

Learning Attention for Object Tracking with Adversarial Learning Network

Xu Cheng (✉ xcheng@nuist.edu.cn)

Nanjing University of Information Science and Technology <https://orcid.org/0000-0003-2355-9010>

Chen Song

Nanjing University of Information Science and Technology

Yongxiang Gu

Nanjing University of Information Science and Technology

Beijing Chen

Nanjing University of Information Science and Technology

Research

Keywords: Surveillance, Deep Learning, Object Tracking, Generative Adversarial Learning

Posted Date: February 28th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-15512/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on November 11th, 2020. See the published version at <https://doi.org/10.1186/s13640-020-00535-1>.

Abstract

Artificial intelligence has been widely studied on solving intelligent surveillance analysis and security problems in recent years. Although many multimedia security approaches have been proposed by using deep learning network model, there are still some challenges on their performances which deserve in-depth research. On one hand, high computational complexity of current deep learning methods makes it hard to be applied to real-time scenario. On the other hand, it is difficult to obtain the specific features of a video by fine-tuning the network online with the object state of the first frame, which fails to capture rich appearance variations of the object. To solve above two issues, in this paper, an effective object tracking method with learning attention is proposed to achieve the object localization and reduce the training time in adversarial learning framework. First, a prediction network is designed to track the object in video sequences. The object positions of the first ten frames are employed to fine-tune prediction network, which can fully mine a specific features of an object. Second, the prediction network is integrated into the generative adversarial network framework, which randomly generates masks to capture object appearance variations via adaptively dropout input features. Third, we present a spatial attention mechanism to improve the tracking performance. The proposed network can identify the mask that maintains the most robust features of the objects over a long temporal span. Extensive experiments on two large-scale benchmarks demonstrate that the proposed algorithm performs favorably against state-of-the-art methods.

1. Introduction

With rapid growth and use of multimedia signal processing and Internet technology, a number of multimedia security issues have also emerged correspondingly in recent years, such as intelligent analysis for surveillance, copy-move forgery in digital images and videos, and biometric spoofing. Meanwhile, artificial intelligence has been widely studied on solving a variety of difficult problems using deep learning network model, such as convolution neural networks for steganalysis and forensics, and generative adversarial networks for coverless steganography.

Intelligent analysis for surveillance technology has been a hot topic in multimedia security community. It is the basis of advanced video processing tasks such as follow-up steganography [1], data hiding [2], JPEG compressed [3] and object recognition [4], and is a necessary prerequisite for implementing high-level intelligent behavior analysis. Object tracking is one of the fundamental tasks in intelligent surveillance technology. The aim is to localize the object in a video sequence with a bounding box of the object at the first frame of the video. Although object tracking has made great progress in recent years and some effective algorithms have been proposed to solve challenging problems in specific scenarios, there are still exist many issues such as occlusion, illumination changes, etc. So the topic deserves in-depth investigation, which is important in both academia and industry. Figure 1 shows intelligent surveillance applications.

Currently, most of the trackers based on deep learning network model use a large-scale benchmark datasets [5, 6] to train the network offline and the sample of the first frame is used to fine-tune the parameters of network online. However, training deep network model online is challenging due to the limited training samples which cannot capture the diversity of the object appearance variations, and offline pre-training is very time consuming. In addition, the quality of the image is also important factor for the training process. The existing tracking methods utilize the sampling scheme around the object state to obtain the training samples, such as KCF [15], DLT [29], MCPF [53], etc. However, the positive samples extracted from each frame are highly overlapped and they fail to capture rich appearance variations. In this work, we use the first ten frames of a video sequence to fine-tune the parameters of deep network model. But manual annotating the object positions from the first ten frames are always impractical and time consuming. To solve the above-mentioned problems, we exploit the advantages of the pre-trained network on a large-scale benchmark datasets to predict the object position. The parameters of deep network have been pre-trained on large-scale datasets. Then we use the prediction network to track the object in the video sequences and automatically obtain the object positions of the first ten frames. The generative adversarial network has great advantages in augmenting training samples. Furthermore, the positive and negative samples are extracted from the first ten frames. Therefore, the positive and negative samples are obtained to fine-tune the generative adversarial network online. The proposed tracking algorithm can capture the changes of the object appearance in the video sequences. During the tracking, the generative model is used to occlude the image features through a randomly generated mask to enhance the diversity of positive samples. The discriminative model employs its discriminative performance to identify the object. These features are robust enough to address the challenges of object tracking. Adversarial learning network can identify masks that retain the most robust features of the object appearance.

We summarize the main contributions of this work as follows:

- (1) We propose an end-to-end the prediction network model for object tracking to improve the tracking accuracy and the computational complexity of training process, which can jointly train the prediction network and spatial attention model in generative adversarial network framework.
- (2) An effective spatial attention mechanism is developed, which can adaptively generate the response maps. The feature representations are employed to online tracking process to alleviate over-fitting. In addition, the positive and negative samples are augmented in the feature space to capture a variety of appearance changes over a temporal span by using generative adversarial network.
- (3) We conduct the extensive experiments on two popular benchmarks, which demonstrate the proposed object tracking method with learning attention significantly outperforms state-of-the-art methods.

The rest of the paper is organized as follows. In Sect. 2, we review related work of existing object tracking algorithms. Section 3 briefly introduces the generative adversarial network. In Sect. 4, we introduce our object tracking approach for intelligent surveillance analysis. The experimental settings are presented in

Sect. 5. In Sect. 6, we present experimental results and discussion in two tracking benchmarks. Finally, Sect. 7 concludes this paper.

2. Related Work

Object tracking is one of the fundamental tasks in computer vision and has been extensively studied over the last decade. There are extensive surveys of object tracking in the literature [7, 8, 9]. In this section, we mainly discuss the representative object trackers based on correlation filters and deep learning.

Correlation filter based tracking: Before the emergence of the object tracking algorithms based on correlation filters, most tracking algorithms used particle filter framework for object tracking, such as Kalman filter [10] and particle filter [11]. However, the disadvantage of these methods is that the number of particles limits the tracking speed. In recent years, correlation filters (CF) have been widely used in numerous applications such as object detection and recognition. It transfers operations into the Fourier domain as element-wise multiplication. Correlation filters [12, 13, 14, 15, 16] have attracted considerable attention due to its computational efficiency and competitive performance. Correlation filters trackers regress all the circular shifted versions of the input features to a Gaussian function. Bolme et al. [12] propose a minimum output sum of squared error (MOSSE) tracker for object tracking on grayscale images, which encodes object appearance through an adaptive correlation filter by optimizing the output sum of squared error. MOSSE can achieve several hundreds of frames per second. In 2012, Heriques et al. [17] propose the CSK algorithm based on the improvement of MOSSE, which solves the problem of a small number of training samples in the object tracking process through the cyclic matrix, and further improves the tracking accuracy of the algorithm by using the kernel technique. However, the above two algorithms adopt simple grayscale features, which are easily disturbed by the external environment, resulting in inaccurate tracking results. They are further improved by the kernelized correlation filters (KCF) [15] with HOG features in a Fourier domain. KCF performs well in OTB50 [4] benchmark in terms of tracking speed and accuracy. In addition, scale change is also a common problem in object tracking. In [18], the DSST tracker learns adaptive multi-scale correlation filters using HOG features to handle the scale and translation changes of the objects. The SRDCF method [19] reduces the boundary effect problem by weighting the weight space of CF. However, its optimization is complicated and the tracking speed is slow. To improve its weakness, the CSR-DCCF method [20] adds feature channel and space stability constraints based on SRDCF, and uses the augmented Lagrangian scheme to facilitate fast FFT solution, which greatly improves the tracking accuracy and speed. The C-COT method [21] proposes a strategy for training continuous convolution filters, which facilitates the integration of multi-scale CNN features and achieves sub-pixel level tracking accuracy. However, the tracking model framework still adopts the SRDCF method and the computational complexity is high. ECO [22] method further proposes a factorization convolution scheme to reduce the computational complexity of the tracking model. The UPDT [23] proposes a novel adaptive fusion approach that leverages the complementary properties of deep and shallow features to improve both robustness and accuracy. In [24], the STRCF introduces temporal regularization to SRDCF with single sample. The formulation of this method can not only serve as a reasonable approximation to SRDCF with multiple training samples, but also provide a more robust

appearance model than SRDCF in the case of large appearance variations. Although these approaches achieve a satisfied performance in some constrained scenarios, they have an inherent limitation that they resort to low level hand-crafted features, which are vulnerable in dynamic situations including illumination changes, occlusion, deformations, background clutter etc. Inspired by the success of deep learning model in object detection, recognition and classification tasks, researchers have started to focus on combining of deep learning and correlation filter. In HCF [25] and HDT [26], deep feature is used to extract the object features instead of handcrafted features. It is worth noting that CFNet [27] and DCFNet [28] achieve an end-to-end representation learning.

