# Unsupervised Deep Representation Learning for Real-Time Tracking

**Ning Wang · Wengang Zhou · Yibing Song · Chao Ma · Wei Liu · Houqiang Li**

**Abstract** The advancement of visual tracking has continuously been brought by deep learning models. Typically, supervised learning is employed to train these models with expensive labeled data. In order to reduce the workload of manual annotations and learn to track arbitrary objects, we propose an unsupervised learning method for visual tracking. The motivation of our unsupervised learning is that a robust tracker should be effective in bidirectional tracking. Specifically, the tracker is able to forward localize a target object in successive frames and backtrace to its initial position in the first frame. Based on such a motivation, in the training process, we measure the consistency between forward and backward trajectories to learn a robust tracker from scratch merely using unlabeled videos. We build our framework on a Siamese correlation filter network, and propose a multi-frame validation scheme and a cost-sensitive loss to facilitate unsupervised learning. Without bells and whistles, the proposed unsupervised tracker achieves the baseline accuracy as classic fully supervised trackers while achieving a real-time speed. Furthermore, our unsupervised framework exhibits a potential in leveraging more unlabeled or weakly labeled data to further improve the tracking accuracy.

**Keywords** Visual tracking · Unsupervised learning · Correlation filter · Siamese network

Ning Wang
The CAS Key Laboratory of GIPAS, University of Science and Technology of China, Hefei, China.
E-mail: wn6149@mail.ustc.edu.cn

Wengang Zhou
The CAS Key Laboratory of GIPAS, University of Science and Technology of China, Hefei, China.
Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, Hefei, China.
E-mail: zhwg@ustc.edu.cn

Yibing Song
Tencent AI Lab, Shenzhen, China.
E-mail: yibingsong.cv@gmail.com

Chao Ma
The MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai, China.
E-mail: chaoma@sjtu.edu.cn

Wei Liu
Tencent AI Lab, Shenzhen, China.
E-mail: wl2223@columbia.edu

Houqiang Li
The CAS Key Laboratory of GIPAS, University of Science and Technology of China, Hefei, China.
Institute of Artificial Intelligence, Hefei Comprehensive National Science Center, Hefei, China.
E-mail: lihq@ustc.edu.cn

Corresponding Authors: Wengang Zhou and Houqiang Li

## 1 Introduction

Visual object tracking is a fundamental task in computer vision with numerous applications including video surveillance, autonomous driving, augmented reality, and human-computer interactions. It aims to localize a moving object annotated at the initial frame with a bounding box. Recently, deep models have improved the tracking accuracies by strengthening the feature representations [41, 13, 9] or optimizing networks end-to-end [1, 32, 46, 58]. These models are offline pretrained with full supervision, which requires a large number of annotated ground-truth labels during the training stage. Manual annotations are always expensive and time-consuming, whereas a huge number of unlabeled videos are readily available on the Internet. On the other hand, visual tracking differs from other recognition tasks (e.g., object detection, image classification) in the sense that object labels vary according to target initializations on the first frame. The extensive and uncertain labeling process for supervised learning raises our interest to develop an alterna-

**Fig. 1** Visual tracking via supervised and unsupervised learnings. Supervised learning requires ground-truth labels for individual frames in the training videos, while our proposed unsupervised learning is free of any labels by measuring the trajectory consistency between forward and backward tracking.

tive learning scheme by using unlabeled video sequences in the wild.

In this paper, we propose an unsupervised learning approach for visual tracking. Instead of using off-the-shelf deep models, we train the visual tracking network from scratch. The intuition of unsupervised learning resides on the bidirectional motion analysis in video sequences. Tracking an object can be executed in both the forward and backward ways. Initially, given the bounding box annotation of a target object in the first frame, we can track the target object forward in the subsequent frames. When tracking backward, we use the predicted location in the last frame as the initial target bounding box, and track it backward towards the first frame. Ideally, the estimated bounding box location in the first frame is identical with the given one in the forward pass. In this work, we measure the difference between the forward and backward target trajectories and formulate it as a loss function. We use the computed loss to train our network in a self-supervised manner[1], as shown in Fig. 1. By repeatedly tracking forward and backward, our model learns to locate target objects in consecutive frames without any supervision.

The proposed unsupervised training aims to learn a generic feature representation instead of strictly focusing on tracking a complete object. In the first frame, we initialize a bounding box that covers the informative local region with high image entropy. The bounding box may contain arbitrary image content and may not cover an entire object. Then, our tracking network learns to track the bounding box region in the training video sequences. Our unsupervised annotation shares similarity with the part-based [36] and edge-based [34] tracking methods that track the subregions of a target object. We expect our tracker not only to concentrate on the shape of a complete object, but also to track any part

of it. The bounding box initialization by image entropy gets rid of the manual annotation on the first frame and thus ensures the whole learning process unsupervised.

We employ unsupervised learning under the Siamese correlation filter framework. The training steps consist of forward tracking and backward verification. A limitation of the forward and backward consistency measurement is that the target trajectory in the forward pass may coincide with that in the backward pass although the tracker loses the target. The consistency loss function fails to penalize this situation because the predicted target region can still backtrace to the initial position on the first frame regardless of losing the target. In addition, challenges such as heavy occlusion or out-of-view in training videos will degrade the CNN feature representation capability. To tackle these issues, we introduce a multi-frame validation scheme and a cost-sensitive loss to facilitate unsupervised training. If the tracker loses the target, the trajectories predicted from the forward and backward directions are unlikely to be consistent when more frames are used in the training stage. Besides, we propose a new cost-sensitive loss to alleviate the impact of the noisy samples during unsupervised learning. The training samples containing background texture will be excluded by the image entropy measurement. Based on the multi-frame validation and sample selection strategies discussed above, our network training is stabilized.

We evaluate our method on the challenging benchmark datasets including OTB-2013 [69], OTB-2015 [70], Temple-Color [35], VOT2016 [25], VOT2017/2018 [24], LaSOT [16], and TrackingNet [45]. Extensive experimental results indicate that without bells and whistles, the proposed unsupervised tracker is even comparable with the baseline configuration of fully supervised trackers [1,58,64]. When integrated with an adaptive online model update [11,9], the proposed tracker shows state-of-the-art performance. It is worth mentioning that our tracker trained via unsupervised learning achieves comparable performance with that via supervised learning when only limited or noisy labels are available. In addition, we demonstrate the potential of our tracker to further boost the accuracy by using more unlabeled data. A complete analysis of various training configurations is given in Section 4.2.

