An Enhanced Multiview Person Tracking Model with YOLOv5s Combined Hungarian Matching Operations

Prof Prashant V. Ingole Department of Information Technology Prof. Ram Meghe Institute of Technology and Research, Sant Gadge Baba Amravati University, Amravati Ashish D. Thete Post Graduate Scholar Department of Information Technology, Prof. Ram Meghe Institute of Technology and Research, Sant Gadge Baba Amravati University, Amravati

Abstract: A rising demand for strong multi-camera person tracking in surveillance, retail analytics, and smart cities is being emphasized with real-time accurate and robust identity association across views affected in the process. Most of the current solutions suffer high computation costs or low generalization to appearance changes, and many of them have a limited scalability to MultiView configurations. Furthermore, many methodologies adopt complicated architectures that are unsuitable for real-time use or for the edge environment. In concern to these, present work lightweight, modularize the pipeline for cross-view person detection, re-identification, and tracking with front-top synchronized videos. The system comprises four key modules optimized for speed, accuracy, and interpretability: at First, YOLOv5s trains per frame detection for fast, compact object identifier trained on the COCO dataset. Speedaccuracy trade-off enables real-time inference. Active feature-extracting vector cores, OSNet x0.25 extract 512-dimension feature vectors, while input is using resized person crops of 128×256. Efficiently cross-view appearance representation is maintained by the model while this section extracts scaled feature vectors. Next, identity association is done using cosine similarity of feature vector closeness and then optimally hooking globally by Hungarian algorithm. An empirical similarity threshold (0.7-0.75) filters out spurious matches, while achieving a mean matching accuracy of ~88%. Centroid-based tracking algorithms measure movements per ID in the top view, with dynamic distance thresholds that preserve temporal continuity and minimize ID switches. The system achieves >92% trajectory completeness with real-time operations, making it suited for indoor surveillance, retail behavior analysis, and crosscamera tracking. Modular design allows deployment on edge devices and generalization to diverse environments with minimal adaptations.

Keywords: MultiView Tracking, YOLOv5s, Person Re Identification, OSNet, Hungarian Algorithm, Scenarios

1. INTRODUCTION

The increasing deployment of intelligent surveillance systems in public squares, transport terminals, and retail environments has made the need for robust MultiView person tracking absolutely important Effective person tracking under these conditions needs that detection, association of identities, and trajectory estimation works very accurately even in the presence of significant viewpoint and appearance changes. Almost all single View tracking approaches fail in a few conditions of occlusions or different lighting conditions and perspective distortions in individuals' activity. Also, existing multi-camera tracking systems are either computationally heavy or use largescale models not suitable for real-time deployment [1, 2, 3] and have error propagation from tightly coupled modules. In light of the above, present work proposes such an improved modular pipeline that shall combine efficiency, accuracy, and scalability for person detection, appearance-based re-identification, and cross-view trajectory estimations. The system employs YOLOv5s, a lightweight convolutional neural network well-known for real-time object detection, to detect persons in each frame. In addition, compact and discriminative ReID embeddings are generated using OSNet x0.25, an omni-scale feature learning architecture. The cosine similarity, a metric that is robust to scale and illumination, compares those embeddings and is followed by Hungarian matching, a global optimization algorithm that guarantees one-to-one identity association across views during processing operations. A signature feature of the model is the integration of a centroid-based tracking algorithm that estimates top view movement trajectories, along which temporal continuity is even tighter because of low ID switch rates. The entire pipeline is then tuned for effective deployment in edge environments while providing very high matching accuracy in real time scenarios. This research brings scalable and interpretable solutions that bridge the existing gaps between computational efficiency and accuracy of MultiView identity association to define a new horizon in the practical surveillance and analysis systems for behavior sets.

