

An Analytical Review on Some Recent Advances in Deep Learning Object Tracking Approaches

K. Nani Kumar[#], M. James Stephen^{*}, P.V.G.D. Prasad Reddy⁺

[#]Department of CS & SE, Andhra University, Visakhapatnam, AP, INDIA.

^{*}Professor, Department of CSE, WISTM Engineering College, Visakhapatnam.

⁺Sr. Professor, Department of CS & SE, Andhra University, Visakhapatnam, AP, INDIA.

Abstract:- Visual Object tracking in real world, real time application scenarios is a complex problem, therefore, it remains a most active area of research in computer vision. This paper presents a detailed review on some of the recent advances in Deep Learning Based Visual Object Tracking Approaches from a wide variety of algorithms often cited in research literature. Some of the recent advances in various aspects have been extensively investigated and critically analyzed. Out of this study, some key issues with available research are presented. This paper concludes with some interesting future potential research directions.

Keywords: Deep Learning, Neural Networks, CNN, Machine Learning, Pre-training, Online learning, Visual Tracking

INTRODUCTION

Tracking is the problem of generating an inference about the motion of an object given a sequence of images. Good solutions of this problem have a variety of applications. Visual object tracking has been extensively studied in computer vision. Visual Tracking has been widely used in various real world, real-time applications including Advance Driver Assistance Solutions, Medical, Military, Video Surveillance, Traffic Control, Navigation, Robotics, Augmented Reality, sports to name a few.



Fig 1: Some Visual Tracking Applications

Basic traditional visual tracking methods utilize various frameworks like Discriminative Correlation Filters (DCF) [24]–[28], silhouette tracking [28, 29], Kernel tracking [30]–[32], point tracking [33], and so forth – these methods were not able to provide satisfactory results in unconstrained environments.

After careful analysis from the wide variety of algorithms cited in the literature, this paper focuses on some potential

and well performed Deep Learning based visual tracking methods. The trackers include from diverse classification methods like, Network Architecture, Network Exploitation, Training, Network Objective, Network Output and Correlation Filter Exploitation.

Although Convolutional Neural Networks (CNN) have been used in Deep Learning (DL)based methods from the research literature, some of other network architectures were also proposed to improve the efficiency and robustness of visual trackers in recent years. The CNN-based visual trackers are classified under three categories, robust target representation, Balancing training data and Computational complexity problem. Some of the advantages of CNN based methods utilized are parameter sharing, sparse interactions, and dominant representations.

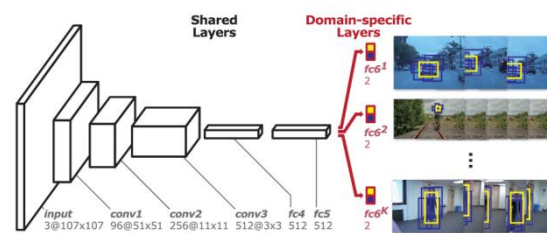


Fig 3: The architecture of Multi-Domain Network (MDNet)

Fig. 3 consists of shared layers and K branches of domain-specific layers. Yellow and blue bounding boxes denote the positive and negative samples in each domain, respectively[4].

To overcome the limitations of pre-trained deep CNNs and take full advantage of end-to-end learning for real-time applications, Siamese Neural Network (SNN) has evolved, some of the SNN based approaches in recent years were proposed to achieve real-time speed. These proposed SNN based methods utilize combination of discriminative target representation, adapting target appearance variation and balancing training data.

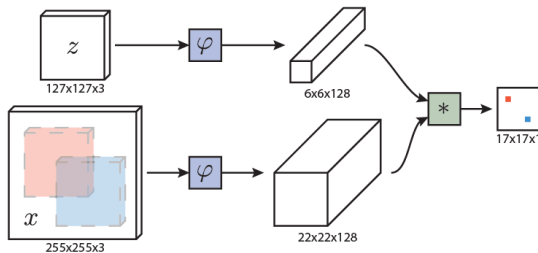


Fig 4: Fully-convolutional Siamese architecture

The Siamese architecture is fully convolutional with respect to the search image x . The output is a scalar-valued score map whose dimension depends on the size of the search image. This enables the similarity function to be computed for all translated sub-windows within the search image in one evaluation. In this example, the red and blue pixels in the score map contain the similarities for the corresponding sub-windows [21].