In summary, the contributions of this work are threefold:

- We propose an unsupervised learning method on the Siamese correlation filter network. The unsupervised learning consists of forward and backward trackings to measure the trajectory consistency for network training.
- We propose a multi-frame validation scheme to enlarge the trajectory inconsistency when the tracker loses the target. In addition, we propose a cost-sensitive loss and an entropy selection metric to reduce the contributions from easy samples in the training process.

---

[1] In this paper, we do not distinguish between the terms *unsupervised* and *self-supervised*, as both refer to learning without ground-truth annotations.

– The extensive experiments carried out on seven standard benchmarks show the favorable performance of the proposed tracker. We provide an in-depth analysis of our unsupervised representation and reveal the potential of unsupervised learning in visual tracking.

In the remainder of this paper, we describe the related work in Section 2, the proposed method in Section 3, and the experiments in Section 4. Finally, we conclude the paper in Section 5.

## 2 Related Work

In this section, we perform a literature review on deep tracking methods, forward-backward motion analysis, and unsupervised representation learning.

### 2.1 Deep Visual Tracking

Deep models have influenced visual tracking mainly from two perspectives. The first one is to provide a discriminative CNN feature representation by using off-the-shelf backbones (e.g., VGG [52,3]), while the second one is to formulate a complete tracking network for end-to-end training and predictions. The discriminative correlation filters (DCFs) [2, 19,10,21,42,55,38] handle the visual tracking task by solving a ridge regression task using densely sampled candidates. While being integrated with discriminative CNN features, the remaining operations (e.g., regression solver, online update) are kept still in the DCF trackers [13,33,63, 9]. On the other hand, the end-to-end learning network can be categorized as classification and regression based networks. The classification networks [46,54,22,48] incrementally train a binary classifier to differentiate the target and background distractors. The regression networks [53,37,40] use CNN layers to regress CNN features of the search region to a response map for accurate localization. These end-to-end learning networks need online update and inevitably increase the computational burden.

Recently, the Siamese network has received huge investigations because of its efficiency in online prediction. The SiamFC tracker [1] uses a cross-correlation layer to measure feature similarity between the template patch and search patches. The fully convolutional nature of SiamFC efficiently predicts the target response without redundancy. By incorporating the region proposal network (RPN) [50], the SiamRPN methods [32,80] achieve state-of-the-art performance while running at 160 FPS. Other improvements based on Siamese networks include ensemble learning [18], dynamic memory [71], attention modulation [65], capacity increments [76], and reinforcement learning [20,15]. By integrating the correlation filter, the Siamese correlation filter
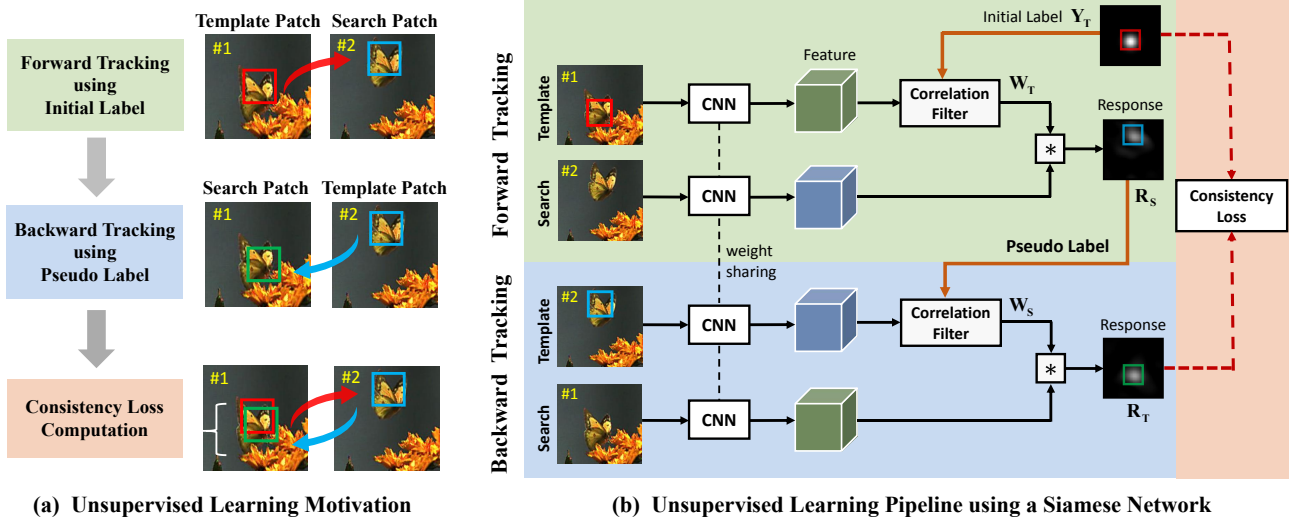
network [58,64] achieves favorable performance even with an extremely lightweight model. Different from the above deep trackers that train a CNN model in a supervised manner or directly use off-the-shelf deep models, we propose to learn a tracking network from scratch using unlabeled videos via unsupervised learning.

### 2.2 Forward-Backward Analysis

The forward and backward strategy has been investigated in motion analysis scenarios. Meister et al. [43] combined the forward-backward consistency estimation and pixel construction to learn optical flows. Wang et al. [67] leveraged the cycle-consistency across multiple steps temporally to learn feature representations for different tasks. The differentiable tracker in [67] is deliberately designed to be weak for feature representation learning. In contrast, aiming at robust visual tracking, we adopt a strong tracking baseline (Siamese correlation filter network), which is not fully-differentiable in the trajectory loop due to the pseudo labeling. However, by repeating forward tracking and backward verification, we incrementally promote the tracking network by pseudo-labeling based self-training. The forward-backward consistency check is also applied in image alignment [79,78] and depth estimation [73,77]. In the visual tracking community, the forward-backward consistency is mainly used for the output reliability or uncertainty measurement. The tracking-learning-detection (TLD) [23] uses the Kanade-Lucas-Tomasi (KLT) tracker [56] to perform forward-backward matching to detect tracking failures. Lee et al. [30] proposed to select the reliable base tracker by comparing the geometric similarity, cyclic weight, and appearance consistency between a pair of forward-backward trajectories. However, these methods rely on empirical metrics to identify the target trajectories. In addition, repeatedly performing forward and backward trackings brings in a heavy computational cost for online tracking and largely hurts the real-time performance. Differently, in TrackingNet [45], forward-backward analysis is used for evaluating the tracking performance and annotating a sparsely labeled dataset such as Youtube-BoundingBox [49] to obtain the per-frame object bounding box labels. In this work, we target at visual tracking but revisit the forward-backward scheme from a different view, i.e., training a deep visual tracker in under unsupervised manner.

