# Improving Multi-Person Pose Tracking with A Confidence Network

Zehua Fu†, Wenhang Zuo†, Zhenghui Hu, Qingjie Liu*, *Member, IEEE*, Yunhong Wang, *Fellow, IEEE*

arXiv:2310.18920v1 [cs.CV] 29 Oct 2023

*Abstract*—

Human pose estimation and tracking are fundamental tasks for understanding human behaviors in videos. Existing top-down framework-based methods usually perform three-stage tasks: human detection, pose estimation and tracking. Although promising results have been achieved, these methods rely heavily on high-performance detectors and may fail to track persons who are occluded or miss-detected. To overcome these problems, in this paper, we develop a novel keypoint confidence network and a tracking pipeline to improve human detection and pose estimation in top-down approaches. Specifically, the keypoint confidence network is designed to determine whether each keypoint is occluded, and it is incorporated into the pose estimation module. In the tracking pipeline, we propose the Bbox-revision module to reduce missing detection and the ID-retrieve module to correct lost trajectories, improving the performance of the detection stage. Experimental results show that our approach is universal in human detection and pose estimation, achieving state-of-the-art performance on both PoseTrack 2017 and 2018 datasets.

*Index Terms*—Multi-person Pose Tracking, Keypoint confidence network, Pose estimation, Bbox-revision.
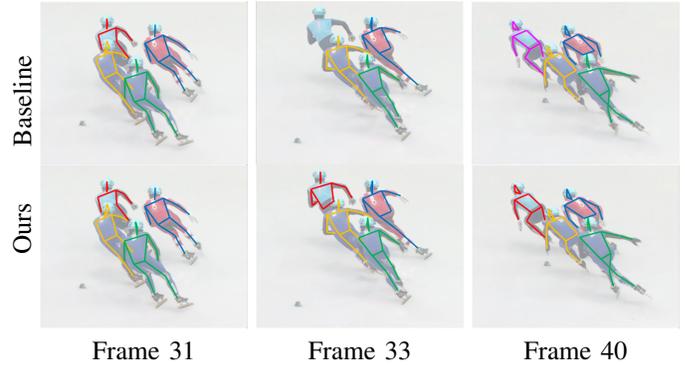


Fig. 1: The Baseline is a top-down multi-person pose tracking method based on HRNet. The baseline may cause the detector failure in occlusion (Frame 33), which results in loss of tracking trajectory (Frame 40). Although using the same detector as the baseline, our approach overcomes this limitation through the Bbox-revision module in the pose tracking pipeline.

## I. INTRODUCTION

Multi-person pose tracking, which intends to detect the body joints of all persons in the video frames and output the pose trajectories over time consistently, is a fundamental task for human-centered video understanding. It is widely used in human-computer interaction, video surveillance, and action recognition [1]–[4], etc. Multi-person pose tracking tasks serve as integral components in multimedia applications, facilitating advances in domains such as video analysis, gesture recognition, gaming and interactive experiences. These tasks allow precise interpretation and assessment of human movements, thus increasing user interaction, immersion, and comprehensiveness in multimedia contexts. The multi-person pose tracking can be categorized into bottom-up [5] and top-down [6]–[9] methods. The former one estimates all joints in each frame and associates each person's joints over time in a spatio-temporal optimization manner without detecting human bounding boxes. While the latter pipeline first detects the bounding box of each person and then estimates his/her

† These authors contribute equally to this work and should be considered co-first authors. * Corresponding author.

Wenhang Zuo, Qingjie Liu and Yunhong Yang are with the State Key Laboratory of Virtual Reality Technology and Systems, School of Computer Science and Engineering, Beihang University, Beijing 100191, China. E-mail: zuowenhang@gmail.com, {qingjie.liu, yhwang}@buaa.edu.cn.

Zehua Fu and Zhenghui Hu are with Hangzhou Innovation Institute, Behang University, Hangzhou 310051, China. Email: {zehua_fu, zhenghuihu2021}@buaa.edu.cn.

pose. In the final stage, a multi-person tracking methodology is applied to retrieve the trajectories of joints. Thanks to the advancement of object detectors in recent years [10]–[19], the top-down pipeline has made great progress and has surpassed the bottom-up pipeline to become the mainstream. However, there are two barriers that prevent these methods from being perfect: occlusion and fast motion (e.g., Figure 1 ). Top-down methods filter keypoints based on heatmaps that are predicted by pose estimators optimized for images instead of video frames. The estimators suffer from motion blurs, and thus it is hard to produce accurate keypoints. Furthermore, occlusions between adjacent persons may fool estimators into making wrong predictions. Finally, occlusions and motion blurs may also degrade person detectors, eventually leading to tracking failure.

In this paper, we attack these two problems with a novel confidence estimation and temporal correction strategy. Specifically, we design a confidence network measuring the visibility of the keypoints in addition to the location probabilities in the heatmap. Then, we build an online tracking pipeline to perform multi-person pose tracking, which consists of three modules, including an association module, an ID-retrieve module, and a Bbox-revision module. Similarly to previous work [20], the association module assigns a unique identification number to each pose in frames. We follow [20] and employ the Hungarian algorithm to achieve this. However, due to challenges such as motion blurs and occlusions, there exist non-matched poses

in both sets. We solve this problem using the following two modules. The ID-retrieve module assigns the ID of a person in the current frame who has no matched ID in the previous frame but may be matched in history. The Bbox-revision module serves as an assistant to the detector. It helps to repair missing detections by leveraging the motion of persons. We also introduce a pose filtering operation to improve the overlapping poses caused by error detections.

Research conducted by Yang et al. [21] addresses the same issue as our proposed Bbox-revision module, namely the missed detection in the current frame on challenging scenes such as occlusion and fast motion. Both work [21] and ours aggregate the pose of the current frame $t$ from the estimator and the predictor. Work [21] predicts the pose of the current frame $t$ from $n$ historical poses using GNN (Graph Neural Network), while ours predicts the pose of the current frame $t$ from the $t-1$ pose using optical flow. Considering that poses in the current frame and historical frames may not be correct, there are two main differences between work [21] and ours. First, we filter the wrong detections, namely bounding boxes without objects with the score calculated from the proposed keypoint confidence. Second, we proposed a novel similarity matrix based on our keypoint confidence for one-to-one mapping to remove redundant detections. Thus, we improve the tracking on missed detection with a novel keypoint confidence. Besides, we also propose an ID-retrieve module to correct lost trajectories, which further improves the tracking performance.

To summarize, the main contributions of this study are as follows:

- We design a keypoint confidence network to measure the visibilities of keypoints, which is helpful to improve the pose estimator's performance on occluded joints; we also propose a tracking pipeline that corrects lost trajectories and miss-detections.
- The keypoint confidence network and pose tracking pipeline are universal and can be used with different pose estimation networks and human detectors.
- The proposed method achieves MOTAs of 69.2%, 72.2%, and 63.5% on the 2018 validation set, the 2017 validation set, and the 2017 test set of the PoseTrack dataset, respectively, achieving state-of-the-art performance.

## II. RELATED WORK

We briefly review the following two related topics, including multi-person pose estimation and multi-person pose tracking.

### A. Multi-person Pose Estimation in Image

Human pose estimation, which aims to locate the keypoint of the human body in images, can be classified into single-person [22]–[25] and multi-person pose estimation [26]–[31]. The multi-person pose estimation is more realistic and challenging. It has received increased attention and made significant progress in recent years [24], [25], [32]–[43]. Multi-person pose estimation is generally classified into top-down methods [27], [31], [44]–[47] and bottom-up methods [30], [48]–[53].

The two most important components of top-down methods are the human detector and the single-person pose estimator. Top-down methods first detect persons and generate the person bounding boxes. Then, pose estimation is conducted to detect the human pose for each bounding box. Unlike top-down methods, bottom-up methods do not rely on human detectors. Bottom-up methods first detect all body joints of every person and then group them by some fitting algorithms to form human poses. In this work, we mainly focus on the top-down method. To verify the effectiveness of our method, we will conduct experiments on Hourglass [24], SimpleBaselinet [47] and HRNet [31].

### B. Multi-person Pose Tracking

Recently, multi-person pose tracking has received significant attention since the topic was first introduced by the PoseTrack dataset. Pose estimation in images can be extended to pose tracking in the video by running independently on each frame and then using data association to correctly link the continuous trajectory of each person over time. Bottom-up methods [54], [55] estimate all joints in each frame and associate the joints in a spatio-temporal optimization manner without detecting human bounding boxes. For example, Spatio-Temporal Affinity Fields (STAF) [54] build upon Part Affinity Fields representation [26] and propose an architecture that can encode and predict spatio-temporal affinity fields across a video sequence.