Deep learning based tracking: Features representations are important for object tracking. Experiments prove that the features designed manually (e.g., haar-like features, histogram, HOG features, etc.) are not necessarily suitable for all the objects. Deep learning in the tracking uses the adaptive selection scheme of object features instead of hand-crafted features. The popular trend is to design the deep network structures and pre-train them in order to learn object-specific features. The main problem of deep learning in the tracking is the lack of training data. The object tracking only provides the object information of the first frame as training data. In this case, it is difficult to train the deep network model with small data. In [29], Wang proposes the idea of off-line pre-training deep network model and online fine-tuning tracking model called DLT tracker, which greatly solves the problem of insufficient training samples in the tracking. SO-DLT [30] continues the strategy of DLT, and also greatly improves the problems of DLT by using CNN as a network model for extraction deep features and classifications. In FCNT [31], authors analyze the performance of CNN features pre-trained on ImageNet [32] and design the subsequent network structure based on the analysis results. In [33], the authors utilize a two-layer convolution neural network to learn hierarchical features from auxiliary video sequences, which takes the complex object motion and object appearance changes into account. In recent years, the Siamese Networks [34, 35] have received more and more attention due to its two stream identical structure. In the offline training phase, a matching score function is trained through the structure of the Siamese Network. Then the matching score function is used to determine the similarity between the current object candidate state and the object template of the first frame during the tracking, which improves the tracking efficiency. In recent years, the generative adversarial network (GAN) has been widely used in many fields, such as object detection [36] and intelligent recommendation system [37]. GAN is first proposed by Ian Goodfellow in 2014, which was originally used to generate realistic-looking images [38]. The main idea behind GAN is to have two competitive neural network models. One takes noise as input and generates samples called a generator. Another model is called discriminator which receives samples from the generator and tries to discriminate true object between two sources. The generator and the discriminator are trained simultaneously by competing with each other. Realizing that there is a huge difference between image classification and tracking, TADA [39] identifies the importance of each convolutional filter and selects the object-aware features based on activations for the object representation. Then the object-aware features are integrated into a Siamese matching network for visual tracking. Different from the existing approaches using extensive annotated data for supervised learning, UDT [40] is trained on large-scale unlabeled videos in an unsupervised manner. UDT tracker achieves the baseline accuracy of fully supervised trackers which

require complete and accurate labels during the training. In this work, we apply adversarial learning to augment training samples in the feature space to capture appearance variations in temporal domain. In addition, we can exploit robust features over the long temporal span instead of the discriminative features in individual frames.

Attention mechanisms: Attention mechanisms were first introduced in neuroscience area [41]. They have spread to image classification, multi-object tracking, etc. DAVT [42] employs a discriminative spatial attention scheme for visual tracking. CSR-DCF [43] utilizes color histograms to construct a foreground spatial map in the correlation filter framework, which learns the attention via an end-to-end deep network model. ACFN [44] chooses a subset from the associated correlation filters as an attention mechanism for visual tracking.

3. Generative Adversarial Network Learning Introduction

The generative adversarial network [38] has been widely used in object detection and semantic segmentation. The generative adversarial networks mainly consist of the generative model and discriminative model. The generative model takes a noise as an input, and the discriminative model takes samples from generative model or training data and outputs the classification probability. This learning process can be written as:

$$L = \min_G \max_D E_{x \sim p_{\text{data}(x)}} [\log D(x)] + E_{z \sim p_{\text{noise}(z)}} [\log(1 - D(G(z)))] \quad (1)$$

where G denotes the generative model; D represents the discriminative model, E is the mean operation. x and z are two vectors from two distributions $P_{\text{data}}(x)$ and $P_{\text{noise}}(z)$, respectively. The training process encourages G to fit $P_{\text{data}}(x)$ so that D will not be able to discriminate x from G(z). After the training process, G is removed and only D is kept for inference.

4. Proposed Object Tracking Method

4.1 Motivation

The existing trackers based on deep learning are performed as off-line training and online fine-tuning for surveillance analysis, and only use the object information of the first frame to fine-tune online the learned deep network parameters. However, it is difficult to capture previously unseen object features from one or few examples. In addition, the positive samples in each frame are highly spatially overlapped and they fail to capture rich appearance variations. On the other hand, amount of positive and negative samples for training deep learning network model are nearly impossible to meet up in real world.

In this section, a novel attention generative adversarial network is given at first to describe the overall training architecture. The proposed generative model takes VGG network as input, which is mainly used to capture the object appearance variations of continuous video frames. The discriminative model is

introduced as a supervisor and provides guidance on the advantages of the generated object appearance details. To stabilize the training of the generative adversarial networks, we present the mean squared loss to punish the classification error for each pixel. In order to improve the tracking performance, a novel spatial attention mechanism is developed to adapt the offline learned deep model to online object tracking. The VGG network is used to sense the tracked object and decode the object features into the attention response maps. At last, online tracking is described consisting of model updating and scales. The object attention maps are captured by inputting the object appearance information provided in the first ten frames and remaining video frames into the generative model. The score with maximum response score is regarded as the tracking result. This process will be continued for video frames until the end of the video sequences. Fig. 2 shows the flowchart of the proposed tracker.

4.2 Network architecture

In Fig 2, the generative model of GAN follows the encoder-decoder framework which attempts to encode the input of the object appearance into feature representation, and decode it into corresponding outputs. The discriminative model is a standard convolutional neural network.

The network includes two branches, and in the lower part of the architecture which utilizes the first ten frames of a video sequences as input, which is called the prediction network. We use the prediction network to track the one to ten frames of a video sequence to obtain the object position of each frame. These features extracted from the predicted object location will be used to fine-tune the fully connected layers of the network located in the top half of the architecture. The object feature of each frame is taken as input of the network from the frame 11 to the end of video. The weight masks are applied to adaptively dropout input features. Adversarial learning identifies the weight mask that maintains the most robust features over a long temporal span while removing the discriminative features from individual frames.

It is worthy to note that the deep learning model is initialized with the weights of a VGG-16 model pre-trained on the ImageNet benchmark for object classification. Most of deep learning based trackers use this offline learned network, and then utilize the first frame to fine-tune the network parameters during the tracking. However, it is difficult to obtain the object specific feature of a video by training the deep network model only with the sample of the first frame. On the other hand, if the deep learning network is fine-tuned by using the first n frames of a video, manually labeling the object position will be expensive and impractical. Therefore, a prediction network is introduced into deep learning framework, which can automatically predict the position of the object in the video sequence. The network structure is shown in Fig.2, which has three convolution layers and two fully connected layers. The architecture of the prediction network is depicted in the lower part of Fig.2. We directly use a VGG-M [45] model pre-trained in the classification task from ImageNet [32], and the parameters of the convolution layers is fixed and only the fully connected layers is fine-tuned online. The cross-entropy loss is adopted for fine-tuning network parameters online. The prediction network is optimized by minimizing the cross-entropy loss function with SGD as follows:

$$\arg \min \frac{1}{N} \sum_{i=1}^N \left(-\sum_{j=1}^2 p(j) \log(q(j)) \right) \quad (2)$$

where p and q denote training samples and corresponding labels, respectively; N is the number of training samples.

The object features are extracted from the convolution layer and fed to the fully connected layer for classification. Fig.3 reports the foreground response maps predicted by using different VGG feature maps. Finally, the sample with the highest response score in each frame is regarded as the tracking result. This prediction network is interpreted as a generative network in generative adversarial network framework and the samples drawn from the predicted location will be used to fine-tune the fully connected layers of the generative model.

The discriminative model is employed to make the generative model produce attention response map that is robust to occlusion, deformation, and background clutter, etc. In this work, the attention response map and corresponding RGB frame of a video sequence are considered as the input of discriminative model.

