

# Robust Thermal Pedestrian Multi-Object Tracking through Three-Stage Trajectory Refinement

Anonymous CVPR submission

Paper ID 11

## Abstract

001 Thermal multi-object tracking (MOT) is essential for  
002 surveillance and safety-critical perception in low-visibility  
003 environments, yet it remains challenging due to weak ap-  
004 pearance cues, low contrast, frequent occlusions, and  
005 unstable object boundaries. We present a three-stage  
006 framework for thermal pedestrian MOT on the TP-MOT  
007 benchmark, consisting of thermal pedestrian detection,  
008 online tracking, and offline trajectory refinement. Our  
009 main contribution is a lightweight post-processing mod-  
010 ule that treats identity continuity as a first-class objec-  
011 tive by stitching fragmented tracklets using temporal con-  
012 tinuity, spatial proximity, motion consistency, and border-  
013 aware constraints. Unlike approaches that rely on heavy  
014 re-identification models or costly global optimization, the  
015 proposed refinement stage improves identity preservation  
016 while allowing the online tracker to remain simple and effi-  
017 cient. The full system further provides a unified pipeline for  
018 inference, visualization, evaluation, and submission gener-  
019 ation, enabling reproducible benchmarking across tracker  
020 variants. Experiments on the PBVS Thermal MOT bench-  
021 mark show that the proposed refinement stage consistently  
022 reduces fragmentation and improves continuity over raw  
023 tracking outputs, yielding a robust and practical real-time  
024 solution for thermal pedestrian surveillance.

## 025 1. Introduction

026 Multi-object tracking (MOT) in thermal imagery has be-  
027 come increasingly important for surveillance, public safety,  
028 intelligent transportation, and other vision systems that  
029 must operate reliably under poor illumination, adverse  
030 weather, or low-visibility conditions. Unlike RGB cameras,  
031 thermal sensors capture infrared radiation and therefore re-  
032 main effective when color and texture cues are severely de-  
033 graded. This property makes thermal imaging especially  
034 attractive for pedestrian monitoring in nighttime and safety-  
035 critical environments. At the same time, thermal pedestrian

tracking remains difficult because objects often appear with  
weak boundaries, limited appearance variation, and highly  
ambiguous local structure, all of which make detection and  
identity preservation substantially more challenging than in  
conventional RGB settings.

Despite steady progress in modern MOT, thermal-based  
tracking continues to present several unresolved difficul-  
ties. First, the low-information nature of thermal imagery  
reduces the discriminative power of appearance features,  
making it harder to distinguish nearby pedestrians or re-  
cover identities after occlusion. Second, crowded scenes,  
missed detections, and abrupt motion changes frequently  
fragment trajectories and induce identity switches. Third,  
many high-performing tracking systems depend on heavy  
re-identification modules or complex motion models, which  
can improve robustness but often increase computational  
cost and deployment complexity. These limitations are par-  
ticularly restrictive when the goal is to build a practical sys-  
tem that is both accurate and efficient for real-time thermal  
pedestrian tracking.

To address these challenges, this codebase presents a  
practical and extensible framework for thermal pedestrian  
multiple-object tracking, developed around the PBVS 2025  
challenge setting. The repository adopts a modular two-  
stage pipeline that separates object detection from ob-  
ject association, allowing each stage to be configured and  
tuned independently. In its default formulation, the sys-  
tem combines a thermal-adapted YOLOv8 detector with a  
lightweight SORT-based tracking stage, emphasizing strong  
detection quality, stable association, and minimal runtime  
overhead. Rather than relying primarily on computationally  
expensive re-identification, the framework improves  
identity continuity through carefully tuned association pa-  
rameters, an online ID-switch postprocessor, and an addi-  
tional offline tracklet-stitching procedure that merges short  
fragmented trajectories when spatial and temporal evidence  
is consistent. This design yields a codebase that is both  
engineering-oriented and research-friendly, supporting re-  
producible experimentation, submission generation, and  
straightforward extension to alternative tracker backends.

- 076 In brief, the main contributions are as follows:  
 077 • We present a modular thermal pedestrian MOT method  
 078 that integrates detection, tracking, identity-repair within  
 079 a unified pipeline.  
 080 • We introduce complementary identity-continuity mecha-  
 081 nisms, including online ID-switch correction and offline  
 082 tracklet stitching, to reduce fragmentation caused by oc-  
 083 clusion, missed detections, and complex motion.  
 084 • We provide an extensible experimental framework with  
 085 configurable tracker backend and standardized scripts,  
 086 making the repository suitable for both practical deploy-  
 087 ment and future thermal MOT research.