Top-down methods [2], [9], [21], [56] are based on top-down pose estimation and exploiting spatio-temporal context for tracking. CombDet [6] extended HRNet [31] from 2D to 3D to build a tracking pipeline to alleviate missed detection, but the approach is limited by the size of clip lengths. TKMRNet [7] proposes refinement networks to improve pose precision and design new keypoint similarity metrics in the association module. Existing top-down methods treat pose estimation and pose tracking as two relatively independent tasks designed and optimized separately. Thus, the pose tracking network directly utilizes the estimated pose from the pose estimation network. In this work, we try to boost the tracking performance by leveraging the additional information output from the pose estimation network, i.e., the proposed keypoint confidence.

## III. METHODOLOGY

In this section, we introduce the Keypoint Confidence Network (KCN) and the Pose Tracking Pipeline (PTP) for multi-person pose tracking. Our method works by first detecting all persons in the current frame and then estimating the keypoints and their confidence for each person by KCN. In the tracking stage, we use two modules to correct the tracking: we retrieve the lost trajectories in previous frames using the ID-retrieve module and revise the bounding boxes of persons with the aid of optical flows using the Bbox-revision module. In what follows, we introduce each component of the proposed method in detail.
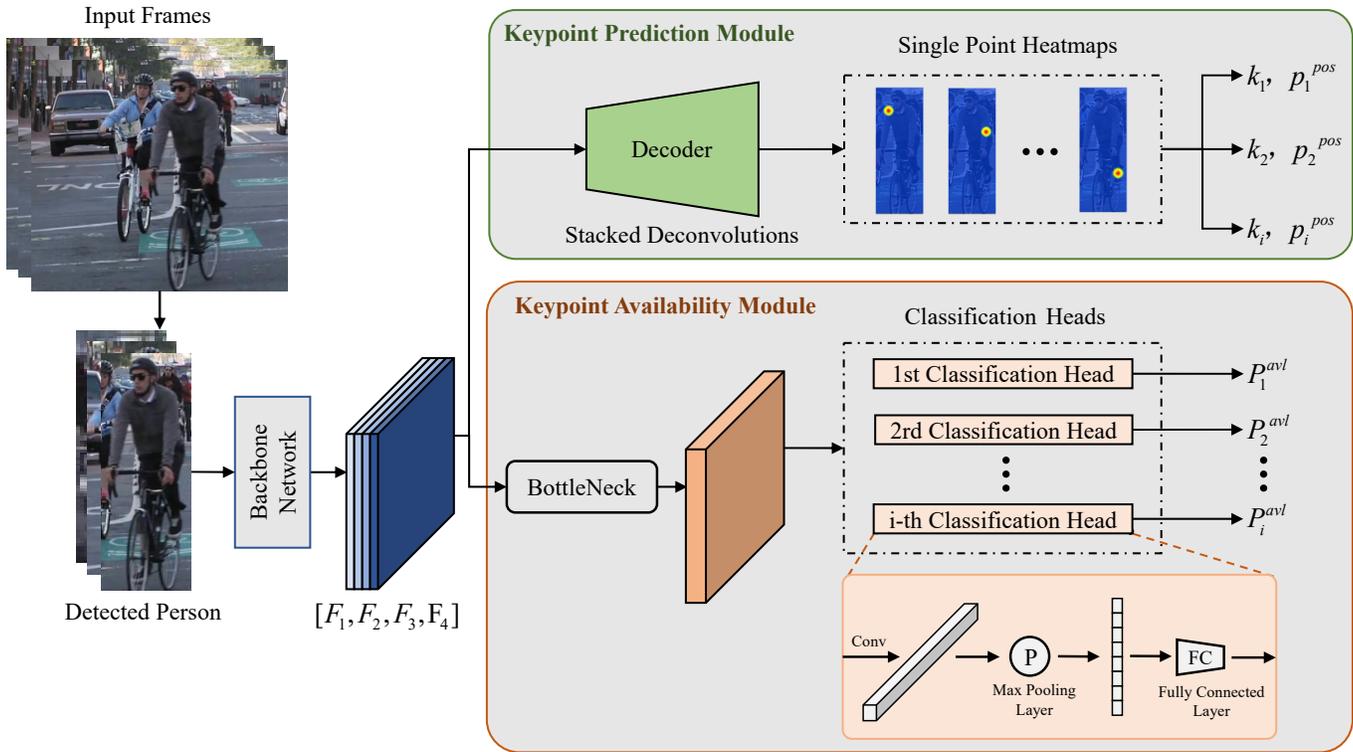
Fig. 2: Illustrating the architecture of the Keypoint Confidence Network (KCN). It consists of a keypoint prediction module and a keypoint availability module. $[F_1, F_2, F_3, F_4]$ is the high-resolution representation after multiple-scale fusion with exchange units.

## A. Keypoint Confidence Network

In the process of pose tracking, interactive occlusions can lead to the unavailability of some keypoints, necessitating filtering. Existing research employs keypoint location probability, namely, the probability of each keypoint at its current location, as the basis for filtering. However, this approach may lead to several issues: on the one hand, occluded keypoints that are mistakenly assigned to other individuals maintain high location probabilities and are erroneously retained; on the other hand, under conditions of frame blurring, location probabilities for each keypoint are generally low, causing inaccurate filtering. The occurrence of these errors complicates the setting of keypoint thresholds, undermines keypoint recognition accuracy, and directly impacts tracking performance. To address these issues, we exploit the global spatial structure information of the human body for modeling the availability probability of each keypoint. And then introduce a Keypoint Confidence Network (KCN). By combining keypoint location probability and availability probability, we effectively alleviate the aforementioned issues and subsequently improve tracking performance.

The architecture of our proposed KCN is shown in Figure 2. Our proposed KCN consists of a backbone for extracting features and two parallel branches for pose estimation: the Keypoint Prediction Module (KPM) for predicting keypoint location as well as the location probability and the Keypoint Availability Module (KAM) for estimating the probability of keypoint availability. We use HRNet [31] as the backbone for

feature extraction.

**Keypoint Prediction Module.** The keypoint prediction module consists of three $3 \times 3$ deconvolution layers and generates $K$ heatmaps, where $K$ is the number of keypoints for each person. For each heatmap $M$, we obtain the keypoint location $l$ from the highest response in the heatmap and use its response value as keypoint location probability $p^{loc}$. Following [6], we convert point-wise annotations into heatmaps using 2D Gaussian convolutions and consider them as ground truths of keypoints during training. The loss function of this module is defined as

$$L = \frac{1}{KWH} \sum_{k}^{K} \sum_{i}^{W} \sum_{j}^{H} \|M_{kij} - G_{kij}\|_2^2 \qquad (1)$$

where $W$ and $H$ represent the width and height of heatmaps. $M$ and $G$ represent the predicted heatmaps and the ground truth maps.