4.3 Training

In our work, mean squared error (MSE) is utilized to measure the difference between estimated attention response map and ground truth map. Given an image I , and its dimension is $N = W \times H$. The mean squared loss can be formulated as:

$$L_{MSE} = \frac{1}{N} \sum_{j=1}^N (S_j - \hat{S}_j)^2 \quad (3)$$

where S and \hat{S} denote the attention response maps and its corresponding ground truth, respectively.

However, mean squared loss function focuses on pixel-level features, and learned deep network can produce a coarse attention response maps. Therefore, training the network with the adversarial loss can be further improved the tracking performance. We iteratively train G and D , and the adversarial loss function is written as:

$$\begin{aligned} L_{AL} = & \min_G \max_D E_{(C,M) \sim P(C,M)} [\log D(M \cdot C)] \\ & + E_{C \sim P(C)} [\log(1 - D(G(C) \cdot C))] + \lambda E_{(C,M) \sim P(C,M)} \|G(C) - M\|^2 \end{aligned} \quad (4)$$

where C is the input image feature; $G(C)$ is the mask generated by the G network; M is the actual mask identifying the discriminative feature. The dot is the dropout operation on the feature C . As described in Eq.(4), G is used to predict a weight mask $G(C)$ which operates on the extracted features. The mask is randomly initialized at the beginning and each mask represents a specific type of appearance variation. Through the adversarial learning process, G will gradually identify the mask that degrades the performance of classifier.

In each iteration of the training process, object features of the input frames are extracted from convolutional layers and fed into G network to obtain the predicted mask m^* . Then obtained deep features are multiplied by the predicted mask m^* and sent into D network. We keep the labels unchanged and train D through supervised learning method. D is trained to discriminate features from individual frames relying on more robust features over a long temporal span. Thus, it avoids the overfitting issue. G is used to predict different masks according to different input deep features. It enables D to focus on the temporal robust features without discriminative feature interference from single frame. Given an input image, multiple output features based on several random masks are created. Diversified features are performed through the dropout operation, which are sent to D for classification, and we choose the one with the highest loss. The corresponding mask of the selected feature is effective in decreasing the impact of the discriminative features. We set this mask as M in equation (4) and update G accordingly.

Finally, we combine the MSE loss with adversarial loss to obtain more stable and fast convergence for GAN model. The final loss function for the adversarial training can be formulated as:

$$L_{GAN} = L_{AL}(D(C, G(C)), 1) + \lambda L_{MSE} \quad (5)$$

where λ is a trade-off parameter, we experimentally set it as 1/20 in our implementation.

4.4 Spatial Attention

Attention from the training samples can be captured to share a common attention. In practical sceneries, some attention maps are obtained by the initialization of matrix of ones. They are too restrictive to constrain all samples and the object to share a single deep network structure. Therefore, we propose a spatial attention scheme to model attention response map. The proposed attention mechanism can capture the general features and distinct the object from the background in the video. It can encode the global information of the object and has a low computational load. The output of attention module is passed through a global pooling layer to produce a channel-wise descriptor. Then three fully connected (FC) layers are added, in which learned for each channel by a self-gating mechanism based on channel dependence. This is followed by reweighting the original feature maps to generate the output of attention module. The cosine similarity criterion is utilized to measure the similarity between current frame features $\varphi_t(p)$ and the features $\varphi_{t-1}(p)$ extracted from $t-1$ frame.

$$w_t(p) = \text{SoftMax}\left(\frac{\varphi_t(p) \cdot \varphi_{t-1}(p)}{\|\varphi_t(p)\| \|\varphi_{t-1}(p)\|}\right) \quad (6)$$

If the current frame features is close to the features of the last frame, it is prone to the foreground object and assigned with a larger weight, otherwise, a smaller weight is assigned to background pixel.

4.5 Online Tracking

In this subsection, we illustrate how our tracker works for visual object tracking. We involve the generative model during the training and remove it in the tracking stage.

We first draw the samples from the first ten frames of a video sequence to fine-tune generative model online. Then, we track the object in all videos. Given an input frame, we generate multiple candidate proposals and extract their deep features. Deep features of the candidate proposals are fed into the classifier to obtain the probability scores. During the online update, we employ these training samples jointly train the generative model and the discriminative model. The object tracking result is obtained by finding the maximum response score in the attention map.

Model updating plays a critical role in object tracking, and most of trackers update their appearance model in each frame or at a fixed interval. However, this updating strategy may introduce background information into the object appearance model when the tracking result is inaccurate due to occlusion or illumination variations. In this paper, model updating is performed when the number of iteration or maximum value of response map are satisfied.

To handle the scale change, we follow the approach in [18] and use patch pyramid with the scale factors. The proposed object tracking algorithm can be summarized as follows.

Algorithm 1 The Proposed Tracking Algorithm

1. **Input:** Initial object bounding box (x_0, y_0, w_0, h_0); video frame sequence I_i ($i=1,2,3\dots$);
2. **Output:** Estimated object state (x_i, y_i, w_i, h_i);
3. **Initialization:**
Use the prediction network to get the object position of the first ten frames and fine-tune the generative adversarial network with the obtained object position;
4. **Tracking process:**
 5. **for** $i = 11 : \text{end}$
 6. Sample the object candidate states in frame i according to the previous object position ($x_{i-1}, y_{i-1}, w_{i-1}, h_{i-1}$);
 7. Extract deep features of these candidates using convolution layers;
 8. Feed extracted deep features into the discriminative model for classification and predict the possibility score of each candidate;
 9. Estimate new object position (x_i, y_i, w_i, h_i) by the candidate with the highest score S ;
 10. **if** $\max(S) > 0$ **or** $i \% 10 == 0$ **then**
 11. Extract features from the object location of the successfully tracked frame;
 12. Use adversarial learning to update the discriminative model by Eq.(5);
 13. **end**
 14. **end**

5. Experiments

In this section, we introduce the implementation details of the proposed tracking algorithm. We then compare our tracker with state-of-the-art trackers on two benchmarks for performance evaluation.

In this work, the first three convolution layers from the VGG-M model are utilized as feature extraction network. The network is pre-trained on a large-scale benchmark datasets. During the adversarial learning, both G and D are learned by the SGD scheme. The learning rate for training G and D are set to 10^{-3} and 10^{-4} , respectively. During the tracking, we draw 256 candidate samples around the object location of each frame for classification. The masks are set randomly and the resolution of each mask is the same as that of the input features. We update the generative adversarial network using 10 iterations in every 10 frames or the response score of tracked result is less than a predefined threshold. Our experiments are performed on a workstation by using MatConvNet toolbox [46] with E5 2.4 GHz CPU and Quadro K2200 GPU.

Benchmarks: We conduct the experiments on two standard benchmarks: OTB-2013 [4] and OTB-2015 [5]. Video sequences are defined with bounding box annotations. These datasets cover various challenging aspects in visual tracking task, such as fast motion, background clutter, deformation, occlusion, illumination variations, low resolution, etc. The performance of all the trackers can be well tested by using two benchmarks.

Evaluation metrics: We follow the standard evaluation metrics [5] from two benchmarks. For the OTB-2013 and OTB-2015 benchmarks, we use the one-pass evaluation (OPE) with precision and area-under-the-curve (AUC) success rate criteria. The precision metric measures the rate of frame locations within a certain threshold distance from those of the ground truth. The threshold distance is set to 20 for all the trackers. The success rate criterion measures the overlap ratio between the predicted bounding box and the ground truth bounding box.

6. Results And Discussion

6.1 Quantitative Evaluation

We perform quantitative evaluation on two benchmark datasets. The experimental results of the proposed tracking algorithm are reported as follows.

OTB-2013 Benchmark: We use the OTB-2013 benchmark to confirm that our tracker is on par with the state-of-the-art trackers. The trackers that we compared including the 29 trackers from the OTB-2013 benchmark and other state-of-the-art trackers including KCF [15], MUSTer [16], DSST [18], SRDCF [19], CCOT [21], ECO [22], HCF [25], HDT [26], CFNet [27], FCNT [31], SiameFC [34], SINT [35], MDNet [47], LCT [48], VITAL [49], CREST [50], TCDL [51], Staple [52], MCPF [53], DLS-SVM [54], CNN-SVM [55], GOTURN [56], SRDCFdecon [57], DeepSRDCF [58], SCT [59] and ADNet [60].