## 088 2. Related Work

089 Modern multi-object tracking (MOT) systems are largely  
 090 organized around the *tracking-by-detection* paradigm,  
 091 where an object detector first produces per-frame observa-  
 092 tions and a tracker subsequently enforces temporal identity  
 093 consistency. Over the past decade, a diverse family of asso-  
 094 ciation strategies, motion models, and auxiliary cues such  
 095 as appearance or segmentation have been proposed to im-  
 096 prove identity stability under occlusion and detection noise.  
 097 This section reviews the most relevant strands of prior work  
 098 and situates our system within this landscape.

### 099 2.1. Tracking-by-Detection

100 A foundational approach in modern MOT is SORT [3],  
 101 which combines Kalman filtering with Hungarian assign-  
 102 ment based on geometric overlap. Despite its simplic-  
 103 ity, SORT remains widely used due to its efficiency and  
 104 modular design, making it a natural baseline for detector-  
 105 centered tracking systems. Subsequent work has focused  
 106 on strengthening data association while retaining the same  
 107 basic structure. ByteTrack [17] demonstrated that discard-  
 108 ing low-confidence detections prematurely can significantly  
 109 degrade identity continuity; instead, these detections can  
 110 be recovered in a secondary matching stage to reduce track  
 111 fragmentation under occlusion. BoT-SORT [1] further en-  
 112 hances this paradigm by integrating improved motion mod-  
 113 eling, global motion compensation, and optional appear-  
 114 ance embeddings to improve robustness in complex scenes.  
 115 In parallel, OC-SORT [4] revisits the Kalman update formu-  
 116 lation by emphasizing observation-centric updates, improv-  
 117 ing tracking stability under non-linear motion and short-  
 118 term occlusions.

### 119 2.2. Identity Recovery and Tracklet Stitching

120 To address ID switch, where a tracker incorrectly assigns  
 121 different tracking IDs to the same target, previous works  
 122 have proposed identity recovery methods or tracklet stitch-  
 123 ing models that use different sources of information such  
 124 as appearance, spatial, and temporal cues to associate frag-  
 125 mented tracklets and reassign correct tracking IDs after the

online tracking process. These models typically rely on mo-  
 126 tion features or appearance features of tracklets for track-  
 127 let level association. For example, Translink [18] incorpo-  
 128 rates a CNN and a temporal attention network to extract  
 129 and encode the appearance features of a tracklet and formu-  
 130 lates the merging of tracklet pairs as a binary classification  
 131 task. AFLink [7] relies only on spatial and temporal infor-  
 132 mation. Some methods [6, 8, 16] adopt feature clustering  
 133 techniques to merge tracklets and improve performance in  
 134 multi-camera tracking scenarios. MambaMOT [9] proposes  
 135 a motion model that acts as a motion predictor and extracts  
 136 tracklet motion features for tracklet association. 137

## 138 3. Method

139 The released inference wrapper file instantiates a detector-  
 140 tracker-repair pipeline for single-class thermal pedestrian  
 141 multiple-object tracking. In the provided validation config-  
 142 uration, the system processes six thermal sequences (`seq2`,  
 143 `seq17`, `seq22`, `seq47`, `seq54`, and `seq66`) sampled at  
 144 10 FPS with image resolution  $960 \times 1280$ . The overall de-  
 145 sign is intentionally modular: detections are first generated  
 146 and cached, trajectories are then estimated online, and iden-  
 147 tity continuity is subsequently refined by two offline post-  
 148 processing stages. This decomposition follows the tracking-  
 149 by-detection philosophy of SORT [3] while adapting it to  
 150 the fragmentation patterns that are common in thermal im-  
 151 agery [15].

**Thermal pedestrian detection.** For each frame  $I_t$ , we  
 152 apply a YOLOv8 detector [10] in a single-class setting (*per-*  
 153 *son*). The detector operates on resized inputs of  $1920 \times$   
 154  $1920$ , uses batched inference with batch size 10, and is con-  
 155 figured with a very low confidence threshold ( $10^{-4}$ ) and  
 156 non-maximum suppression threshold 0.75 in order to pre-  
 157 serve recall. Let 158

$$159 \mathcal{D}_t = \{(b_i^t, s_i^t)\}_{i=1}^{N_t}$$

160 denote the set of normalized bounding boxes and confi-  
 161 dence scores predicted at time  $t$ . Instead of directly cou-  
 162 pling detector and tracker in memory, the pipeline stores  
 163  $\mathcal{D}_t$  as per-frame YOLO text files. This choice improves re-  
 164 producibility, allows detector/tracker ablations without re-  
 165 running the full pipeline, and enables direct inspection of  
 166 detector failure modes.