**Keypoint Availability Module.** The keypoint availability module is composed of a bottleneck layer [57] and $K$ classifier heads to obtain the keypoint availability probability $p^{avl}$. The bottleneck layer consists of three convolution layers with sizes of $1 \times 1$, $3 \times 3$, and $1 \times 1$, respectively. Each layer has a channel size of 384. The classifier heads predict the availability probabilities of keypoints. For each head, there are a $1 \times 1$ convolution layer, a global max-pooling layer, and a fully connected layer followed by a softmax layer. The convolution
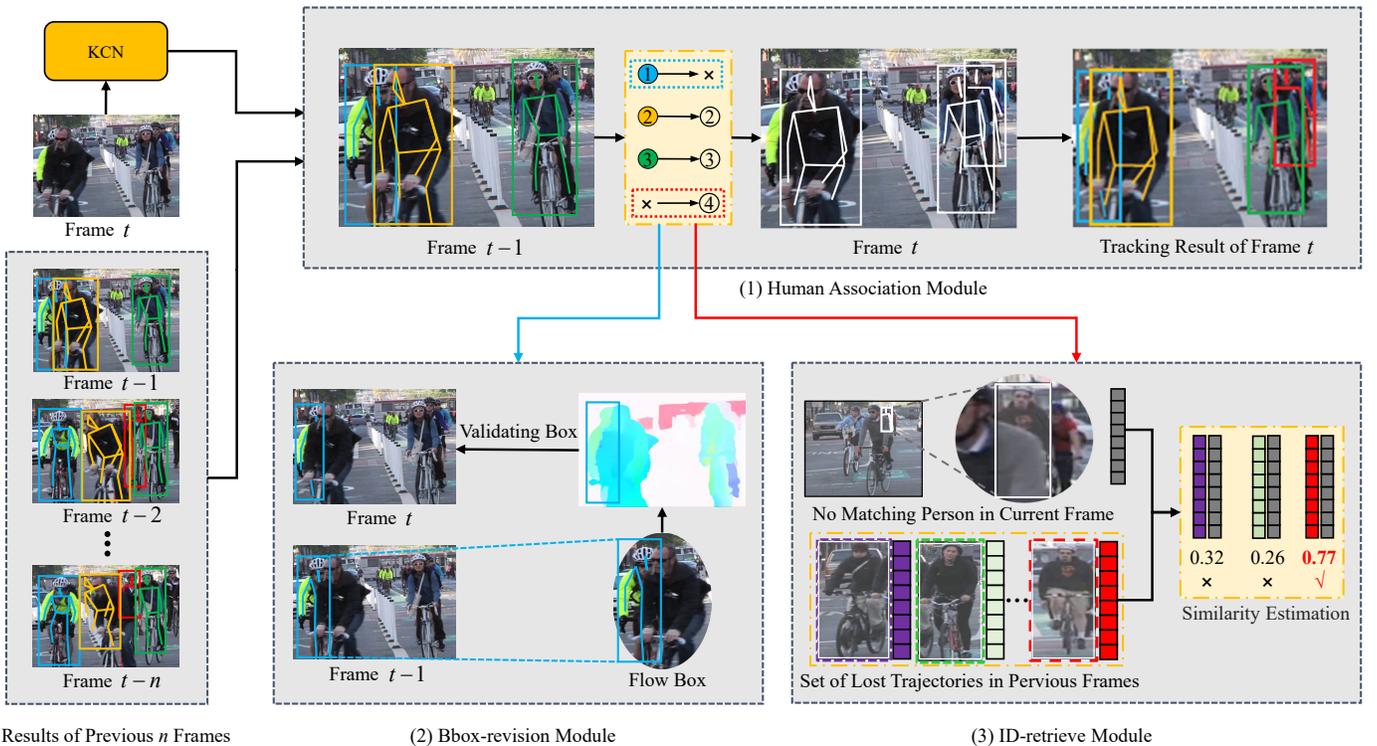
Fig. 3: Illustration of the proposed pose tracking pipeline. In the first stage, we detect persons and estimate their poses using our keypoint confidence network. Then, in the tracking stage, (1) performs the identity association between frames and (2) generates the bounding boxes from the previous frame for the unmatched trajectories. Finally, (3) identifies a person in the current frame who has no matched ID in the previous frame but may be matched in history.

layer has 384 channels as input and outputs 512 channels. The fully connected layer has 512 input units and 2 output units.

During training, we use an $\alpha$-balanced variant of the focal loss [58] $L_f$ for optimization

$$L_f = \begin{cases} -\alpha(1-p)^\gamma \log(p) & , y = 1 \\ -(1-\alpha)p^\gamma \log(1-p) & , y = 0 \end{cases} \quad (2)$$

where $p$ is the keypoint availability probability of $p^{loc}$, and $y$ is the label of keypoint availability. We use the default settings where $\gamma = 2$, $\alpha = 0.25$.

**Keypoint Confidence.** Through KCN, we obtain both keypoint location probability $p^{loc}$ and keypoint availability probability $p^{avl}$. By combining them, we can calculate the keypoint confidence $p^{conf}$ of the $i$-th body joint as follows

$$p_i^{conf} = p_i^{avl} \times p_i^{loc} \quad (3)$$

### B. Pose Tracking Pipeline

The illustration of our proposed pose tracking pipeline is shown in Figure 3. In the tracking stage, we first use a human association module to match persons. Two things may happen: 1) a person is lost in previous frames but shows up again in the current frame; 2) the detector may miss detecting persons due to occlusion or fast motion. We propose the following two modules, the ID-retrieve module and the Bbox-revision module, to solve the above two challenges. The ID-retrieve module retrieves lost IDs using a person re-identification technique. We use the Bbox-revision module to

generate bounding boxes that are missed by the detector in the current frame. The details of each module are given below.

**Human Association Module.** During tracking, we first analyze the detection and pose estimation results in the current frame and then associate them with the existing trajectories by assigning a unique identification number. The Human Association Module in our approach addresses the weighted bipartite graph matching problem, aiming to establish an optimal one-to-one correspondence between the trajectory from the previous frame, and the detection and pose estimation results from the current frame. To accomplish this, we utilized the Hungarian algorithm, which requires the construction of a weight matrix. Specifically, to measure the similarity of a person across frames, we employ the Object Keypoint Similarity (OKS) [59], an evaluation metric in multi-person pose tracking commonly used in the COCO keypoints challenge and PoseTrack challenge. Consider both pose similarity and distance, ensuring only poses that are similar and close poses receive higher weights.

**Bbox-revision Module.** The Bbox-revision module is used to reduce missing detection in the current frame. As shown in Figure 3 (2), the person detected in the previous frame failed to be detected in the current frame, resulting in the unmatched trajectory. To address this issue, we first use the optical flow sub-module to generate bounding boxes from the previous frame and verify their correctness by a pose filtering sub-module. The verified bounding boxes are reserved and identified with the ID of the corresponding unmatched

trajectories.

For the optical flow sub-module, we apply the RAFT (Recurrent All-Pairs Field Transforms) [60] method to get the offset of each pixel in the previous frame. Thus, the position of each pixel in the unmatched pose trajectory in the current frame is provided and the minimum bounding rectangle is used as the bounding box.

After bounding boxes are generated, the pose filtering sub-module is used to verify their correctness in two steps. First, KCN is applied to estimate human poses, and the average keypoint confidence is used as the score of each bounding box for filtering. Then, we propose a novel similarity metric that evaluates the overlap between the bounding boxes generated by optical flow and those detected in the current frame. This approach addresses the challenge of sub-region overlapping of bounding boxes generated by optical flow. The optical flow sub-module may produce wrong trajectories due to the disappearance of the person or missing detection in the previous stage. Therefore, in the scene where the human is interlaced and occluded, the Non-Maximum Suppression (NMS) algorithm based on the Intersection over Union (IoU) and OKS cannot filter out the overlapping box, which ultimately affects the final performance. To address this issue, we propose a new similarity metric that calculates the overlap of poses by focusing on the shared keypoint regions within both boxes, allowing a more efficient filtering process for overlapping bounding boxes. It is defined as

$$\delta\left(IoU_{p,q} > 0.1\right) \frac{\sum_i e^{-\frac{d_i^2}{2s^2 k_i^2}} \delta\left(p_i^{\mathrm{conf}} > \theta\right)\delta\left(q_i^{\mathrm{conf}} > \theta\right)}{\sum_i \delta\left(p_i^{\mathrm{conf}} > \theta\right)\delta\left(q_i^{\mathrm{conf}} > \theta\right)} \quad (4)$$

where the $d_i^2$ is the euclidean distance between person $p$ and person $q$. $\delta$ is an indicator function. $q_i^{conf}$ and $p_i^{conf}$ are the $i$-th keypoint confidences of person $p$ and person $q$. $s$ is the object scale, and $k_i$ is a per-keypoint constant that controls falloff. $\theta$ is the confidence threshold. $IoU_{p,q}$ computes the IoU metric between person $p$ and person $q$, dealing with special cases where two persons do not overlap.

**ID-retrieve Module.** The OKS-based human association may fail when occlusion and persons walk out of sight. As shown in Figure 3 (3), to remedy such matching failures, we propose the ID-retrieve module to improve the robustness of pose tracking.

For the input of this module, we reuse the features of the KCN backbone to perform matching followed by adaptive average pooling operations to obtain compact pedestrian features. Meanwhile, we maintain a feature set for lost persons. For a new detection in the current frame, we first perform feature matching with historical persons. If the similarity score reaches the threshold, we assign its historical person ID; otherwise, we assign a new ID.

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metrics

We evaluate the proposed method on PoseTrack, which is a large-scale benchmark for multi-person pose estimation and pose tracking in videos. It contains several video sequences
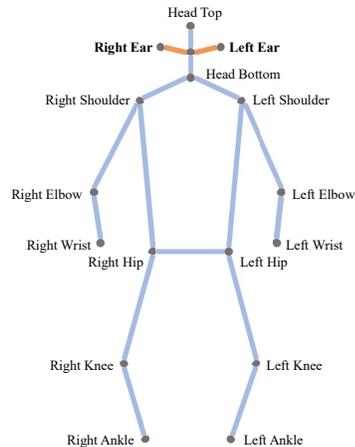


Fig. 4: Illustration of body joints in PoseTrack 2017 and PoseTrack 2018 datasets.

