the scene where the human body is occluded, and the third row shows the scene with a dense crowd.

### 4.3.3. Results on the PoseTrack2021 Dataset.

Our model is evaluated using the validation set from Pose-Track2021. A total of 7 methods are compared in the table **??**, including Tracktor++ w. poses[4, 9], CorrTrack[9, 27], CorrTrack w. ReID[9, 27], Tracktor++ w. Corr.[4, 9] , DCPose[21], FAMIPose[22] and our VEPE.

The difference between PoseTrack21 and Pose-

Track2018 is the addition of more challenging pose test scenarios such as small figures and figures in crowds. From the table, we can find that in the more challenging PoseTrack21, our model has nearly improved compared to DCpose by 0.3mAP, which shows that our model still has good robustness in complex scenes. Specific visualization examples are shown in Figure.5

## 4.4. Ablation Experiments

We conducted ablation studies focusing on the contribution of each component in the VEPE framework. We also examined the STPE and Instance Consistency Mechanism for further discussion. All experiments were performed on the PoseTrack 2017 validation set.

### 4.4.1. Study on components of VEPE

Table.2 presents a summary of the impact of various design components on the posetrack17 dataset. The baseline model refers to the model with only the spatial phase of training, which has a mAP of 77.2. Finally, all components were used to reach a mAP of 83.0, which is an improvement of 5.8 mAP with respect to the baseline model, with the Spatio-Temporal Pose Encoder (STPE) module showing the largest improvement.

Table 2. Ablation of different components in VEPE.

| Method | STPE | ICM | STDME | STPD | Mean |
|--------|------|-----|-------|------|------|
| Baseline | | | | | 77.2 |
| (a) | ✓ | | | | 80.7 |
| (b) | ✓ | ✓ | | | 82.1 |
| (c) | ✓ | ✓ | ✓ | | **82.5** |
| (d) | ✓ | ✓ | ✓ | ✓ | **83.0** |

### 4.4.2. Study on Pose Query Selection

In dense scenes, instance frequently exhibit considerable overlap and similarity due to the static characteristics common to human bodies. To tackle this redundancy and improve the efficacy of the spatio-temporal pose encoder (STPE), we have developed a module specifically for selecting pose queries. This module enhances query quality by employing a confidence threshold, appropriately set by the pose instance evaluator, to filter out queries. The optimal threshold is established through rigorous experimental analysis. The effects of different confidence thresholds on the selection of pose queries are detailed in the following table3.

### 4.4.3. Study on Instance Consistency Mechanism

The instance consistency mechanism aims to improve the accuracy of the temporal pose encoder in pose query recognition and tracking and reduce the mismatch problem in

Table 3. Effect of Different Threshold Settings on Pose Query Selection

| threshold | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|-----------|-----|-----|-----|-----|-----|
| mAP(%) | 80.4 | 80.5 | 80.7 | 80.4 | 80.1 |

cross-frame pose query matching by enhancing cross-frame instance query consistency and differentiation. Briefly, we introduce instance queries that correspond one-to-one with the pose query and ensure that the instance query feature space representations of the same instance in different frames have similarity. The tracking of instances is realized by calculating the similarity between instance queries, which guides the pose query for cross-frame matching and improves the accuracy of the temporal pose encoder in instance identification and tracking.

As can be seen from Table.2, the inclusion of the Instance Consistency Mechanism (ICM) improves the model accuracy by 1.4%, which proves the effectiveness of this mechanism. Figure.6 illustrates the tracking effect of instance queries under the instance consistency mechanism, using two example frames. The heatmaps in the lower-right corner show the similarity between all instance queries, allowing for the identification of the same individual across frames. For instance, the similarity heatmap in the first row shows that ID-0 and ID-14 represent the same human individual and they have the highest similarity.
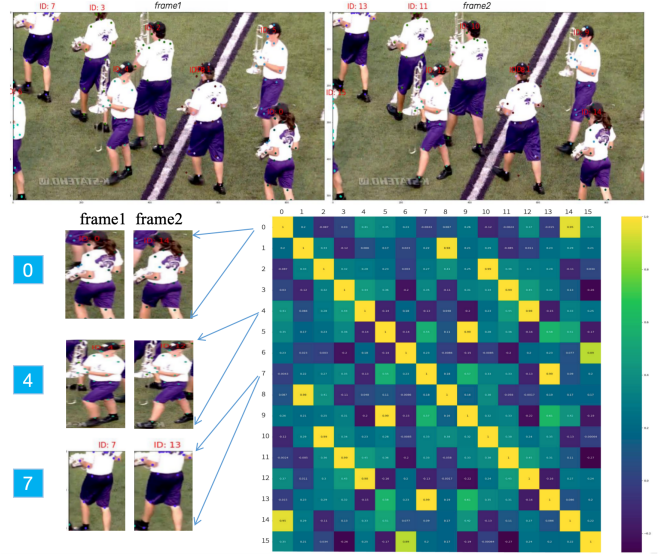


