

Leveraging Tacit Information Embedded in CNN Layers for Visual Tracking

Kourosh Meshgi¹, Maryam Sadat Mirzaei¹, and Shigeyuki Oba²

¹ RIKEN Center for Advanced Intelligence Project (AIP), Tokyo, Japan

² Graduate School of Informatics, Kyoto University, Kyoto, Japan

Abstract. Different layers in CNNs provide not only different levels of abstraction for describing the objects in the input but also encode various implicit information about them. The activation patterns of different features contain valuable information about the stream of incoming images: spatial relations, temporal patterns, and co-occurrence of spatial and spatiotemporal (ST) features. The studies in visual tracking literature, so far, utilized only one of the CNN layers, a pre-fixed combination of them, or an ensemble of trackers built upon individual layers. In this study, we employ an adaptive combination of several CNN layers in a single DCF tracker to address variations of the target appearances and propose the use of style statistics on both spatial and temporal properties of the target, directly extracted from CNN layers for visual tracking. Experiments demonstrate that using the additional implicit data of CNNs significantly improves the performance of the tracker. Results demonstrate the effectiveness of using style similarity and activation consistency regularization in improving its localization and scale accuracy.

1 Introduction

Discovering new architectures for deep learning and analyzing their properties, have resulted in a rapid expansion in computer vision, along with other domains. Among these architectures, convolutional neural networks (CNNs) have played a critical role to capture the statistics and semantics of natural images. CNNs are widely used in different computer vision applications since they are able to effectively learn complicated mappings while using minimal domain knowledge [1]. Deep learning has been introduced to visual tracking in different forms, mostly to provide features for established trackers based on correlation filters [2], particle filters [3,4] and detector-based trackers [5]. Although deep features extracted by fully-connected layers of CNN are shown to be adequately generic to be used for a variety of computer vision tasks, including tracking [6]. Further studies revealed that convolutional layers are even more discriminative, semantically meaningful and capable of learning structural information [7]. However, the direct use of CNNs to perform tracking is complicated because of the need to re-train the classifier with the stream of target appearances during tracking, the diffusion of background information into template [8], and “catastrophic forgetting” of previous appearance in the face of extreme deformations and occlusions [9].

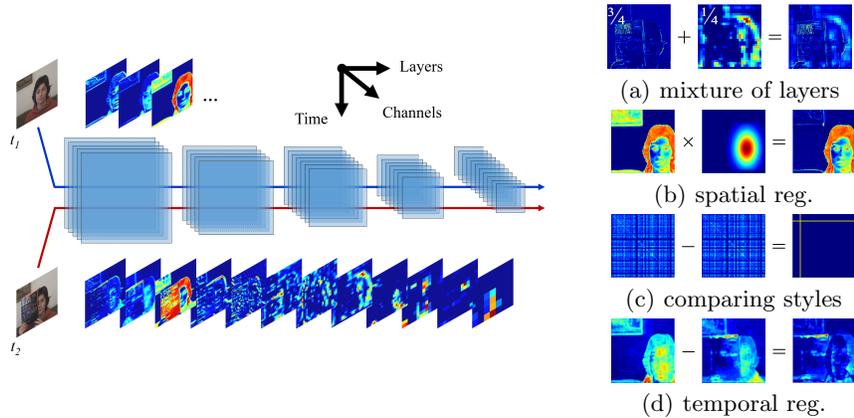


Fig. 1. When presented with a sequence of images, CNN neurons are activated in a particular way, involving information about the spatial, semantic, style and transformations of the target. **(a)** Combining information from different layers balances the amount of spatial and semantic information, **(b)** spatial weighting of the response would discard most of the background distraction, **(c)** changes of the co-activations of neurons (measured by Gram matrices) for different subsequent images indicate style changes of the target (the plot is exaggeratedly enhanced for visibility), and **(d)** changes in the activations of the neurons themselves signals the appearance transformations and pose changes (in shallower layers) or alteration of semantic contents.

Early studies in the use of deep learning in tracking utilized features from autoencoders [10, 11] and fully-connected layers of pre-trained (CNN-based) object detector [12], but later the layers were used to serve as features balancing the abstraction level needed to localize the target [13], provide ST relationship between the target and its background [14], combine spatial and temporal features [15, 16], and generate probability maps for the tracker [17]. Recently trackers employ other deep learning approaches such as R-CNN for tracking-by-segmentation [18], Siamese Networks for template similarity measuring [19–22], GANs to augment positive samples [23, 24], and CNN-embeddings for self-supervised image coloring and tracking [25]. However, the tacit information in pre-trained CNNs including information between layers, within layers, and activation patterns across the time axis are underutilized (Fig. 1).

Different layers of CNNs provide different levels of abstraction [26], and it is known that using multiple layers of CNNs can benefit other tasks such as image classification [27]. Such information was used as a coarse-to-fine sifting mechanism in [13], as a tree-like pre-defined combination of layers with fixed weights and selective updating [9], as the features for different trackers in an ensemble tracker [28], or in a summation over all layers as uniform features [29]. However, direct adaptive fusion of these features in a single tracker that can address different object appearances is still missing in the literature.

CNN stores information not only between layers but within layers in different channels, each representing a different filter. These filters may respond to dif-

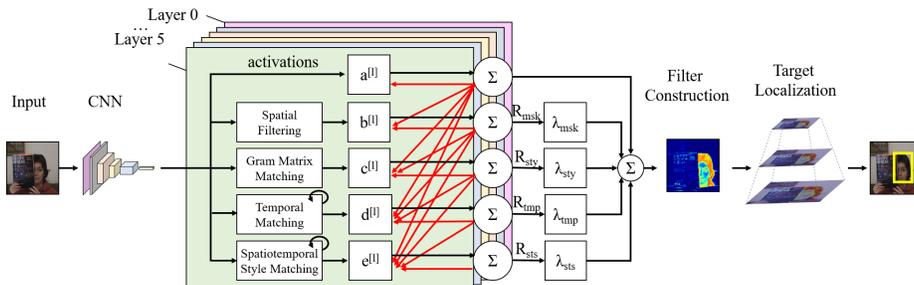


Fig. 2. Schematic of the proposed tracker. Given the input image, the activations of different layers of CNN are processed independently and its spatial irregularities, style mismatches, temporal inconsistencies and ST pattern changes compared to the template are calculated. These errors are then adaptively fused with that of other layers to form the regularization terms R_x . The final filter is then constructed and used in the baseline DCF tracker that uses multi-scale matching to find the position and scale of the target in the input image. The weights of different error terms are then updated in reverse proportion of their contribution in total error of each level.

ferent visual stimuli. The shallower layers have a Gabor-like filter response [30] whereas in the deeper layers, they respond to angles, color patterns, simple shapes, and gradually highly complex stimuli like faces [26]. The co-occurring patterns within a layer, activate two or more different channels of the layer. Such co-incidental activations are often called style in the context of neural style transfer (NST), and different approaches are proposed to measure the similarity between two styles [31,32]. The loss functions for NST problem can serve as the similarity index for two objects (e.g., in the context of image retrieval [33]).