We evaluate all the trackers on 50 video sequences using the one-pass evaluation with distance precision and overlap success metrics. Fig.5 shows the tracking results from all compared trackers. We only show

the top ten trackers for presentation clarity. The number listed in the legend indicates the AUC overlap success rate and precision score at 20 pixels. Overall, it clearly illustrates that our tracking method outperforms the state-of-the-art trackers significantly in both evaluation measures. The OTB-2013 dataset has 11 attributes (e.g., background clutter, occlusion, deformation, scale variation, illumination variation, etc.) to describe the different challenges in the tracking. These attributes are useful for analyzing the performance of trackers in different aspects. Fig.6 shows the results of different tracking algorithms on eight main challenging attributes. It demonstrates that our tracker can effectively handle the challenges and achieve leading performance. The proposed method performs favorably against the state-of-the-art trackers when evaluating with eight challenging factors.

OTB-2015 Benchmark: For more detailed analysis, we also compare our tracker with the state-of-the-art trackers on the OTB-2015 benchmark. Fig.7 shows that the proposed tracker performs well. Although the ECO tracker has achieved a good performance, the proposed tracker uses the samples of the first ten frames to train deep network model, so both precision and success rate are leading.

6.2 Qualitative Evaluation

In Fig.8, we qualitatively report the results of other four state-of-the-art trackers (such as, CNN-SVM, CCOT, MDNet, ECO) and the proposed tracker on 12 challenging video sequences. In the most of the video sequences, CNN-SVM is unable to locate the object position due to the limited performance of SVM classifier. MDNet improves CNN-SVM through an end-to-end CNN network formulation, and it performs well on deformation (Trans), low resolution (Skiing) and fast motion (Diving). However, it does not perform well in handling out-of-plane rotation (Ironman) and occlusion (Human4). The correlation filter based trackers such as CCOT and ECO use deep features for visual object tracking, but they fails to exploit more sophisticated deeper architectures. They perform well in handling occlusion (Human4, Box) and deformation (Trans). However, the tracked object drifts when it undergoes heavy occlusions (Bird1). Overall, our tracker captures the appearance variations of the object by fine-tuning the network and the adversarial learning scheme enhances the discriminative ability of the classifier. Therefore, the proposed tracker performs well in estimation both the scale and position of the object on these challenging video sequences. The proposed tracker performs favorably against state-of-the-art.

Moreover, the accurate prediction of spatial attention response map is a key factor in our tracker. Thanks to the utilization of mean squared loss and adversarial loss, the predicted attention response maps are robust for most challenging cases. Even if the attention is not precise, the tracking results will not be largely affected in our experiments.

6.3 Feature comparison

We compare the feature effects of different layers of deep learning network model on the OTB-2015 benchmark, which is shown in Table 1. We can see that the combination of the features extracted from conv3 and conv4 layers achieves the best results, which verifies the rationality of the feature selection strategy of the proposed tracking algorithm. The best results are in bold.

Table 1 Results of different features

ConvLayer	2	3	4	3, 4	2, 3, 4
Precision	0.624	0.705	0.764	0.821	0.795
Success rate	0.502	0.524	0.557	0.587	0.566

6.4 Failure cases

Although the proposed tracking algorithm can achieve a satisfied performance, a few failure cases occur when object suffers from the long-term occlusions. When the tracked object reappears and becomes very small, the proposed tracking method fails to follow the object due to the limited pixels and appearance variations, which can result in poor tracking performance. A feature selection implementation strategy using the feature from conv2 is able to track the object, because the features of conv2 layer have higher resolution than the features from deeper layers. For the Biker sequence, the object suddenly moves violently beyond the search area of the proposed tracking method. Many single object trackers aren't able to cope with this challenge problem in this sequence.

7 Conclusion

In this paper, we propose an effective object tracking method with learning attention. We design a prediction network which is pre-trained off-line and used to predict the object positions of a video sequence. The object positions of the first ten frames are employed to fine-tune prediction network for obtaining rich appearance variations. The positive and negative samples are also augmented. Furthermore, these object locations are captured to mine the domain-specific information through fine-tuning the adversarial generative network. We adaptively use dropout to mine the discriminative features which are originally diminished during the training process. The adaptive dropout is achieved via adversarial learning to find discriminative features according to different inputs. In addition, we present a spatial attention mechanism to improve the tracking performance. Compared with the state-of-the-art, the proposed tracking method achieves outstanding performance in two large public tracking benchmarks. Further research directions include applying the spatial attention into multi-modal applications.

Abbreviations

DL: Deep learning; CF: Correlation filter; CNN: Convolutional neural network; AUC: Area under curve; MOSSE: Minimum output sum of squared error; HOG: Histogram of oriented gradient; OPE: One pass evaluation.

Declarations

Acknowledgments

The authors would like to thank the anonymous reviewers for useful and constructive comments that help improve the quality of this paper.

Authors' contributions

All authors took part in the discussion of the work described in this paper. All authors read and approved the final manuscript.

Funding

This work is supported in part by the National Natural Science Foundation of China (Grant No. 61802058); in part by the International Cooperation and Exchange of the National Natural Science Foundation of China (Grant No. 61911530397); in part by the Equipment Advance Research Foundation Project of China (Grant No. 61403120106); in part by the Startup Foundation for Introducing Talent of Nanjing University of Information Science and Technology (Grant No. 2018r057); in part by the Project funded by the China Postdoctoral Science Foundation (Grant No. 2019M651650); in part by the Open Project Program of the State Key Lab of CAD&CG (Grant No. A1919), Zhejiang University, and the PAPD fund.

Availability of data and materials

We used publicly available dataset in order to illustrate and test our methods. The OTB dataset can be found in <http://cvlab.hanyang.ac.kr/trackerbenchmark/datasets.html>.

Competing interests

The authors declare that they have no competing interests.

Author details

¹School of Computer and Software, Nanjing University of Information Science and Technology, 210044, China

²Engineering Research Center of Digital Forensics, Ministry of Education, Nanjing, 210044, China.

References

1. Zhang, X. Luo, Y. Guo, et al. Multiple Robustness Enhancements for Image Adaptive Steganography in Lossy Channels. *IEEE Transactions on Circuits and Systems for Video Technology*, Published online (Early Access), (2019). DOI:10.1109/TCSVT.2019.2923980
2. Qin, W. Zhang, F. Cao, X. Zhang, and C. Chang, Separable Reversible Data Hiding in Encrypted Images via Adaptive Embedding Strategy with Block Selection. *Signal Processing*, 153: 109–122, (2018).