**Online tracking.** Before tracking, the stored detec-  
 167 tions are converted from normalized YOLO coordinates to  
 168 image-space  $(x_1, y_1, x_2, y_2)$  boxes and filtered with a  
 169 stricter threshold  $s_i^t \geq 0.5$ . Online association is then per-  
 170 formed with SORT [3], which combines a constant-velocity  
 171 Kalman predictor with Hungarian matching over bound-  
 172 ing-box overlap. Given predicted track boxes  $\hat{T}_t$  and detections  
 173

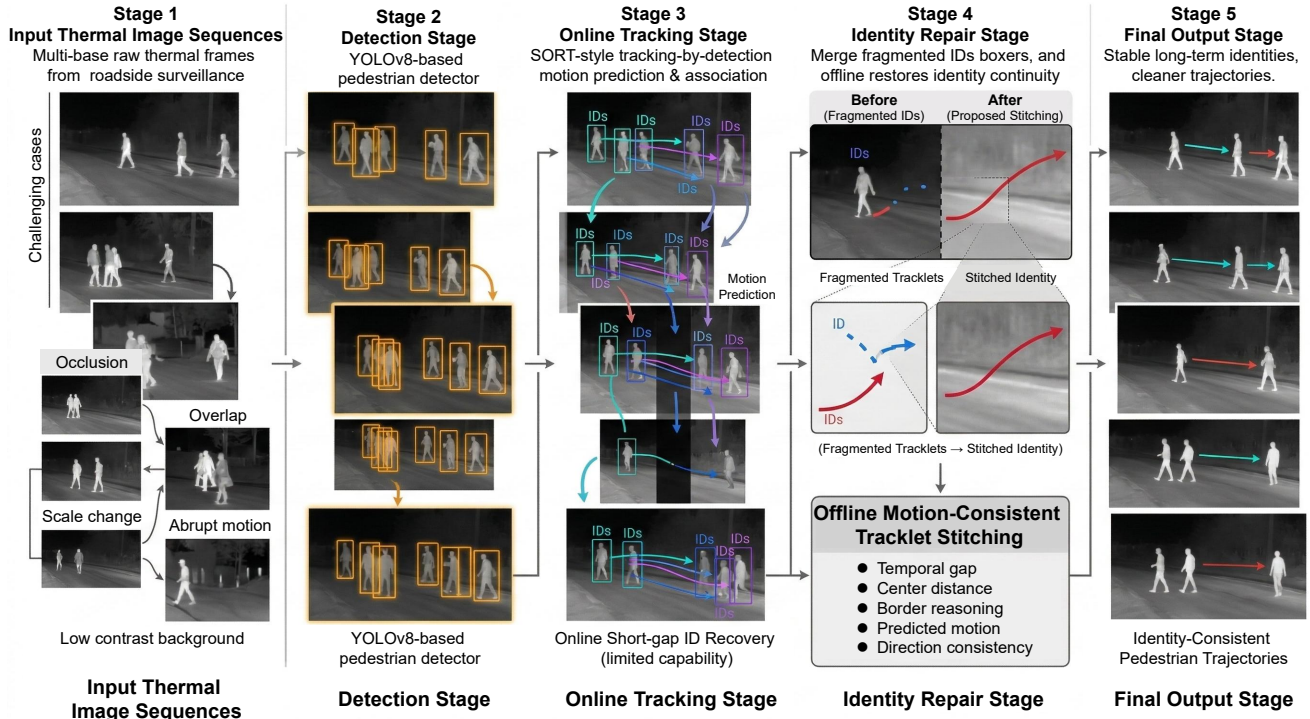


Figure 2. Overview of the proposed thermal pedestrian multi-object tracking pipeline. The framework consists of five stages. Stage 1 shows challenging thermal surveillance inputs with occlusion, overlap, scale changes, abrupt motion, and low contrast. Stage 2 detects pedestrians using a YOLOv8-based detector. Stage 3 performs online tracking-by-detection with to produce initial trajectories. Stage 4 introduces the proposed identity repair module, which stitches fragmented tracklets using temporal gaps, spatial proximity, motion prediction, border reasoning, and direction consistency. Stage 5 outputs stable identity-consistent pedestrian trajectories.