Most of the current CNN-based techniques use architectures with 2D convolutions to achieve different invariances to the variations of the images. Meanwhile, the invariance to transformations in time is of paramount importance for video analysis [34]. Modeling temporal information in CNNs has been tackled by applying CNNs on optical flow images [35], reformulating R-CNNs to exploit the temporal context of the images [36] or by the use of separate information pathways for spatial and temporal pathways [37,38]. Motion-based CNNs typically outperform CNN representations learned from images for tasks dealing with dynamic target, e.g. action recognition [34]. In these approaches, a CNN is applied on 3-channel optical flow image [39], and different layers of such network provide different variances toward speed and the direction of the target’s motion [40]. In visual tracking, deep motion features provide promising results [41]. However, this requires the tracker to fuse the information from two different CNN networks (temporal+spatial) [38], and their inconsistency hinders a meaningful layer-wise fusion and only the last layers of temporal CNN are used for tracking [41].

Contributions: We propose an adaptive fusion of different CNN layers in a tracker to combine high spatial resolution of earlier layers for precise localization, and semantic information of deeper layers to handle large appearance changes

and alleviate model drift. We utilize the tacit information between each layer’s channels at several timepoints, i.e., the style, to guide the tracker. To our best knowledge, this is the first use of within-layer information of CNNs especially spatial and ST co-activations for tracking. We also introduced temporal constraint helps to better preserve target’s temporal consistency, suppress jitter, and promote scale adaptation.

- (i) We propose an intuitive adaptive weight adjustment to tune the effect of components (spatial, background avoiding, co-incident, temporal, and ST) both within and between CNN layers. Temporal and style regs are typically complementary: target changes are punished by a reduction in activation and style similarity, big changes are punished by spatial reg and style keep track of the target. Employing multiple layers not only gives different realization of details-semantic trade-off [13], but also provides richer statistics compared to one layer. We incorporate different regularization terms on CNN layers using only one feed-forward pass on an arbitrary pre-trained CNN, with no change to its architecture or CNN block design (e.g. ST block as in [42]) and no additional computation (e.g., compute optical flow).
- (ii) We introduce a Gram-based style matching in our tracker to capture style alterations of the target. The style matching exploits the co-activation patterns of the layers and improves the localization accuracy of the target and provides complementary information to the baseline which relies on spatial matching.
- (iii) We introduce the temporal coherence regularization to the tracker by monitoring activations of each layer through the course of tracking, to enhance tracker’s movement speed and direction invariance, adaptation to different degrees of changes in the target appearance, stability, and scale tuning.
- (iv) Our system is tested on various public datasets (OTB50 & 100, LaSOT, VOT2015 & 2018, and UAV123), and considering the simplicity of our baseline (SRDCF [8]) we obtained results on par with many sophisticated trackers. The results shed light on the hidden information within CNN layers, that is the goal of this paper. The layer-fusion, style-matching, and temporal regularization component of the proposed tracker is further shown to advance the baseline significantly and outperformed state-of-the-art trackers.

It should be noted that our proposed method differs with HDT [28], CCOT [43] and ECO [44] that also integrates the multi-layer deep features by considering them in the continuous domain. Here, we employed a pre-trained CNN for object detection task, simplified the need to deal with different dimensions of convolutional layers by maintaining the consistency between layers while isolating them to make an ensemble of features (rather than an ensemble of trackers in [28]). Additionally, we proposed style and temporal regularization that can be incorporated into Conjugate Gradient optimization of C-COT [29], ADMM optimization of BACF [45] and Gauss-Newton optimization of ECO [44].

2 Method

We propose a tracker that adaptively fuse different convolutional layers of a pre-trained CNN into a DCF tracker. The adaptation is inspired by weight tuning of the different tracker components as well as the scale-aware update of [46]. As Figure 2 illustrates, we incorporated spatial regularization (a multi-layer extension of [2]), co-incidental regularization (which matches style between the candidate patch and historical templates), temporal regularization (that ensures a smooth alteration of the activations in normal condition), and ST regularization that captures the change patterns of spatial features.

2.1 Discriminative Correlation Filter Tracker

The DCF framework utilizes the properties of circular correlation to efficiently train and apply a classifier in a sliding window fashion [47]. The resulting classifier is a correlation filter applied to the image/feature channels, similar to conv layers in a CNN. Different from CNNs, the DCF is trained by solving a linear least-squares problem using Fast Fourier Transform (FFT) effectively.

A set of example patches x_τ are sampled at each frame $\tau = 1, \dots, t$ to train the discriminative correlation filter f_t , where t denotes the current frame. The patches are all of the same size (conveniently, the input size of the CNN) centered at the estimated target location in each frame. We define feature x_τ^k as the output of channel k at a convolutional layer in the CNN. With this notion, the tracking problem is reduced to learn a filter f_t^k for each channel k , that minimizes the L^2 -error between the responses S_{f_t} on samples x_τ and the desired filter form y_k :

$$\epsilon = \sum_{\tau=1}^t \alpha_\tau \|S(f_t, x_\tau) - y_\tau\|^2 + \lambda \|f_t\|^2 \quad (1)$$

where $S(f_t, x_\tau) = f_t \star x_\tau$ in which the \star denotes circular correlation generalized to multichannel signals by computing inner products. The desired correlation output y_τ is set to a Gaussian function with the peak placed at the target center location [48]. A weight parameter λ controls the impact of the filter size regularization term, while the weights α_τ determine the impact of each sample.

To find an approximate solution of eq(1), we use the online update rule of [46]. At frame t , the numerator g_t and denominator \hat{h}_t of the discrete Fourier transformed (DFT) filter \hat{f}_t are updated as,

$$\hat{g}_t^k = (1 - \gamma)\hat{g}_{t-1}^k + \gamma \overline{\hat{y}_t} \cdot \hat{x}_t^k \quad (2a)$$

$$\hat{h}_t^k = (1 - \gamma)\hat{h}_{t-1}^k + \gamma \left(\sum_{k'=1}^{n_C} \overline{\hat{x}_t^{k'}} \cdot \hat{x}_t^{k'} + \lambda \right) \quad (2b)$$

in which the ‘hat’ denotes the 2D DFT, the ‘bar’ denotes complex conjugation, ‘ \cdot ’ denotes pointwise multiplication, $\gamma \in [0, 1]$ is learning rate and n_C is the number of channels. Next, the filter is constructed by a point-wise division $\hat{f}_t^k = \hat{h}_t^k / \hat{g}_t^k$.

To locate the target at frame t , a sample patch s_t is extracted at the previous location. The filter is applied by computing the correlation scores in the Fourier domain $\mathcal{F}^{-1}\left\{\sum_{k'=1}^{n_C} \overline{\hat{f}_{t-1}^{k'}} \cdot \hat{s}_t^{k'}\right\}$, in which \mathcal{F}^{-1} denotes the inverse DFT.

2.2 Incorporating Information of CNN Layers

Here, we extend the DCF tracker formulation to accept a linear combination of multiple CNN layers $l \in \mathcal{L}$ with dimensions $n_W^{[l]} \times n_H^{[l]} \times n_C^{[l]}$. We embed spatial focus, style consistency, temporal coherency, and ST style preserving terms as regularizations over the minimization problem.

$$\epsilon = \sum_{\tau=1}^t \alpha_{\tau} \left(\sum_{l \in \mathcal{L}} a_t^{[l]} \|S(f_t^{[l]}, x_{\tau}) - y_{\tau}\|^2 + \sum_{x \in \{\text{msk}, \text{sty}, \text{tmp}, \text{sts}\}} \lambda_x R_x(f_t, x_{\tau}) \right) \quad (3)$$

where the desired filter form for all layers $l \in \mathcal{L}$, $\mathcal{A}_t = \{a_t^{[l]}\}$ is the activation importance of the layers l , and $\Lambda = \{\lambda_{\text{msk}}, \lambda_{\text{sty}}, \lambda_{\text{tmp}}, \lambda_{\text{sts}}\}$ are the regularization weights for tracker components.