3. Wang, H. Wang, J. Li, X. Luo, Y. Shi, S. Jha, Detecting double JPEG compressed color images with the same quantization matrix in spherical coordinates. *IEEE Transactions on Circuits and Systems for Video Technology*, DOI: 10.1109/TCSVT.2019.2922309.
4. Wu, C. Luo, Q. Zhang, J. Zhou, H. Yang and Y. Li, Text Detection and Recognition for Natural Scene Images Using Deep Convolutional Neural Networks. *Computers, Materials & Continua*. 61(1): 289-300, (2019).
5. Wu, J. Lim, and M.-H. Yang. Online object tracking: A benchmark, in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, Portland, OR, USA, Jun. (2013), pp. 2411–2418.
6. Wu, J. Lim, and M. Yang, Object tracking benchmark, *IEEE Transactions on Pattern Analysis and Machine Intelligence*. 37(9): 1834–1848, (2015).
7. Yilmaz, O. Javed, M. Shah, Object tracking: A survey. *ACM computing surveys (CSUR)*. 38(4): 13, (2006).
8. Wang, T. Li, X. Luo, Y. Shi, S. Jha, Identifying Computer Generated Images Based on Quaternion Central Moments in Color Quaternion Wavelet Domain. *IEEE Transactions on Circuits and Systems for Video Technology*. 29(9): 2775-2785, (2018).
9. W. M. Smeulders, D. M. Chu, R. Cucchiara, et al. Visual tracking: An experimental survey. *IEEE transactions on pattern analysis and machine intelligence*. 36(7): 1442-1468, (2014).
10. Ababsa, Robust extended kalman filtering for camera pose tracking using 2D to 3D lines correspondences. *IEEE/ASME International Conference on Advanced Intelligent Mechatronics*, (2009), 1834-1838.
11. S. Arulampalam, S. Maskell, N. Gordon, et al., A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *IEEE Transactions on Signal Processing*. 50(2): 174-188, (2002).
12. S. Bolme, J. R. Beveridge, B. A. Draper, et al. Visual object tracking using adaptive correlation filters. *IEEE Conference on Computer Vision and Pattern Recognition*, (2010): 2544-2550.
13. Cheng, Y. Zhang, L. Zhou, Y. Zheng. Visual tracking via auto-encoder pair correlation filter. *IEEE Transactions on Industrial Electronics*. 67(4): 3288-3297, (2020).
14. Cheng, YATA: Yet another Proposal for Traffic Analysis and Anomaly Detection, *Computers, Materials & Continua*. 60(3): 1171-1187, (2019).
15. F. Henriques, R. Caseiro, P. Martins, et al. High-speed tracking with kernelized correlation filters. *IEEE transactions on pattern analysis and machine intelligence*, (2015), 37(3): 583-596.
16. Hong, Z. Chen, C. Wang, et al. Multi-store tracker (MUSTer): A cognitive psychology inspired approach to object tracking. *IEEE Conference on Computer Vision and Pattern Recognition*. (2015): 749-758.
17. F. Henriques, R. Caseiro, P. Martins, et al. Exploiting the circulant structure of tracking-by-detection with kernels. *European conference on computer vision*. (2012): 702-715.
18. Danelljan, G. Häger, F. Khan, et al. Accurate scale estimation for robust visual tracking. *British Machine Vision Conference*, Nottingham. (2014): 1-5.

19. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg, Learning spatially regularized correlation filters for visual tracking. IEEE International Conference on Computer Vision, (2015), 4310–4318.
20. Lukezic, T. Vojir, L. C. Zajc, et al. Discriminative Correlation Filter with Channel and Spatial Reliability. IEEE Conference on Computer Vision and Pattern Recognition. (2017).
21. Danelljan, A. Robinson, F. Khan, and M. Felsberg. Beyond correlation filters: Learning continuous convolution operators for visual tracking. European Conference on Computer Vision, (2016): 472-488.
22. [Danelljan](#), G. Bhat, F. S. Khan, M. Felsberg, ECO: Efficient Convolution Operators for Tracking, IEEE Conference on Computer Vision and Pattern Recognition, (2017): 6931-6939.
23. Bhat, J. Johnander, M. Danelljan, et al. Unveiling the power of deep tracking. European Conference on Computer Vision. (2018): 483-498.
24. Li, C. Tian, W. Zuo, et al. Learning spatial-temporal regularized correlation filters for visual tracking. IEEE Conference on Computer Vision and Pattern Recognition. (2018): 4904-4913.
25. Ma, J. B. Huang, X. Yang, et al. Hierarchical convolutional features for visual tracking. IEEE international conference on computer vision. (2015): 3074-3082.
26. Qi, S. Zhang, L. Qin, et al. Hedged deep tracking. IEEE conference on computer vision and pattern recognition. (2016): 4303-4311.
27. Valmadre, L. Bertinetto, J. Henriques, et al. End-to-end representation learning for correlation filter based tracking. IEEE Conference on Computer Vision and Pattern Recognition. (2017): 2805-2813.
28. Wang, J. Gao, J. Xing, M. Zhang, and W. Hu. Dcfnet: Discriminant correlation filters network for visual tracking. arXiv preprint arXiv:1704.04057, 2017.
29. Wang and D.-Y. Yeung. Learning a deep compact image representation for visual tracking. The Annual Conference on Neural Information Processing Systems, (2013): 809–817.
30. Wang, S. Li, A. Gupta, et al., Transferring rich feature hierarchies for robust visual tracking. arXiv preprint arXiv:1501.04587, (2015).
31. Wang, W. Ouyang, X. Wang, et al., Visual tracking with fully convolutional networks, IEEE International Conference on Computer Vision, (2015): 3119-3127.
32. Deng, W. Dong, R. Socher, et al. ImageNet: A large-scale hierarchical image database. IEEE conference on computer vision and pattern recognition. (2009): 248-255.
33. Zhang, Q. Liu, Y. Wu, et al., Robust visual tracking via convolutional networks without training. IEEE Transactions on Image Processing. 25(4): 1779-1792, (2016).
34. Bertinetto, J. Valmadre, J. F. Henriques, et al. Fully-convolutional siamese networks for object tracking. European conference on computer vision. (2016): 850-865.
35. Tao, E. Gavves, A. W. M. Smeulders, Siamese instance search for tracking. IEEE conference on computer vision and pattern recognition. (2016): 1420-1429.
36. Wang, A. Shrivastava, A. Gupta, A-fast-rcnn: Hard positive generation via adversary for object detection. IEEE Conference on Computer Vision and Pattern Recognition. (2017): 2606-2615.

37. Jiang, J. Chen, Y. Jiang, Y. Xu, Y. Wang, L. Tan and G. Liang, A New Time-Aware Collaborative Filtering Intelligent Recommendation System. *Computers, Materials & Continua*. 61(2): 849-859, (2019).
38. Goodfellow, J. Pouget-Abadie, M. Mirza, et al. Generative adversarial nets. *Advances in neural information processing systems*. (2014): 2672-2680.
39. Li, C. Ma, B. Wu, et al. Target-Aware Deep Tracking. *IEEE Conference on Computer Vision and Pattern Recognition*, (2019).
40. Wang, Y. Song, C. Ma, W. Zhou, et al. Unsupervised Deep Tracking. *IEEE Conference on Computer Vision and Pattern Recognition*, (2019).
41. A. Olshausen, C. H. Anderson, and D. C. Van Essen. A neurobiological model of visual attention and invariant pattern recognition based on dynamic routing of information. *Journal of Neuroscience*, 13(11):4700–4719, (1993).
42. Fan, Y. Wu, and S. Dai. Discriminative spatial attention for robust tracking. In *European Conference on Computer Vision*. 480–493, (2010).
43. Lukezic, T. Vojir, L. Cehovin Zajc, J. Matas, and M. Kristan. Discriminative correlation filter with channel and spatial reliability. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
44. Choi, H. J. Chang, S. Yun, T. Fischer, and Y. Demiris. Attentional correlation filter network for adaptive visual tracking. *IEEE Conference on Computer Vision and Pattern Recognition*. pp: 4807–4816, (2017).
45. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, (2014).
46. Vedaldi, K. Lenc, Matconvnet: Convolutional neural networks for MATLAB. The 23rd ACM international conference on Multimedia. ACM, (2015): 689-692.
47. Nam, B. Han, Learning multi-domain convolutional neural networks for visual tracking. *IEEE Conference on Computer Vision and Pattern Recognition*. (2016): 4293-4302.
48. Ma, X. Yang, C. Zhang, et al. Long-term correlation tracking. *IEEE conference on computer vision and pattern recognition*. (2015): 5388-5396.
49. Song, C. Ma, X. Wu, et al. VITAL: Visual tracking via adversarial learning. *IEEE Conference on Computer Vision and Pattern Recognition*. (2018): 8990-8999.
50. Song, C. Ma, L. Gong, et al. CREST: Convolutional residual learning for visual tracking. *IEEE International Conference on Computer Vision*. (2017): 2555-2564.
51. Cheng, Y. Zhang, J. Cui, et al. Object Tracking via Temporal Consistency Dictionary Learning. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 47(4): 628-638, (2017).
52. Bertinetto, J. Valmadre, S. Golodetz, et al. Staple: Complementary learners for real-time tracking. *IEEE conference on computer vision and pattern recognition*. (2016): 1401-1409.
53. Zhang, C. Xu, M. H Yang. Multi-task correlation particle filter for robust object tracking. *IEEE Conference on Computer Vision and Pattern Recognition*. (2017): 4335-4343.