### 2.3 Unsupervised Representation Learning

Our tracking framework relates to unsupervised representation learning. Learning feature representations from raw videos under an unsupervised manner has gained increasing attention in recent years. These approaches typically

**(a) Unsupervised Learning Motivation**   **(b) Unsupervised Learning Pipeline using a Siamese Network**

**Fig. 2** An overview of unsupervised learning in deep tracking. We show our motivation in (a) that we track forward and backward to compute the consistency loss for network training. The detailed training procedure is shown in (b), where unsupervised learning is integrated into a Siamese correlation filter network. In the testing stage, we only track forward to predict the target location.

design ingenious techniques to explore and utilize the free supervision in images or videos. In [31], the feature representation is learned by shuffling the video frames and then sorting them again to achieve self-supervised training. The multi-layer autoencoder on large-scale unlabeled data has been explored in [29]. Vondrick et al. [59] proposed to anticipate the visual representation of frames in the future. In [60], Vondrick et al. colorized gray-scale videos by copying colors from a reference frame to learn a CNN model. Wang and Gupta [66] used the KCF tracker [19] to pre-process the raw videos, and then selected a pair of tracked images together with another random patch for learning CNNs using a ranking loss. Our method differs from [66] significantly in two aspects. First, we integrate the tracking algorithm into unsupervised training instead of merely utilizing an off-the-shelf tracker as the data pre-processing tool. Second, our unsupervised framework is coupled with a tracking objective function, so the learned feature representation is effective in characterizing the generic target objects.

In the visual tracking community, unsupervised learning has rarely been touched. According to our knowledge, the only related but different approach is the autoencoder based method [62]. However, the encoder-decoder is a general unsupervised framework [47], whereas our unsupervised method is specially designed for the tracking task. Since the visual objects or scenes in videos typically change smoothly, the motion information of the objects in a forward-backward trajectory loop provides free yet informative self-supervision signals for unsupervised learning, which is naturally suitable for the motion-related visual tracking task.

## 3 Proposed Method

The motivation of our unsupervised learning is shown in Fig. 2(a). We first select a content-rich local region as the target object. Given this initialized bounding box label, we track forward to predict its location in the subsequent frames. Then, we reverse the sequence and take the predicted bounding box in the last frame as the pseudo label for backward verification. The predicted bounding box in the first frame via backward tracking is ideally identical to the original bounding box. We measure the difference between the forward and backward trajectories using the consistency loss to train the network. Fig. 2(b) shows an overview of our unsupervised Siamese correlation filter network.

In the following, we first revisit the correlation filter as well as the Siamese network. In Section 3.2, we present our unsupervised learning prototype for an intuitive understanding. In Section 3.3, we improve our prototype to facilitate unsupervised training. Finally, training details and online tracking are elaborated in Sections 3.4 and 3.5, respectively.

### 3.1 Revisiting Correlation Tracking

The Discriminative Correlation Filters (DCFs) [2,19] regress the circularly shifted versions of the input features of a search patch to a soft target response map for target localization. When training a DCF, we select a template patch $\mathbf{X}$ with the corresponding ground-truth label $\mathbf{Y}$, which is Gaussian-shaped with the peak localized at the target position. The size of the template patch is usually larger than that of the target. Fig. 2 shows an example of the template patch, where there are both target and background contents.

The filter $\mathbf{W}$ can be learned by solving the following ridge regression problem:

$$\min_{\mathbf{W}} \|\mathbf{W} * \mathbf{X} - \mathbf{Y}\|_2^2 + \lambda \|\mathbf{W}\|_2^2, \tag{1}$$

where $\lambda$ is a regularization parameter and $*$ denotes the circular convolution. Eq. 1 can be efficiently calculated in the Fourier domain [2,10,19] and the DCF can be computed by

$$\mathbf{W} = \mathscr{F}^{-1} \left( \frac{\mathscr{F}(\mathbf{X}) \odot \mathscr{F}^{\star}(\mathbf{Y})}{\mathscr{F}^{\star}(\mathbf{X}) \odot \mathscr{F}(\mathbf{X}) + \lambda} \right), \tag{2}$$

where $\odot$ is the element-wise product, $\mathscr{F}(\cdot)$ is the Discrete Fourier Transform (DFT), $\mathscr{F}^{-1}(\cdot)$ is the inverse DFT, and $\star$ denotes the complex-conjugate operation. In each subsequent frame, given a search patch $\mathbf{Z}$, its corresponding response map $\mathbf{R}$ can be computed in the Fourier domain:

$$\mathbf{R} = \mathbf{W} * \mathbf{Z} = \mathscr{F}^{-1} \left( \mathscr{F}^{\star}(\mathbf{W}) \odot \mathscr{F}(\mathbf{Z}) \right). \tag{3}$$

The above DCF framework starts from learning the target template's correlation filter (i.e., $\mathbf{W}$) using the template patch and then convolves it with a search patch $\mathbf{Z}$ to generate the response. Recently, the Siamese correlation filter network [58,64] embeds the DCF in the Siamese framework and constructs two shared-weight branches to extract feature representations, as shown in Fig. 2(b). The first one is the template branch which takes a template patch $\mathbf{X}$ as input and extracts its features to further generate a target template filter via DCF. The second one is the search branch which takes a search patch $\mathbf{Z}$ as input for feature extraction. The template filter is then convolved with the CNN features of the search patch to generate the response map. The advantage of the Siamese DCF network is that both the feature extraction CNN and correlation filter are formulated into an end-to-end framework, so the learned features are more related to the visual tracking scenarios.

