# Commentary: Hierarchical Clustering and Refinement for Generalized Multi-Camera Person Tracking

Zeal Liang

z5325156

## 1. Introduction

In this paper the authors propose a new framework for Multi-Target Multi-Camera Tracking (MTMCT). This framework has great potential in surveillance scenarios. However, the task is extremely challenging due to factors such as different viewing angles, severe occlusions and lighting variations. The main goal of the authors is to address these challenges in order to improve the accuracy and robustness of MTMCT.

There are many issues in the field of MTMCT. First, with the development of virtual surveillance scenarios, applying MTMCT to diverse surveillance scenarios is increasingly worth exploring. Second, target tracking within a single camera can only capture short-term trajectories, is not robust enough to associate a target with a single trajectory, and is susceptible to changes in viewing angle, illumination, and image resolution. In public places, the same identity tends to appear and disappear multiple times under the same camera, which results in multiple trajectories for the same identity under the same camera. Therefore, the authors propose a hierarchical clustering and refinement framework in the hope of reducing false matches through progressive clustering and refinement strategies [1].

There are several important aspects to addressing these issues. First, in real-world surveillance applications, it is critical to ensure accuracy and consistency in tracking people across different cameras, which can improve the effectiveness of public safety systems. Second, with the rise of virtual surveillance scenarios, extending MTMCT to virtual environments can provide an integrated solution for both real and virtual world surveillance. By introducing hierarchical clustering and refinement, the authors are aiming to overcome the problems in the existing methods and thus bring a positive impact on the development of the MTMCT field.

If these issues are effectively addressed, there will be multiple benefits. Surveillance system operators would have access to more accurate and consistent personnel tracking information, leading to improved emergency response and security. In practical applications, it can also provide richer information for crime detection, incident analysis, and so on. In addition, the successful implementation of this method will change the current practice approach in the field of MTMCT and provide new ideas and directions for future research and applications.

## 2. Methods

The authors' proposed framework contains four main components: person detection, person re-recognition, single-camera tracking, and multi-camera tracking. The overall process can be summarized as follows:(1) obtaining a person bounding box from each camera view using a person detector; (2) extracting pedestrian re-recognition features from each bounding box by using a pre-trained re-recognition model; (3) generating single-camera trajectory segments for each camera using a single-camera tracking model; (4) clustering the trajectory segment features to correlate intra-camera trajectory segments; (5) Apply clustering methods to associate inter-camera trajectory segments; (6) Refine inter-camera trajectory segments using appearance, spatio-temporal and face constraints [1].

Person detection is the first step in multi-camera person tracking, so choosing a reliable detector is crucial. The authors have chosen to use YOLOv5 as a detector [1]. It is a highly accurate and efficient target detection algorithm for a variety of application scenarios. However, the paper does not discuss possible drawbacks, such as the fact that YOLOv5 may be poor at detecting small-sized targets, which may affect tracking accuracy in some cases. Therefore, I suggest that the authors should try several different people detection algorithms and give comparisons to make sure that the most suitable method has been chosen.

Person re-identification plays a crucial role in multi-camera person tracking. The authors chose the TransReID-SSL based person re-identification model, which is pre-trained on the LUperson dataset and has the ability to extract robust and domain invariant ReID features. By fine-tuning the model on multiple datasets, the authors further improved its performance. However, this method may also be affected by mismatches during the re-recognition process, especially in cross-camera tracking where ID switching problems may occur. As well, the paper does not explicitly state the domain differences between these datasets and how the domain bias can be addressed.

The authors use ByteTrack as a single-camera tracking algorithm that is able to associate each detection frame with a unique identity by considering motion information and visual similarity. In multi-camera tracking, the authors use clustering methods to associate intra-camera tracklets and employ appearance, spatio-temporal, and face constraints to refine cross-camera tracklets [1].However, I noticed that the paper does not discuss whether the ByteTrack algorithm is able to maintain its efficiency and accuracy even in dense scenarios, so more validation and analysis may be needed.

For multi-camera person association, the authors used the K-means clustering algorithm to group the aggregated track segments. However, K-means clustering is very sensitive to noise and outliers and may require parameter tuning based on data distribution. The authors mention a weighting strategy to solve some of the problems, but do not provide a specific

method for selecting the weights.

