

Optimizing Long-Term Player Tracking and Identification in NAO Robot Soccer by fusing Game-state and External Video

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Monitoring a fleet of robots requires stable long-term tracking with
2 re-identification, which is yet an unsolved challenge in many scenarios. One ap-
3 plication of this is the analysis of autonomous robotic soccer games at RoboCup.
4 Tracking in these games requires handling of identically looking players, strong
5 occlusions, and non-professional video recordings, but also offers state informa-
6 tion estimated by the robots. In order to make effective use of the information
7 coming from the robot sensors, we propose a robust tracking and identification
8 pipeline. It fuses external non-calibrated camera data with the robots’ internal
9 states using quadratic optimization for tracklet matching. The approach in this
10 work is validated using game recordings from previous RoboCup World Cups.

11 1 Introduction

12 Robust tracking with stable object identification is a crucial step towards extracting game statistics
13 and improving gameplay in many team sports. While this is usually approached using an exter-
14 nal camera only, our application in understanding soccer games played by humanoid robots allows
15 us to fuse this information with measurements from robot-mounted sensors. In this work, we fo-
16 cus on the RoboCup Standard Platform League (SPL) where humanoid NAO robots compete fully
17 autonomously. Game analytics in this setting can offer objective feedback on the algorithms’ per-
18 formance to the teams and help to improve the gameplay.

19 Our problem differs in multiple ways from the well-known tracking and identification problem in
20 game analytics: RoboCup games are recorded with non-professional uncalibrated camera equip-
21 ment, robots look identical except for their jerseys, jersey numbers are too small to be detected
22 reliably, and human referees often occlude the camera. To handle these challenges, we propose a
23 long-term tracking pipeline consisting of the following modules:

- 24 1. Camera calibration, to estimate camera distortion, intrinsics, and pose relative to the field.
- 25 2. Short-term object tracking based on Tracktor [1] and trained on our data to generate tracklets.
- 26 3. Optimization-based long-term tracking and player identification by fusing cues from an external
27 camera and the robot sensor data.

28 2 Related Work

29 **Multi Object Tracking (MOT)** describes the tracking of all objects belonging to a given set of
30 categories [2]. In joint tracking and detection approaches, the object detector is a fundamental part of
31 the tracking pipeline [1, 3, 4]. We use Tracktor [1], which follows this paradigm as a building block
32 in our pipeline. Another category of trackers uses detections provided by a separate object detector,
33 followed by solving a data association problem. This framework includes fully deep learning based
34 methods [5, 6, 7] as well as optimization based approaches [8, 9, 10].

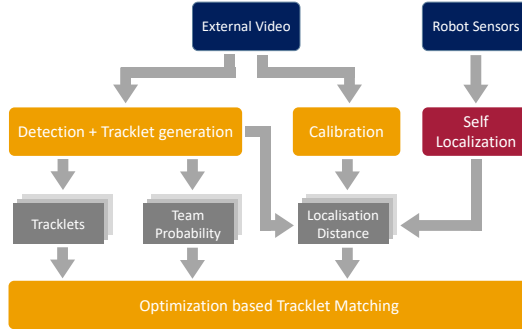


Figure 1: Overview of the proposed approach.

35 **Game Analytics** covers the tracking and identification of the players in videos [11] as a key chal-
 36 lenge, where MOT is an important component. This can further be aided by robust team detec-
 37 tion [11, 12]. A key aspect of our work, usually not addressed in player tracking and identification
 38 is the integration of player-mounted sensor data, as it is not easily applicable to human players.

39 **Camera Calibration** Tracking and identity assignment requires accurate camera intrinsics and ex-
 40 trinsics. Standard calibration processes generally use point correspondances, robot’s motion or cal-
 41 ibration patterns to provide accurate intrinsics [13, 14, 15], which cannot be used in our case due to
 42 poor point correspondances. Therefore, our approach utilizes a technique proposed by Alvarez et
 43 al. [16], which minimizes an energy objective based on rectifying lines present on the field.

44 3 Method

45 In this section, we detail our pipeline for consistent player tracking and identification/ Figure 1
 46 provides an overview of the key components and information flow. First, the camera is calibrated
 47 to compute its intrinsic and extrinsic parameters with respect to the known soccer field, which is
 48 required only once for each video sequence. Then, the tracklets, the team color for each player and
 49 all relevant robot state information is extracted. Subsequently, each tracklet is associated with a
 50 specific robot player by optimizing a binary quadratic program as described in Section 3.5.