2.3 Regularizing the Filter

We embed five different regularizations to push the resulting filter toward ideal form given the features of the tracker. To localize the effect of features a mask reg R_{msk} is used, to penalize the style mismatches between target and the template, co-incidental reg R_{sty} is employed, to push temporal consistency of the target and smoothness of tracking, the temporal reg R_{tmp} is proposed, and to punish abrupt ST changes, the ST style reg R_{sts} is introduced.

Mask Component To address the boundary problems induced by the periodic assumption [8] and minimizing the effect of background [49] we use mask regularization to penalize filter coefficients located further from object’s center:

$$R_{\text{msk}}(f_t, x_{\tau}) = \sum_{l \in \mathcal{L}} b_t^{[l]} \sum_{k=1}^{n_C^{[l]}} \|\mathbf{w} \cdot f_t^{k, [l]}\|^2 \quad (4)$$

in which $\mathbf{w} : \{1, \dots, n_W^{[l]}\} \times \{1, \dots, n_H^{[l]}\} \rightarrow [0, 1]$ is the spatial penalty function, and $\mathcal{B}_t = \{b_t^{[l]}\}$ is the spatial importances of the layers. We use Tikhonov reg. similar to [8] as \mathbf{w} smoothly increase with distance from the target center.

Co-incidental Component CNN activations entails spatial information of the target but may suffer from extreme target deformations, missing information in the observation (e.g. due to partial occlusions) and complex transformations (e.g. out-of-plane rotations). On the other hand, Gram-based description of a

CNN layer encodes the second order statistics of the set of CNN filter responses and tosses spatial arrangements [50]. Although, this property may lead to some unsatisfying results in NST domain, it is desired in the context of visual tracking as a complement for raw activations.

$$R_{\text{sty}}(f_t, x_\tau) = c_{\text{norm}} \sum_{l \in \mathcal{L}} c_t^{[l]} \|G^{[l]}(f_t) - G^{[l]}(f_\tau)\|_F^2 \quad (5)$$

where $\mathcal{C}_t = \{c_t^{[l]}\}$ is the layers' co-incidental importance, $c_{\text{norm}} = \sum_{l \in \mathcal{L}} c_t^{[l]} = (2n_H^{[l]}n_W^{[l]}n_C^{[l]})^{-2}$ as normalizing constant and $\|\cdot\|_F^2$ is the Frobenios norm operator and $G^{[l]}(\cdot)$ is the cross-covariance of activations in layer l , the Gram matrix:

$$G_{kk'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_W^{[l]}} q_{ijk}^{[l]} q_{ijk'}^{[l]}, \quad (6)$$

where $q_{ijk}^{[l]}$ is the activation of neuron in (i, j) of channel k in layer l . The Gram matrix captures the correlation of activations across different channels of layer l , indicating the frequency of co-activation of features of a layer. It is a second-degree statistics of the network activations, that captures co-occurrences between neurons, known as ‘‘style’’ in spatial domain. While network activations reconstruct the image based on the features in each layer, style information encodes input patterns, which in lowest form is considered as the texture, known to be useful for tracking [51]. The patterns of deeper layers contain higher levels of co-occurrences, such as the relation of the body-parts and shape-color.

Temporal Component This term is devised to ensure the smoothness of activation alterations of CNNs, which means to see the same features in the same positions of the input images, and punish big alterations in the target appearance, which may happen due to misplaced sampling window. Another benefit of this term is to prefer bounding boxes which include all of the features and therefore improve the scale adaptation ($\mathcal{D}_t = \{d_t^{[l]}\}$):

$$R_{\text{tmp}}(f_t, x_\tau) = \sum_{l \in \mathcal{L}} d_t^{[l]} \|S(f_t^{[l]}, x_\tau) - S(f_t^{[l]}, x_{\tau-1})\|^2 \quad (7)$$

Spatiotemporal Style Component To capture the style of target’s ST changes, the style of the spatial patterns in consecutive frames is compared. It promotes the motion smoothness of the spatial features, and monitors the style in which each features evolve throughout the video ($\mathcal{E}_t = \{e_t^{[l]}\}$):

$$R_{\text{sts}}(f_t, x_\tau) = \sum_{l \in \mathcal{L}} e_t^{[l]} \|G^{[l]}(f_\tau) - G^{[l]}(f_{\tau-1})\|_F^2 \quad (8)$$

Model Update Extending filter update equations (eq(2)) to handle multiple layers is not trivial. It is the importance weights $\mathcal{A}_t, \dots, \mathcal{E}_t$, that provides high degree of flexibility for the visual tracker to tailor its use of different layers (i.e., their activations, styles, spatial, temporal, and spatitemporal coherences) to the target. As such, we use a simple yet effective way of adjusting the weights, considering the effect of the layer they represent among all the layers. Here, we denote $\tilde{z}_t^{[l]}$ as the portion of error in t caused by layer l among all layers \mathcal{L} :

$$z_{t+1}^{[l]} = 1 - \frac{\eta + \tilde{z}_t^{[l]}}{\eta + \sum_{l' \in \mathcal{L}} \tilde{z}_t^{[l']}} \quad , \text{ where } z_t \in \{a_t, b_t, c_t, d_t, e_t\} \quad (9)$$

in which η is a small constant and error terms $\tilde{z}_t^{[l]}$ are defined as follows:

$$\tilde{a}_t^{[l]} = \|S(f_t^{[l]}, x_t) - y_t\|^2 \quad (10a)$$

$$\tilde{b}_t^{[l]} = a_t^{[l]} \tilde{a}_t^{[l]} + \sum_{k=1}^{n_c^{[l]}} \|w \cdot f_t^{k,[l]}\|^2 \quad (10b)$$

$$\tilde{c}_t^{[l]} = a_t^{[l]} \tilde{a}_t^{[l]} + b_t^{[l]} \tilde{b}_t^{[l]} + \|G^{[l]}(f_t) - G^{[l]}(f_{t-1})\|_F^2 \quad (10c)$$

$$\tilde{d}_t^{[l]} = a_t^{[l]} \tilde{a}_t^{[l]} + b_t^{[l]} \tilde{b}_t^{[l]} + c_t^{[l]} \tilde{c}_t^{[l]} + \|S(f_t^{[l]}, x_t) - S(f_{t-1}^{[l]}, x_{t-1})\|^2 \quad (10d)$$

$$\tilde{e}_t^{[l]} = a_t^{[l]} \tilde{a}_t^{[l]} + b_t^{[l]} \tilde{b}_t^{[l]} + c_t^{[l]} \tilde{c}_t^{[l]} + d_t^{[l]} \tilde{d}_t^{[l]} + \|G^{[l]}(f_t) - G^{[l]}(f_{t-1})\|_F^2 \quad (10e)$$

In the update phase, first, \tilde{a}_t is calculated that represents the reconstruction error. Plugged into eq(9) (which is inspired by AdaBoost), a_{t+1} is obtained. a_t for layer l is the weight of reconstruction error of this layer compared to the other layers, which is weighted by its importance \tilde{a}_t . Next, the weighted reconstruction error is added to the raw mask error to give the \tilde{b}_t . This is, in turn, used to calculate the weight of the mask error in this layer. This process is repeated for coincidental error, temporal component, and ST component. The errors of each layer are also accumulated to update the weight of the next. Hence, the network won't rely on the style information of a layer with large reconstruction error, etc. The same holds for ST co-occurrences.

