

Object Detection Using CNN

¹N. Ramesh, ²P. Rakshitha Reddy, ³K. Mahesh, ⁴L. Vishnu, ⁵Dr. Lakshmaiah,

^{1,2,3,4} U.G.Scholar, Department of ECE, Sri Indu College Of Engineering & Technology, Ibrahimpatnam, Hyderabad.

⁵Research Guide, Department of ECE, Sri Indu College Of Engineering & Technology, Ibrahimpatnam, Hyderabad.

Abstract:-This paper reviews recent advances in Deep Learning-based Visual Object Tracking Approaches from various algorithms, highlighting key issues with existing research and suggesting future potential directions. It provides a comprehensive review of recent advancements in this complex area of computer vision, highlighting the importance of real-world application scenarios in this field. This paper reviews recent advancements in Deep Learning-based Visual Object Tracking Approaches, analyzing various algorithms. It highlights key issues in existing research and suggests future research directions, concluding with a comprehensive review of recent developments.

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

INTRODUCTION

Tracking is the process of inferring an object's motion from images. Visual object tracking, a field of computer vision, has applications in various fields like driver assistance, medical, military, video surveillance, traffic control, navigation, robotics, augmented reality, and sports. Visual object tracking, extensively studied in computer vision, is utilized in real-world applications such as Advance Driver Assistance Solutions, Medical, Military, Video Surveillance, Traffic Control, Navigation, Robotics, Augmented Reality, and sports.

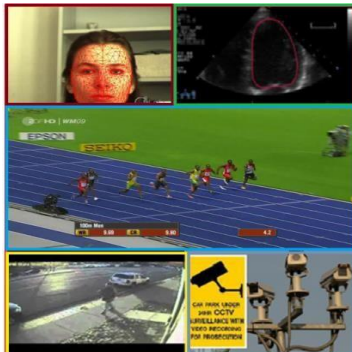


Fig1: Some Visual Tracking Applications

Basic traditional visual tracking methods utilize various frameworks like Discriminative Correlation Filters (DCF), silhouette tracking, Kernel tracking, point tracking, 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 the 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.

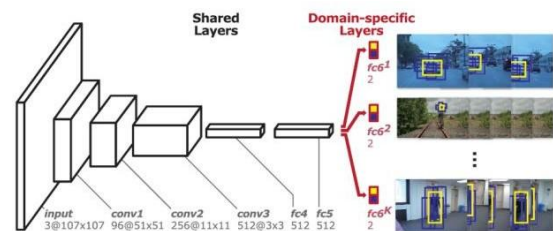


Fig3: The architecture of Multi-Domain Network (MDNet)

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.

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.

Siamese Neural Network (SNN) Evolution for Real-Time Applications

- Overcomes limitations of pre-trained deep CNNs.
- Utilizes discriminative target representation, adapting target appearance variation.
- Balances training data for real-time speed.

These proposed SNN based methods utilize combination of discriminative target representation, adapting target appearance variation and balancing training data.

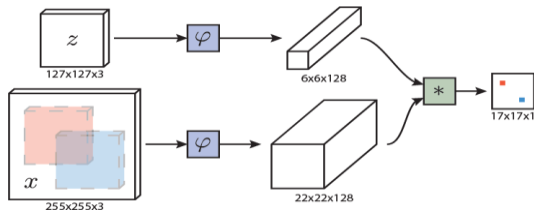


Fig4: Fully-convolutional Siamese architecture

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 movements 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.

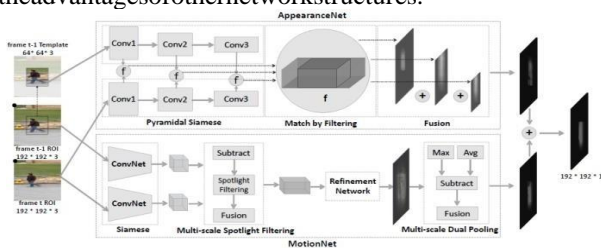


Fig5: Deep collaborative tracking network [19].

The majority of the basic deep learning based visual tracking method 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 represented and then described and discuss some experimental comparisons of the presented visual tracking methods. Finally, some of the research issues are represented with future directions at the conclusion.

SOME RECENT METHODS IN THE LITERATURE

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 representation 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.

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 (15fps).

Described a factorized convolution operator that dramatically reduces the number of parameters in the DCF model. Designed a compact generative model of the training samplespace 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. Siamese Neural Network (SNN) has evolved to overcome pre-trained deep CNN limitations and optimize end-to-end learning for real-time applications. Recent approaches use discriminative target representation, adapting target appearance variation, and balancing training data to achieve real-time speed. Deep analysis of Siamese trackers and proved that when using dee

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 multi-stage 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), enable 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. 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 100fps.

[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 susceptible to the non-uniform distributions of the feature map, and can better suppress the background regions. Extensive experiments have been 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 improvement over SRDCF by tracking accuracy.

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.

SiamMask, a simple approach that enables fully-convolutional Siamese tracker 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.

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 benchmark sets show that this tracker performs favorably against many state-of-the-art algorithms, with real-time performance of 28 fps.

This paper presented a systematic study on the factors of backbone network 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 benchmark baseline datasets.

ANALYTICAL REVIEW

According to the in-depth analysis, the deep visual object tracking methods that are consistent in performance on standard visual tracking datasets are [2], [4], [10], [3], [1], [5], [11], [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 lack in speed (≤ 1 Frames Per Second) because of huge computational complexity, for real-time applications these methods will not be suitable.

According to the analysis, though [2] outperformed 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) presents in a scene.

Based on systematic study of reveal 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 and online adaptation of Convolutional Neural Network (CNN) models. This method responded well to the deformation and in-plane rotation visual challenges.

The methods [14], [13] and [9] are the DCF based and describe on fusing the HOG with deep off-the-shelf features to improve the consistency of the results. Though these methods show 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], [5], [3], [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] 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] 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 motions 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 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 that deep 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 may not 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

The study reveals that the most challenging attributes for Deep Learning-based visual tracking methods are OCC, DEF, OV, SV, and FM. To improve robustness, multiple complementary features from efficient methods should be exploited. Integrating both offline and online training methods can lead to more robust visual trackers. To achieve success rates, more efficient combination of network architectures and intelligent methods should be employed.

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] 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.
- [14] 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.