54. Ning, J. Yang, S. Jiang, et al. Object tracking via dual linear structured SVM and explicit feature map. IEEE conference on computer vision and pattern recognition. (2016): 4266-4274.
55. Hong, T. You, S. Kwak, et al. Online tracking by learning discriminative saliency map with convolutional neural network. International conference on machine learning. (2015): 597-606.
56. Held, S. Thrun, S. Savarese. Learning to track at 100 fps with deep regression networks. European Conference on Computer Vision. (2016): 749-765.
57. Danelljan, G. Hager, K. F. Shabaz, et al. Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking. IEEE Conference on Computer Vision and Pattern Recognition. (2016): 1430-1438.
58. Danelljan, G. Hager, K. F. Shabaz, et al. Convolutional features for correlation filter based visual tracking. IEEE International Conference on Computer Vision Workshops. (2015): 58-66.
59. Choi, J. Chang, J. Jeong, et al. Visual tracking using attention-modulated disintegration and integration. IEEE Conference on Computer Vision and Pattern Recognition. (2016): 4321-4330.
60. Yun, J. Choi, Y. Yoo, et al. Action-decision networks for visual tracking with deep reinforcement learning. IEEE Conference on Computer Vision and Pattern Recognition. (2017): 2711-2720.

Figures

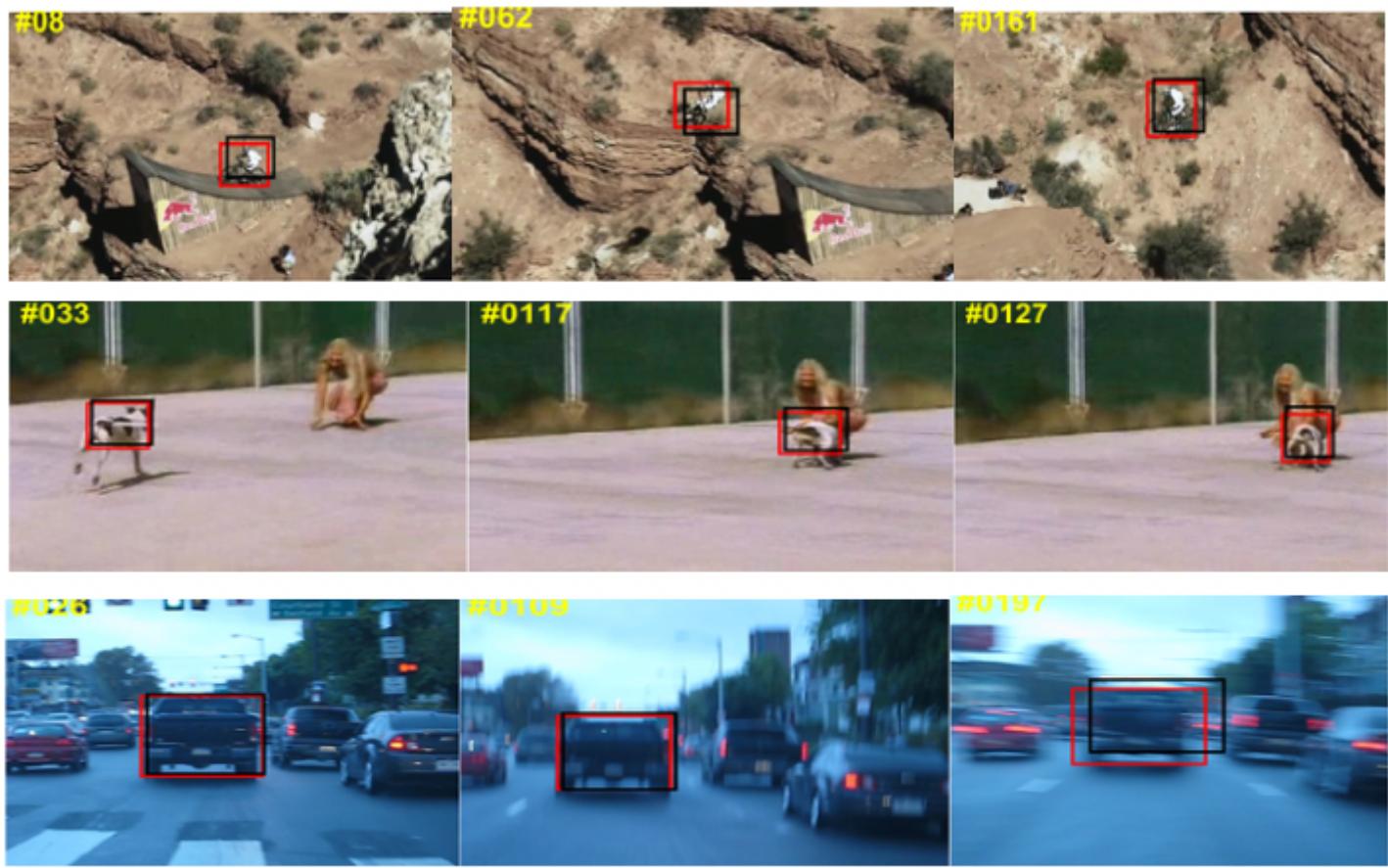


Figure 1

Object tracking technique for intelligent surveillance analysis, tracking results of the proposed tracker (Red) and the correlated VITAL tracker [34] (Black) on the OTB2015 benchmark.