Optimization and Target Localization Following [8], we used the Gauss-Seidel iterative approach to compute filter coefficients. The cost can be effectively minimized in Fourier domain due to the sparsity of DFT coefficients after regularizations. Image patch with the minimum error of eq(3) is a target candidate and target scale is estimated by applying the filter at multiple resolutions. The maximum filter response corresponds to the target's location and scale.

2.4 Implementation Details

We used VGG19 network consisting of 16 convolutional and 5 max-pooling layers, as implemented in MatConvNet [52] and pre-trained on the ImageNet dataset for the image classification. To be consistent with [2] and [31], we used

the conv layers after the pooling . We also added the input as Layer 0 which enables the tracker to benefit from NCC tracking, hence $\mathcal{L} = \{\text{input}, \text{conv1_1}, \text{conv2_1}, \text{conv3_1}, \text{conv4_1}, \text{conv5_1}\}$. In our implementation, $\sum_{l \in \mathcal{L}} a_t^{[l]}, \dots, e_t^{[l]} = 1$, regularization weights λ are determined with cross-validation on YouTubeBB [53] and are fixed for all the experiments, others parameters are similar to [2].

3 Experiments

We take a step-by-step approach to first prove that adding co-incidental and temporal regularization to the baseline improves the tracking performance, and then show that combining multiple layers can improve the tracking performance significantly. We also show that the regularization based on activation, style, and temporal coherence is helpful only if proper importance parameters are selected for different layers. Then we discuss the effect of different regularization terms on the performance of the tracker in different scenarios. Finally, we compare our proposed algorithm with the state-of-the-art and discuss its merits and demerits. For this comparison, we used success and precision plots and PASCAL metric ($IoU > 0.50$) over OTB50 [54]. For each frame, the area of the estimated box divided by the area of the annotation is defined as *scale ratio*, and its average and standard deviation represents the scale adaptation and jitteriness of a tracker.

For the comparison with latest trackers, we use OTB100 [55], UAV123 [56] and LaSot [57] with success and precision indexes and VOT2015 [58] and VOT2018 [59] using accuracy, robustness, and expected average overlap (EAO).³ We have developed our tracker with Matlab using MatConvNet and C++ and on a Nvidia RTX2080 GPU, we achieved the speed of 53.8 fps.

3.1 The Effects of Regularization on Single Layer

In this experiment, we study the effect of proposed regularizations on different CNN layers, used as the features in our baseline tracker, the single layer Deep-DCF using eq(1). Mask regularization (eq(4)) as MR, proposed co-incidental (CR, eq(5)), temporal (TR, eq(7)) and ST (SR, (8)) are then progressively added to the baseline tracker to highlight their contribution in the overall tracker performance (all importance weights are fixed to 1).

Layer-wise Analysis: Table 1 shows that the activations of features in the shallower layer of CNN generally yields better tracking compared to the deeper layers, except L5 which according to [2] contains high-level object-specific features. Shallower layers encodes more spatial information while accommodating a certain degree of invariance in target matching process. Contrarily, deeper layers ignore the perturbations of the appearance and perform semantic matching.

Mask Reg: Results shows that the use of mask regularization for tracking improve the tracking performance around 2.1-3.3%, where shallower layers benefit more from such regularization.

³ More info: <http://ishiilab.jp/member/meshgi-k/ccnt.html>.

Table 1. The effectiveness of regularizations with single layer of CNN with success rate $IoU > 0.50$. Here, we benchmarked baseline (B) with mask (MR), co-incidental (CR), temporal (TR), and ST style (SR) regularizations on OTB50 [54].

Layer (l)	L0	L1	L2	L3	L4	L5
B	46.2	62.3	57.4	53.9	52.9	56.3
B + MR	49.5	65.1	60.0	56.4	55.0	58.5
B + CR	45.1	61.9	57.5	57.4	55.1	58.4
B + TR	45.8	62.2	57.9	55.6	54.0	57.8
B + MR + CR	46.2	62.7	59.0	55.9	54.8	58.3
B + MR+ CR + TR	47.5	64.3	60.1	57.3	57.5	62.8
B + MR+ CR + TR + SR	48.1	64.7	60.1	58.3	58.2	62.0

Style Reg: The style information (CR) generally improves the tracking, especially in deeper layers which the activations are not enough to localize the target. However, when applied to shallower layers, especially input image, the style information may be misleading for the tracker which is aligned with the observation of Gatys et al. [31].

Temporal Reg: Deeper layers enjoys temporal regularization more. This is due to the fact that changes in activations in deeper layers signals semantic changes in the appearance, such as misalignment of the window to the target, partial or full occlusions or unaccounted target transformations such as out-of-plane rotations. In contrary, the changes in shallower layers come from minor changes in the low-level features (e.g. edges) that is abundant in the real-world scenarios, and using only this regularization for shallow layers is not recommended.

Spatiotemporal Reg: Using this regularization on top of temporal regularization, often improves the accuracy of the tracking since non-linear temporal styles of the features cannot be always handled using temporal reg.

All Regularizations: The combination of MR and CR terms, especially helps the tracking using deeper layers and starting from L2 it outperforms both MR and CR regularizations. The combination of all regularization terms proved to be useful for all layers, improving tracking accuracy by 2-6% compared to baseline.

Feature Interconnections: Feature interconnections can be divided into (i) spatial-coincidental (when two features have high filter responses at the same time), (ii) spatial-temporal (a feature activates in a frame), (iii) spatial-ST style (a feature coactivates with a motion feature), (iv) style-temporal (coupled features turns on/off simultaneously), (v) style-ST style (coupled features moves similarly), temporal-ST style (a features starts/stops moving). The features are designed to capture different aspects of the object’s appearance and motion; they are sometimes overlapping. Such redundancy improves tracking, with more complex features improving semantics of tracking, and low-level features improving the accuracy of the localization.

3.2 Employing Multiple Layers of CNN

To investigate different components of the proposed tracker, we prepared several ablated versions and compared them in Table 2. Three settings have been

considered for the importance weights: *uniform weights*, *random weights*, and *optimized weights based on the model update* (eq(9)). The random initial weights were generated for each time t (summed up to 1), and the experiment was repeated five times and averaged. By adding each regularization term, the speed of the tracker degrades, therefore, we added the ratio of the custom tracker speed to the baseline (first row) in the last column. It should be noted that when the spatial coefficient $b_t^{[l]}$ are zero, the L2 norm of all filter responses (of all layers) is used to regularize. *Uniform weighting* keeps reg. weights fixed and equal during tracking, *random weighting* assigns random weights to different components of each layer and our *proposed* AdaBoost-inspired weighting penalizes components proportional to their contribution in the previous mistakes.

Table 2. The effect of using multiple CNN layers with various importance weight strategies. This is based on the success rate ($IoU > 0.50$) on OTB50. Last column presents the speed of the ablated trackers (+ model update) compared to baseline (%).