2. REVIEW OF EXISTING MODELS FOR MULTIVIEW PERSON ANALYSIS

The growth of multicamera person ReIDs and tracking domains can be easily explained as they cater to the recent demands for intelligent surveillance and behavior analysis systems. Several works have attempted to improve reliability, transferability, and scalability of ReID systems, especially when the variables are viewpoint, light, occlusion, and intercamera spatial separation. The review by Li et al. of typical deep-learning-based Person ReID models are discussed to target more exact situations, including the types of variation (e.g., spatial difference for parameterized models seems to provide stable results) and bias necessitating the use of invariance descriptors and view-regularizing tracking models [1]. The kernel of any deep learning model development is in deep feature extraction, which has had the attention of certain recent developments and progress using multi-attention mechanisms [2] and omni-scale feature learning [8], claiming to strengthen identity discrimination. Wang et al. [8] specifically designed Open Set Network (OSNet) employing varied receptive fields on part paralleled convolutional streams of images to absorb omni-scale interactions, proving

effective in constrained excited environments. The architecture has sent off ripples to the design of lightweight systems for real-time situation—this was respected as a recent choice in the model. Accelerated deployment on hardware was demonstrated by Tapuhi et al. [3] by integrating multiple people ReID pipelining on the AI-focused edge accelerator, Hailo-8, which again adds weight to the vitality of computational efficiency. Similarly, Elgendy et al. [2] and Gu et al. [15] came up with methodologies enhancing robustness in dynamic industrial scenes where open-set recognition and never-seen identities are big challenges. Their works exalt generalized model-rules, and this is why the Law proposed certain elements-beyond-threshold matching and a modular deployment scheme came into the picture.

A suite of multi-target/multi-object-tracking methods, built on deep learning, are nothing but intensifying. Li et al. [4], for instance, floated a hybrid mechanism in the direction of detection-tracing, while Ma et al. [5] focused on putting the 3D spatial information to accommodate an ongoing initialization mode for Object Tracking in MultiView scenarios. The current work relates to [5] in suggesting, rather conceptually, additional facts about the use of top View for spatiotemporal anchoring. Model operation has been completed: That, particularly, represents an endeavor into the temporal side of color-matching theory. Zhao et al. This spatial-temporal relations-aware attention mechanism is brought in, enhancing the model's capacity to keep identities very consistent over time. Zhang et al. [10] developed the proxy anchor loss for deep metric learning, which enables the handling of cases that have very high intra-class variance and suppressed inter-class similarity in multi-cam setting. These investigations clearly assert that using learned similarity is crucial; however, the implementable version presented in this work favors deterministic cosine similarities for simple interpretability. Many recent benchmarks and surveys greatly contribute to the normalization of evaluation criteria, thereby fostering reproducibility. For example, Baiju [11] and Dijkinga [12] concentrated on scrutinizing the state-of-the-art multi-camera ReID techniques, particularly exploiting attention-based fusion and feature aggregation and hierarchical comparison pipelines, somehow shaping the proposed model into the architectural separation of the detection, embedding, and association modules. Of course, their ideas have given a lot of value to the idea of modularity and the need for a system that could easily produce interpretable outputs. The need for robust datasets to apply any scientific research cannot be overemphasized. Yang et al. [9] and Gu et al. [15] envisaged open-set ReID issues that occur in operating environments; hence projects should address models that accept unknown or rarely seen identities in their process. Although the adopted model neither explicitly addresses open-set scenarios, it nevertheless gives consideration to an identity-stop mechanism with a similarity threshold and is based on empirically-attuned rejection logic fueling the filtering of low-confidence matches. The combination of detection and embedding comes with Chen et al.'s unified ReID-detection framework YOLO ReIDNet. Unified systems may be very closely aligned with each other, but in reality, they tend to make it hard to separate for any other specialized utilization in different environments. To this end, the proposed system prefers loosely coupled modules, YOLOv5s for detection can be combined with OSNet for learning while using Hungarian matching for assignment, enabling the independent optimization of each stage and realtime updates to individual components.

In conclusion, Khan et al. detailed scale-based implementation of CNNs for MultiView tracking, which

brought into being a dynamic adaptation of receptive fields across various spatial perspectives [7]. This matches OSNet very well with its omni-scale feature learning, validating its suitability, as the ReID backbone, for the proposed architecture. Overall, the current literature builds a portfolio of varied strategies from attention mechanisms to metric learning, and spatiotemporal modeling, with the possibility of observing edge-aware deployment and benchmark designs. The proposed model seeks to build on these aspects with lightweight, interpretable, and highly performative parts, each in full sync with MultiView synchronization events, not only for modular analysis but also for estimating real-time trajectories, thus allowing the community to share in the pools. With its scalable solutions, this model can claim to nest somewhere between the strictness of academia and pragmatism other industrial applications often bring.