Figure 6. Instance tracking between cross-frame instance queries

# 5. Conclusion

In this paper, we propose a simple video-based end-to-end multi-person pose estimation framework, VEPE, which uses the Transformer to naturally associate instance targets in a video as a sequence-to-sequence task and performs end-to-end training without any post-processing. The VEPE learns the temporal context of pose instances and visual cues across frames through the Spatio-Temporal Pose Encoder (STPE) and the Spatio-Temporal Deformable Memory Encoder (STDME), and finally, Spatio-Temporal Pose Decoder (STPD) cascades to update the pose results of the keyframes by utilizing the spatio-temporal cues from Spatio-Temporal Pose Encoder and Spatio-Temporal Deformable Memory Encoder. Furthermore, to reduce the mismatch problem during the cross-frame pose query matching process, an instance consistency mechanism is introduced, which aims to enhance the consistency and discrepancy of the cross-frame instance query and realize the instance tracking function, which in turn guides the pose query to perform accurate cross-frame matching. Extensive experiments on the Posetrack dataset demonstrate that VEPE surpasses most two-stage models in performance and triples inference efficiency.

## References

[1] Mykhaylo Andriluka, Umar Iqbal, Eldar Insafutdinov, Leonid Pishchulin, Anton Milan, Juergen Gall, and Bernt Schiele. Posetrack: A benchmark for human pose estimation and tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5167–5176, 2018. 7

[2] Bruno Artacho and Andreas Savakis. Unipose: Unified human pose estimation in single images and videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7035–7044, 2020. 2, 3

[3] Qian Bao, Wu Liu, Yuhao Cheng, Boyan Zhou, and Tao Mei. Pose-guided tracking-by-detection: Robust multi-person pose tracking. *IEEE Transactions on Multimedia*, 23: 161–175, 2020. 8

[4] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 941–951, 2019. 8

[5] Gedas Bertasius, Christoph Feichtenhofer, Du Tran, Jianbo Shi, and Lorenzo Torresani. Learning temporal pose estimation from sparsely-labeled videos. In *Advances in Neural Information Processing Systems*, pages 3027–3038, 2019. 3, 8

[6] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020. 3

[7] Xiao Chu, Wei Yang, Wanli Ouyang, Cheng Ma, Alan L Yuille, and Xiaogang Wang. Multi-context attention for human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1831–1840, 2017. 1

[8] Andreas Doering, Umar Iqbal, and Juergen Gall. Joint flow: Temporal flow fields for multi person tracking. *arXiv preprint arXiv:1805.04596*, 2018. 8

[9] Andreas Doering, Di Chen, Shanshan Zhang, Bernt Schiele, and Juergen Gall. Posetrack21: A dataset for person search, multi-object tracking and multi-person pose tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20963–20972, 2022. 7, 8

[10] Hao-Shu Fang, Shuqin Xie, Yu-Wing Tai, and Cewu Lu. Rmpe: Regional multi-person pose estimation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2334–2343, 2017. 1, 8

[11] Rohit Girdhar, Georgia Gkioxari, Lorenzo Torresani, Manohar Paluri, and Du Tran. Detect-and-track: Efficient pose estimation in videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 350–359, 2018. 8

[12] Hengkai Guo, Tang Tang, Guozhong Luo, Riwei Chen, Yongchen Lu, and Linfu Wen. Multi-domain pose network for multi-person pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 8

[13] Jihye Hwang, Jieun Lee, Sungheon Park, and Nojun Kwak. Pose estimator and tracker using temporal flow maps for limbs. In *2019 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2019. 8

[14] Umar Iqbal, Anton Milan, and Juergen Gall. Posetrack: Joint multi-person pose estimation and tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2011–2020, 2017. 7

[15] Tao Jiang, Peng Lu, Li Zhang, Ningsheng Ma, Rui Han, Chengqi Lyu, Yining Li, and Kai Chen. Rtmpose: Real-time multi-person pose estimation based on mmpose. *arXiv preprint arXiv:2303.07399*, 2023. 1

[16] Sheng Jin, Wentao Liu, Wanli Ouyang, and Chen Qian. Multi-person articulated tracking with spatial and temporal embeddings. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5664–5673, 2019. 8

[17] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 7

[18] Yanjie Li, Sen Yang, Peidong Liu, Shoukui Zhang, Yunxiao Wang, Zhicheng Wang, Wankou Yang, and Shu-Tao Xia. Simcc: A simple coordinate classification perspective for human pose estimation. In *European Conference on Computer Vision*, pages 89–106. Springer, 2022. 1

[19] Huan Liu, Qiang Chen, Zichang Tan, Jiang-Jiang Liu, Jian Wang, Xiangbo Su, Xiaolong Li, Kun Yao, Junyu Han, Errui Ding, et al. Group pose: A simple baseline for end-to-end multi-person pose estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15029–15038, 2023. 3

[20] Zhenguang Liu, Shuang Wu, Shuyuan Jin, Qi Liu, Shijian Lu, Roger Zimmermann, and Li Cheng. Towards natural and

accurate future motion prediction of humans and animals. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10004–10012, 2019. 1