51 **Data and Application:** We consider RoboCup SPL matches between humanoid NAO robots, us-
 52 ing a dataset comprised of 8 annotated 5000-frame sequences recorded by wide-angle cameras at
 53 30 FPS. The videos were recorded at RoboCup 2019 and 2022, together with the corresponding
 54 team communication and game controller logs. Annotations include the bounding box, jersey color
 55 and number for each active player and frame. Object detection and classification models and the
 56 optimizer hyperparameters are trained on five sequences, and evaluated on the remaining three.

57 3.1 Camera Calibration and Extrinsic

58 We assume a static camera over the sequence and compute the median of each pixel over all frames
 59 to remove moving objects and obtain a clear view of the field. Then the wide-angle image is undis-
 60 torted, by estimating the distortion using [16] on detections of field line candidates. We apply the
 61 SOLD2 [17] line detector on the undistorted image and obtain intersection points on the field which
 62 can then be matched to known field 3D coordinated. The camera pose is computed using P3P [18].

63 3.2 Multi Object Tracker

64 To generate tracklets, we use Tractor [1] with Faster-RCNN [19] and a ResNet-50 backbone. We
 65 initialize the model with MS-COCO [20] pretrained weights and fine-tune it on the 5 training se-
 66 quences of our dataset. The tracklets are robust in easy tracking scenarios without occlusions, but
 67 do not cover a whole video. For further processing, each tracklet is projected to field coordinates



Figure 2: Visualization of robots identified by the tracker. The tracking result is represented by bounding boxes and IDs at their top. Ground truth positions are represented by green crosses and corresponding green IDs.

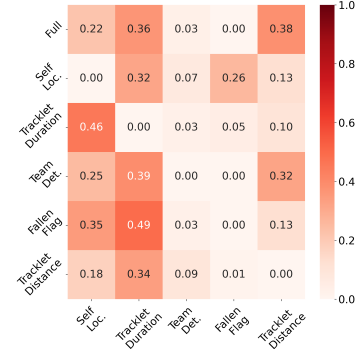


Figure 3: Ablation weights.

68 using the estimated camera pose. Subsequently, each projected tracklet is smoothed using a Kalman
 69 filter with a constant velocity model.

70 3.3 Jersey color detection

71 In the SPL, 9 distinct jersey colors are used. The team colors for a match provide a strong signal
 72 to associate tracklets to players from either team. We detect colors using a VGG16 network that
 73 assigns a score for each of the team colors used by the two teams in a game.

74 3.4 Robot States

75 We furthermore use the robots' states for our matching problem. These include information from
 76 the robot sensors as well as the game state:

- 77 – Self Localization: The robots calculate their relative position on the field based on their views.
- 78 – Fallen Robots: The robots use the IMU sensor information to determine if they have fallen.
- 79 – Penalized Robots: The robots removed from the field is a constraint in the optimization problem.

80 3.5 Global optimization

81 Occlusions and distractors cause Tracktor to split the ideally long tracks into a large number of
 82 shorter tracklets. Therefore, we frame the long-term tracking problem as an assignment of tracklets
 83 to a fixed number of player tracks similar to [10] as a constrained quadratic binary optimization
 84 problem. We denote the index set of player tracks $I = \{1, \dots, N\}$ (with $N = 10$) and generated
 85 tracklets $J = \{1, \dots, M\}$. The objective is to minimize:

$$H(x) = \sum_{i \in I} \sum_{j \in J} x_{i,j} (O_u + \sum_{l \in L} w^l c_{i,j}^l) + \sum_{i \in I} \sum_{j \in J} \sum_{k \in J} x_{i,j} x_{i,k} (\sum_{p \in P} w^p c_{j,k}^p) \quad (1)$$

86 where $x_{i,j} \in \{0, 1\}$ are binary optimization variables, with $x_{i,j} = 1$ meaning tracklet j is assigned
 87 to track i , L and P the sets of unary and pairwise cost functions with costs $c_{i,j}^l$ and $c_{i,j,k}^p$. w^l and
 88 w^p are used to weight different cost terms. The offsets are negative to penalize the trivial solution
 89 of assigning nothing ($x_{i,j} = 0 \forall i, j$). Two constraints are imposed to ensure feasible matchings:

- 90 1. One tracklet can only be assigned to one track.
- 91 2. Temporally overlapping tracklets cannot be merged to the same track.

92 3.6 Cost terms

93 Different cost terms control the assignment of tracklets to tracks. We use the following terms.

- 94 • **Duration:** Penalizes short tracklets, as these are often spurious tracklets.
- 95 • **Self-localization:** The distance between the position estimated by a robot from its camera
 96 and the position estimated for a tracklet from the external camera.