Visual Object Tracking involves both spatial and temporal information of video frames, over the recent years. Recurrent Neural Network (RNN) architecture based methods are proposed to consider motion or movement simultaneously. Because of tedious training and a numerous number of parameters, the number of RNN-based methods is limited comparatively in the available literature.

Couple of recent methods in literature utilized Generative Adversarial Network (GAN) architecture to address the imbalance distribution of training samples and also to deal self-learning problem of visual tracking.

The trend in describing custom networks is also seen, which is a combination of multiple above mentioned network architectures like CNN, SNN, RNN and GAN to mainly tackle the limitations of ordinary methods by exploiting the advantages of other network structures.

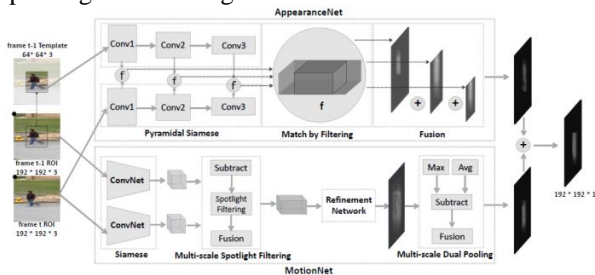


Fig 5: Deep collaborative tracking network [19].

The majority of the basic deep learning based visual tracking methods exploits end-to-end learning with train or re-trains a DNN by applying gradient based optimization algorithms.

The main contributions of this review are summarized as below:

- In depth study and analysis of Deep Learning based trackers recent advances in various aspects.

- Summarize experimental results of trackers cited in research literature publications.
- Describe present research issues of visual tracking research and present some interesting potential future research directions.

The rest of this paper is organized as follows: At first, various deep visual tracking methods from the available literature are presented and then described and discussed some experimental comparisons of the presented visual tracking methods. Finally, some of the research issues are presented with future directions at the conclusion.

SOME RECENT METHODS IN THE LITERATURE

[L. Bertinetto, J. Valmadre., 2016] [21] proposed an alternative approach that focuses on learning strong embedding's in an offline phase. Their experiments show that deep embedding's provide a naturally rich source of features for online trackers, and enable simplistic test-time strategies to perform well. They described an algorithm with a novel fully-convolutional Siamese network trained end-to-end on the ILSVRC15 dataset for object detection in video.

[H. Nam and B. Han. 2016] [4] proposed a multi-domain learning framework based on CNNs, which separates domain-independent information from domain-specific one, to capture shared representations effectively. They have successfully implemented a framework to learn domain-specific information adaptively.

[D. Held, S. Thrun, and S. Savarese 2016] [6] at a high level, they feed frames of a video into a neural network and the network successively outputs the location of the tracked object in each frame. They trained the tracker entirely offline with video sequences and images. Their tracker learns offline a generic relationship between an object's appearance and its motion, allowing network to track novel objects at real-time (100fps) speeds.

[S. Yun, J. J. Y. Choi., 2016] [15] [16] proposed a tracker which is controlled by action-decision network (ADNet), pursues the target object by sequential actions iteratively trained by deep reinforcement learning. They proposed a tracker which is designed to achieve a light computation as well as satisfactory tracking accuracy in both location and scale. In this method the deep network to control actions is pre-trained using various training sequences and fine-tuned during tracking for online adaptation to target and background changes. This tracker achieved a real-time speed (15 fps).

[M. Danelljan, G. Bhat., 2017] [18] presented a novel formulation which addresses the issues of state-of-the-art DCF trackers. Described a factorized convolution operator that dramatically reduces the number of parameters in the DCF model. Designed a compact generative model of the training sample space that effectively reduces the number of samples in the learning, while maintaining their diversity. Comprehensive experiments clearly demonstrate that this approach concurrently improves both tracking performance and speed.

[B. Li, W. Wu., 2018] [1] Described a deep analysis of Siamese trackers and proved that when using deep

networks the decrease in accuracy comes from the destroying of the strict translation invariance. Proposed a sampling strategy to break the spatial invariance restriction which successfully trains Siamese tracker driven by a ResNet architecture.

[Y. Song, C. Ma., 2018] [2] proposed to use a generative adversarial network (GAN) to augment positive samples in the feature space to capture a variety of appearance changes over a temporal span. Demonstrated to use higher-order cost sensitive loss to mine hard negative samples to handle class imbalance. Conducted extensive validation of method on benchmark datasets with large-scale sequences.