Model Update	uniform	random	proposed	speed (%)
B ($\mathcal{B}_t = \mathcal{C}_t = \mathcal{D}_t = \mathcal{E}_t = 0$)	66.8	64.4	79.2	100.0
B + MR ($\mathcal{C}_t = \mathcal{D}_t = \mathcal{E}_t = 0$)	67.3	66.4	81.7	95.2
B + CR ($\mathcal{B}_t = \mathcal{D}_t = \mathcal{E}_t = 0$)	69.1	69.9	82.8	83.1
B + TR ($\mathcal{B}_t = \mathcal{C}_t = \mathcal{E}_t = 0$)	67.3	67.0	81.1	98.4
B + MR + CR ($\mathcal{D}_t = \mathcal{E}_t = 0$)	68.3	72.6	85.9	80.8
B + MR + CR + TR ($\mathcal{E}_t = 0$)	69.0	73.0	86.5	78.0
B + MR + CR + TR + SR	69.2	73.3	86.9	78.7

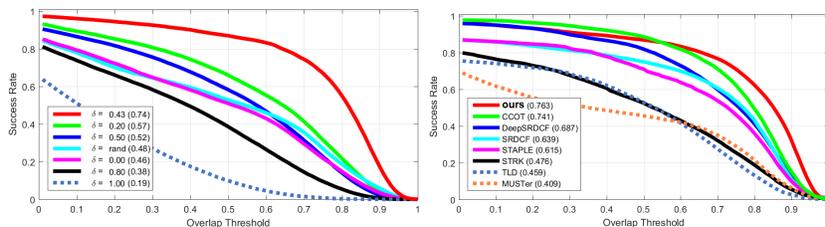
Comparing Model Update Schemes: Table 2 shows that with the use of proposed model update, different components of the tracker may collaborate to improve the overall performance of the tracker when combining different layers. However, uniform weights for all parameters (equal to $|\mathcal{L}|^{-1}$) cannot provide much improvement compared to the single layer DeepDCF, especially when compared to the L1 column of Table 1. Interestingly, random weights outperform uniform weights when applied to style regularization, which shows that not all layers contain equivalently useful co-incidental information.

Multiple Layers: By comparing each row of the Table 2 with the corresponding row of Table 1, the effect of combining different layers is evident. Comparing the first rows shows the advantage of combining layers without using any regularization. Uniform weights for the layers raise the performance only by 4.5% (all layers vs. L1), whereas the proposed fusion can boost the combination performance up to 16.9%. This is a recurring pattern for other rows that show the benefit of the layer combination for activations, as well as regularization terms.

While our method can be seen as a feature selection/weight tuning, it is crucial to see the tuning procedure as a layer-wise adaptation. In each layer, the effect of different regularization term is determined by its contribution in the loss term. This calculation is isolated from other layers. Features of each layer compete with each other to better represent the target, but cooperate with each other to deliver the best overall representation that can be obtained from

Table 3. Scale adaptation obtained by proposed regularizations on OTB-100 measured by the mean of estimate-to-real scale ratio and its standard deviation (jitter).

Tracker	B	B+MR	B+CR	B+TR	B+MR+TR	B+MR+CR+TR	ALL
Avg.Ratio	92.2	93.1	93.3	93.8	93.3	94.2	94.7
Jitter	8.17	7.13	5.81	2.66	5.11	2.40	2.35

**Fig. 3.** (left) The activation vs. style trade-off for the custom tracker on OTB50. While $\delta \rightarrow 1$ puts too much emphasize on the style, $\delta = 0$ overemphasizes on the activations. (middle) Performance comparison of trackers on OTB100 using success plot.

that particular layer. Additionally, to use different types of features, we utilize the combination of different layers to balance the detail-semantic trade-off in different tracking scenarios; therefore layers' importance should be adaptively adjusted.

Applying Different Reg: Similar to the case of single layers, regularization multiple layers is also effective. In case of uniform weights, using CR outperforms MR+CR which indicates that without proper weight adjustment, different regularization cannot be effectively stacked in the tracker. Therefore, it is expected that the proposed adaptive weight can handle this case, as table shows.

3.3 Scale Adaptation

The proposed style and temporal regs, tend to discard candidates with mismatching scale due to style and continuity inconsistencies. Additionally, temporal reg tend to reduce the jittering of the position and scale. Table 3 demonstrates the proposed tracker with multi-layers of CNN, adaptive weights and different regs.

3.4 Activation vs. Style

As seen in NST literature [31,32], various amount of focus on the content image and style image yields different outcomes. In tracking, however, the accuracy provides a measure to balance this focus. We conducted an experiment to see the effect of the regularization weights λ_{sty} on the tracking performance. Hence, we set $\lambda_{sty} = \delta$ while disabling spatial and temporal regularizers. Figure 3-left depicts the success plot for several δ and the optimal value δ^* (via annealing and cross-validation on OTB-50, with proposed model update for layers in \mathcal{L}).

This figure also depicts the performance of the obtained tracker with various values of δ . The results reveal that when the tracker ignores the style information ($\delta = 0$, the base multi-layer tracker) its performance is better than when it ignores activations ($\delta = 1$) since the style information is not suitable in isolation for tracking. The values between these two extremes work better by enhancing the activations with style information. However, finding a sweetspot is difficult and scenario-dependent, e.g., for textureless object [60] more spatial and semantic information is required, whereas textured objects benefit from style feedback.

3.5 Preliminary Analysis

We compared our tracker with TLD [61], STRUKK [62], MEEM [63], MUSTer [64], STPL [65], CMT [66], SRDCF [8], dSRDCF [67] and CCOT [29].

Table 4. Quantitative evaluation of trackers (top) using average success on OTB50 [54] for different tracking challenges; (middle) success and precision rates on OTB100 [55], estimated-on-real scale ratio and jitter; (bottom) robustness and accuracy on VOT2015 [58]. The first, second and third best methods are shown in color.

	TLD	STRUKK	MEEM	MUSTer	STAPLE	SRDCF	dSRDCF	CCOT	Ours
OTB50	Illumination	0.48	0.53	0.62	0.73	0.68	0.70	0.71	0.75 0.80
	Deformation	0.38	0.51	0.62	0.69	0.70	0.67	0.67	0.69 0.78
	Occlusion	0.46	0.50	0.61	0.69	0.69	0.70	0.71	0.76 0.79
	Scale Changes	0.49	0.51	0.58	0.71	0.68	0.71	0.75	0.76 0.82
	In-plane Rot.	0.50	0.54	0.58	0.69	0.69	0.70	0.73	0.72 0.80
	Out-of-plane Rot.	0.48	0.53	0.62	0.70	0.67	0.69	0.70	0.74 0.81
	Out-of-view	0.54	0.52	0.68	0.73	0.62	0.66	0.66	0.79 0.81
	Low Resolution	0.36	0.33	0.43	0.50	0.47	0.58	0.61	0.70 0.74
	Background Clutter	0.39	0.52	0.67	0.72	0.67	0.70	0.71	0.70 0.78
	Fast Motion	0.45	0.52	0.65	0.65	0.56	0.63	0.67	0.72 0.78
	Motion Blur	0.41	0.47	0.63	0.65	0.61	0.69	0.70	0.72 0.78
Average Success	0.49	0.55	0.62	0.72	0.69	0.70	0.71	0.75 0.80	
OTB100	Average Success	0.46	0.48	0.65	0.57	0.62	0.64	0.69	0.74 0.76
	Average Precision	0.58	0.59	0.62	0.74	0.73	0.71	0.81	0.85 0.85
	<i>IoU</i> > 0.5	0.52	0.52	0.62	0.65	0.71	0.75	0.78	0.88 0.86
	Average Scale	116.4	134.7	112.1	-	110.8	88.5	101.8	94.0 93.7
Jitter	8.2	8.7	8.2	-	5.9	4.1	4.9	3.8 2.3	
VOT	Accuracy	-	0.47	0.50	0.52	0.53	0.56	0.53	0.54 0.76
	Robustness	-	1.26	1.85	2.00	1.35	1.24	1.05	0.82 0.65

Attribute Analysis: We use partial subsets of OTB50 [54] with a distinguishing attribute to evaluate the tracker performance under different situations. Table 4 shows the superior performance of the algorithm, especially in handling deformations (by adaptive fusion of deep and shallow layers of CNN) and background clutter (by spatial and style reg.) and motion (by temporal reg.). Figure 4 demonstrates the performance of the tracker on several challenging scenarios. **OTB100:** Figure 3 (right) and Table 4 presents the success and precision plots of our tracker along with others. Data shows that proposed algorithm has superior performance, less jitter, and comparable localization and scale adaptation.