Full	Self Loc.	Tracklet Duration	Team Det.	Fallen Flag	Penalized Flag	Tracklet Distance	Time Frame	30	150	300	900	1800	3600	5400
88.1	15.4	51.1	76.5	86.2	76.3	83.3	MPIR	42.5	42.3	40.1	43.8	38.4	38.5	38.5

Table 1: Tracking performance and ablation study. Results are provided in percent MPIR.

Table 2: DeepSORT Performance with different re-identification time.

- 97 • **Jersey color detection:** The score for the color-based team detection throughout a tracklet.
- 98 • **Global trajectory continuity:** The pairwise spatial distance between the end of a tracklet
- 99 and the start of a new temporally close tracklet.

100 3.7 Reference Method: DeepSORT

101 To fulfill the task of long-term tracking and player identification under strong occlusions we aug-
 102 ment the DeepSORT tracker [21, 22] by a greedy tracklet matching algorithm that matches any new
 103 tracklet to the spatially closest inactive track. Constraints are applied to prevent temporarily overlap-
 104 ping tracklets from being assigned to the same track. To provide a strong baseline, we provide this
 105 approach with oracle information; the total number of robots that are present in a sequence, which
 106 defines the maximum number of tracks as well as the ground truth ID of the first tracklet for each
 107 robot. Finally, the best re-identification time for DeepSORT is selected on the testset.

108 4 Results

109 We evaluate our approach over a test set containing 3 video sequences of 3 minutes. Each video
 110 covers a different game, thus testing our approach under different conditions.

111 Table 1 shows the Mean Player Identification Recall (MPIR), the ratio of times each player has been
 112 identified correctly. The first column shows our full approach. Subsequent columns show ablations,
 113 each feature removed separately with the cost term optimized using PSO [23] for each scenario.

114 With all features we achieve 88.1% MPIR. Removing the robot self localization has the strongest
 115 impact with 15.4% MPIR, while removing the fallen robot flag results in the least performance drop.
 116 The self-localization is an important feature since it provides information about the position of the
 117 robot. The fallen robot flag is unreliable, as it relies on the robot’s IMU and a heuristic to detect
 118 whether the robot has fallen in the video.

119 Table 2 shows the performance of our DeepSORT baseline over different reidentification times. The
 120 best performance as achieved with 900 frames which corresponds to 30s of video. In this case,
 121 the performance is 43.8% MPIR, compared to 88.1% MPIR with our method using all available
 122 features.

123 We further analyze each feature’s importance through its weighting, where a higher weight indicates
 124 a more important feature. Figure 3 shows the importance of the features in the different ablations.
 125 Strong weights are assigned to the self-localization and tracklet duration. Removing these features
 126 shows that weighting is redistributed: the noisy fallen robot events are incorporated when no self-
 127 localization information is available, as it can provide unique information about a robot’s ID.

128 5 Conclusion

129 In this work, we presented a sensor fusion based method for tracking multiple similar humanoid
 130 robots. We utilize information from both visual data and their own sensors by combining tracklets
 131 using a quadratic optimization technique. The method allows automated tracking of robot players
 132 over a long time on a stationary video sequence. Open points that we will investigate in the future
 133 include the evaluation in more complex environments, the interpolation of tracks during occlusions
 134 as well as the extraction of high-level game statistics.