[H. Fan and H. Ling 2018] [5] presented a multistage tracking framework, the Siamese Cascaded RPN (CRPN), to solve the problem of class imbalance by performing hard negative sampling. They designed a novel feature transfer block (FTB), enables to fuse the high-level features into low-level RPN, which further improves its discriminative power to deal with complex background, resulting in better performance of C-RPN. They have conducted extensive experiments on six benchmarks

[C. Sun, D. Wang., 2018] [7] in this model, they proposed and developed a spatial-aware KRR model by introducing a cross-patch similarity kernel. They have implemented a model with both regression coefficients and patch reliability, which enables a model to be robust to the unreliable patches. The regression coefficient and similarity weight vectors are simultaneously optimized via an end-to-end neural network.

[J. Choi, H. J. Chang., 2018] [8] proposed a visual tracking framework based on context-aware deep feature compression using multiple auto-encoders. In this model they introduced a context-aware scheme which includes expert auto encoders specializing in one context, and a context-aware network which is able to select the best expert auto-encoder for a specific tracking target. Conducted experiments lead to the compelling finding that this framework achieved a high-speed tracking ability of over 100 fps.

[G. Bhat, J. Johnander., 2018] [9] analyzed the influential characteristics of deep and shallow features for visual tracking. They systematically studied the impact of a variety of data augmentation techniques. Deep investigation of the accuracy-robustness trade-off in the discriminative learning of the target model has done. They proposed a fusion strategy to combine the deep and shallow appearance models leveraging their complementary characteristics. Experiments are performed on four challenging datasets.

[S. Pu, Y. Song., 2018] [10] Described and implemented a reciprocal learning algorithm to exploit visual attention within the tracking by detection framework. In this method, for temporal robust features they introduced attention maps as regularization terms coupled with the classification loss to train deep classifiers. They also conducted experiments on benchmark datasets.

[Y. Zhang, L. Wang., 2018] [11] proposed a local pattern detection scheme, which can automatically identify discriminative local parts of target objects. To achieve more accurate tracking results, they implemented the

message passing process via differentiable operations, and reformulate it through a neural network module.

[Z. Zhu, Q. Wang., 2018] [12] Distractor-aware Siamese Region Proposal Networks (DaSiamRPN) method, Analyzed the features used in conventional Siamese trackers in detail and they found that the imbalance of the non-semantic background and semantic distractor in the training data is the main obstacle for the learning. Proposed a framework to learn distractor-aware features in the off-line training, and explicitly suppress distractors during the inference of online tracking.

[C. Sun, D. Wang., 2018] [13] described a model which includes both discrimination and reliability information using the correlation filter framework. They have introduced local response consistency constraint to ensure that different sub-regions of the base filter have similar importance. In this model to depict the importance of each sub-region in the filter (i.e. reliability learning) the reliability weight map is exploited. This tracker is insusceptible to the non-uniform distributions of the feature map, and can better suppress the background regions. Extensive experiments have conducted to show the superiority of the algorithm compared and this tracker achieves remarkable tracking performance on the OTB-2013, OTB-2015 and VOT-2016 benchmarks.

[F. Li, C. Tian., 2018] [14] presented a spatial-temporal regularized correlation filters (STRCF) model by incorporating both spatial and temporal regularization into the DCF framework. The proposed STRCF serves as an approximation of SRDCF with multiple training samples. They have implemented ADMM algorithm for solving STRCF efficiently, where each sub-problem has the closed form solution. This STRCF Model with hand-crafted feature can run in real-time, achieves notable improvements over SRDCF by tracking accuracy.

[X. Jiang, X. Zhen., 2018] [17] presented Deep Collaborate Tracking Network (DCTN) a unified framework that jointly encodes both appearance and motion information for generic object tracking. They described establishing a unified tracking framework of a two-stream network that can fully capture complementary motion and appearance information with an end-to-end learning architecture; Described a motion net (MotionNet) to fulfill end-to-end trainable motion detection and an appearance net (AppearanceNet) for multi-scale appearance matching to achieve object localization.

[Q. Wang, L. Zhang., 2018] [23] introduced SiamMask, a simple approach that enables fully-convolutional Siamese trackers to produce class-agnostic binary segmentation masks of the target object. They proposed two variants of SiamMask are initialized with a simple bounding box, operate online, run in real-time and do not require any adaptation to the test sequence.