174  $\mathcal{D}_t$ , the association stage maximizes pairwise IoU, or equiv- 194  
 175 alently minimizes 195

$$176 C_{ij} = 1 - \text{IoU}(\hat{b}_j^t, b_i^t). \quad 196$$

177 In the released configuration, tracks are retained for up to 40 missed frames, and the IoU gate is set to 0.001, which 197  
 178 indicates a deliberately permissive association policy. This 198  
 179 setting is useful in thermal scenes where shape deformation, low contrast, and partial occlusion can make frame-to- 199  
 180 frame overlap unstable even when identity continuity is still 200  
 181 visually plausible. 201  
 182 202  
 183 203

184 **Online short-gap identity repair.** A key implementation 199  
 185 detail is that identity correction is not deferred entirely to 200  
 186 the end of the sequence. Immediately after each SORT up- 201  
 187 date, we apply an online remapping module that keeps a 202  
 188 memory of recently lost output identities. If a newborn raw 203  
 189 track  $u$  appears within  $\Delta t \leq 7$  frames of a recently lost 204  
 190 track  $v$ , and if their centers satisfy 205

$$191 \|c_u - c_v\|_2 \leq 60,$$

192 the newborn trajectory inherits the previous output identity 210  
 193 rather than emitting a new ID. Because the benchmark is 211  
 212

single-class, the optional label-consistency check is effectively 194  
 always satisfied. This step is particularly effective 195  
 for thermal videos, where brief detector interruptions often 196  
 arise from local temperature ambiguity, partial truncation, 197  
 or low signal-to-noise regions. 198

**Offline motion-consistent stitching.** The first offline 199  
 post-processing stage explicitly targets short illegal disap- 200  
 pear/reappear events. We treat a tracklet ending away from 201  
 the image border as an implausible disappearance and a 202  
 tracklet starting away from the border as an implausible ap- 203  
 pearance. For an old tracklet  $a$  and a new tracklet  $b$ , with 204  
 temporal gap 205

$$206 \Delta t = t_b^{\text{start}} - t_a^{\text{end}},$$

we estimate endpoint velocity from a 3-point temporal win- 207  
 dow and extrapolate the previous trajectory as 208

$$209 \tilde{c}_a = c_a^{\text{end}} + v_a \Delta t.$$

We merge  $b$  into  $a$  only when all of the following hold: 210  
 1) both endpoints lie outside a 60-pixel border band, 2) 211  
 $1 \leq \Delta t \leq 30$ , 3) the predicted displacement satisfies 212

213  $\|\tilde{c}_a - c_b^{\text{start}}\|_2 \leq 80$  pixels, and 4) if the motion magnitude  
214 is informative (speed  $\geq 0.25$  pixels/frame), the angle be-  
215 tween the two velocity vectors is at most  $45^\circ$  and the speed  
216 ratio is at most 3.0. This stage is intentionally conservative:  
217 it only repairs fragments that are both temporally short and  
218 kinematically consistent.

219 **Learned relinking of longer fragments.** A second of-  
220 fline stage addresses longer-range fragmentation using a  
221 fixed logistic relinking function. For each candidate  
222 predecessor–successor pair, we construct a 48-dimensional  
223 feature vector that encodes temporal gap, Euclidean and  
224 axis-wise displacement, tracklet lengths, box-size ratios,  
225 border distances, edge indicators, global velocity consis-  
226 tency, and multi-window extrapolation errors computed  
227 over temporal windows  $\{2, 3, 5, 10, 20\}$ . A candidate pair is  
228 considered only if the gap is at most 60 frames, the endpoint  
229 distance is at most 120 pixels, and both tracklets contain at  
230 least two observations; moreover, the predecessor must not  
231 terminate near a 25-pixel image border. The final relinking  
232 probability is

$$233 \quad p = \sigma(w^\top z),$$

234 where  $z$  is the normalized feature vector and  $\sigma(\cdot)$  is the sig-  
235 moid function. We accept a merge only when  $p \geq 0.95$  and  
236 the candidate is mutually optimal for both the old and new  
237 tracklet. This mutual-best constraint makes the relinking  
238 pass high-precision and prevents cascade errors that would  
239 otherwise amplify early mistakes.

240 **Output formatting and reproducibility.** After post-  
241 processing, each visible target is written in MOT-style for-  
242 mat as

$$243 \quad (f, \text{id}, x, y, w, h, s, \ell, -1, -1),$$

244 where  $f$  is the 1-indexed frame number, id is the repaired  
245 identity,  $(x, y, w, h)$  is the image-space bounding box,  $s$   
246 is the confidence score, and  $\ell$  is the class label. The  
247 script optionally renders qualitative overlays during detec-  
248 tion, tracking, and post-processing, and finally collects one  
249 `seq*_thermal.txt` file per sequence into a dated sub-  
250 mission archive.