## 3.2 Unsupervised Learning Prototype

Given two consecutive frames $P_1$ and $P_2$, we crop the template and search patches from them, respectively. By conducting forward tracking and backward verification, the proposed framework does not require additional supervision. The location difference between the initial bounding box and the predicted bounding box in $P_1$ will formulate a consistency loss. We utilize this loss to train the network without ground-truth annotations.

### 3.2.1 Forward Tracking

Following the previous approaches [58,64], we build a Siamese correlation filter network to track the initialized

bounding box region in frame $P_1$. After generating the template patch $\mathbf{T}$ from the first frame $P_1$, we compute the corresponding template filter $\mathbf{W_T}$ as follows:

$$\mathbf{W_T} = \mathscr{F}^{-1} \left( \frac{\mathscr{F}(\varphi_\theta(\mathbf{T})) \odot \mathscr{F}^{\star}(\mathbf{Y_T})}{\mathscr{F}^{\star}(\varphi_\theta(\mathbf{T})) \odot \mathscr{F}(\varphi_\theta(\mathbf{T})) + \lambda} \right), \tag{4}$$

where $\varphi_\theta(\cdot)$ denotes the CNN feature extraction operation with trainable network parameters $\theta$, and $\mathbf{Y_T}$ is the label of the template patch $\mathbf{T}$. This label is a Gaussian response centered at the initialized bounding box center. Once we obtain the learned template filter $\mathbf{W_T}$, the response map of a search patch $\mathbf{S}$ from frame $P_2$ can be computed by

$$\mathbf{R_S} = \mathscr{F}^{-1}(\mathscr{F}^{\star}(\mathbf{W_T}) \odot \mathscr{F}(\varphi_\theta(\mathbf{S}))). \tag{5}$$

If the ground-truth Gaussian label of patch $\mathbf{S}$ is available, the network $\varphi_\theta(\cdot)$ can be trained by computing the $L_2$ distance between $\mathbf{R_S}$ and the ground-truth label. Different from the supervised framework, in the following, we present how to train the network without requiring labels by exploiting backward trajectory verification.

### 3.2.2 Backward Tracking

After generating the response map $\mathbf{R_S}$ for frame $P_2$, we create a pseudo Gaussian label centered at its maximum value, which is denoted by $\mathbf{Y_S}$. In backward tracking, we switch the role between the search patch and the template patch. By treating $\mathbf{S}$ as the template patch, we generate a template filter $\mathbf{W_S}$ using the pseudo label $\mathbf{Y_S}$. The template filter $\mathbf{W_S}$ can be learned using Eq. 4 by replacing $\mathbf{T}$ with $\mathbf{S}$ and replacing $\mathbf{Y_T}$ with $\mathbf{Y_S}$, as follows:

$$\mathbf{W_S} = \mathscr{F}^{-1} \left( \frac{\mathscr{F}(\varphi_\theta(\mathbf{S})) \odot \mathscr{F}^{\star}(\mathbf{Y_S})}{\mathscr{F}^{\star}(\varphi_\theta(\mathbf{S})) \odot \mathscr{F}(\varphi_\theta(\mathbf{S})) + \lambda} \right). \tag{6}$$

Then, we generate the response map $\mathbf{R_T}$ of the template patch through Eq. 5 by replacing $\mathbf{W_T}$ with $\mathbf{W_S}$ and replacing $\mathbf{S}$ with $\mathbf{T}$, as shown in Eq. 7.

$$\mathbf{R_T} = \mathscr{F}^{-1}(\mathscr{F}^{\star}(\mathbf{W_S}) \odot \mathscr{F}(\varphi_\theta(\mathbf{T}))). \tag{7}$$

Note that we only use one Siamese correlation filter network for executing forward and backward trackings. The network parameters $\theta$ are fixed during the tracking steps.

### 3.2.3 Consistency Loss Computation

After forward and backward tracking, we obtain the response map $\mathbf{R_T}$. Ideally, $\mathbf{R_T}$ should be a Gaussian label with the peak located at the initialized target position. In other words, $\mathbf{R_T}$ should be as similar as the originally given label $\mathbf{Y_T}$. Therefore, the representation network $\varphi_\theta(\cdot)$ can be trained under an unsupervised manner by minimizing the reconstruction error as follows:

$$\mathcal{L}_{\text{un}} = \|\mathbf{R_T} - \mathbf{Y_T}\|_2^2. \tag{8}$$

**Forward Stage:** Data labeling using the tracking model



**Backward Stage:** Tracking model update using labeled data

**Fig. 3** The intuition of pseudo-labeling based self-training. We use the same network for both forward and backward predictions. The forward stage generates a pseudo label for the search patch. The backward stage updates the tracking network using training pairs via loss back-propagation. During training iterations, the response map of the template gradually approaches the initial label via self supervision.

Our unsupervised learning can be viewed as an incremental self-training process that iteratively predicts labels and updates the model to steadily improve the tracking capability. Fig. 3 shows the intuition, where we use the same network for both forward and backward predictions. In the forward tracking, we generate a pseudo label $\mathbf{Y_S}$ for the search patch $\mathbf{S}$. Then we treat generated $\mathbf{Y_S}$ as the label of $\mathbf{S}$ and create a corresponding sample. Using these labeled training pairs (i.e., with initial or pseudo labels), we can update the Siamese correlation filter network in a similar way to supervised learning. During loss back-propagation, we follow the Siamese correlation filter methods [64,74] to update the network:

$$\frac{\partial \mathcal{L}_{\mathrm{un}}}{\partial \varphi_\theta(\mathbf{T})} = \mathscr{F}^{-1}\left(\frac{\partial \mathcal{L}_{\mathrm{un}}}{\partial \left(\mathscr{F}\left(\varphi_\theta(\mathbf{T})\right)\right)^\star} + \left(\frac{\partial \mathcal{L}_{\mathrm{un}}}{\partial \left(\mathscr{F}\left(\varphi_\theta(\mathbf{T})\right)\right)}\right)^\star\right),$$

$$\frac{\partial \mathcal{L}_{\mathrm{un}}}{\partial \varphi_\theta(\mathbf{S})} = \mathscr{F}^{-1}\left(\frac{\partial \mathcal{L}_{\mathrm{un}}}{\partial \left(\mathscr{F}\left(\varphi_\theta(\mathbf{S})\right)\right)^\star}\right).$$

$$(9)$$

The above unsupervised training process is based on the forward-backward consistency between two frames, which is summarized by Algorithm 1. In the next section, we extend this prototype framework to consider multiple frames for better network training.