[K. Dai, D. Wang., 2019] [22] presented a novel Adaptive Spatially Regularized Correlation Filters (ASRCF) model to simultaneously optimize the filter coefficients and the spatial regularization weight. Their tracker effectively and efficiently estimates both location and scale with two Correlation Filter models: one exploits complicated features for accurate localization; and the other exploits shallow

features for fast scale estimation. Conducted extensive experiments on five recent benchmarksshow that this tracker performsfavorably against many state-of-the-art algorithms, with real-time performance of 28fps.

[Z. Zhang and H. Peng 2019] [3] presented a systematic study on the factors of backbone networks that affect tracking accuracy, and provides architectural design guidelines for the Siamese tracking framework. They have found and described receptive field size, network padding and stride are crucial factors. Based on no padding, residual units they have designed new deeper and wider network architectures for Siamese trackers. Conducted multiple experiments on five benchmarks baseline datasets.

ANALYTICAL REVIEW

According to the in depth analysis, the deep visual object tracking methods that are consistent in performance on some standard visual tracking datasets are[2], [4], [10], [23], [3], [1], [5], [11], [22], [12], [9], [7], [14], and [13].These methods performed well in terms of precision success measures and accuracy robustness on some famous standard Data Sets.The research results are promising when considered individual visual challenging attributes, but the results are not encouraging when considered all the visual challenging attributes simultaneously.

The methods [2], [4], and [10] take advantage from both Offline and Online Training of Deep Neural Network (DNN), However these methods lacking in speed (≤ 1 Frames Per Second) because of huge computational complexity, for real time applications these methods will not be suitable.

According the analysis, though [2] out performed in couple of visual challenging scenarios(deformation (DEF),in-plane rotation (IPR), out-of-plane rotation (OPR)) but lacking in robustness when significant scale variation (SV) present in a scene.

Based on systematic study of [15, 16] reveals that this method failed to follow the abrupt movement of the target and the proposed actions could not adapt to the sudden aspect ratio change.

Though [1] performed well in terms of accuracy, Overall system speed is far less comparative to real time. In[7] utilizes kernelized ridge regression (KRR) to concentrate on reliable regions of the target, consideration of rotation information andonline adaptation of Convolution Neural Network (CNN) models.This method responded well to the deformationand in-plane rotationvisual challenges.

The methods [14], [23], [13] and [9] are the DCF based and describe on fusing the HOG with deep off the self-features to improve the consistency of the results. Though these methodsshow the competitive performance but very much suffer from the limitations of the computational complexity of appearance variation and deep features.

According to the research results[1], [22], [5], [3], [18], [7], [13], [9], and [14] on some standard data sets have failed in scenarios that consist of simultaneous multiple critical visual attributes.

In depth study of [9], [14], [13], [7], and [23]describes exploiting on deep off-the-shelf features and take the

advantage of DCF framework to address some of the challenging visual attributes.

The methods [3], [1], [5], [22] and [12] are based on the fast SiamRPN method and exploit on one-shot detection task to solve some of the visual tracking problems.Based on research results of [4] performs well on deformation, low resolution and fast motion scenarios.

In [6]When there is a size change or no variation, the tracker performs slightly worse when using the previous frame. Under a large size change, the corresponding appearance change is too drastic for network to perform an accurate comparison between the previous frame and the current frame. The tracker is acting as a local generic object detector in such a case.

In [5], [11], and [12] methods exploit the shallow AlexNet as their backbone network. In [11] describes to decrease the sensitivity of SNN-based (Siamese Network) methods specifically for non-rigid appearance change and partial occlusion (POC) attributes, this method detects contextual information of local patterns and their relationships and matches them by a Siamese network in real-time speed.

In the method [22]present three-branch architecture to estimate the target location by a rotated Bounding Box,which includes the binary mask of the target.The most failure reasons of SiamMask are the motion blur (MB) and out-of-view (OV) attributes that produce erroneous target masks.

RESEARCH CHALLENGES

Despite rapid considerable advancements that are emerged in Deep Visual Tracking, the mentioned trackers from research literature are still unable to handle the real-world challenges efficiently.

Studies on references from literature revealed thatdeep learning based methods are still not reliable for real-world applications as they are lacking in intelligent situation understanding with real time speed.

It is understood that despite decades of research still have the problems to simultaneously handle challenging scenarios which significantly consists of visual attributes such as OCC(Occlusion), OV (Out-of-View), DEF(Deformation), SV(Scale Variation) and FM(Fast Motion).