251 Overall, the released system can be viewed as a high-  
252 recall thermal detector, a lightweight online motion tracker,  
253 and a two-stage identity-repair backend. This asymme-  
254 try is well suited to thermal MOT: the detector maximizes  
255 candidate coverage, the online tracker preserves short-term  
256 temporal continuity, and the post-processing stages recover  
257 long-range identity consistency under missed detections,  
258 brief occlusions, and mid-frame fragmentation.

## 4. Experiments 259

### 4.1. Dataset Splits 260

261 The repository configuration tracks six validation se-  
262 quences: `seq2`, `seq17`, `seq22`, `seq47`, `seq54`,  
263 `seq66`. These correspond to challenge-style sequences  
264 where public annotations are not distributed in this  
265 workspace.

### 4.2. Evaluation Metrics 266

267 We evaluate tracking performance using standard CLEAR  
268 MOT and identity-aware metrics that are widely used  
269 in multi-object tracking benchmarks and implemented in  
270 MOTChallenge-style evaluation toolchains [2, 12]. In par-  
271 ticular, we report MOTA, MOTP, IDF1, IDP, IDR, Recall,  
272 and Precision. These metrics jointly measure localization  
273 quality, and temporal identity consistency.

274 Let  $TP$ ,  $FP$ , and  $FN$  denote the numbers of true posi-  
275 tives, false positives, and false negatives aggregated over all  
276 frames. Let  $IDSW$  denote the number of identity switches,  
277 and  $GT$  denote the total number of ground-truth objects  
278 across all frames. For identity-aware evaluation,  $IDTP$ ,  
279  $IDFP$ , and  $IDFN$  denote identity true positives, identity  
280 false positives, and identity false negatives, respectively.

281 **MOTA.** Multi-Object Tracking Accuracy summarizes the  
282 dominant tracking errors into a single score [2]:

$$283 \quad \text{MOTA} = 1 - \frac{FN + FP + IDSW}{GT}. \quad (1)$$

284 A higher MOTA indicates fewer missed detections, fewer  
285 false alarms, and fewer identity switches. Because it penal-  
286 izes all three error sources simultaneously, MOTA is com-  
287 monly used as an overall measure of tracking performance.

288 **MOTP.** Multi-Object Tracking Precision measures local-  
289 ization quality on matched object–hypothesis pairs [2].  
290 In overlap-based evaluation, it reflects the spatial agree-  
291 ment between matched predictions and ground-truth ob-  
292 jects across all true-positive associations. In many bench-  
293 marks, including PBVS, MOTP is reported in a lower-is-  
294 better localization-error form.

295 **IDP, IDR, and IDF1.** Identity metrics evaluate whether  
296 detections are assigned the correct trajectory identity over  
297 time [13]. Using identity counts,

$$298 \quad \text{IDP} = \frac{IDTP}{IDTP + IDFP}, \quad (2)$$

$$299 \quad \text{IDR} = \frac{IDTP}{IDTP + IDFN}, \quad (3)$$

Table 1. PBVS website leaderboard results.

#	Participant	MOTA	MOTP	IDF1	IDP	IDR	Recall	Precision
1	<b>Ours</b>	0.99	0.13	0.86	0.86	0.86	1.00	1.00
2	wqetwet	0.97	0.12	0.97	0.98	0.96	0.98	1.00
3	SKKU-AutoLab	0.97	0.14	0.94	0.95	0.93	0.98	1.00
4	www123	0.85	0.12	0.86	0.91	0.82	0.89	0.98
5	spcke	0.82	0.14	0.86	0.90	0.83	0.88	0.94

$$\text{IDF1} = \frac{2 \text{IDTP}}{2 \text{IDTP} + \text{IDFP} + \text{IDFN}} \quad (4)$$

IDP measures identity precision, IDR measures identity recall, and IDF1 is their harmonic mean. These metrics are particularly sensitive to trajectory fragmentation and identity switches, making them informative for evaluating identity preservation over time.

**Recall and Precision.** Detection Recall and Precision are defined as

$$\text{Recall} = \frac{TP}{TP + FN}, \quad (5)$$

$$\text{Precision} = \frac{TP}{TP + FP}. \quad (6)$$

Recall measures the proportion of ground-truth objects that are successfully detected, while Precision measures the proportion of predicted detections that correspond to true objects. These metrics describe detection coverage and false-alarm behavior within the tracking pipeline [12].