TABLE I: The description of PoseTrack 2017 Dataset.

|  | Num Poses | Num Trajectories | Num Videos |
|---|---|---|---|
| Training set | 61178 | 2437 | 292 |
| Validation set | 18996 | 695 | 50 |
| Test set | 73471 | 2334 | 208 |
| Total | **153615** | **5446** | **550** |

with a lot of annotated poses with various activities being performed. PoseTrack has the 2017 and 2018 versions of this benchmark. Each dataset has a publicly available training set and validation set, as well as an evaluation server for benchmarking on a held-out test set. PoseTrack 2017 annotates 15 body parts for each body pose, including the head, nose, neck, shoulders, elbows, wrists, hips, knees, and ankles; while PoseTrack 2018 annotates two more ears, as shown in Figure 4. For our experiments, only the original 15 body joints are used for both training and inference. We conduct experiments on both PoseTrack 2017 and PoseTrack 2018. Especially, PoseTrack 2017 includes 250 videos for training, 50 videos for validation, and 214 videos for tests, as shown in Table I. PoseTrack 2018 is expanded on the basis of PoseTrack2017, including 593 videos for training, 170 videos for validation, and 375 videos for the test. This is more than double the amount of data from PoseTrack 2017. In the training set, the videos are densely annotated 30 frames from the center of frames. In the validation and testing set, the videos are densely annotated 30 frames of the middle, and afterward annotated every fourth frame. In addition, we use the COCO dataset to pre-train the multi-person pose estimation model used in our experiments.

TABLE II: The Pose Tracking performance in MOTA (%) with the different confidence threshold on PoseTrack validation datasets.

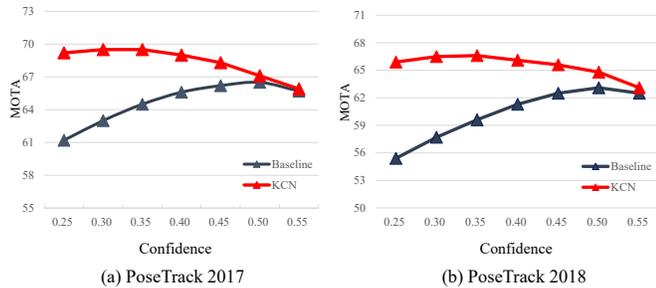| Dataset | Method | Threshold Value | | | | | | |
|---|---|---|---|---|---|---|---|---|
|  |  | 0.25 | 0.30 | 0.35 | 0.40 | 0.45 | 0.50 | 0.55 |
| PoseTrack 2017 | Baseline | 61.2 | 63.0 | 64.5 | 65.6 | 66.2 | 66.5 | 65.3 |
|  | Ours | 69.2 | 69.5 | 69.5 | 69.0 | 68.3 | 67.1 | 65.4 |
| PoseTrack 2018 | Baseline | 55.4 | 57.7 | 59.6 | 61.3 | 62.5 | 63.1 | 62.5 |
|  | Ours | 65.9 | 66.5 | 66.6 | 66.1 | 65.6 | 64.8 | 63.1 |

(a) PoseTrack 2017    (b) PoseTrack 2018

Fig. 5: The impact of the different confidence thresholds on PoseTrack validation datasets. The numbers in the figure refer to MOTA (%).

TABLE III: Effectiveness analysis of KCN by pose tracking performance (MOTA) on PoseTrack validation datasets.

| Dataset | Method | Head | Shou | Elb | Wri | Hip | Knee | Ankl | Total |
|---|---|---|---|---|---|---|---|---|---|
| PoseTrack 2017 | Baseline | 73.6 | 76.4 | 61.1 | 55 | 63.6 | 64.8 | 52.7 | 64.5 |
| | Ours | 76.5 | 77.8 | 70.0 | 63.2 | 66.2 | 68.9 | 60.5 | 69.5(**+5.0**) |
| PoseTrack 2018 | Baseline | 68.0 | 71.7 | 55.7 | 51.6 | 58.5 | 58.2 | 49.4 | 59.6 |
| | Ours | 69.8 | 74.7 | 69.7 | 62.3 | 63.9 | 64.7 | 59.2 | 66.6(**+7.0**) |

The performance of our method is evaluated from two aspects. We evaluate the accuracy of multi-person pose estimation using the mean Average Precision (mAP). We use Multiple Object Tracking Accuracy (MOTA) to evaluate the performance of trackers at keeping trajectories, including the performance of false positives, misses, and mismatches. MOTA is the main evaluation metric for PoseTrack Benchmark, which is defined as

$$\text{MOTA} = 1 - \frac{\sum_t \left( FP_t + FN_t + IDS_t \right)}{\sum_t g_t} \qquad (5)$$

where the subscript $t$ refers that current values are computed at the $t$-th frame. $FP_t$, $FN_t$, and $IDS_t$ denote the number of false positives, missed targets, and identity switches, respectively, at time $t$. $g_t$ stands for ground truth.

### B. Implementation Details

In the human detection stage, we use an HTC detector [61] to detect all the person instances in frames and extract crops of size 384×288 around detected person instances as input to our proposed keypoint confidence network. We use the pre-trained model trained on the COCO dataset in mmdetection [62]. Note that no additional fine-tuning of the detectors on the PoseTrack dataset is performed. In Non-Maximum Suppresion (NMS) operation for human detection, we changed the metric from IoU to OKS metric and set the threshold to 0.6.

TABLE IV: Keypoint statistical analysis of wrong detection and missing detection performed on the validation set of PoseTrack 2018.

| Dataset | Method | Wrong Detection | Missing Detection | Total |
|---|---|---|---|---|
| PoseTrack 2017 | Baseline | 22249 | 42099 | 64348 |
| | Ours | 19940(**-10.4%**) | 38270(**-9.1%**) | 58210(**-9.5%**) |
| PoseTrack 2018 | Baseline | 74353 | 97297 | 171650 |
| | Ours | 66058(**-11.6%**) | 88565(**-9.0%**) | 154623(**-9.9%**) |

In the pose estimation stage, our keypoint confidence network used HRNet [31] as the backbone. In the training phase, we first train the keypoint prediction module and then attach the point confidence module to refine the two modules together. Different from other methods, we combine the training set and the validation set of COCO dataset and the training set of the PoseTrack dataset for model training. The COCO dataset and the PoseTrack dataset do not agree on two joint types. The PoseTrack dataset contains the Head top and Head bottom joint points, while the COCO dataset contains the left ear and right ear joint points. Since the two different types of joint points between these two datasets are similar in position and both belong to the joint points of the face, we directly use the visibility of the ear instead of the visibility of the head for training. In particular, we use the default parameters of the HRNet in the pose estimation task for training. When training the point confidence module, we trained a total of 20 epochs to fine-tune the model.

To train the ID-retrieve module of the proposed pose tracking pipeline, we built a dataset based on PoseTrack 2018, which contains 119656 images with 4613 person tags. We use the Euclidean distance to measure the similarity between two features. When the similarity is less than the threshold value (e.g. 100), we consider it to be a person.

### C. Ablation Study

In this section, we provide ablation experiments and analysis to demonstrate the effectiveness of our proposed framework and the strength of the key components. For clarity, we denote the Keypoint Confidence Network as KCN and the Pose Tracking Pipeline as PTP.

***Analysis of Confidence Threshold.*** For top-down methods, the keypoint confidence threshold significantly impacts the final permanence. For a fair comparison, our baseline network only replaces the HRNet-based keypoint confidence network with an HRNet-based pose estimation method. As shown in Table II, with the confidence threshold ranging from 0.25 to 0.55, we compare the multi-person pose tracking performance variations of our KCN and the baselines on PoseTrack 2017 and PoseTrack 2018, respectively. We also visualize the result in Figure 5. As can be observed, the performance variation curves of KCN are much flatter between the two, implying that our KCN reduces the impact of the confidence threshold compared to our baseline. Meanwhile, our KCN outperforms the baseline method stably as the threshold changes. We observe that the performance of KCN begins to decline when the confidence threshold is greater than 0.35. It may be caused by the mistaken filtration of correct keypoints. The confidence threshold is set to 0.35 for the rest of the experiments.