### 3.3 Enhancement for Unsupervised Learning

The proposed unsupervised learning method constructs the objective function based on the consistency between $\mathbf{R_T}$ and $\mathbf{Y_T}$. In practice, the tracker may deviate from the target in the forward tracking but still return to the original position during the backward process. However, the proposed

---

**Algorithm 1:** Unsupervised training prototype

**Input:** Unlabeled videos.
**Output:** Pretrained tracking network $\varphi_\theta(\cdot)$.
1  Crop the patches (i.e., $\mathbf{T}$ and $\mathbf{S}$) from the raw videos;
2  Initialize the CNN model $\varphi_\theta(\cdot)$ with random weights $\theta$;
3  **for** *each training epoch* **do**
4      **for** *each training pair* **do**
5          Obtain $\varphi_\theta(\mathbf{T})$ and $\varphi_\theta(\mathbf{S})$;
6          // Forward Trajectory
7          Construct $\mathbf{W_T}$ using $\varphi_\theta(\mathbf{T})$ and $\mathbf{Y_T}$ (Eq. 4);
8          Compute $\mathbf{R_S}$ using $\mathbf{W_T}$ (Eq. 5) and obtain the pseudo label of $\mathbf{S}$;
9          // Backward Trajectory
10         Construct $\mathbf{W_S}$ using $\varphi_\theta(\mathbf{S})$ and $\mathbf{Y_S}$ (Eq. 6);
11         Compute the response map $\mathbf{R_T}$ of $\mathbf{T}$ (Eq. 7);
12         // Calculate Consistency Loss
13         Compute the consistency loss of $\mathbf{Y_T}$ and $\mathbf{R_T}$ (Eq. 8);
14     **end**
15     Update network $\varphi_\theta(\cdot)$ using the computed loss;
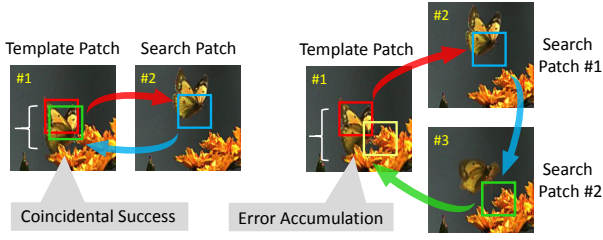16 **end**

---

loss function does not penalize this deviation because of the consistent trajectories. Meanwhile, the raw videos may contain textureless or occluded training samples that deteriorate the unsupervised learning process. In this section, we propose a multi-frame validation scheme and a cost-sensitive loss to tackle these two limitations.

#### 3.3.1 Multi-frame Validation

We propose a multi-frame validation approach to enlarge the trajectory inconsistency when the tracker loses the target. Our intuition is to incorporate more frames during training to reduce the limitation that the erroneous localization in the subsequent frame successfully backtraces to the initial position in the first frame. In this way, the reconstruction error in Eq. 8 will effectively capture the inconsistent trajectory. As shown in Fig. 3, adding more frames in the forward stage further challenges the model tracking capability.

Our unsupervised learning prototype can be easily extended to multiple frames. To build a trajectory cycle using three frames, we can involve another frame $P_3$ which is the subsequent frame after $P_2$. We crop a search patch $\mathbf{S}_1$ from $P_2$ and another search patch $\mathbf{S}_2$ from $P_3$. If the generated response map $\mathbf{R}_{\mathbf{S}_1}$ is different from its corresponding ground-truth response, the difference tends to become larger in the next frame $P_3$. As a result, the inconsistency is more likely to appear in backward tracking, and the generated response map $\mathbf{R_T}$ is more likely to differ from $\mathbf{Y_T}$, as shown in Fig. 4. By involving more search patches during forward and backward trackings, the proposed consistency loss will be more effective to penalize the inaccurate localizations.

We can further extend the number of frames utilized for multi-frame validation. The length of trajectory will increase as shown in Fig. 5. The limitation of consistent trajectory

**Fig. 4** Single frame validation and multi-frame validation. The inaccurate localization in single frame validation may not be captured as shown on the left. By involving more frames as shown on the right, we accumulate the localization errors to break the prediction consistency during forward and backward trackings.



**Fig. 5** An overview of multi-frame trajectory consistency. We denote $\mathbf{T}$ as a template and $\mathbf{S}$ as a search patch, respectively. Our unsupervised training prototype is shown in (a), where only two frames are involved. Using more frames as shown in (b) and (c), we can gradually improve the training performance to overcome consistent trajectories when losing the target.

when losing the target is more unlikely to affect the training process. Let $\mathbf{R}_{(\mathbf{S}_k \to \mathbf{T})}$ denote the response map of the template $\mathbf{T}$, which is generated (or tracked) by the DCF trained using the $k$-th search patch $\mathbf{S}_k$. The corresponding consistency loss function can be computed as follows:

$$\mathcal{L}_k = \left\| \mathbf{R}_{(\mathbf{S}_k \to \mathbf{T})} - \mathbf{Y}_{\mathbf{T}} \right\|_2^2 . \tag{10}$$

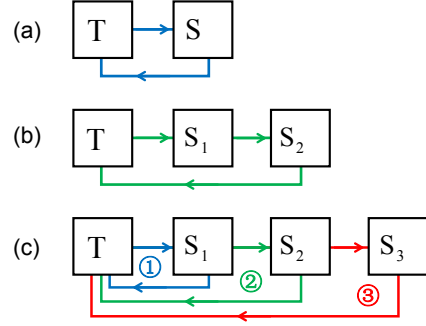Considering different trajectory cycles, the multi-frame consistency loss can be computed by

$$\mathcal{L}_{\text{un}} = \sum_{k=1}^{M} \mathcal{L}_k, \tag{11}$$

where $k$ is the index of the search parch. Taking Fig. 5(c) as an example, the final consistency objective contains three losses (i.e., $M = 3$ in Eq. 11), which are denoted by the blue, green, and red cycles in Fig. 5(c), respectively.