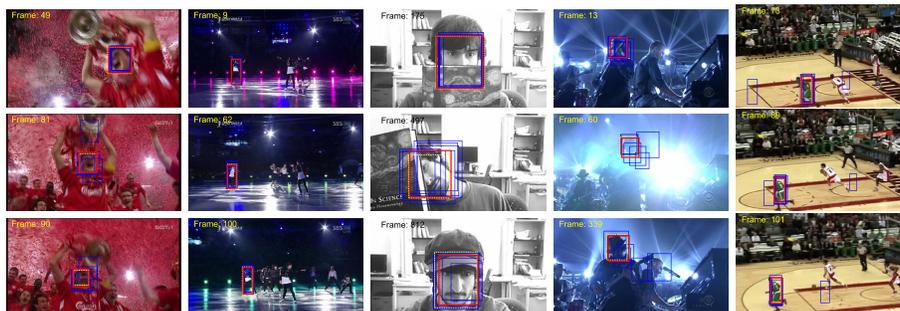


Fig. 4. Example tracking results on Soccer, Skating1, FaceOcc1, Shaking, and Basketball with severe occlusion, noise and illumination changes, scaling and 3D rotations, and clutter. (Red: proposed tracker, Blue: other trackers, Yellow: GT).

Table 5. Evaluation of deep trackers on OTB100 [55] using success rate and precision.

	ECO	ATOM	VITAL	HDT	YCNN	MDNet	dSTRCF	STResCF	CBCW	SiFC	SiRPN	SiRPN++	SINT++	Ours
Avg. Succ	0.69	0.66	0.68	0.65	0.60	0.67	0.68	0.59	0.61	0.59	0.63	0.69	0.57	0.76
Avg. Prec	0.91	-	0.91	0.84	0.84	0.90	-	0.83	0.81	0.78	0.85	0.91	0.76	0.85
$IoU > \frac{1}{2}$	0.74	0.86	0.79	0.68	0.74	0.85	0.77	0.76	0.76	0.76	0.80	0.83	0.78	0.86

VOT2015: Table 4 also shows superior performance in terms of accuracy coupled with decent robustness.

3.6 Comparison with State-of-the-Art

Deep Trackers: We compared our tracker against recent deep trackers on OTB100, including ECO [44], ATOM [68], VITAL [24], HDT [69], YCNN [16], MDNet [70], dSTRCF [71], STResCF [42], CBCW [72], SiamFC [20], SiamRPN [21], SiamRPN++ [22], SINT++ [23], and DiMP [73]. Table 5 shows that although our proposed tracker has some issues in accurate localization, it has a superior overall performance and success rate in handling various scenarios.

The experiments revealed that the proposed tracker is resistant to target abrupt or extensive target changes. Temporal and ST features in our method monitor the inconsistency in target appearance and style, co-occurrence features in different levels of abstraction provide different levels of robustness to target changes (from low-frequency features to the high-level features such as object part relations). The dynamic weighting enables the tracker to have the flexibility to resort to more abstract feature when the target undergoes drastic changes, and ST features handle abnormalities such as temporal occlusion and deformations.

Recent Public Datasets: Our method is compared with recent state-of-the-art methods in VOT2018 [59], UAV123 [56], LaSOT [57], GOT-10K [74] and TrackingNet [75] datasets. In phase I of LaSOT evaluation, our tracker is trained on our own data and tested on all 1400 training video sequences of LaSOT. In phase II, the training data is limited to the given 1120 training videos and

Table 6. Evaluation on VOT2018 by the means of EAO, robustness and accuracy.

	STURCK	MEEM	STAPLE	SRDCF	CCOT	SiamFC	ECO	SiamRPN	SiamRPN++	ATOM	Ours
EAO	0.097	0.192	0.169	0.119	0.267	0.188	0.280	0.383	0.414	0.401	0.408
Accuracy	0.418	0.463	0.530	0.490	0.494	0.503	0.484	0.586	0.600	0.590	0.586
Robustness	1.297	0.534	0.688	0.974	0.318	0.585	0.276	0.276	0.234	0.204	0.281

Table 7. Evaluation on LaSOT with protocol I (testing on all videos) and protocol II (training on given videos and testing on the rest). We get better results with dataset’s own videos as training due to large training set and matching domain.

	STAPLE	SRDCF	SiamFC	SINT	MDNet	ECO	BACF	VITAL	ATOM	SiamRPN++	DiMP	Ours
(I) Accuracy	0.266	0.271	0.358	0.339	0.413	0.340	0.277	0.412	0.515	0.496	0.596	0.521
(I) Robustness	0.231	0.227	0.341	0.229	0.374	0.298	0.239	0.372	-	-	-	0.411
(II) Accuracy	0.243	0.245	0.336	0.314	0.397	0.324	0.259	0.390	-	-	-	0.507
(II) Robustness	0.239	0.219	0.339	0.295	0.373	0.301	0.239	0.360	-	-	-	0.499

Table 8. Evaluation on UAV123 by success rate and precision. Our algorithm is having difficulty with small/ low resolution targets.

	TLD	STRUCK	MEEM	STAPLE	SRDCF	MUSTer	ECO	ATOM	SiamRPN	SiamRPN++	DiMP	Ours
Success	0.283	0.387	0.398	0.453	0.473	0.517	0.399	0.650	0.527	0.613	0.653	0.651
Precision	0.439	0.578	0.627	-	0.676	-	0.591	-	0.748	0.807	-	0.833

Table 9. Benchmarking on TrackingNet and GOT-10k

	ECODa	Siam-RPN	ATOM	SiamRPN++	DiMP	SiamMask	D3S	SiamFC++	SiamRCNN	ours		
T-Net	Prec.	0.492	0.591	0.648	0.694	0.687	0.733	-	0.705	0.800	0.711	
	N-Prec.	0.618	0.733	0.771	0.800	0.801	0.664	-	0.800	0.854	0.810	
	Success	0.554	0.638	0.703	0.733	0.740	0.778	-	0.754	0.812	0.752	
GOT-10k	AO	0.316	0.417	0.556	0.518	0.611	0.514	0.597	0.595	-	0.601	
	SR	0.750	0.111	0.149	0.402	0.325	0.492	0.366	0.462	0.479	-	0.479
CCOT	SR	0.5	0.309	0.461	0.635	0.618	0.717	0.587	0.676	0.695	-	0.685

tested on the rest. The results are better than SiamRPN++ in VOT2018 and UAV123, and VITAL in LaSOT, despite using a pre-trained CNN. This method benefits from multi-layer fusion, adaptive model update, and various regularization. Comparing the results of the benchmark with ATOM, SiamRPN++ and DiMP showed that just using convolutional layers and using the underlying features is not enough to perform a high-level tracking. Having a deeper network (ResNet-18 in ATOM and ResNet-50 in DiMP) and having auxiliary branches (region proposal in SiamRPN, IOU prediction in ATOM, and model prediction in DiMP) are two main differences between our method and the SotA. However, the proposed method offers insight about extra features to be used in tracking, along with the activation such as cooccurrence features that are embedded in all CNNs ready to be exploited. Further, good performance on specific datasets such as UAV123 demonstrates that these features can support different object types and contexts.