References

- [1] P. Bergmann, T. Meinhardt, and L. Leal-Taixe. Tracking Without Bells and Whistles. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 941–951, Seoul, Korea (South), Oct. 2019. IEEE. ISBN 978-1-72814-803-8. doi:10.1109/ICCV.2019.00103.
- [2] P. Dendorfer, H. Rezatofighi, A. Milan, J. Shi, D. Cremers, I. Reid, S. Roth, K. Schindler, and L. Leal-Taixé. MOT20: A benchmark for multi object tracking in crowded scenes. *arXiv:2003.09003 [cs]*, Mar. 2020.
- [3] S. Karthik, A. Prabhu, and V. Gandhi. Simple Unsupervised Multi-Object Tracking. *arXiv:2006.02609 [cs]*, June 2020.
- [4] Z. Zheng, X. Yang, Z. Yu, L. Zheng, Y. Yang, and J. Kautz. Joint discriminative and generative learning for person re-identification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [5] G. Braso and L. Leal-Taixe. Learning a Neural Solver for Multiple Object Tracking. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6246–6256, Seattle, WA, USA, June 2020. IEEE. ISBN 978-1-72817-168-5. doi:10.1109/CVPR42600.2020.00628.
- [6] X. Zhou, V. Koltun, and P. Krähenbühl. Tracking Objects as Points. In A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, editors, *Computer Vision – ECCV 2020*, Lecture Notes in Computer Science, pages 474–490, Cham, 2020. Springer International Publishing. ISBN 978-3-030-58548-8. doi:10.1007/978-3-030-58548-8_28.
- [7] J.-N. Zaech, D. Dai, A. Liniger, M. Danelljan, and L. V. Gool. Learnable online graph representations for 3d multi-object tracking. In *IEEE International Conference on Robotics and Automation, ICRA*, 2022. URL <https://arxiv.org/abs/2104.11747>.
- [8] A. Hornakova, R. Henschel, B. Rosenhahn, and P. Swoboda. Lifted Disjoint Paths with Application in Multiple Object Tracking. In *Proceedings of the 37th International Conference on Machine Learning*, pages 4364–4375. PMLR, Nov. 2020.
- [9] S. Tang, M. Andriluka, B. Andres, and B. Schiele. Multiple People Tracking by Lifted Multicut and Person Re-identification. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3701–3710, Honolulu, HI, July 2017. IEEE. ISBN 978-1-5386-0457-1. doi:10.1109/CVPR.2017.394.
- [10] J.-N. Zaech, A. Liniger, M. Danelljan, D. Dai, and L. Van Gool. Adiabatic quantum computing for multi object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8811–8822, June 2022.
- [11] T. Yamamoto, H. Kataoka, M. Hayashi, Y. Aoki, K. Oshima, and M. Tanabiki. Multiple players tracking and identification using group detection and player number recognition in sports video. In *IECON 2013-39th Annual Conference of the IEEE Industrial Electronics Society*, pages 2442–2446. IEEE, 2013.
- [12] S. Gerke, K. Muller, and R. Schafer. Soccer jersey number recognition using convolutional neural networks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 17–24, 2015.
- [13] Z. Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000.
- [14] M. Pollefeys, L. Van Gool, and A. Oosterlinck. The modulus constraint: a new constraint self-calibration. In *Proceedings of 13th International Conference on Pattern Recognition*, volume 1, pages 349–353. IEEE, 1996.

- 180 [15] T. Probst, K.-K. Maninis, A. Chhatkuli, M. Ourak, E. Vander Poorten, and L. Van Gool. Auto-
181 matic tool landmark detection for stereo vision in robot-assisted retinal surgery. *IEEE Robotics*
182 *and Automation Letters*, 3(1):612–619, 2017.
- 183 [16] L. Alvarez, L. Gómez, and J. R. Sendra. An algebraic approach to lens distortion by line
184 rectification. *Journal of Mathematical Imaging and Vision*, 35(1):36–50, 2009.
- 185 [17] R. Pautrat, J.-T. Lin, V. Larsson, M. R. Oswald, and M. Pollefeys. Sold2: Self-supervised
186 occlusion-aware line description and detection. In *Proceedings of the IEEE/CVF Conference*
187 *on Computer Vision and Pattern Recognition (CVPR)*, pages 11368–11378, June 2021.
- 188 [18] L. Kneip, D. Scaramuzza, and R. Siegwart. A novel parametrization of the perspective-three-
189 point problem for a direct computation of absolute camera position and orientation. In *CVPR*
190 *2011*, pages 2969–2976. IEEE, 2011.
- 191 [19] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards Real-Time Object Detection
192 with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelli-*
193 *gence*, 39(6):1137–1149, June 2017. ISSN 0162-8828, 2160-9292. doi:10.1109/TPAMI.2016.
194 2577031.
- 195 [20] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan,
196 C. L. Zitnick, and P. Dollár. Microsoft COCO: Common Objects in Context. *arXiv:1405.0312*
197 *[cs]*, May 2014.
- 198 [21] N. Wojke, A. Bewley, and D. Paulus. Simple online and realtime tracking with a deep asso-
199 ciation metric. In *2017 IEEE International Conference on Image Processing (ICIP)*, pages
200 3645–3649. IEEE, 2017. doi:10.1109/ICIP.2017.8296962.
- 201 [22] N. Wojke and A. Bewley. Deep cosine metric learning for person re-identification. In *2018*
202 *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 748–756. IEEE,
203 2018. doi:10.1109/WACV.2018.00087.
- 204 [23] R. Poli, J. Kennedy, and T. M. Blackwell. Particle swarm optimization. *Swarm Intelligence*, 1:
205 33–57, 1995.