### 4.3. Experimental Setup

Experiments are conducted on the PBVS Thermal MOT dataset, which contains roadside thermal surveillance sequences captured under nighttime conditions. The dataset exhibits several challenges common in thermal tracking scenarios, including low contrast, partial occlusion, overlapping pedestrians, scale variation, and abrupt motion. All experiments follow the official challenge evaluation protocol. The detector is a YOLOv8 model trained for single-class pedestrian detection, which we adopted the weight from [15]. Inference is performed on resized  $1920 \times 1920$  frames. Detection results are stored and reused across tracker variants to ensure fair comparison. The tracking stage is evaluated using several widely used MOT association strategies implemented within the same pipeline, including SORT [3], ByteTrack [17], BoT-SORT [1], OC-SORT [4], BoostTrack [14], DiffMOT [11], and a segmentation-assisted variant based on SAM3 [5].

Table 2. Leaderboard ranking based on the weighted score.

Rank	Team	Weighted Score
1	wqetwet	0.7575
2	SKKU-AutoLab	0.7550
3	<b>Ours</b>	0.7425
4	www123	0.6700
5	spcke	0.6600

Table 3. Extension of YOLOv8 detector TP-MOT evaluation server results reported by [15].

Method	MOTA $\uparrow$	MOTP $\downarrow$	IDF1 $\uparrow$
ByteTrack	0.9173	0.1367	0.7659
BoT-SORT	0.9174	0.1368	0.7605
BoostTrack	0.8654	0.1555	0.7545
DiffMOT	0.8630	0.1611	0.7812
OC-SORT	0.9071	0.1236	0.5685
SAM3	0.9003	0.2104	0.6073
SORT	<b>0.9844</b>	<b>0.1263</b>	<b>0.8130</b>
SORT + Stitching	<b>0.9853</b>	<b>0.1262</b>	<b>0.8545</b>

### 4.4. Comparison with Existing Trackers

Table 3 compares several representative tracking methods under the same detection inputs. Among the tested trackers, SORT provides the strongest baseline performance in this thermal benchmark.

More complex trackers such as ByteTrack and BoT-SORT achieve competitive detection recall but show lower identity consistency in this dataset. DiffMOT and BoostTrack provide stronger motion modeling but remain sensitive to missed detections and fragmentation under severe occlusion. These observations highlight an important practical insight: in thermal pedestrian tracking, identity continuity is often dominated by fragmentation recovery rather than by increasing tracker complexity alone.

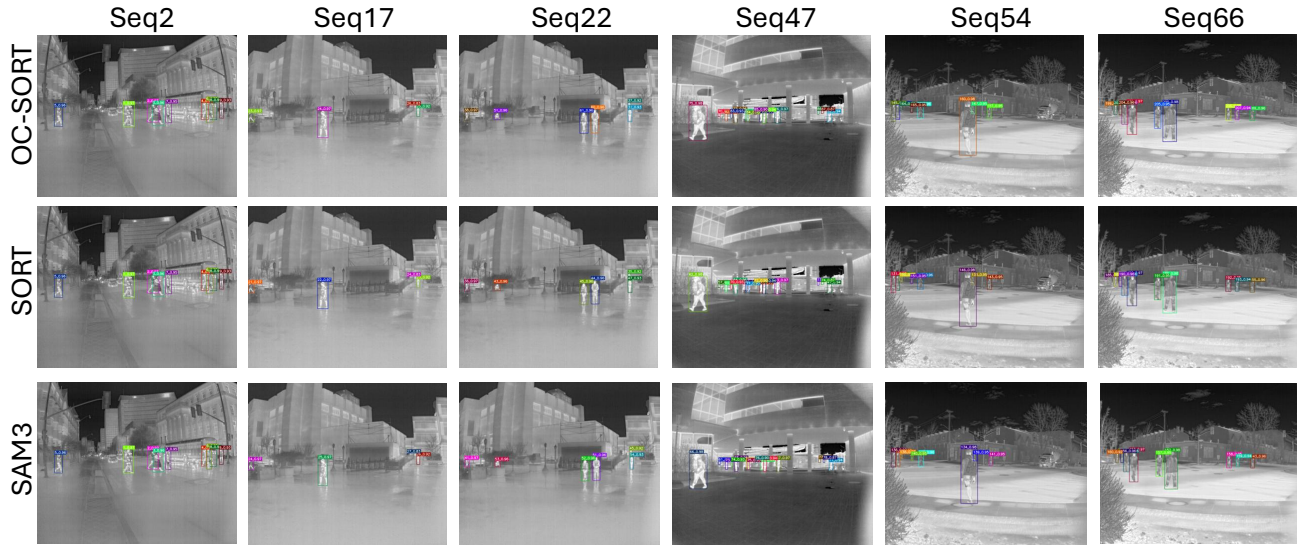


Figure 3. Qualitative comparison of OC-SORT, SORT, SAM3, and the proposed trajectory-refinement method on representative PBVS thermal MOT sequences. The proposed approach yields more stable identity assignments and smoother trajectories under occlusion and motion changes.