4 Conclusion

We proposed a tracker that exploits various CNN statistics including activations, spatial data, co-occurrences within a layer, and temporal changes and patterns between time slices. It adaptively fuses several CNN layers to negate the demerits of each layer with merits of others. It outperformed recent trackers in various experiments, promoting the use of spatial and temporal style in tracking. Our regularizations can be used with other CNN-based methods.

References

1. Li, H., Li, Y., Porikli, F.: Deeptrack: Learning discriminative feature representations online for robust visual tracking. *IEEE TIP* **25** (2016) 1834–1848 [1](#)
2. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Convolutional features for correlation filter based visual tracking. In: *Proceedings of the IEEE International Conference on Computer Vision Workshops*. (2015) 58–66 [1](#), [5](#), [8](#), [9](#)
3. Wang, L., Liu, T., Wang, G., Chan, K.L., Yang, Q.: Video tracking using learned hierarchical features. *IEEE TIP* **24** (2015) 1424–1435 [1](#)
4. Zhang, T., Xu, C., Yang, M.H.: Multi-task correlation particle filter for robust object tracking. In: *CVPR*. Volume 1. (2017) 3 [1](#)
5. Hong, S., You, T., Kwak, S., Han, B.: Online tracking by learning discriminative saliency map with convolutional neural network. In: *International Conference on Machine Learning*. (2015) 597–606 [1](#)
6. Sharif Razavian, A., Azizpour, H., Sullivan, J., Carlsson, S.: Cnn features off-the-shelf: an astounding baseline for recognition. In: *CVPRw*. (2014) 806–813 [1](#)
7. Cimpoi, M., Maji, S., Vedaldi, A.: Deep convolutional filter banks for texture recognition and segmentation. *arXiv preprint arXiv:1411.6836* (2014) [1](#)
8. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Learning spatially regularized correlation filters for visual tracking. In: *ICCV’15*. (2015) 4310–4318 [1](#), [4](#), [6](#), [8](#), [13](#)
9. Nam, H., Baek, M., Han, B.: Modeling and propagating cnns in a tree structure for visual tracking. *arXiv preprint arXiv:1608.07242* (2016) [1](#), [2](#)
10. Wang, N., Yeung, D.Y.: Learning a deep compact image representation for visual tracking. In: *NIPS*. (2013) 809–817 [2](#)
11. Zhou, X., Xie, L., Zhang, P., Zhang, Y.: An ensemble of deep neural networks for object tracking. In: *Image Processing (ICIP), 2014 IEEE International Conference on, IEEE* (2014) 843–847 [2](#)
12. Fan, J., Xu, W., Wu, Y., Gong, Y.: Human tracking using convolutional neural networks. *IEEE Transactions on Neural Networks* **21** (2010) 1610–1623 [2](#)
13. Ma, C., Huang, J.B., Yang, X., Yang, M.H.: Hierarchical convolutional features for visual tracking. In: *ICCV*. (2015) 3074–3082 [2](#), [4](#)
14. Zhang, K., Liu, Q., Wu, Y., Yang, M.: Robust visual tracking via convolutional networks without training. *IEEE TIP* **25** (2016) 1779–1792 [2](#)
15. Zhu, Z., Huang, G., Zou, W., Du, D., Huang, C.: Uct: learning unified convolutional networks for real-time visual tracking. In: *ICCVw*. (2017) 1973–1982 [2](#)
16. Chen, K., Tao, W.: Once for all: a two-flow convolutional neural network for visual tracking. *IEEE CSVT* (2018) 1–1 [2](#), [14](#)
17. Wang, N., Li, S., Gupta, A., Yeung, D.Y.: Transferring rich feature hierarchies for robust visual tracking. *arXiv preprint arXiv:1501.04587* (2015) [2](#)

18. Drayer, B., Brox, T.: Object detection, tracking, and motion segmentation for object-level video segmentation. arXiv preprint arXiv:1608.03066 (2016) [2](#)
19. Tao, R., Gavves, E., Smeulders, A.W.: Siamese instance search for tracking. In: CVPR. (2016) 1420–1429 [2](#)
20. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.: Fully-convolutional siamese networks for object tracking. In: ECCV, Springer (2016) 850–865 [2](#), [14](#)
21. Li, B., Yan, J., Wu, W., Zhu, Z., Hu, X.: High performance visual tracking with siamese region proposal network. In: CVPR'18. (2018) 8971–8980 [2](#), [14](#)
22. Li, B., Wu, W., Wang, Q., Zhang, F., Xing, J., Yan, J.: Siamrpn++: Evolution of siamese visual tracking with very deep networks. In: CVPR'19. (2019) 4282–4291 [2](#), [14](#)
23. Wang, X., Li, C., Luo, B., Tang, J.: Sint++: Robust visual tracking via adversarial positive instance generation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 4864–4873 [2](#), [14](#)
24. Song, Y., Ma, C., Wu, X., Gong, L., Bao, L., Zuo, W., Shen, C., Rynson, L., Yang, M.H.: Vital: Visual tracking via adversarial learning. In: CVPR. (2018) [2](#), [14](#)
25. Vondrick, C., Shrivastava, A., Fathi, A., Guadarrama, S., Murphy, K.: Tracking emerges by colorizing videos. In: ECCV. (2018) [2](#)
26. Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: ECCV, Springer (2014) 818–833 [2](#), [3](#)
27. Liu, L., Shen, C., van den Hengel, A.: The treasure beneath convolutional layers: Cross-convolutional-layer pooling for image classification. In: CVPR. (2015) 4749–4757 [2](#)
28. Qi, Y., Zhang, S., Qin, L., Huang, Q., Yao, H., Lim, J., Yang, M.H.: Hedging deep features for visual tracking. PAMI (2018) [2](#), [4](#)
29. Danelljan, M., Robinson, A., Khan, F.S., Felsberg, M.: Beyond correlation filters: Learning continuous convolution operators for visual tracking. In: ECCV, Springer (2016) 472–488 [2](#), [4](#), [13](#)
30. Bovik, A.C., Clark, M., Geisler, W.S.: Multichannel texture analysis using localized spatial filters. PAMI (1990) 55–73 [3](#)
31. Gatys, L.A., Ecker, A.S., Bethge, M.: Image style transfer using convolutional neural networks. In: CVPR. (2016) 2414–2423 [3](#), [8](#), [10](#), [12](#)
32. Johnson, J., Alahi, A., Fei-Fei, L.: Perceptual losses for real-time style transfer and super-resolution. In: ECCV, Springer (2016) 694–711 [3](#), [12](#)
33. Matsuo, S., Yanai, K.: Cnn-based style vector for style image retrieval. In: Proceedings of the 2016 ACM on International Conference on Multimedia Retrieval, ACM (2016) 309–312 [3](#)
34. Varol, G., Laptev, I., Schmid, C.: Long-term temporal convolutions for action recognition. PAMI **40** (2018) 1510–1517 [3](#)
35. Gkioxari, G., Malik, J.: Finding action tubes. In: CVPR. (2015) 759–768 [3](#)
36. Chao, Y.W., Vijayanarasimhan, S., Seybold, B., Ross, D.A., Deng, J., Sukthankar, R.: Rethinking the faster r-cnn architecture for temporal action localization. In: CVPR. (2018) 1130–1139 [3](#)
37. Simonyan, K., Zisserman, A.: Two-stream convolutional networks for action recognition in videos. In: NIPS. (2014) 568–576 [3](#)
38. Zhu, Z., Wu, W., Zou, W., Yan, J.: End-to-end flow correlation tracking with spatial-temporal attention. CVPR **42** (2017) 20 [3](#)