### 3.3.2 Cost-sensitive Loss

We initialize a bounding box region as a training sample in the first frame during unsupervised training. The image content within this bounding box region may contain arbitrary or partial objects. Fig. 6 shows an overview of these regions. To alleviate the background interference, we propose a cost-sensitive loss to effectively exclude noisy samples for network training. For simplicity, we use three consecutive frames as an example to illustrate sample selection, which can be naturally extended to more frames. The pipeline of using three frames is shown in Fig. 5(b).

During unsupervised learning, we construct multiple training triples from video sequences. For a trajectory containing three frames, each training triple consists of one initialized template patch $\mathbf{T}$ in frame $P_1$ and two search patches $\mathbf{S}_1$ and $\mathbf{S}_2$ in the subsequent frames $P_2$ and $P_3$, respectively. We use several triples to form a training batch for Siamese network learning. In practice, we find that some training triples with extremely high losses prevent network training from convergence. To reduce these outlier effects in

pseudo-labeling based self-training, we exclude 10% of the whole training triples which contain the highest loss values. Their losses can be computed using Eq. 10. To this end, we assign a binary weight $\mathbf{A}_{\text{drop}}^i$ to each training triple. All these weights constitute a vector $\mathbf{A}_{\text{drop}}$, where 10% of its elements are 0 and the others are 1.

In addition to the outlier training pairs, the raw videos include meaningless image patches, where there are textureless backgrounds or still objects. In these patches, the objects (e.g., sky, grass, or tree) do not contain big movements. We assign a motion weight vector $\mathbf{A}_{\text{motion}}$ to all the training pairs to increase the large motion effect for network learning. Each element $\mathbf{A}_{\text{motion}}^i$ within this vector can be computed by

$$\mathbf{A}_{\text{motion}}^i = \left\| \mathbf{R}_{\mathbf{S}_1}^i - \mathbf{Y}_{\mathbf{T}}^i \right\|_2^2 + \left\| \mathbf{R}_{\mathbf{S}_2}^i - \mathbf{Y}_{\mathbf{S}_1}^i \right\|_2^2 , \tag{12}$$

where $\mathbf{R}_{\mathbf{S}_1}^i$ and $\mathbf{R}_{\mathbf{S}_2}^i$ are the response maps in the $i$-th training pair, and $\mathbf{Y}_{\mathbf{T}}^i$ and $\mathbf{Y}_{\mathbf{S}_1}^i$ are the corresponding initial (or pseudo) labels. Eq. 12 calculates the target motion difference from frame $P_1$ to $P_2$ and $P_2$ to $P_3$. When the value of $\mathbf{A}_{\text{motion}}^i$ is large, the target object undergoes fast motion in this trajectory. On the other hand, the large value of $\mathbf{A}_{\text{motion}}^i$ represents the hard training pair which the network should pay more attention to. We normalize the motion weight and the binary weight as follows:

$$\mathbf{A}_{\text{norm}}^i = \frac{\mathbf{A}_{\text{drop}}^i \cdot \mathbf{A}_{\text{motion}}^i}{\sum_{i=1}^{N} \mathbf{A}_{\text{drop}}^i \cdot \mathbf{A}_{\text{motion}}^i}, \tag{13}$$

where $N$ is number of the training pairs in a mini-batch. The sample weight $\mathbf{A}_{\text{norm}}^i$ serves as a scalar that reweighs the training data without gradient back-propagation.

The final unsupervised loss for the case of Fig. 5(b) in a mini-batch is computed as:

$$\mathcal{L}_{\text{3-frame}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{A}_{\text{norm}}^i \cdot \left\| \mathbf{R}_{(\mathbf{S}_2 \to \mathbf{T})}^i - \mathbf{Y}_{\mathbf{T}}^i \right\|_2^2 . \tag{14}$$

We can naturally extend Eq. 14 to the following by using more frames to construct trajectories of different lengths, as illustrated by the toy example of Fig. 5(c). Combining with Eq. 11, we compute the final unsupervised loss function using $M$ subsequent frames as:

$$\mathcal{L}_{\text{final}} = \frac{1}{N} \sum_{k=1}^{M} \sum_{i=1}^{N} \mathbf{A}_{\text{norm}}^i \cdot \mathcal{L}_k^i, \tag{15}$$

where $\mathcal{L}_k^i = \left\| \mathbf{R}_{(\mathbf{S}_k \to \mathbf{T})}^i - \mathbf{Y}_{\mathbf{T}}^i \right\|_2^2$ is similar to that in Eq. 10 but with the index $i$ for differentiating different samples in a mini-batch.

### 3.4 Unsupervised Training Details

**Network Structure.** We follow the DCFNet [64] to use a shallow Siamese network consisting of two convolutional layers for tracking. This shallow structure is demonstrated effective in CFNet [58] to integrate the DCF formulation. The filter sizes of these convolutional layers are $3 \times 3 \times 3 \times 32$ and $3 \times 3 \times 32 \times 32$, respectively. Besides, a local response normalization (LRN) layer is employed at the end of convolutional layers following [64]. This lightweight structure enables efficient forward inferences for online tracking.

**Training Data.** We choose ILSVRC 2015 [51] as our training data, which is the same dataset employed by existing supervised trackers. In the data pre-processing step, supervised approaches [1,58,64] require per-frame labels. Besides, the frames will be removed, where the target object is occluded, partially out-of-view, or in an irregular shape (e.g., snake). The data pre-precessing for the supervised approaches is time-consuming with human labor. In contrast, our method does not rely on manually annotated labels for data pre-processing.

In our approach, for the first frame in a raw video, we crop overlapped small patches ($5 \times 5$ in total) by sliding windows as shown in Fig. 7. Then, we compute the image entropy of each image patch. Image entropy effectively measures the content variance of an image patch. When an image patch only contains the unitary texture (e.g., the sky), the entropy of this patch approaches 0. When an image patch contains textured content, the entropy will become higher. We select the cropped image patch containing the highest image entropy. This image patch initializes the KCF [19] tracker for localization in the subsequent frames. Then, we crop a larger image patch with a padding of 2 times of the target size following DCFNet [64], which is further resized to $125 \times 125$ as the input of our network. Fig. 6 exhibits some examples of the cropped patches. We randomly choose 4 cropped patches from the continuous 10 frames in a video to form a training trajectory, and one of them is defined as the template and the rest as search patches. This is based on the



**Fig. 6** Examples of the cropped image patches from ILSVRC 2015 [51]. Most of these samples contain meaningful objects, while some samples are less meaningful (e.g., last row).