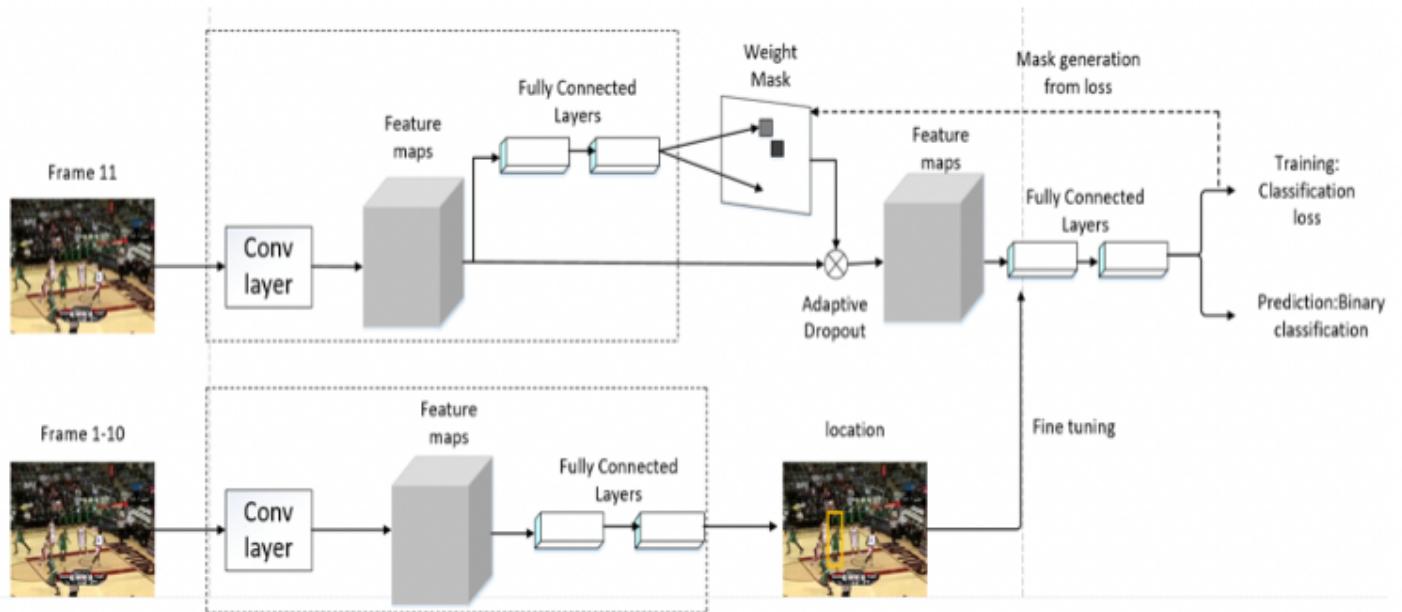


Figure 2

The architecture of the proposed object tracking algorithm

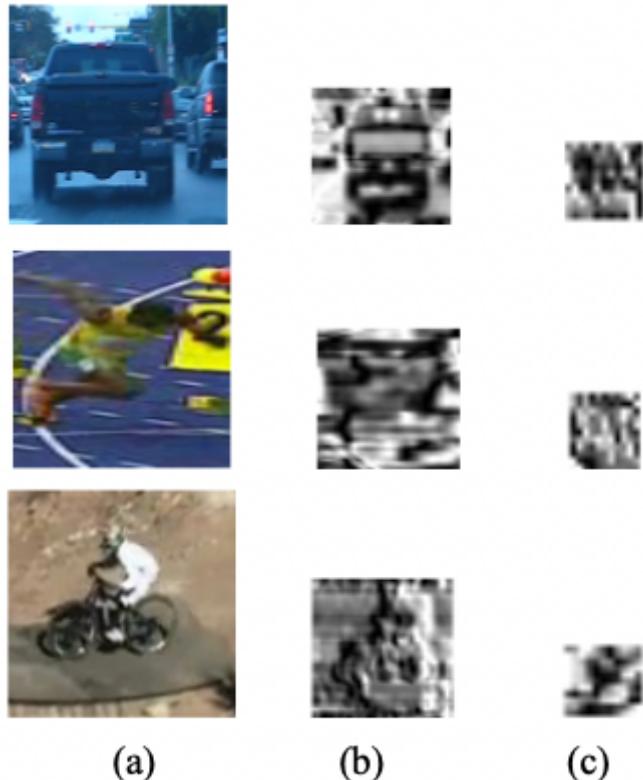


Figure 3

Foreground response maps predict using different VGG feature maps. (a) Input image. (b) Using Conv5-1 feature only. (c) Using the concatenation of Conv5-2 and Conv5-3 feature with the proposed adversarial learning with attention.

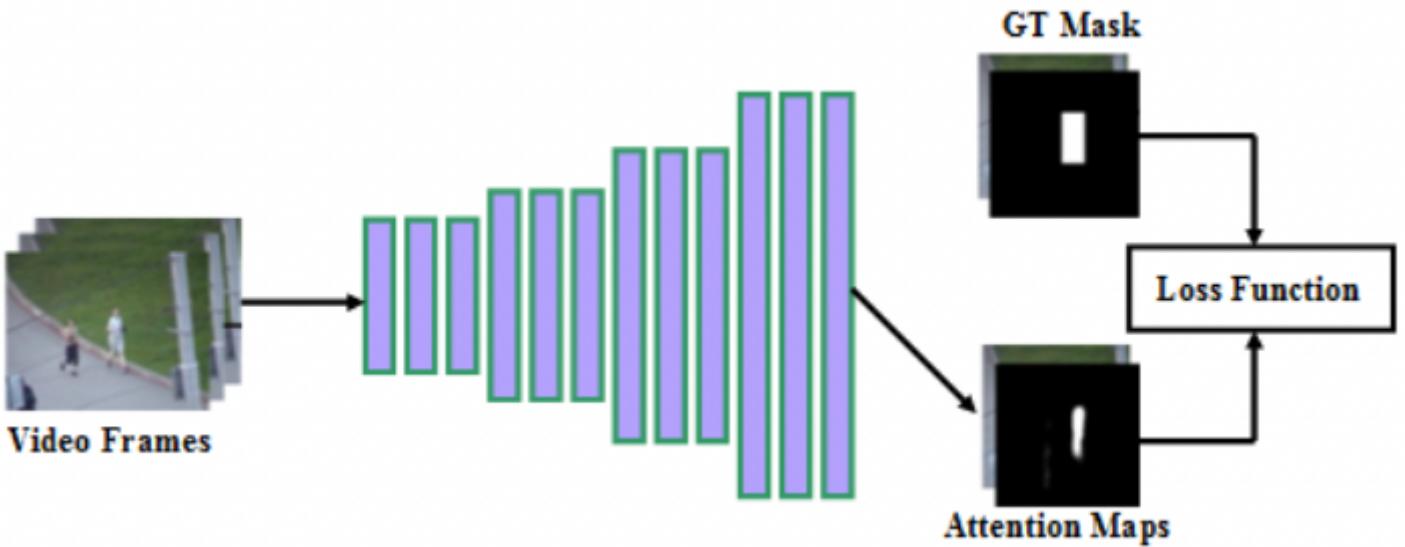


Figure 4

Spatial attention response map

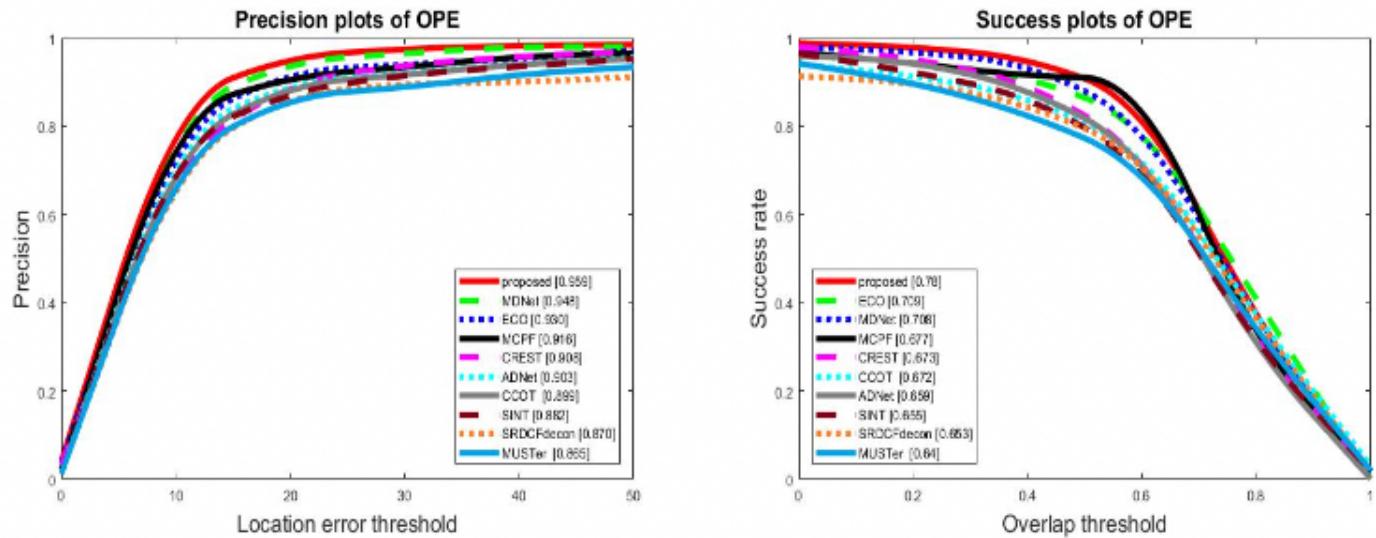


Figure 5

Precision and success plots on the OTB-2013 dataset using the one-pass evaluation. The number in the legend indicates the average precision scores at 20 pixels and the AUC scores.