39. Dosovitskiy, A., Fischer, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., Van Der Smagt, P., Cremers, D., Brox, T.: FlowNet: Learning optical flow with convolutional networks. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 2758–2766 [3](#)
40. Feichtenhofer, C., Pinz, A., Wildes, R.P., Zisserman, A.: What have we learned from deep representations for action recognition? *connections* **19** (2018) 29 [3](#)
41. Gladh, S., Danelljan, M., Khan, F.S., Felsberg, M.: Deep motion features for visual tracking. In: ICPR, IEEE (2016) 1243–1248 [3](#)
42. Zhu, Z., et al.: STResNet_cf tracker: The deep spatiotemporal features learning for correlation filter based robust visual object tracking. *IEEE Access* **7** (2019) [4](#), [14](#)
43. Danelljan, M., Robinson, A., Khan, F.S., Felsberg, M.: Beyond correlation filters: Learning continuous convolution operators for visual tracking. In: European Conference on Computer Vision, Springer (2016) 472–488 [4](#)
44. Danelljan, M., Bhat, G., Khan, F.S., Felsberg, M.: Eco: Efficient convolution operators for tracking. In: CVPR. (2017) [4](#), [14](#)
45. Galoogahi, H.K., Fagg, A., Lucey, S.: Learning background-aware correlation filters for visual tracking. *arXiv preprint arXiv:1703.04590* (2017) [4](#)
46. Danelljan, M., Häger, G., Khan, F., Felsberg, M.: Accurate scale estimation for robust visual tracking. In: BMVC, BMVA Press (2014) [5](#)
47. Henriques, J.F., Caseiro, R., Martins, P., Batista, J.: Exploiting the circulant structure of tracking-by-detection with kernels. In: ECCV’12, Springer (2012) 702–715 [5](#)
48. Bolme, D.S., Beveridge, J.R., Draper, B.A., Lui, Y.M.: Visual object tracking using adaptive correlation filters. In: CVPR’10, IEEE (2010) 2544–2550 [5](#)
49. Jepson, A.D., Fleet, D.J., El-Maraghi, T.F.: Robust online appearance models for visual tracking. *PAMI* (2003) [6](#)
50. Berger, G., Memisevic, R.: Incorporating long-range consistency in cnn-based texture generation. *ICLR* (2017) [7](#)
51. Wiyatno, R.R., Xu, A.: Physical adversarial textures that fool visual object tracking. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 4822–4831 [7](#)
52. Vedaldi, A., Lenc, K.: Matconvnet: Convolutional neural networks for matlab. In: Proceedings of the 23rd ACM international conference on Multimedia, ACM (2015) 689–692 [8](#)
53. Real, E., Shlens, J., Mazzocchi, S., Pan, X., Vanhoucke, V.: Youtube-boundingboxes: A large high-precision human-annotated data set for object detection in video. In: CVPR’17. (2017) 5296–5305 [9](#)
54. Wu, Y., Lim, J., Yang, M.H.: Online object tracking: A benchmark. In: CVPR’13, IEEE (2013) 2411–2418 [9](#), [10](#), [13](#)
55. Wu, Y., Lim, J., Yang, M.H.: Object tracking benchmark. *PAMI* (2015) [9](#), [13](#), [14](#)
56. Mueller, M., et al.: A benchmark and simulator for uav tracking. In: ECCV’16, Springer (2016) [9](#), [14](#)
57. Fan, H., et al.: LaSOT: A high-quality benchmark for large-scale single object tracking. *CVPR’19* (2019) [9](#), [14](#)
58. Kristan, M., Matas, J., Leonardis, A., Felsberg, M., Cehovin, L., Fernández, G., Vojir, T., Hager, G., Nebehay, G., Pflugfelder, R.: The visual object tracking vot2015 challenge results. In: ICCVw’15. (2015) 1–23 [9](#), [13](#)
59. Kristan, M., et al.: The sixth visual object tracking vot2018 challenge results. In: ECCV’18. (2018) [9](#), [14](#)

60. Choi, C., Christensen, H.I.: 3d textureless object detection and tracking: An edge-based approach. In: 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, IEEE (2012) 3877–3884 [13](#)
61. Kalal, Z., Mikolajczyk, K., Matas, J.: Tracking-learning-detection. PAMI **34** (2012) 1409–1422 [13](#)
62. Hare, S., Saffari, A., Torr, P.H.: Struck: Structured output tracking with kernels. In: ICCV’11. (2011) [13](#)
63. Zhang, J., Ma, S., Sclaroff, S.: Meem: Robust tracking via multiple experts using entropy minimization. In: ECCV. (2014) [13](#)
64. Hong, Z., Chen, Z., Wang, C., Mei, X., Prokhorov, D., Tao, D.: Multi-store tracker (muster): A cognitive psychology inspired approach to object tracking. In: CVPR’15. (2015) 749–758 [13](#)
65. Bertinetto, L., Valmadre, J., Golodetz, S., Miksik, O., Torr, P.H.: Staple: Complementary learners for real-time tracking. In: CVPR. (2016) 1401–1409 [13](#)
66. Meshgi, K., Oba, S., Ishii, S.: Active discriminative tracking using collective memory. (In: MVA’17) [13](#)
67. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Convolutional features for correlation filter based visual tracking. In: ICCVw. (2015) 58–66 [13](#)
68. Danelljan, M., Bhat, G., Khan, F.S., Felsberg, M.: Atom: Accurate tracking by overlap maximization. In: CVPR’19. (2019) 4660–4669 [14](#)
69. Qi, Y., Zhang, S., Qin, L., Yao, H., Huang, Q., Lim, J., Yang, M.H.: Hedged deep tracking. In: CVPR. (2016) 4303–4311 [14](#)
70. Nam, H., Han, B.: Learning multi-domain convolutional neural networks for visual tracking. In: CVPR. (2016) 4293–4302 [14](#)
71. Li, F., et al.: Learning spatial-temporal regularized correlation filters for visual tracking. In: CVPR’18. (2018) [14](#)
72. Zhou, Y., et al.: Efficient correlation tracking via center-biased spatial regularization. IEEE TIP **27** (2018) 6159–6173 [14](#)
73. Bhat, G., Danelljan, M., Gool, L.V., Timofte, R.: Learning discriminative model prediction for tracking. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 6182–6191 [14](#)
74. Huang, L., Zhao, X., Huang, K.: Got-10k: A large high-diversity benchmark for generic object tracking in the wild. IEEE TPAMI (2019) [14](#)
75. Muller, M., Bibi, A., Giancola, S., Alsubaihi, S., Ghanem, B.: Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In: ECCV’2018. (2018) [14](#)