***Analysis of Keypoint Confidence Network.*** In this part, we first evaluate the effectiveness of the proposed Keypoint Confidence Network (KCN) on the validation set of PoseTrack datasets. Our baseline network is the same as in the previous experiment. As shown in Table III, we compare our KCN with the baseline on the multi-person pose estimation and tracking task, where the performance is evaluated as MOTA and all joints are counted. As can be seen, after applying KCN,

TABLE V: Generalization of KCN across three top-down multi-person pose estimation networks by pose tracking performance in MOTA (%).

| Dataset | Method | KCN | Head | Shou | Elb | Wri | Hip | Knee | Ankl | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| PoseTrack 2017 | Hourglass | | 71.6 | 71.8 | 57.6 | 49.2 | 52.6 | 54.2 | 43.5 | 58.2 |
| | Hourglass | ✓ | 73.8 | 73.3 | 62.6 | 53.1 | 54.2 | 55.6 | 49.9 | 61.2(**+3.0**) |
| | Pose_ResNet | | 72.5 | 71.8 | 57.7 | 48.4 | 58.7 | 60.3 | 49.2 | 60.6 |
| | Pose_ResNet | ✓ | 73.5 | 72.6 | 63.4 | 53.1 | 59.3 | 62.6 | 55.5 | 63.6(**+3.0**) |
| | HRNet | | 73.6 | 76.4 | 61.1 | 55.0 | 63.6 | 64.8 | 52.7 | 64.5 |
| | HRNet | ✓ | 76.5 | 77.8 | 70.0 | 63.2 | 66.2 | 68.9 | 60.5 | 69.5(**+5.0**) |
| PoseTrack 2018 | Hourglass | | 65.6 | 66.5 | 52.6 | 45.5 | 50.9 | 51.0 | 42.4 | 54.3 |
| | Hourglass | ✓ | 67.7 | 68.0 | 58.9 | 50.1 | 53.3 | 53.6 | 48.5 | 57.9(**+3.6**) |
| | Pose_ResNet | | 66.4 | 68.2 | 52.4 | 46.1 | 56.3 | 55.5 | 47.0 | 56.7 |
| | Pose_ResNet | ✓ | 67.9 | 69.3 | 61.6 | 52.3 | 59.2 | 59.8 | 54.4 | 61.1(**+4.4**) |
| | HRNet | | 68.0 | 71.7 | 55.7 | 51.6 | 58.5 | 58.2 | 49.4 | 59.6 |
| | HRNet | ✓ | 69.8 | 74.7 | 69.7 | 62.3 | 63.9 | 64.7 | 59.2 | 66.6(**+7.0**) |

TABLE VI: Generalization of the proposed pose tracking pipeline with different backbones on PoseTrack 2018 validation set. Performance is evaluated by mAP (%) and MOTA (%) for pose estimation and pose tracking in the multi-person pose estimation and tracking task, respectively.

| Detector | ID-retrieve | Bbox-revision | mAP | MOTA |
|---|---|---|---|---|
| HTC [61] | | | 76.6 | 66.6 |
| | ✓ | | 76.6 | 67.3(**+0.7**) |
| | | ✓ | 78.1(**+1.5**) | 68.5(**+1.9**) |
| | ✓ | ✓ | 78.1(**+1.5**) | 69.2(**+2.6**) |
| YOLOv5 [63] | | | 76.3 | 66.9 |
| | ✓ | | 76.3 | 67.7(**+0.8**) |
| | | ✓ | 77.1(**+0.8**) | 68.0(**+1.1**) |
| | ✓ | ✓ | 77.1(**+0.8**) | 68.8(**+1.9**) |

TABLE VII: Effectiveness analysis of each component in Pose Tracking Pipeline by pose estimation performance in mAP (%) and pose tracking performance in MOTA (%) in multi-person pose estimation and tracking task on the validation set of PoseTrack datasets.

| Dataset | ID-retrieve | Bbox-revision | mAP | MOTA |
|---|---|---|---|---|
| PoseTrack 2017 | | | 80.0 | 64.5 |
| | ✓ | | 80.0 | 65.4(**+0.9**) |
| | | ✓ | 81.8(**+1.8**) | 67.4(**+2.9**) |
| | ✓ | ✓ | 81.8(**+1.8**) | 67.8(**+3.3**) |
| PoseTrack 2018 | | | 78.3 | 59.6 |
| | ✓ | | 78.3 | 60.5(**+0.9**) |
| | | ✓ | 80.1(**+1.8**) | 62.4(**+2.8**) |
| | ✓ | ✓ | 80.1(**+1.8**) | 62.7(**+3.1**) |

the overall MOTA metrics improve significantly by 5.0% and 7.0% on PoseTrack 2017 and PoseTrack 2018, respectively. Besides, our KCN outperforms the baseline model on all joints, especially for joints at the elbow and ankle on Pose-Track 2018 dataset, where the MOTA metrics improved by 14.0% and 9.8%, respectively. We also observed that on more challenging tracking areas, such as elbows, ankles, and wrists, the performance gets impressive improvement with our KCN. To further demonstrate the effectiveness of KCN, we count the keypoint numbers of wrong detection and missing detection on PoseTrack 2018 dataset, as shown in Table IV. Compared to the baseline method, our KCN shows an 11.6% reduction in false detection and a 9.0% reduction in missed detection, for a total reduction of 9.9%. We believe the reason is that

the baseline method only uses location probability, which may cause failure when filtering keypoints. For example, obscured keypoints may get high location probabilities as they are incorrectly labeled to other persons, thus they will be incorrectly detected; in the case of frame blurring, keypoints will get low location probabilities, thus they will be incorrectly filtered, resulting in missed detection.

We also verify the generalisability of KCN with three different top-down multi-person pose estimation networks, namely HRNet [31], Hourglass [24] and Pose_ResNet [47] on both PoseTrack 2017 and PoseTrack 2018 validation sets, as shown in Table V. The experimental results indicate that our method can stably improve the performance with different pose estimation networks and has strong generality. Our KCN can be plugged into top-down framework-based multi-person pose tracking methods to improve the performance of multi-person pose tracking.

***Analysis of Pose Tracking Pipeline.*** With two detectors, HTC [61] and YOLOv5 [63], built on top of our proposed pose estimation network KCN, we first evaluate the different components of our pose tracking pipeline and quantify how much each of them contributes to the final performance of our proposed method on the validation set of PoseTrack 2018, as shown in Table VI. As we can see, the ID-retrieve module does not enhance performance in multi-person pose estimation, while the Bbox-revision module effectively boosts the performance in both multi-person pose estimation and pose tracking tasks. We also observe that the ID-retrieve module improves comparable performance in MOTA with both detectors. Meanwhile, the improvement brought by Bbox-revision with HTC detector is significantly surpassing that with YOLOv5 detector in both mAP and MOTA. The reason may be that YOLOv5 detector has more missing detections than HTC detector.

Besides, we also evaluate the effectiveness of each component in PTP with our baseline pose estimation network. As shown in Table VII, similarly, both ID-retrieve and Bbox-revision modules boost the pose tracking performance on PoseTrack 2017 and PoseTrack 2018 validation sets. This amply demonstrates the effectiveness of our proposed ID-retrieve module and Bbox-revision.

***Effectiveness of KCN and PTP.*** In this part, we evaluate the contribution of two components, the Keypoints Confi-

TABLE VIII: Ablation study of components KCN and PTP on pose tracking. Trackers are evaluated by MOTA metric (%).

| Dataset | KCN | PTP | Head | Shou | Elb | Wri | Hip | Knee | Ankl | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| PoseTrack 2017 | | | 73.6 | 76.4 | 61.1 | 55 | 63.6 | 64.8 | 52.7 | 64.5 |
| | ✓ | | 76.5 | 77.8 | 70.0 | 63.2 | 66.2 | 68.9 | 60.5 | 69.5(+**5.0**) |
| | | ✓ | 76.7 | 79.7 | 64.8 | 59.4 | 67.0 | 67.3 | 55.5 | 67.8(+**3.4**) |
| | ✓ | ✓ | 79.5 | 81.2 | 72.8 | 66 | 69.5 | 70.8 | 61.8 | 72.2(+**7.7**) |
| PoseTrack 2018 | | | 68.0 | 71.7 | 55.7 | 51.6 | 58.5 | 58.2 | 49.4 | 59.6 |
| | ✓ | | 69.8 | 74.7 | 69.7 | 62.3 | 63.9 | 64.7 | 59.2 | 66.6(+**7.0**) |
| | | ✓ | 70.3 | 74.5 | 59.6 | 55.8 | 61.3 | 61.2 | 52.6 | 62.7(+**3.1**) |
| | ✓ | ✓ | 72.2 | 77.4 | 72.4 | 64.7 | 66.1 | 67.2 | 61.3 | 69.2(+**9.6**) |