3. PROPOSED MODEL DESIGN ANALYSIS

The proposed model incorporates a sum of carefully selected modules to achieve MultiView person detection, cross-View ReIdentification, identity association, and trajectory estimation in real-time environments. The design philosophy focuses on modularity, computational efficiency, and robustness against inter-View appearance chang.



Figure 1. Two camera view one for face identification and another for tracking on ceiling

Each module in the system complements others by tackling the different stages of the pipeline, thus allowing the optimization of performance while setting architectural interpretability apart. The operation begins, as shown in Fig. 2, with synchronized video frames from two spatially distinct viewpoints on the front view and top view. Each frame enters a YOLOv5s-based person detector for the process. YOLOv5s uses a CSPDarknet53 backbone and PANet neck structure, producing dense feature maps with convolutional layers to predict bounding boxes B =

 $\{bi\}$ where each $bi \in \mathbb{R}^5$ consists of a class label, centre coordinates, and box dimensions. Once a score for every bounding box is calculated, a Non-Maximum Suppression (NMS) is utilized to delete boxes with lower scores in the process. Scoring via "(3.1)",

$$bi = argmax\{bj \in B, IoU(bi, bj) < \theta\}(\sigma(bj))...(3.1)$$

Where $\sigma(bj)$ is the object confidence score and θ is the IoU threshold typically set to 0.5 for this process. Hence, during the process, only confidently detected, non-overlapping detections are retained. Each detected bounding box is cropped and resized to a standard resolution of 128×256 and then passed to the OSNet x0.25 feature extractor. OSNet encapsulates multiple convolutional streams, producing omni-scale representations $fi \in \mathbb{R}^{512}$, where each vector is L2 Normalized Via "(3.2)",

$$\hat{f}i = \frac{fi}{\parallel fi \parallel} \dots (3.2)$$

Which guarantees that comparisons between features are scale invariant and robust to brightness and contrast changes in the views in process. To associate identities between views, a similarity matrix $S \in \mathbb{R}^{\vee}\{M \times N\}$ is constructed between front view features $\{fi^{\circ}F\}$ and top view features $\{fj^{\circ}T\}$ using cosine similarity. The similarity between a front view feature $fi^{\circ}F$ and a top view feature $fj^{\circ}T$ is expressed Via "(3.3)",

$$Sij = \frac{\langle \hat{f}i'F, \hat{f}j'T \rangle}{\|\hat{f}i'F\|^{2} \cdot \|\hat{f}j'T\|^{2}} = \hat{f}i'F \cdot \hat{f}j'T \dots (3.3)$$

To resolve for optimal identity assignments, the Hungarian algorithm would thus be applied to negative similarity matrix -S for the fact that this algorithm minimizes cost sets. Let the assignment matrix be $A \in \{0,1\}$ ' { $M \times N$ } and in that case, the assignment is defined Via "(3.4)",

$$A = argmin\{A \in \mathcal{A}\}\sum Aij(-Sij)...(3.4)$$

Where \mathcal{A} , refers to the set of valid one-to-one assignment matrices. For the purpose of reducing falsepositive matches, those matches with $Sij < \tau$ are rejected, where the threshold τ has been empirically set anywhere between 0.7 and 0.75 to provide the best trade-off between precision and recall. Trajectories are estimated by obtaining centroids $C_t k = (xc'k(t), yc'k(t))$ for every top view detection at timestamp 't'. The Euclidean norms measure the distance between centroids across frames. A match between two centroids C(t, i) and C(t+1,j) is accepted if the Identity Represented Via equation "(3.5)" is fulfilled,