Some tracking approaches perform well in specific video/image scenarios and Standard Data Sets.While applied to other cases, however, they maynot produce satisfying results.

Though some of the methods presented performs well in a challenging scenarios but they are not robust enough to handle the diversity of situations.

Based on the review it is understood that maintaining accuracy in numerous situations, robustness to visual variances and computational efficiency all at once is an existing research challenge.

Based on thorough analysis computational complexity and memory usage is the biggest research challenge in Deep Learning Based Visual Tracking even to address single challenging visual attribute.

All these issues restrict further development of the Visual Tracking research and its applications in real-world, real-time systems. Recently, attempts to deal with some of these issues have been made, for example, the Benchmark new data sets provides a large set of testing video sequences, standard baseline evaluation tools, new methods of evaluation etc. This is likely to advance the further studies and developments of Visual Object Tracking techniques.

CONCLUSION AND FUTURE DIRECTIONS

According to the study, the following observations have been made:

- The most difficult attributes for Deep Learning based visual tracking methods are OCC (occlusion), DEF (deformation), OV (Out of View), SV (scale variation) and FM (fast motion).
- Multiple complementary features exploitation from different efficient Deep Learning based tracking methods may improve robustness of the model.
- Integration of both offline and online training methods may lead to more robust visual trackers.
- Improving tracking efficiency in handling simultaneously challenging visual attributes.
- To handle success rate, robustness and efficiency need more exploitation on integrating multiple network architectures and the efficient combination of different methods.
- Work towards more intelligent methods to reduce the computational complexity that is a necessity in real-time applications to achieve real time performance.