TABLE IX: Comparison with state-of-the-art methods on multi-person pose estimation (with keypoint filtering) on the Validation sets of PoseTrack dataset in terms of mAP metric (%). '-' indicates that the result is not provided in the referred paper.

| Dataset | Method | Head | Shou | Elb | Wri | Hip | Knee | Ankl | Total |
|---|---|---|---|---|---|---|---|---|---|
| PoseTrack 2017 | BUTD [64] | 79.1 | 77.3 | 69.9 | 58.3 | 66.2 | 63.5 | 54.9 | 67.8 |
| | RPAF [65] | 83.8 | 84.9 | 76.2 | 64 | 72.2 | 64.5 | 56.6 | 72.6 |
| | ArtTrack [5] | 78.7 | 76.2 | 70.4 | 62.3 | 68.1 | 66.7 | 58.4 | 68.7 |
| | PoseFlow [66] | 66.7 | 73.3 | 68.3 | 61.1 | 67.5 | 67.0 | 61.3 | 66.5 |
| | STAF [54] | - | - | - | 65.0 | - | - | 62.7 | 72.6 |
| | ST-Embed [55] | 83.8 | 81.6 | 77.1 | 70.0 | 77.4 | 74.5 | 70.8 | 77 |
| | DAT [20] | 67.5 | 70.2 | 62.0 | 51.7 | 60.7 | 58.7 | 49.8 | 60.6 |
| | FlowTrack [47] | 81.7 | 83.4 | 80.0 | 72.4 | 75.3 | 74.8 | 67.1 | 76.9 |
| | TKMRNet [7] | 85.3 | 88.2 | 79.5 | 71.6 | 76.9 | 76.9 | **73.1** | 79.5 |
| | LDGNNTrack [21] | **88.4** | **88.4** | **82.0** | 74.5 | **79.1** | 78.3 | **73.1** | **81.1** |
| | Ours | 86.6 | 87 | 80.1 | **75.5** | 77.3 | **78.6** | 71.6 | 80.0 |
| PoseTrack 2018 | STAF [54] | - | - | - | 64.7 | - | - | 62 | 70.4 |
| | TML++ [67] | - | - | - | - | - | - | - | 74.6 |
| | TKMRNet [7] | - | - | - | - | - | - | - | 76.7 |
| | LightTrack [68] | - | - | - | - | - | - | - | 77.3 |
| | LDGNNTrack [21] | 80.6 | 84.5 | 80.6 | 74.4 | 75.0 | 76.7 | 71.9 | 77.9 |
| | SKCTrack [69] | - | - | - | - | - | - | - | **79.2** |
| | Ours | 80.6 | 85.3 | 80.7 | 74.3 | 76.1 | 76.7 | 71.9 | 78.1 |

TABLE X: Comparison with state-of-the-art methods on pure multi-person pose estimation (without keypoint filtering) on the validation sets of PoseTrack dataset in terms of mAP (%).

| Dataset | Method | Head | Shou | Elb | Wri | Hip | Knee | Ankl | Total |
|---|---|---|---|---|---|---|---|---|---|
| PoseTrack 2017 | SKCTrack [69] | 86.1 | 87.0 | 83.4 | 76.4 | 77.3 | 79.2 | 73.3 | 80.8 |
| | PoseWarper [70] | 81.4 | 88.3 | 83.9 | 78.0 | 82.4 | 80.5 | 73.6 | 81.2 |
| | CombDet [6] | 89.4 | 89.7 | 85.5 | 79.5 | 82.4 | 80.8 | 76.4 | 83.8 |
| | LDGNNTrack [21] | **90.9** | 90.7 | 86.0 | 79.2 | **83.8** | **82.7** | **78.0** | **84.9** |
| | Ours | 89.5 | **90.9** | **87.6** | **81.8** | 81.1 | 82.6 | 76.1 | 84.6 |
| PoseTrack 2018 | SKCTrack [69] | **86.0** | 78.3 | 84.8 | 78.3 | 79.1 | 81.1 | 75.6 | 82.0 |
| | KeyTrack [9] | 84.1 | 87.2 | 85.3 | 79.2 | 77.1 | 80.6 | 76.5 | 81.6 |
| | CombDet [6] | 84.9 | 87.4 | 84.8 | 79.2 | 77.6 | 79.7 | 75.3 | 81.5 |
| | LDGNNTrack [21] | 85.1 | 87.7 | 85.3 | 80.0 | **81.1** | **81.6** | **77.2** | 82.7 |
| | Ours | 85.1 | **88.9** | **86.4** | **80.7** | 80.9 | 81.5 | 77.0 | **83.1** |

dence Network (KCN) and Pose Tracking Pipeline (PTP), to the final performance of our method on both PoseTrack 2017 and 2018 datasets, as shown in Table VIII. It can be observed that the results improve substantially on both two datasets with our proposed KCN and PTP. In particular, KCN delivers significantly higher performance improvement compared to PTP. The model incorporating KCN and PTP obtain 12.0%∼16.1% performance gains over our baseline, which proves the effectiveness of our design.

### D. Comparison with State-of-the-art Methods

We compare our proposed method with the state-of-the-art methods in the multi-person pose estimation and tracking task on PoseTrack 2017 and PoseTrack 2018 datasets.

***Multi-person Pose Estimation.*** Table IX and Table X and Table XI show the comparisons between our proposed method and existing methods on multi-person pose estimation

task on the validation sets of PoseTrack 2017 and PoseTrack 2018 datasets. Table IX shows the results of multi-person pose estimation with filtering the low confidence keypoints for pose tracking. It can be observed that our approach outperforms the most competitive top-down approach, TKMRNet, by 0.5 mAP and outperforms the best bottom-up approach, ST-Embed, by 3.0 mAP on PoseTrack 2017 validation set. On PoseTrack 2018, our method achieves comparable results with the state-of-the-art method TKMRNet. Table X shows the results of multi-person pose estimation in videos where we evaluate the poses without filtering the low confidence keypoints. Our method can well recover the human instances that are missed by the detector in videos for multi-person pose estimation. It can be seen that our approach achieves the second best performance on PoseTrack 2017 and the best performance on PoseTrack 2018. Our method outperforms the most competitive approach, LDGNNTrack, by 0.4 mAP on PoseTrack 2018

TABLE XI: Comparison with state-of-the-art methods on multi-person pose tracking on the Validation sets of PoseTrack dataset in terms of MOTA (%). '-' indicates that the result is not provided in the referred paper.

| Dataset | Method | Head | Shou | Elb | Wri | Hip | Knee | Ankl | Total |
|---|---|---|---|---|---|---|---|---|---|
| PoseTrack 2017 | BUTD [64] | 71.5 | 70.3 | 56.3 | 45.1 | 55.5 | 50.8 | 37.5 | 56.4 |
| | ST-Embed [55] | 78.7 | 79.2 | 71.2 | 61.1 | 74.5 | 69.7 | 64.5 | 71.8 |
| | DAT [20] | 61.7 | 65.5 | 57.3 | 45.7 | 54.3 | 53.1 | 45.7 | 55.2 |
| | FlowTrack [47] | 73.9 | 75.9 | 63.7 | 56.1 | 65.5 | 65.1 | 53.5 | 65.4 |
| | PGPT [2] | - | - | - | - | - | - | - | 67.1 |
| | SKCTrack [69] | - | - | - | - | - | - | - | 68.3 |
| | CombDet [6] | 80.5 | 80.9 | 71.6 | 63.8 | 70.1 | 68.2 | 62.0 | 71.6 |
| | TKMRNet [7] | 81.0 | 82.9 | 69.8 | 63.6 | 72 | 71.1 | 60.8 | 72.2 |
| | LDGNNTrack [21] | **82.0** | **83.1** | **73.4** | 63.5 | **72.3** | **71.3** | **63.5** | **73.4** |
| | Ours | 79.5 | 81.2 | 72.8 | **66.0** | 69.5 | 70.8 | 61.8 | 72.2 |
| PoseTrack 2018 | STAF [54] | - | - | - | - | - | - | - | 60.9 |
| | TML++ [67] | **76.0** | 76.9 | 66.1 | 56.4 | 65.1 | 61.6 | 52.4 | 65.7 |
| | PT_CPN++ [8] | 68.8 | 73.5 | 65.6 | 61.2 | 54.9 | 64.6 | 56.7 | 64.0 |
| | LightTrack [68] | - | - | - | - | - | - | - | 64.9 |
| | KeyTrack [9] | - | - | - | - | - | - | - | 66.6 |
| | PGPT [2] | 75.4 | 77.3 | 69.4 | 71.5 | 65.8 | **67.2** | 59.0 | 68.4 |
| | CombDet [6] | 74.2 | 76.4 | 71.2 | 64.1 | 64.5 | 65.8 | **61.9** | 68.7 |
| | TKMRNet [7] | - | - | - | - | - | - | - | 68.9 |
| | SKCTrack [69] | - | - | - | - | - | - | - | 69.1 |
| | LDGNNTrack [21] | 74.3 | 77.6 | 71.4 | 64.3 | 65.6 | 66.7 | 61.7 | **69.2** |
| | Ours | 72.8 | **77.7** | **72.4** | **64.8** | **66.3** | **67.2** | 61.1 | **69.2** |