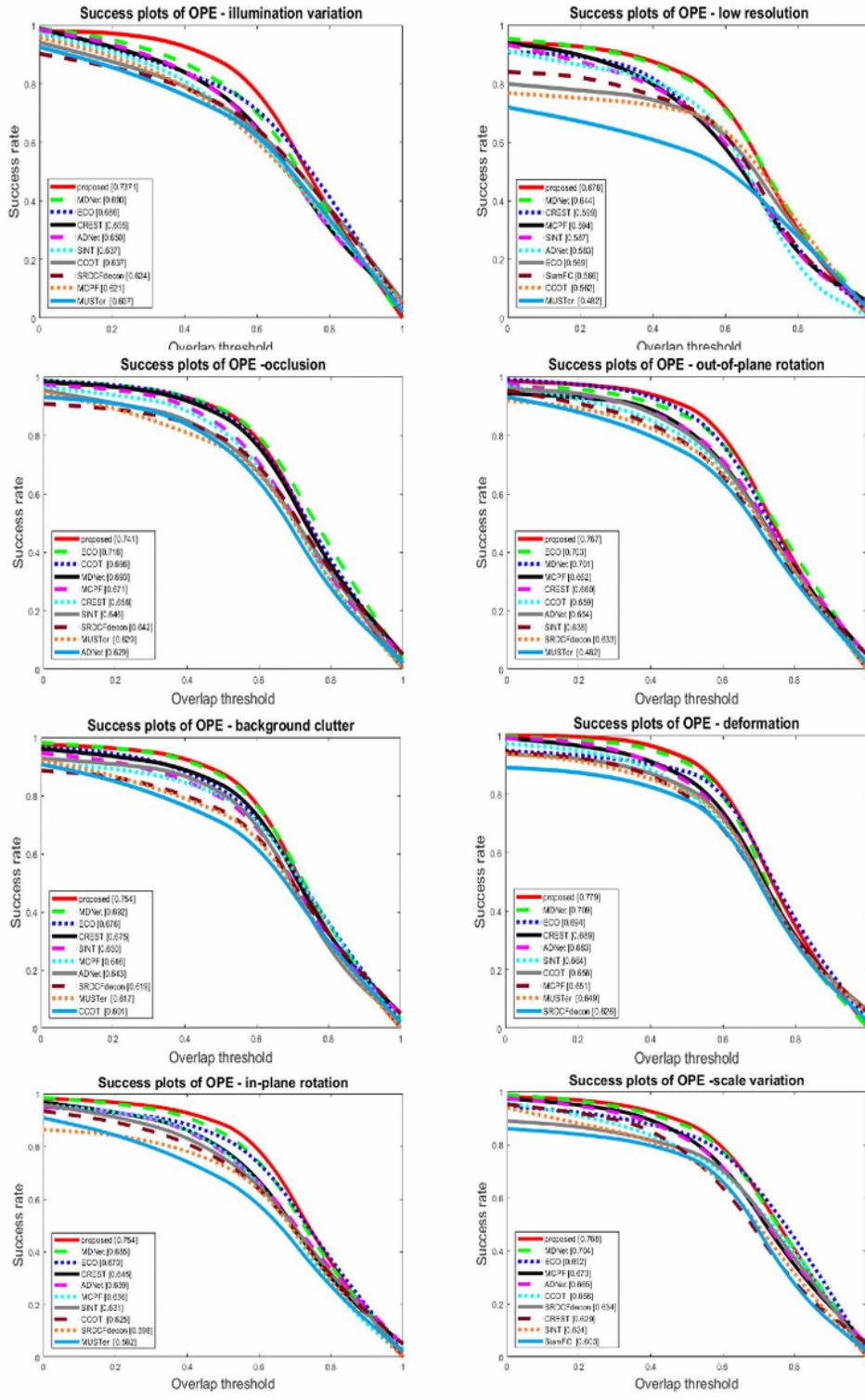


Figure 6

The success rate plots of eight challenge attributes: illumination variation, deformation, in-plane rotation, out-of-plane rotation, background clutter, occlusion, scale variation and low resolution. The legend contains the AUC score for each attribute.

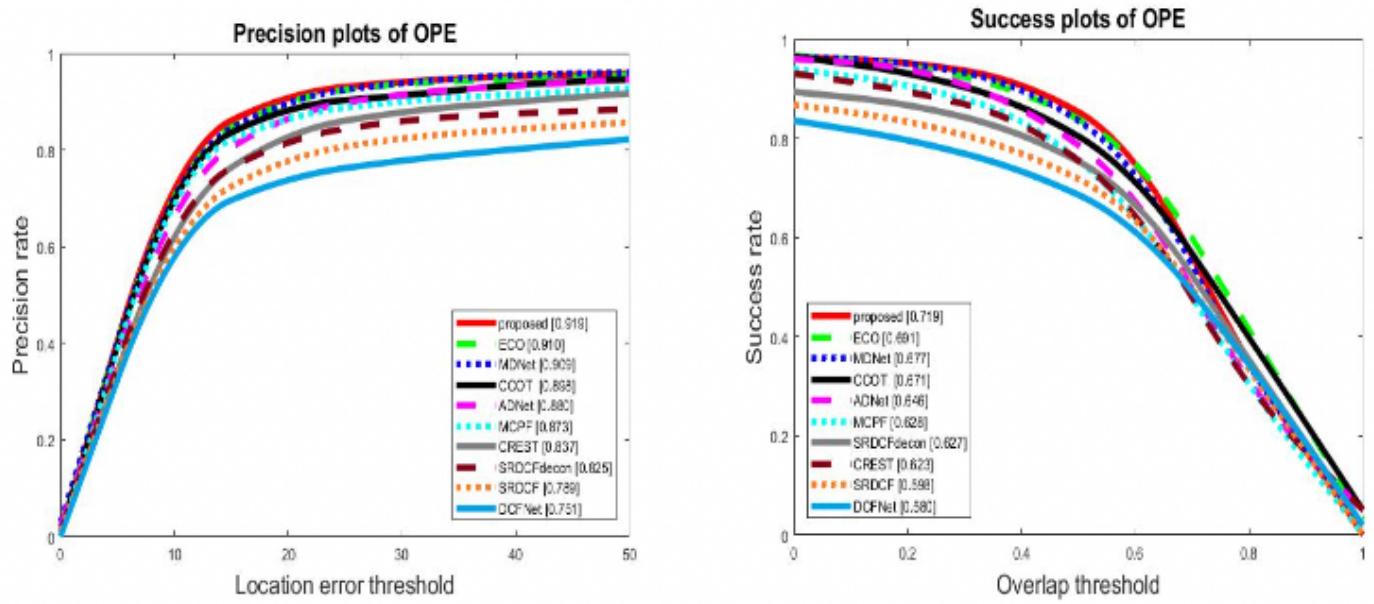


Figure 7

Precision and success rate plots on the OTB-2015 benchmark by using the one-pass evaluation. The number in the legend indicates the average distance precision scores at 20 pixels and AUC success scores.



Figure 8

Qualitative evaluation of the proposed tracker, CNN-SVM, CCOT, MDNet, ECO on 12 challenging sequences.

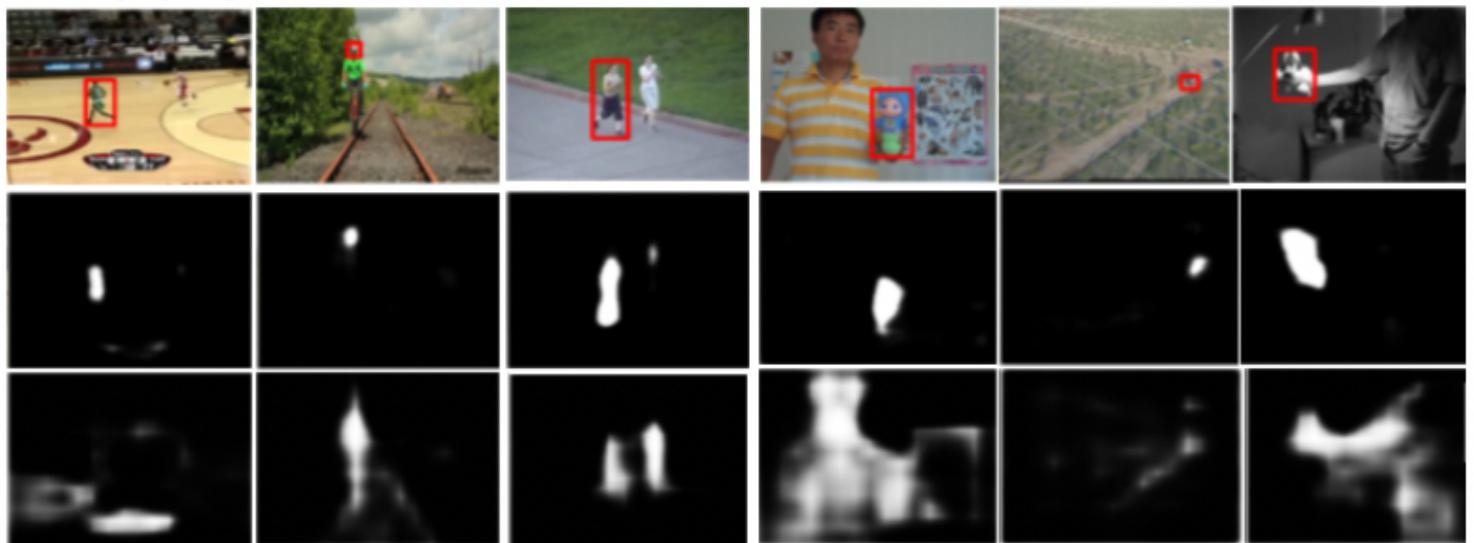


Figure 9

The maps generated by the proposed spatial attention scheme (middle row) and saliency detection algorithm Deep Saliency (bottom row).