REFERENCES

- [1] B. Li, W. Wu, Q. Wang, F. Zhang, J. Xing, and J. Yan, "SiamRPN++: Evolution of Siamese visual tracking with very deep networks," 2018. [Online]. Available: <http://arxiv.org/abs/1812.11703>
- [2] Y. Song, C. Ma, X. Wu, L. Gong, L. Bao, W. Zuo, C. Shen, R. W. Lau, and M. H. Yang, "VITAL: Visual tracking via adversarial learning," in Proc. IEEE CVPR, 2018, pp. 8990–8999.
- [3] Z. Zhang and H. Peng, "Deeper and wider Siamese networks for real-time visual tracking," 2019. [Online]. Available: <http://arxiv.org/abs/1901.01660>
- [4] H. Nam and B. Han, "Learning multi-domain convolutional neural networks for visual tracking," in Proc. IEEE CVPR, 2016, pp. 4293–4302.
- [5] H. Fan and H. Ling, "Siamese cascaded region proposal networks for real-time visual tracking," 2018. [Online]. Available: <http://arxiv.org/abs/1812.06148>
- [6] D. Held, S. Thrun, and S. Savarese, "Learning to track at 100 FPS with deep regression networks," in Proc. ECCV, 2016, pp. 749–765.
- [7] C. Sun, D. Wang, H. Lu, and M. Yang, "Learning spatial-aware regressions for visual tracking," in Proc. IEEE CVPR, 2018, pp. 8962–8970.
- [8] J. Choi, H. J. Chang, T. Fischer, S. Yun, K. Lee, J. Jeong, Y. Demiris, and J. Y. Choi, "Context-aware deep feature compression for high-speed visual tracking," in Proc. IEEE CVPR, 2018, pp. 479–488.
- [9] G. Bhat, J. Johnander, M. Danelljan, F. S. Khan, and M. Felsberg, "Unveiling the power of deep tracking," in Proc. ECCV, 2018, pp. 493–509.
- [10] S. Pu, Y. Song, C. Ma, H. Zhang, and M. H. Yang, "Deep attentive tracking via reciprocal learning," in Proc. NIPS, 2018, pp. 1931–1941.
- [11] Y. Zhang, L. Wang, J. Qi, D. Wang, M. Feng, and H. Lu, "Structured Siamese network for real-time visual tracking," in Proc. ECCV, 2018, pp. 355–370.
- [12] Z. Zhu, Q. Wang, B. Li, W. Wu, J. Yan, and W. Hu, "Distractor aware Siamese networks for visual object tracking," in Proc. ECCV, vol. 11213 LNCS, 2018, pp. 103–119.
- [13] C. Sun, D. Wang, H. Lu, and M. H. Yang, "Correlation tracking via joint discrimination and reliability learning," in Proc. IEEE CVPR, 2018, pp. 489–497.
- [14] F. Li, C. Tian, W. Zuo, L. Zhang, and M. H. Yang, "Learning spatial-temporal regularized correlation filters for visual tracking," in Proc. IEEE CVPR, 2018, pp. 4904–4913.
- [15] S. Yun, J. J. Y. Choi, Y. Yoo, K. Yun, and J. J. Y. Choi, "Action decision networks for visual tracking with deep reinforcement learning," in Proc. IEEE CVPR, 2016, pp. 2–6.
- [16] S. Yun, J. Choi, Y. Yoo, K. Yun, and J. Y. Choi, "Action-driven visual object tracking with deep reinforcement learning," IEEE Trans. Neural Network. Learn. Syst., vol. 29, no. 6, pp. 2239–2252, 2018.
- [17] X. Jiang, X. Zhen, B. Zhang, J. Yang, and X. Cao, "Deep collaborative tracking networks," in Proc. BMVC, 2018, p. 87.
- [18] M. Danelljan, G. Bhat, F. Shahbaz Khan, and M. Felsberg, "ECO: Efficient convolution operators for tracking," in Proc. IEEE CVPR, 2017, pp. 6931–6939.
- [19] Xiaolong Jiang, Xiantong Zhen, Baochang Zhang, Jian Yang, Xianbin Cao "Deep Collaborative Tracking Networks"
- [20] Seyed Mojtaba Marvasti-Zadeh, Li Cheng, Hossein Ghanei-Yakhdan, and Shohreh Kasaei "Deep Learning for Visual Tracking: A Comprehensive Survey".
- [21] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr, "Fully-convolutional Siamese networks for object tracking," in Proc. ECCV, 2016, pp. 850–865.
- [22] K. Dai, D. Wang, H. Lu, C. Sun, and J. Li, "Visual tracking via adaptive spatially-regularized correlation filters," in Proc. CVPR, 2019, pp. 4670–4679.
- [23] Q. Wang, L. Zhang, L. Bertinetto, W. Hu, and P. H. S. Torr, "Fast online object tracking and segmentation: A unifying approach," 2018. [Online]. Available: <http://arxiv.org/abs/1812.05050>
- [24] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui, "Visual object tracking using adaptive correlation filters," in Proc. IEEE CVPR, 2010, pp. 2544–2550.
- [25] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-speed tracking with kernelized correlation filters," IEEE Trans. Pattern Anal. Mach. Intell., vol. 37, no. 3, pp. 583–596, 2015.
- [26] M. Danelljan, G. Hager, F. S. Khan, and M. Felsberg, "Discriminative Scale Space Tracking," IEEE Trans. Pattern Anal. Mach. Intell., vol. 39, no. 8, pp. 1561–1575, 2017.
- [27] S. M. Marvasti-Zadeh, H. Ghanei-Yakhdan, and S. Kasaei, "Rotation-aware discriminative scale space tracking," in Iranian Conf. Electrical Engineering (ICEE), 2019, pp. 1272–1276.
- [28] G. Boudoukh, I. Leichter, and E. Rivlin, "Visual tracking of object silhouettes," in Proc. ICIP, 2009, pp. 3625–3628.
- [29] C. Xiao and A. Yilmaz, "Efficient tracking with distinctive target colors and silhouette," in Proc. ICPR, 2016, pp. 2728–2733.
- [30] V. Bruni and D. Vitulano, "An improvement of kernel-based object tracking based on human perception," IEEE Trans. Syst., Man, Cybern. Syst., vol. 44, no. 11, pp. 1474–1485, 2014.
- [31] W. Chen, B. Niu, H. Gu, and X. Zhang, "A novel strategy for kernel-based small target tracking against varying illumination with multiple features fusion," in Proc. ICICT, 2018, pp. 135–138.
- [32] D. H. Kim, H. K. Kim, S. J. Lee, W. J. Park, and S. J. Ko, "Kernel-based structural binary pattern tracking," IEEE Trans. Circuits Syst. Video Technol., vol. 24, no. 8, pp. 1288–1300, 2014.
- [33] I. I. Lychkov, A. N. Alfimtsev, and S. A. Sakulin, "Tracking of moving objects with regeneration of object feature points," in Proc. GloSIC, 2018, pp. 1–6.