TABLE XII: Comparison with state-of-the-art multi-person pose estimation and tracking methods on PoseTrack 2017 Test set in terms of AP (%) and MOTA (%). Results are from PoseTrack 2017 Test Leaderboard.

| Method | Wrists AP | Ankles AP | Total AP | Total MOTA |
|---|---|---|---|---|
| JointFlow [47] | 53.1 | 50.4 | 63.4 | 53.1 |
| TML++ [67] | 60.9 | 56.0 | 67.8 | 54.5 |
| FlowTrack [47] | 71.5 | 65.7 | 74.6 | 57.8 |
| HRNet [31] | 72.0 | 67.0 | **75.0** | 57.9 |
| POINet [71] | 69.5 | 67.2 | 72.5 | 58.4 |
| KeyTrack [9] | 71.9 | 65.0 | 74.0 | 61.2 |
| CombDet [6] | 69.8 | 65.9 | 74.1 | **64.1** |
| Ours | **72.7** | **68.5** | 74.9 | 63.6 |

validation set.

Overall, excellence in performance on both PoseTrack 2017 and PoseTrack 2018 validation sets in two tasks, pose estimation with keypoint filtering and multi-person pose estimation without keypoint fully prove the effectiveness of our proposed method.

***Mulit-person Pose Tracking.*** We compare our method with state-of-the-art multi-person pose tracking methods on both PoseTrack validation sets and the test set. Since the PoseTrack 2018 test set is not available yet on the benchmark server, for the test split we only provide the comparison on PoseTrack 2017 dataset. As shown in Table XI and Table XII, our approach outperforms other methods and achieves the best performances on PoseTrack 2018 validation set. On PoseTrack 2017, our approach also achieves an excellent pose tracking performance on both validation set and test set.

## V. Conclusion

In this work, we present a confidence-based novel top-down approach for multi-person pose estimation and tracking task. Specifically, we propose a keypoint confidence network and a pose tracking pipeline. Compared to the previous methods, the proposed keypoint confidence network considers the availability probability when estimating keypoint confidence, while others only use the location probability. We also design a pose tracking pipeline with a Bbox-revision module and an ID-retrieve module to improve the tracking performance. The experimental results show that our approach achieves state-of-the-art performance on both PoseTrack 2017 and 2018 datasets.

## References

[1] M. Liu and J. Yuan, "Recognizing human actions as the evolution of pose estimation maps," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1159–1168.

[2] Q. Bao, W. Liu, Y. Cheng, B. Zhou, and T. Mei, "Pose-guided tracking-by-detection: Robust multi-person pose tracking," *IEEE Transactions on Multimedia*, vol. 23, pp. 161–175, 2020.

[3] Y. Wu, D. Kong, S. Wang, J. Li, and B. Yin, "An unsupervised real-time framework of human pose tracking from range image sequences," *IEEE Transactions on Multimedia*, vol. 22, no. 8, pp. 2177–2190, 2020.

[4] Z. Liu, Z. Lin, X. Wei, and S.-C. Chan, "A new model-based method for multi-view human body tracking and its application to view transfer in image-based rendering," *IEEE Transactions on Multimedia*, vol. 20, no. 6, pp. 1321–1334, 2018.

[5] E. Insafutdinov, M. Andriluka, L. Pishchulin, S. Tang, E. Levinkov, B. Andres, and B. Schiele, "Arttrack: Articulated multi-person tracking in the wild," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6457–6465.

[6] M. Wang, J. Tighe, and D. Modolo, "Combining detection and tracking for human pose estimation in videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 088–11 096.

[7] C. Zhou, Z. Ren, and G. Hua, "Temporal keypoint matching and refinement network for pose estimation and tracking," in *Proceedings of European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, vol. 12367. Springer, 2020, pp. 680–695.

[8] D. Yu, K. Su, J. Sun, and C. Wang, "Multi-person pose estimation for pose tracking with enhanced cascaded pyramid network," in *Proceedings of European Conference on Computer Vision Workshops*, ser. Lecture Notes in Computer Science, vol. 11130, 2018, pp. 221–226.

[9] M. Snower, A. Kadav, F. Lai, and H. P. Graf, "15 keypoints is all you need," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6738–6748.

[10] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *Proceedings of the Europe Conference on Computer Vision*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds., vol. 12346. Springer, 2020, pp. 213–229.

[11] T. Liang, X. Chu, Y. Liu, Y. Wang, Z. Tang, W. Chu, J. Chen, and H. Ling, "Cbnet: A composite backbone network architecture for object detection," *IEEE Transactions on Image Processing*, vol. 31, pp. 6893–6906, 2022.

[12] T.-I. Chen, Y.-C. Liu, H.-T. Su, Y.-C. Chang, Y.-H. Lin, J.-F. Yeh, W.-C. Chen, and W. Hsu, "Dual-awareness attention for few-shot object detection," *IEEE Transactions on Multimedia*, pp. 1–1, 2021, doi: 10.1109/TMM.2021.3125195.

[13] C. Zhang, Z. Li, J. Liu, P. Peng, Q. Ye, S. Lu, T. Huang, and Y. Tian, "Self-guided adaptation: Progressive representation alignment for domain adaptive object detection," *IEEE Transactions on Multimedia*, vol. 24, pp. 2246–2258, 2022.

[14] G. Cheng, J. Han, P. Zhou, and D. Xu, "Learning rotation-invariant and fisher discriminative convolutional neural networks for object detection," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 265–278, 2019.

[15] Z. Zou, Z. Shi, Y. Guo, and J. Ye, "Object detection in 20 years: A survey," 2019, arXiv preprint arXiv:1905.05055.

[16] Z.-Q. Zhao, P. Zheng, S.-T. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 11, pp. 3212–3232, 2019.

[17] T. Kong, F. Sun, H. Liu, Y. Jiang, L. Li, and J. Shi, "Foveabox: Beyond anchor-based object detection," *IEEE Transactions on Image Processing*, vol. 29, pp. 7389–7398, 2020.

[18] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.

[19] M. Tan, R. Pang, and Q. V. Le, "Efficientdet: Scalable and efficient object detection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 10 781–10 790.

[20] R. Girdhar, G. Gkioxari, L. Torresani, M. Paluri, and D. Tran, "Detect-and-track: Efficient pose estimation in videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 350–359.

[21] Y. Yang, Z. Ren, H. Li, C. Zhou, X. Wang, and G. Hua, "Learning dynamics via graph neural networks for human pose estimation and tracking," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2021, pp. 8074–8084.

[22] T. Pfister, J. Charles, and A. Zisserman, "Flowing convnets for human pose estimation in videos," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1913–1921.

[23] A. Toshev and C. Szegedy, "Deeppose: Human pose estimation via deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 1653–1660.

[24] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *Proceedings of Europe Conference on Computer Vision*, ser. Lecture Notes in Computer Science, vol. 9912. Springer, 2016, pp. 483–499.

[25] S.-E. Wei, V. Ramakrishna, T. Kanade, and Y. Sheikh, "Convolutional pose machines," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016, pp. 4724–4732.

[26] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7291–7299.

[27] H.-S. Fang, S. Xie, Y.-W. Tai, and C. Lu, "Rmpe: Regional multi-person pose estimation," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2334–2343.

[28] J. Song, L. Wang, L. Van Gool, and O. 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*, 2017, pp. 4220–4229.

[29] G. Papandreou, T. Zhu, N. Kanazawa, A. Toshev, J. Tompson, C. Bregler, and K. Murphy, "Towards accurate multi-person pose estimation in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4903–4911.

[30] L. Pishchulin, E. Insafutdinov, S. Tang, B. Andres, M. Andriluka, P. V. Gehler, and B. Schiele, "Deepcut: Joint subset partition and labeling for multi person pose estimation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4929–4937.