$\| C(t,i) - C(t+1,j) \|^2 < \lambda D \dots (3.5)$

Where, *D* is the image diagonal length and $\lambda \in [0.05, 0.1]$ is a dynamic threshold factor depending on scene scales. This aids in the continuity of ID tracking while

still allowing the rejection of spurious associations. To smooth the trajectory and reduce jitter, a Gaussian Weighted moving average was executed over the past *n* frames. The smoothed position \tilde{C}_t is computed Via "(3.6)",

$$\tilde{C}_{t} = \left(\frac{1}{Z}\right) \sum wk * C(t-k), wk = exp\left(-\frac{k^{2}}{2\sigma^{2}}\right) \dots (3.6)$$

Where, *Z* is a normalization factor which guarantees $\sum wk = 1$, and σ controls the extent of smoothing in process. A differentiable trajectory energy function is defined for assessing the trajectory smoothness and integrity over temporal instances—formulated Via "(3.7)".

$$E = \int \left(\| \frac{d^2 \tilde{C}_t}{dt^2} \|^2 + \alpha \| \frac{d \tilde{C}_t}{dt} \|^2 \right) dt \dots (3.7)$$

It penalizes both high curvature (acceleration) and velocity spikes in the estimated paths, where the parameter α is a trade-off parameter balancing trajectory stability against responsiveness to motion. The design of this model attempts to equally balance functional precision with architectural simplicity. YOLOv5s detector uses very mild computation for fast and precise localization. OSNet, being scale sensitive and lightweight, thus gives out highly discriminative features, best suited for appearance Invariant ReID applications. The cosine Hungarian identity association thus gives mathematically optimal bipartite matching, while the centroid-based tracker allows real-time trajectory estimation with low latency sets. These schemes, complemented together, thereby form an interpretable and effective MultiView tracking pipeline, suitable for research and deployment in a constrained environment. The next discussion deals with an Iterative Validation evaluation of the proposed model under different scenarios.



Figure 2. Model Architecture of the Proposed Analysis Process

4. COMPARATIVE RESULT ANALYSIS

In this section, we evaluate the proposed MultiView tracking system over various contextual datasets. The proposed model's performance was assessed in a series of experiments controlled and partially controlled against three baseline methods Method [3], Method [8], and Method [15], each representing a traditional method of multi-camera person tracking. These methods all differ in detection architecture, feature encoding, and identity association strategies. The analysis performed, compared on basis of detection accuracy, matching precision, trajectory quality, and runtime efficiency. The experiments were conducted on synchronized dualview video datasets emulating three real surveillance scenarios; indoor retail, airport terminal, and a controlled laboratory setting, with identity ground truths and movement trajectories annotated in each dataset. The front and top view cameras were calibrated and synchronized for frame-level alignment. All models were executed on a workstation powered by an NVIDIA RTX 3090 GPU under PyTorch 1.12.1 on CUDA. Detection was set to work at 30 FPS with the evaluation metrics including Mean Matching Accuracy (MMA), False Match Rate (FMR), Identity Switches (IDS), Trajectory Completeness (TC), Average

Precision (AP), and average runtime per frame (RT) Sets.

Table 1. Person Detection Performance (YOLOv5s vs
Baselines)

Methods	Preci sion (%)	Rec all (%)	AP@0 .5 (%)	AP@0. 75(%)	Runti me (ms/f rame)
Proposed Method	91.2	88.5	90.1	84.3	12.5
Method by Tapuhi, T., Klinger, A., & Sholev, O. [3]	87.5	83.2	85.8	78.4	18.7
Method by Khan, S. [8]	84.1	80.0	81.9	75.0	21.4
Method by Gu, X [15]	88.0	82.5	86.1	79.2	16.3

The proposed YOLOv5s-based detector performs the best in terms of detection accuracy and runtime per frame, thereby making it more suitable for real-world applications. It affords a better trade-off between precision and recall in comparison to other methods with the maximum AP values over the IoU thresholds.

Table 2. Re Identification Accuracy usi	ing OSNet vs
Alternatives	

Methods	Top-1 Accuracy (%)	Top-5 Accuracy (%)	mAP (%)	FMR (%)
Proposed	87.4	94.3	85.1	6.8
Method by Tapuhi, T., Klinger, A., & Sholev, O. [3]	78.9	88.2	76.4	12.5
Method by Khan, S. [8]	80.1	89.7	78.3	11.1
Method by Gu, X [15]	83.7	91.5	81.0	9.4

The OSNet x0.25-based ReID module provides better performance than the deeper and bulkier networks utilized in other approaches. Omni-scale feature learning promotes much better identity discrimination under viewpoint variation while remaining efficient in inference sets.