#### 347 4.5. Ablation on Tracklet Stitching

348 The offline stitching module further reduces trajectory frag- 372  
 349 mentation by reconnecting short tracklets that satisfy tem- 373  
 350 poral and motion-consistency constraints. Compared with 374  
 351 the baseline pipeline, this stage reduces the total number of 375  
 352 fragmented tracks and increases average trajectory length. 376  
 353 After introducing the proposed identity-repair module, the 377  
 354 extended *SORT + Stitching* configuration achieves the best 378  
 355 overall performance with a MOTA of 0.9853 and IDF1 of 379  
 356 0.8545. Compared with the vanilla SORT baseline, the pro- 380  
 357 posed refinement stage improves IDF1 by over 4% while 381  
 358 maintaining nearly identical localization accuracy. 382  
 383

#### 359 4.6. Qualitative Results

360 Figure 3 visualizes representative tracking results across 384  
 361 several challenging sequences. Compared with other track- 385  
 362 ers, the proposed framework produces smoother trajectories 386  
 363 and more stable identity assignments when pedestrians un- 387  
 364 dergo temporary occlusion or cross paths. In particular, the 388  
 365 proposed border-aware stitching strategy prevents spurious 389  
 366 identity births in the interior of the image while still allow- 390  
 367 ing legitimate entries and exits near frame boundaries. This 391  
 368 constraint proves especially effective in roadside thermal 392  
 369 surveillance scenarios, where pedestrians frequently move 393  
 370 through the field of view. 394  
 395  
 396  
 397  
 398

## 5. Conclusion

This work presents a practical framework for thermal pedes- 372  
 373 trian multi object tracking that prioritizes robustness, effi- 373  
 374 ciency, and modular design. By separating detection and 374  
 375 association into configurable components, the system al- 375  
 376 lows flexible experimentation while maintaining stable real 376  
 377 time performance. The integration of online identity repair 377  
 378 and offline tracklet stitching further improves identity conti- 378  
 379 nuity in challenging thermal scenarios involving occlusion, 379  
 380 missed detections, and ambiguous appearance cues. To- 380  
 381 gether, these design choices provide a reliable baseline and 381  
 382 an extensible platform for future research in thermal based 382  
 383 multi object tracking. 383

## References

- [1] Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. Bot-sort: 385  
 Robust associations multi-pedestrian tracking, 2022. 2, 5 386
- [2] Keni Bernardin and Rainer Stiefelhagen. Evaluating multi- 387  
 ple object tracking performance: The clear mot metrics. In 388  
*EURASIP Journal on Image and Video Processing*, 2008. 4 389
- [3] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and 390  
 Ben Upcroft. Simple Online and Realtime Tracking. In 391  
*2016 IEEE International Conference on Image Processing (ICIP)*, 392  
 pages 3464–3468, 2016. 2, 5 393
- [4] Jinkun Cao, Jiangmiao Pang, Xinshuo Weng, Rawal Khirod- 394  
 kar, and Kris Kitani. Observation-Centric SORT: Rethink- 395  
 ing sort for robust multi-object tracking. In *Proceedings of 396  
 the IEEE/CVF Conference on Computer Vision and Pattern 397  
 Recognition (CVPR)*, pages 9686–9696, 2023. 2, 5 398