**Fig. 7** The illustration of training samples generation. The proposed method crops $5 \times 5$ patch candidates in the center region of the initial frame. Then we select the image patch with the highest image entropy. As shown in the right figure, the background patches (e.g., labeled by green and red boxes) have a small image entropy.

assumption that the center located target objects are unlikely to move out of the cropped region in a short span of time. We track the content in the image patch regardless of specific object categories. Although this entropy-based method may not accurately select a target region and the KCF tracker is not robust enough to track the cropped region, this method can well alleviate the meaningless background regions.

### 3.5 Online Object Tracking

After offline unsupervised learning, we perform online tracking in the way of forward tracking as illustrated in Section 3.2. We online update DCF to adapt to the target appearance changes. The DCF update follows a moving average operation shown as follows:

$$\mathbf{W}_t = (1 - \alpha_t)\mathbf{W}_{t-1} + \alpha_t \mathbf{W}, \tag{16}$$

where $\alpha_t \in [0, 1]$ is the linear interpolation coefficient. The target scale is estimated through a patch pyramid with scale factors $\{a^s | a = 1.015, s = \{-1, 0, 1\}\}$ [12]. We name our Tracker as LUDT (i.e., Learning Unsupervised Deep Tracking). Besides, we update our model adaptively via $\alpha_t$ and follow the superior DCF formulation as that in ECO [9]. We name the improved tracker as LUDT+.

We keep the notation of our preliminary tracker UDT and UDT+ [61] in the following experiment section. Our previous UDT uses a 3-frame cycle (Fig. 5(b)) and simply crops the center patch in raw videos. LUDT improves UDT in two aspects: (1) LUDT combines different trajectory cycles as shown in Fig. 5(c), and (2) LUDT utilizes image entropy to select the informative image patches instead of the center crop. LUDT+ and UDT+ improve LUDT and UDT by adopting some online tracking techniques (e.g., adaptive update) proposed in [9], respectively.

## 4 Experiments

In this section, we first analyze the effectiveness of our unsupervised training framework and discuss our network potentials. Then, we compare our tracker LUDT against state-of-the-art trackers on both standard and recently released large-scale benchmarks including OTB-2013 [69], OTB-2015 [70], Temple-Color [35], VOT2016 [25], VOT2017/2018 [24], LaSOT [16], and TrackingNet [45].

### 4.1 Experimental Details

In our experiments, we use the stochastic gradient descent (SGD) with a momentum of 0.9 and a weight decay of 0.005 to train our model. Our unsupervised network is trained for 50 epochs with a learning rate exponentially decaying from $10^{-2}$ to $10^{-5}$ and a mini-batch size of 32. We set the trajectory length as 4. All the experiments are executed on a PC with 4.00GHz Intel Core I7-4790K and NVIDIA GTX 1080Ti GPU. On a single GPU, our LUDT and LUDT+ exhibit about 70 FPS and 55 FPS, respectively[2].

The proposed method is evaluated on seven benchmarks. On the OTB-2013/2015, TempleColor, LaSOT, and TrackingNet datasets, we use one-pass evaluation (OPE) with distance and overlap precision metrics. The distance precision threshold is set as 20 pixels. The overlap success plot uses thresholds ranging from 0 to 1, and the area-under-curve (AUC) is computed to evaluate the overall performance. On the VOT2016 and VOT2017/2018 datasets, we measure the performance using Expected Average Overlap (EAO).

### 4.2 Ablation Experiments and Discussions

#### 4.2.1 Improvements upon UDT

Our preliminary tracker UDT [61] adopts a three-frame validation (i.e., Fig. 5(b)) and the center crop for sample generation. The improvement upon UDT is that we construct

---

[2] The source code is provided at https://github.com/594422814/UDT

**Table 1** Ablation study of Trajectory Enlargement (TE) and RoI Selection (RS). We denote UDT as our preliminary tracker [61]. We integrate TE and RS into UDT during training and report the performance improvement. The evaluation metrics are DP and AUC scores on the OTB-2015 and Temple-Color datasets.

|  | OTB-2015 DP / AUC (%) | Temple-Color DP / AUC (%) |
|---|---|---|
| UDT [61] | 76.0 / 59.4 | 65.8 / 50.7 |
| UDT + TE | 76.5 / 59.8 | 66.7 / 51.2 |
| UDT + RS | 76.5 / 60.0 | 66.9 / 51.3 |
| UDT + TE + RS | 76.9 / 60.2 | 67.1 / 51.5 |

**Table 2** Comparison results of the DCFNet tracking framework with different feature extractors. Random: the randomly initialized feature extractor without pre-training. HOG: adopting HOG [8] without deep features. ED: the backbone network trained via encoder-decoder [62]. The evaluation metrics are DP and AUC scores on OTB-2015.

|  | Random | HOG | ED | LUDT (Ours) |
|---|---|---|---|---|
| DP (%) | 59.1 | 69.2 | 71.6 | 76.9 |
| AUC (%) | 46.9 | 52.1 | 54.5 | 60.2 |

a multi-supervision consistency loss function using more frames. We denote this strategy as Trajectory Enlargement (TE) in Table 1. Meanwhile, we select the RoI (region of interest) from raw videos using the image entropy and KCF tracker, while only the center region is utilized in UDT. We denote RoI Selection as RS in the table. Note that the performance of UDT has been close to that of its supervised configuration and exceeded several supervised trackers. Moreover, under the same training configuration, LUDT steadily improves UDT by using TE and RS during training. LUDT achieves 60.2% and 51.5% in AUC on the OTB-2015 and Temple-Color benchmarks, respectively.