[21] Zhenguang Liu, Haoming Chen, Runyang Feng, Shuang Wu, Shouling Ji, Bailin Yang, and Xun Wang. Deep dual consecutive network for human pose estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 525–534, 2021. 2, 3, 8

[22] Zhenguang Liu, Runyang Feng, Haoming Chen, Shuang Wu, Yixing Gao, Yunjun Gao, and Xiang Wang. Temporal feature alignment and mutual information maximization for video-based human pose estimation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11006–11016, 2022. 3, 8

[23] Yue Luo, Jimmy Ren, Zhouxia Wang, Wenxiu Sun, Jinshan Pan, Jianbo Liu, Jiahao Pang, and Liang Lin. Lstm pose machines. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5207–5215, 2018. 1, 2, 3

[24] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European conference on computer vision*, pages 483–499. Springer, 2016. 1

[25] Tomas Pfister, James Charles, and Andrew Zisserman. Flowing convnets for human pose estimation in videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1913–1921, 2015. 2, 3

[26] Yaadhav Raaj, Haroon Idrees, Gines Hidalgo, and Yaser Sheikh. Efficient online multi-person 2d pose tracking with recurrent spatio-temporal affinity fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4620–4628, 2019. 8

[27] Umer Rafi, Andreas Doering, Bastian Leibe, and Juergen Gall. Self-supervised keypoint correspondences for multi-person pose estimation and tracking in videos. In *European Conference on Computer Vision*, pages 36–52. Springer, 2020. 8

[28] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. 7, 8

[29] Dahu Shi, Xing Wei, Liangqi Li, Ye Ren, and Wenming Tan. End-to-end multi-person pose estimation with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11069–11078, 2022. 2, 3, 7

[30] Jie Song, Limin Wang, Luc Van Gool, and Otmar Hilliges. Thin-slicing network: A deep structured model for pose estimation in videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4220–4229, 2017. 2, 3

[31] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5693–5703, 2019. 7, 8

[32] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. In *Proceedings of*

the IEEE/CVF conference on computer vision and pattern recognition*, pages 14454–14463, 2021. 3

[33] Alexander Toshev and Christian Szegedy. Deeppose: Human pose estimation via deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1653–1660, 2014. 1

[34] 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*, pages 5998–6008, 2017. 5

[35] Fang Wang and Yi Li. Beyond physical connections: Tree models in human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 596–603, 2013. 1

[36] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 1

[37] Manchen Wang, Joseph Tighe, and Davide Modolo. Combining detection and tracking for human pose estimation in videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11088–11096, 2020. 1, 2, 3

[38] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 4724–4732, 2016. 1

[39] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. 1

[40] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *Proceedings of the European conference on computer vision (ECCV)*, pages 466–481, 2018. 8

[41] Yabo Xiao, Kai Su, Xiaojuan Wang, Dongdong Yu, Lei Jin, Mingshu He, and Zehuan Yuan. Querypose: Sparse multi-person pose regression via spatial-aware part-level query. *Advances in Neural Information Processing Systems*, 35: 12464–12477, 2022. 3

[42] Yuliang Xiu, Jiefeng Li, Haoyu Wang, Yinghong Fang, and Cewu Lu. Pose flow: Efficient online pose tracking. *arXiv preprint arXiv:1802.00977*, 2018. 8

[43] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer baselines for human pose estimation. *Advances in Neural Information Processing Systems*, 35:38571–38584, 2022. 1

[44] Sijie Yan, Yuanjun Xiong, and Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *Proceedings of the AAAI conference on artificial intelligence*, 2018. 1

[45] Jie Yang, Ailing Zeng, Shilong Liu, Feng Li, Ruimao Zhang, and Lei Zhang. Explicit box detection unifies end-to-end multi-person pose estimation. *arXiv preprint arXiv:2302.01593*, 2023. 3

[46] Sen Yang, Zhibin Quan, Mu Nie, and Wankou Yang. Transpose: Keypoint localization via transformer. In *Proceedings*

*of the IEEE/CVF International Conference on Computer Vision*, pages 11802–11812, 2021. 1

[47] Yiding Yang, Zhou Ren, Haoxiang Li, Chunluan Zhou, Xinchao Wang, and Gang Hua. 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*, pages 8074–8084, 2021. 8

[48] Jiabin Zhang, Zheng Zhu, Wei Zou, Peng Li, Yanwei Li, Hu Su, and Guan Huang. Fastpose: Towards real-time pose estimation and tracking via scale-normalized multi-task networks. *arXiv preprint arXiv:1908.05593*, 2019. 8

[49] Xiaoqin Zhang, Changcheng Li, Xiaofeng Tong, Weiming Hu, Steve Maybank, and Yimin Zhang. Efficient human pose estimation via parsing a tree structure based human model. In *2009 IEEE 12th International Conference on Computer Vision*, pages 1349–1356. IEEE, 2009. 1

[50] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3, 5