Figure 3. Model's Integrated Result Analysis

Table 3. Cross View Identity Matching (Cosine + Hungarian)

Methods	Mean	ID	Assignmen	Threshol
	Matchi	Switche	t Precision	d Used
	ng	S	(%)	
	Accura			
	cy (%)			
Proposed	88.3	1.3	91.6	0.72
Method	75.4	3.6	81.0	0.65
by				
Tapuhi,				
Т.,				
Klinger,				
A., &				
Sholev, O.				
[3]				
Method	78.9	2.9	84.7	0.68
by Khan,				
S. [8]				
Method	83.5	2.0	88.2	0.70
by Gu, X				
[15]				

Deterministic and accurate identity assignment is ensured by using cosine similarity and the Hungarian algorithm in the proposed model sets. A lower ID switch rate and higher assignment precision show robustness to misdetections and occlusions.

Table 4. Top View Trajectory Estimation Metrics

Methods	Trajector y Complete ness (%)	Average ID Switche s	Temporal Continuit y (s)	Spatia l Drift (px)
Proposed	93.1	1.2	32.5	4.3
Method by Tapuhi, T., Klinger, A., & Sholev, O. [3]	84.5	3.8	21.6	9.7
Method by Khan, S. [8]	86.2	2.9	24.2	7.3
Method by Gu, X [15]	89.0	2.1	27.9	6.0

The centroid tracker with distance resets guarantees highly complete trajectories and sets of long-term temporal stability sets. The proposed model prevents spatial drift and ID fragmentation, especially in top View monitoring contexts.

Scenario	Method	MMA	TC	IDS	RT
		(%)	(%)		(ms)
Indoor	Proposed				
Lab		88.1	93.0	1.0	12.1
	Method by Tapuhi, T., Klinger, A., & Sholev, O. [3]	75.6	83.0	3.5	18.8
	Method by Khan, S. [8]	78.2	86.1	2.8	20.6
	Method by Gu, X [15]	83.0	89.1	2.0	16.4

Table 5. Multi-Scenario Benchmark (Indoor Lab)

The proposed system outperforms all scenarios in accuracy, stability, and speed parameters, thus confirming the model's generalizability under different surveillance settings

Table 6. Ablation Study on Key Components

Configuration	MMA (%)	FMR (%)	IDS	RT (ms)
Full Model	88.3	6.8	1.3	12.5
Without OSNet (Replaced with ResNet-50)	81.2	10.9	2.8	15.6
Without Hungarian (Greedy Matching)	83.1	9.2	2.3	11.7
Without Centroid Tracker (IoU-based)	85.0	7.4	3.5	12.3

The ablation study demonstrates the significance of each module. The OSNet-based feature extractor and Hungarian algorithm bring upon the most extensive performance enhancement. The centroid-based tracker greatly decreases ID switches and enhances temporal consistency when compared to rudimentary IoU-based techniques. Clearly, the proposed model surpasses all baseline methods with regard to all measurement criteria dealing with detection, reidentification, identity association, and tracking. The improvement in performance is apparent in diverse scenarios and constant under real-time restrictions. Each stage can be fine-tuned because of the modularity of the system, which further contributes to an efficient end-to-end deployment architecture set in process.