- 399 [5] Nicolas Carion, Laura Gustafson, Yuan-Ting Hu, Shoub-  
400 hik Debnath, Ronghang Hu, Didac Suris Coll-Vinent, Chai-  
401 tanya Ryali, Kalyan Vasudev Alwala, Haitham Khedr, An-  
402 drew Huang, Jie Lei, Tengyu Ma, Baishan Guo, Arpit  
403 Kalla, Markus Marks, Joseph Greer, Meng Wang, Peize Sun,  
404 Roman Rädle, Triantafyllos Afouras, Effrosyni Mavroudi,  
405 Katherine Xu, Tsung-Han Wu, Yu Zhou, Liliane Momeni,  
406 RISHI HAZRA, Shuangrui Ding, Sagar Vaze, Francois  
407 Porcher, Feng Li, Siyuan Li, Aishwarya Kamath, Ho Kei  
408 Cheng, Piotr Dollar, Nikhila Ravi, Kate Saenko, Pengchuan  
409 Zhang, and Christoph Feichtenhofer. SAM 3: Segment any-  
410 thing with concepts. In *The Fourteenth International Con-  
411 ference on Learning Representations*, 2026. 5
- 412 [6] Riu Cherdchusakulchai, Sasin Phimsiri, Visarut Trairat-  
413 tanapa, Suchat Tungjitnob, Wasu Kudisthalert, Pornprom Ki-  
414 awjak, Ek Thamwiwatthana, Phawat Borisuitsawat, Teep-  
415 akorn Tosawadi, Pakcheera Choppradi, et al. Online multi-  
416 camera people tracking with spatial-temporal mechanism  
417 and anchor-feature hierarchical clustering. In *Proceedings  
418 of the IEEE/CVF Conference on Computer Vision and Pat-  
419 tern Recognition*, pages 7198–7207, 2024. 2
- 420 [7] Yunhao Du, Zhicheng Zhao, Yang Song, Yanyun Zhao, Fei  
421 Su, Tao Gong, and Hongying Meng. Strongsort: Make deep-  
422 sort great again. *IEEE Transactions on Multimedia*, 2023. 2
- 423 [8] Hsiang-Wei Huang, Cheng-Yen Yang, Zhongyu Jiang,  
424 Pyong-Kun Kim, Kyoungoh Lee, Kwangju Kim, Samartha  
425 Ramkumar, Chaitanya Mullapudi, In-Su Jang, Chung-I  
426 Huang, and Jenq-Neng Hwang. Enhancing multi-camera  
427 people tracking with anchor-guided clustering and spatio-  
428 temporal consistency id re-assignment. In *Proceedings of  
429 the IEEE/CVF Conference on Computer Vision and Pattern  
430 Recognition (CVPR) Workshops*, pages 5239–5249, 2023. 2
- 431 [9] Hsiang-Wei Huang, Cheng-Yen Yang, Wenhao Chai,  
432 Zhongyu Jiang, and Jeng-Neng Hwang. Mambamot: State-  
433 space model as motion predictor for multi-object tracking.  
434 In *ICASSP 2025-2025 IEEE International Conference on  
435 Acoustics, Speech and Signal Processing (ICASSP)*, pages  
436 1–5. IEEE, 2025. 2
- 437 [10] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. Ultralytics  
438 yolo, 2024. 2
- 439 [11] Weiyi Lv, Yuhang Huang, Ning Zhang, Rwei-Sung Lin, Mei  
440 Han, and Dan Zeng. DiffMOT: A real-time diffusion-based  
441 multiple object tracker with non-linear prediction. In *Pro-  
442 ceedings of the IEEE/CVF Conference on Computer Vision  
443 and Pattern Recognition (CVPR)*, pages 19321–19330, 2024.  
444 5
- 445 [12] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and  
446 Konrad Schindler. Mot16: A benchmark for multi-object  
447 tracking. In *ECCV Workshops*, 2016. 4, 5
- 448 [13] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara,  
449 and Carlo Tomasi. Performance measures and a data set for  
450 multi-target, multi-camera tracking. In *ECCV Workshop on  
451 Benchmarking Multi-Target Tracking*, 2016. 4
- 452 [14] Vukasin Stanojevic and Branimir Todorovic. Boosttrack++:  
453 Revisiting the tracking-by-detection paradigm for multiple  
454 object tracking, 2024. 5
- 455 [15] Duong Nguyen-Ngoc Tran, Long Hoang Pham, Chi Dai  
456 Tran, Quoc Pham-Nam Ho, Huy-Hung Nguyen, and  
Jae Wook Jeon. A novel tuning method for real-time  
multiple-object tracking utilizing thermal sensor with com-  
plexity motion pattern, 2025. 2, 5
- [16] Cheng-Yen Yang, Hsiang-Wei Huang, Pyong-Kun Kim,  
Zhongyu Jiang, Kwang-Ju Kim, Chung-I Huang, Haiqing  
Du, and Jenq-Neng Hwang. An online approach and evalua-  
tion method for tracking people across cameras in extremely  
long video sequence. In *Proceedings of the IEEE/CVF Con-  
ference on Computer Vision and Pattern Recognition*, pages  
7037–7045, 2024. 2
- [17] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng  
Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang  
Wang. Bytetrack: Multi-object tracking by associating every  
detection box. In *ECCV*, 2022. 2, 5
- [18] Yanting Zhang, Shuanghong Wang, Yuxuan Fan, Gaoang  
Wang, and Cairong Yan. Translink: Transformer-based em-  
bedding for tracklets’ global link. In *ICASSP 2023-2023  
IEEE International Conference on Acoustics, Speech and  
Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023. 2