[31] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5693–5703.

[32] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler, "Efficient object localization using convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, 2015, pp. 648–656.

[33] J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *Advances in Neural Information Processing Systems 27*, 2014, pp. 1799–1807.

[34] A. Bulat and G. Tzimiropoulos, "Human pose estimation via convolutional part heatmap regression," in *Proceedings of European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, vol. 9911. Springer, 2016, pp. 717–732.

[35] J. Carreira, P. Agrawal, K. Fragkiadaki, and J. Malik, "Human pose estimation with iterative error feedback," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition,*, 2016, pp. 4733–4742.

[36] P. Hu and D. Ramanan, "Bottom-up and top-down reasoning with hierarchical rectified gaussians," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5600–5609.

[37] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, and K. Murphy, "Speed/accuracy trade-offs for modern convolutional object detectors," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3296–3297.

[38] U. Rafi, B. Leibe, J. Gall, and I. Kostrikov, "An efficient convolutional network for human pose estimation," in *Proceedings of the British Machine Vision Conference*, 2016.

[39] Y. Xu, J. Zhang, Q. Zhang, and D. Tao, "Vitpose: Simple vision transformer baselines for human pose estimation," *arXiv preprint arXiv:2204.12484*, 2022.

[40] A. Bulat, J. Kossaifi, G. Tzimiropoulos, and M. Pantic, "Toward fast and accurate human pose estimation via soft-gated skip connections," in *Proceedings of the IEEE International Conference on Automatic Face and Gesture Recognition*, 2020, pp. 8–15.

[41] H. Liu, W. Liu, Z. Chi, Y. Wang, Y. Yu, J. Chen, and T. Jin, "Fast human pose estimation in compressed videos," *IEEE Transactions on Multimedia*, pp. 1–1, 2022, doi: 10.1109/TMM.2022.3141888.

[42] G. Kim, H. Kim, K. Kong, J.-W. Song, and S.-J. Kang, "Human body-aware feature extractor using attachable feature corrector for human pose estimation," *IEEE Transactions on Multimedia*, pp. 1–11, 2022.

[43] G. Ning, Z. Zhang, and Z. He, "Knowledge-guided deep fractal neural networks for human pose estimation," *IEEE Transactions on Multimedia*, vol. 20, no. 5, pp. 1246–1259, 2018.

[44] L. Pishchulin, A. Jain, M. Andriluka, T. Thormählen, and B. Schiele, "Articulated people detection and pose estimation: Reshaping the future," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 3178–3185.

[45] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.

[46] Y. Chen, Z. Wang, Y. Peng, Z. Zhang, G. Yu, and J. Sun, "Cascaded pyramid network for multi-person pose estimation," in *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7103–7112.

[47] B. Xiao, H. Wu, and Y. Wei, "Simple baselines for human pose estimation and tracking," in *Proceedings of the Europe Conference on Computer Vision*, 2018, pp. 466–481.

[48] G. Gkioxari, B. Hariharan, R. B. Girshick, and J. Malik, "Using k-poselets for detecting people and localizing their keypoints," in *Proceedings of the IEE Conference on Computer Vision and Pattern Recognition,*, 2014, pp. 3582–3589.

[49] X. Chen and A. L. Yuille, "Parsing occluded people by flexible compositions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3945–3954.

[50] E. Insafutdinov, L. Pishchulin, B. Andres, M. Andriluka, and B. Schiele, "Deepercut: A deeper, stronger, and faster multi-person pose estimation model," in *Proceedings of the European Conference on Computer Vision*, ser. Lecture Notes in Computer Science, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., vol. 9910. Springer, 2016, pp. 34–50.

[51] U. Iqbal and J. Gall, "Multi-person pose estimation with local joint-to-person associations," in *Proceedings of the Europe Conference on Computer Vision Workshops*, ser. Lecture Notes in Computer Science, vol. 9914, 2016, pp. 627–642.

[52] S. Kreiss, L. Bertoni, and A. Alahi, "Pifpaf: Composite fields for human pose estimation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11 977–11 986.

[53] K. Sven, B. Lorenzo, and A. Alexandre, "Openpifpaf: Composite fields for semantic keypoint detection and spatio-temporal association," *IEEE*

*Transactions on Intelligent Transportation Systems*, vol. 23, no. 8, pp. 13 498–13 511, 2022.

[54] Y. Raaj, H. Idrees, G. Hidalgo, and Y. Sheikh, "Efficient online multi-person 2d pose tracking with recurrent spatio-temporal affinity fields," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4620–4628.

[55] S. Jin, W. Liu, W. Ouyang, and C. Qian, "Multi-person articulated tracking with spatial and temporal embeddings," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5664–5673.

[56] H. Guo, T. Tang, G. Luo, R. Chen, Y. Lu, and L. Wen, "Multi-domain pose network for multi-person pose estimation and tracking," in *Proceedings of the Europe Conference on Computer Vision Workshops*, ser. Lecture Notes in Computer Science, vol. 11130. Springer, 2018, pp. 209–216.

[57] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[58] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the International Conference on Computer Vision*, 2017, pp. 2980–2988.

[59] "Object Keypoint Similarity," https://cocodataset.org/#keypoints-eval, accessed: 2023-05-12.

[60] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *Proceedings of the Europe Conference on Computer Vision*, ser. Lecture Notes in Computer Science, vol. 12347. Springer, 2020, pp. 402–419.

[61] K. Chen, J. Pang, J. Wang, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Shi, W. Ouyang *et al.*, "Hybrid task cascade for instance segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4974–4983.

[62] K. Chen, J. Wang, J. Pang, Y. Cao, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Xu, Z. Zhang, D. Cheng, C. Zhu, T. Cheng, Q. Zhao, B. Li, X. Lu, R. Zhu, Y. Wu, J. Dai, J. Wang, J. Shi, W. Ouyang, C. C. Loy, and D. Lin, "MMDetection: Open mmlab detection toolbox and benchmark," *arXiv preprint arXiv:1906.07155*, 2019.

[63] G. Jocher, A. Stoken, J. Borovec, NanoCode012, ChristopherSTAN, L. Changyu, Laughing, tkianai, A. Hogan, lorenzomammana, yxNONG, AlexWang1900, L. Diaconu, Marc, wanghaoyang0106, ml5ah, Doug, F. Ingham, Frederik, Guilhen, Hatovix, J. Poznanski, J. Fang, L. Yu, changyu98, M. Wang, N. Gupta, O. Akhtar, PetrDvoracek, and P. Rai, "ultralytics/yolov5: v3.1 - Bug Fixes and Performance Improvements," Oct. 2020. [Online]. Available: https://doi.org/10.5281/zenodo.4154370

[64] S. Jin, X. Ma, Z. Han, Y. Wu, W. Yang, W. Liu, C. Qian, and W. Ouyang, "Towards multi-person pose tracking: Bottom-up and top-down methods," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, vol. 2, no. 3, 2017, p. 7.

[65] X. Zhu, Y. Jiang, and Z. Luo, "Multi-person pose estimation for posetrack with enhanced part affinity fields," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, vol. 7, 2017, p. 4321.

[66] Y. Xiu, J. Li, H. Wang, Y. Fang, and C. Lu, "Pose flow: Efficient online pose tracking," *arXiv preprint arXiv:1802.00977*, 2018.

[67] J. Hwang, J. Lee, S. Park, and N. Kwak, "Pose estimator and tracker using temporal flow maps for limbs," in *Proceedings of the IEEE International Joint Conference on Neural Networks*, 2019, pp. 1–8.

[68] G. Ning, J. Pei, and H. Huang, "Lighttrack: A generic framework for online top-down human pose tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition workshops*, 2020, pp. 1034–1035.

[69] U. Rafi, A. Doering, B. Leibe, and J. Gall, "Self-supervised keypoint correspondences for multi-person pose estimation and tracking in videos," in *Proceedings of the Europe Conference on Computer Vision*. Springer, 2020, pp. 36–52.

[70] G. Bertasius, C. Feichtenhofer, D. Tran, J. Shi, and L. Torresani, "Learning temporal pose estimation from sparsely-labeled videos," *arXiv preprint arXiv:1906.04016*, 2019.

[71] W. Ruan, W. Liu, Q. Bao, J. Chen, Y. Cheng, and T. Mei, "Poinet: pose-guided ovonic insight network for multi-person pose tracking," in *Proceedings of the ACM International Conference on Multimedia*, 2019, pp. 284–292.