Figure 3. Model's View of Output

5. CONCLUSION & FUTURE SCOPES

This study presents a coherent and modular approach to detection, MultiView person cross View ReIdentification, and top View trajectory tracking from synchronized dual camera video streams. The proposed system combines YOLOv5s for fast detection, OSNet x0.25 for compact and discriminative embedding features, cosine similarity with the Hungarian algorithm to grant robust identity association, and a centroidbased tracker for temporally coherent trajectory estimation. Each of the selected and optimized components was aimed at achieving the highest computational efficiency, identity consistency, and trajectory completeness across spatially distinct viewpoints in real-time. The proposed model consistently outperformed three typical methods (Method [3], Method [8], and Method [15]) across a number of real-world situations. In particular, the detection module yielded an excellent average precision of 90.1% at IoU 0.5 and real-time throughput of 12.5 ms/frame, being better than all the baselines in both accuracy and speed. The OSNet powered ReID module produced a Top-1 accuracy of 87.4% with a mAP of 85.1% and a low false match rate of just 6.8% clearly proving its capability of coping with viewpoint induced appearance variations. In identity association, the model achieved a mean matching accuracy of 88.3% with only 1.3 average ID switches per sequence, showing strong inter-view consistency. The trajectory module also exhibited a trajectory completeness of 93.1% while sustaining lasting temporal continuity of 32 seconds with minute spatial drift. These empirical results solidify the claim for high accuracy of the system with interpretability and efficiency especially desirable in edge-computing or resource-constrained settings. Its modular architecture also allows for flexible adaptation to new sensor setups, extra viewpoints, or different ReID backbones. Future work will keep advancing along three primary scopes. First, incorporating temporal attention-based ReID models could boost the robustness of matching in cases of extreme occlusion or motion blur. Second, extending the model to accommodate asynchronous and/or partly overlapping camera views would broaden its openness to real-world relevance in surveillance systems. Thirdly, integrating a 3D spatial reasoning module using camera calibration data would further boost trajectory accuracy by solving projection and perspective ambiguities. One additional option is to include online learning mechanisms to allow identity models to adapt in real-time to changes in the scene or for never-before-seen subjects. Together, these enhancements aim to narrow the chasm between multi-camera tracking systems pursued in research and systems ready for deployment in surveillance, retail analytics, and public safety sets.

6. REFERENCES

- Li, X., Wu, A., & Zheng, W. (2024). "Person re Identification in special scenes based on deep learning: A comprehensive survey. Mathematics," *12*(16), 2495.
- [2] Elgendy, M., Alrahhal, M., & Sikora, T. (2023). "A multi-attention approach for person re Identification using deep learning. Sensors," 23(7), 3678.
- [3] Tapuhi, T., Klinger, A., & Sholev, O. (2022). "Multicamera multi-person re Identification with Hailo-8," Hailo AI.

- [4] Li, X., et al. (2023). "Multi-target tracking of person based on deep learning. Computer Systems Science and Engineering," 47(2), 2671–2688.
- [5] Ma, L. V., et al. (2024). "Track initialization and re Identification for 3D MultiView multi-object tracking," rXiv:2405.18606.
- [6] Viso.ai. (2024). "Deep learning for person re Identification,"
- [7] Khan, S., et al. (2023). "Scale-driven convolutional neural networks for robust MultiView tracking. IEEE Transactions on Pattern Analysis and Machine Intelligence," 45(3), 1452–1466.
- [8] Wang, G., et al. (2023). "Omni-scale feature learning for real-time person re Identification," IEEE Transactions on Circuits and Systems for Video Technology, 33(8), 4123–4136.
- [9] Yang, J., et al. (2023). "Open-set person re Identification in industrial environments," IEEE Access, 11, 12345– 12356.
- [10] Zhang, Y., et al. (2024). "Proxy anchor loss for deep metric learning in multi-camera tracking. Pattern Recognition," 147, 110023.
- [11] Baiju, N. (2024). "Navigating the maze: Person re Identification across multiple cameras with deep metric learning," IEEE CVPR Workshops.
- [12] Dijkinga, F. J. (2024). "A deep dive into modern multicamera person re Identification techniques," IEEE International Conference on Computer Vision (ICCV).
- [13] Luo, H., et al. (2023). "Spatial-temporal relation-aware global attention for re Identification," Proceedings of the AAAI Conference on Artificial Intelligence, 37(2), 1892–1900.
- [14] Chen, L., et al. (2023). "YOLORe IDNet: A unified framework for multi-camera tracking," IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 9876–9885.
- [15] Gu, X., et al. (2023). Facility-ReID: "A benchmark for industrial open-set re Identification," IEEE International Conference on Robotics and Automation (ICRA).