Overall, although this study proposes a framework for MTMCT, I believe there is still room for improvement in terms of method selection, parameter tuning, and method limitations. The authors could have further explored the rationality of the method selection, provided more experimental and comparative results, as well as discussed the method limitations and possible directions for improvement.

## 3. Results

The authors conduct an extensive experimental analysis of the proposed MTMCT framework, focusing on the use of IDF1, IDP and IDR metrics to assess the effectiveness of the proposed framework. Implementation details, weakening analysis, and comparisons with other models are discussed.

The dataset used for the experiments is comprehensive and includes 1,491 minutes of high-resolution video data from 130 cameras. The authors utilized both real and synthetic data, as well as also using the publicly available Person ReID dataset for model training, further enhancing the robustness of the evaluation.

The quantitative analysis presents a weakening study of the hierarchical clustering and refinement strategies, demonstrating their impact on the overall performance. The results show that the proposed hierarchical clustering strategy (both intra-camera and inter-camera clustering) significantly improves the performance. In addition, various refinement strategies also positively impacted the results, where appearance, spatio-temporal, and trajectory-level refinements improved performance by 2%, 1%, and 2%, respectively. The performance improvements of these strategies are progressive and convincingly demonstrated.

The authors' decision to choose TransReID-SSL as the backbone network is supported by the weakening analysis in different backbone networks. The fine-tuned TransReID-SSL achieves the best results in different backbone networks, contributing to the overall credibility of the paper's findings.

The authors also provide a visualization of the final MTMCT results, enhancing the clarity of the effects of their method. The trajectories visualized across cameras and time illustrate that the model does indeed match trajectories effectively, even across cameras with different viewing angles and occlusions. In addition, the results of the cross-cluster analysis further illustrate the model's ability to manage situations with body overlap and occlusion, highlighting the robustness of the method.

The authors compare the model's performance with other teams, demonstrating that their proposed method achieved a competitive fifth place with an IDF1 score of 0.921 on Track1 of the AIcity2023 challenge. This comparison adds support to the effectiveness of their proposed method.

In conclusion, this paper successfully presents a robust experimental analysis of the proposed MTMCT framework. The validity of the proposed methodology is fully confirmed through rigorous evaluation of multiple metrics, weakening studies, and comparisons with other models. To further enhance the credibility of the results, I suggest that future experiments could explore testing the performance of the framework under a variety of challenging conditions, such as varying light, occlusion, and crowd density. In addition, considering the growing concern for privacy-preserving AI, the paper could discuss how the proposed approach handles sensitive information and complies with privacy regulations.

## 4. Conclusions

Taken together, the MTMCT framework proposed in this paper has achieved remarkable results in many aspects. However, its main advantages and potential shortcomings still need to be comprehensively analyzed in a comprehensive evaluation of the whole work.

First of all, the framework in this paper proposes a series of innovative approaches to the multi-camera tracking problem, such as hierarchical clustering and refinement strategies, as well as the introduction of appearance, spatio-temporal and facial constraints. Through rich experimental analyses, it has been demonstrated that these methods are significantly effective in improving the accuracy and robustness of multi-camera person tracking.

However, in exploring the potential pitfalls of the work, I note a number of aspects that remain to be further explored. For example, the performance of the ByteTrack algorithm in dense scenarios and the detection problem for different target sizes mentioned in the paper need more experimental verification and analysis.

In future research, the following aspects can be considered to further deepen. First, more experimental scenarios can be explored, including the performance performance under different situations such as light change, occlusion and crowd density. In addition, as privacy protection is a growing concern, the authors can also consider how to provide effective tracking while ensuring the handling of sensitive information and compliance with privacy regulations.

In addition, considering the latest technological advances, I suggest incorporating the newly introduced YOLOv8 model into the MTMCT framework in future studies.YOLOv8 is optimized in terms of accuracy and speed, which may bring better performance to the whole system.

Overall, the new MTMCT framework proposed by the authors is indeed of great value and provides some insights for research in this field. However, there is still room for further improvement in terms of rationality of method selection, exploration of parameter tuning, and discussion of method limitations.

# References

[1] Zongyi Li, Runsheng Wang, He Li, Bohao Wei, Yuxuan Shi, Hefei Ling, Jiazhong Chen, Boyuan Liu, Zhongyang Li, Hanqing Zheng; Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops, 2023, pp. 5519-5528