#### 4.2.2 Baseline Performance

To verify the effectiveness of the proposed unsupervised framework, we evaluate our tracker using different feature extractors. As shown in Table 2, without pre-training, the model still exhibits a weak tracking capability, which can be attributed to the discriminating power of the correlation filter. By adopting the empirical HOG representations, the performance is still significantly lower than ours. Furthermore, we leverage the auto-encoder framework [62] to train the backbone network under an unsupervised manner using the same training data. From Table 2, we can observe that our approach is superior to the encoder-decoder in this tracking scenario since our forward-backward based unsupervised training is tightly related to object tracking.

#### 4.2.3 Training Data

We evaluate the tracking performance using different data pre-processing strategies. The results are shown in Table 3. Our unsupervised LUDT method uses the last RoI selection

**Table 3** Evaluation results of our network trained using different data pre-processing strategies. Our LUDT tracker uses RoI selection via image entropy for unsupervised training. The evaluation metrics are DP and AUC scores on the OTB-2015 dataset.

|        | Groundtruth label | Groundtruth label with deviations | Center cropping | RoI Selection via entropy |
|--------|-------------------|-----------------------------------|-----------------|---------------------------|
|        | Full supervision  | Weak supervision                  | Unsupervision   | Unsupervision             |
| DP (%) | 80.6              | 78.9                              | 76.0            | 76.5                      |
| AUC (%)| 62.6              | 61.4                              | 59.4            | 60.0                      |

**Table 4** Evaluation results of our unsupervised model trained using different trajectory lengths. Note that the 4 and 5 frames validations conduct multiple self-supervisions as illustrated in Fig. 5. The evaluation metrics are DP and AUC scores on the OTB-2015 dataset. Compared with 3 frames validation, using more frames further improves the tracking accuracy.

| Frame Number | 2 frames Fig. 5(a) | 3 frames Fig. 5(b) | 4 frames Fig. 5(c) | 5 frames akin to Fig. 5(c) |
|--------------|--------------------|--------------------|--------------------|----------------------------|
| DP (%)       | 73.2               | 76.0               | 76.8               | 76.8                       |
| AUC (%)      | 57.4               | 59.4               | 59.8               | 59.7                       |

**Table 5** A performance study by adding additional training data. Adding more unlabeled data steadily improves the tracking results. The evaluation metrics are DP and AUC scores on the OTB-2015 dataset.

|         | LUDT | Few-shot fine-tune OTB-2015 | More data OxUvA | More data LaSOT |
|---------|------|-----------------------------|-----------------|-----------------|
| DP (%)  | 76.9 | 78.1                        | 77.6            | 78.2            |
| AUC (%) | 60.2 | 61.5                        | 61.4            | 62.0            |

strategy. During the evaluation, we keep the remaining modules fixed in LUDT.

**Comparison with Full Supervision**. Using the same videos (i.e., ILSVRC 2015 [51]), we conduct the supervised training of our network. The supervised learning with ground-truth annotations can be regarded as the upper bound of our unsupervised learning. We observe that the performance gap is small (i.e., 2.6% AUC) between the trackers trained using unsupervised learning (60.0% AUC) and fully supervised learning (62.6% AUC).

**Comparison with Weak Supervision**. On ILSVRC 2015, we add deviations to the ground-truth boxes to crop the training samples. The deviations range from -20 pixels to 20 pixels randomly. The reason for setting sample deviations from the ground-truth bounding boxes is that we aim to simulate the inaccurate object localizations on in-the-wild videos using existing object detection or optical flow approaches. We assume that these deviated samples are predicted by existing methods and then utilized to train our unsupervised network. In Table 3, we observe that our tracker learned by these weakly labeled samples is comparable with the supervised results (61.4% vs. 62.6% AUC). Note that 20 pixels deviations can be achieved with many object localization methods. The comparable performance indicates that our method can be applied to raw videos with weakly or sparsely labeled annotations (e.g., the dataset Youtube-BB [49]). On the other hand, existing object detectors and models are mostly trained by supervised learning. To ensure our method to be fully unsupervised, we use two unsupervised data pre-processing methods: center cropping and RoI selection based on entropy.

**Center Cropping**. In center cropping, we crop the center region of the video frame. Although we crop a fixed region of the image, the image content appears randomly in this region and we denote this operation as center cropping. There may be meaningless content (e.g., textureless sky and ocean) in this region to disturb our unsupervised learning. The tracker learned by center cropping achieves an AUC score of 59.4%.

**RoI Selection**. We use the entropy-based image patch selection as illustrated in Section 3.4. Compared to the center cropping, image-entropy based selection can suppress the meaningless background samples such as sky and grass, and the KCF tracker is able to capture the selected informative region in the subsequent frames. The RoI selection achieves

better performance than center cropping with an AUC score of 60.0%.

### 4.2.4 Trajectory Length

As discussed in Section 3.3.1, trajectory enlargement helps measure the consistency loss when the tracker loses RoI. In Table 4, we show the performance with different trajectory lengths on the OTB-2015 dataset. We use center cropping to generate training samples following UDT for comparison. The prototype of our unsupervised learning is denoted as 2 frames validation. By incorporating the third frame, the learned tracker achieves improvement (i.e., 2.8% DP and 2.0% AUC). The 4 frames validation proposed in this work not only extends the trajectory length but also combines multiple self-supervision constraints, which further improves the accuracy. However, the 5 frames validation seems to be less effective. It may be because the validation with 4 frames already contains adequate self-supervision and effectively measures the consistency loss.

### 4.2.5 Cost-sensitive Loss

On the OTB-2015 dataset, without hard sample reweighing (i.e., $A_{motion}$ in Eq. 13), the performance of our LUDT tracker drops about 1.5% DP and 1% AUC scores. We did not conduct the ablation study of the sample dropout because we observe that the unsupervised training cannot well converge without $A_{drop}$ illustrated in Eq. 14.

### 4.2.6 Unlabled Data Augmentation

**Few-shot Domain Adaptation**. To better fit a new domain such as OTB, we construct a small training set by collecting the first several frames (e.g., 5 frames in our experi-

**Table 6** A performance potential of our unsupervised tracker. When using more data (LaSOT) for network training, the performance is further improved. By incorporating empirical features (HOG), our unsupervised tracker achieves superior results. The performance is evaluated on the OTB-2015 dataset using DP and AUC metrics.
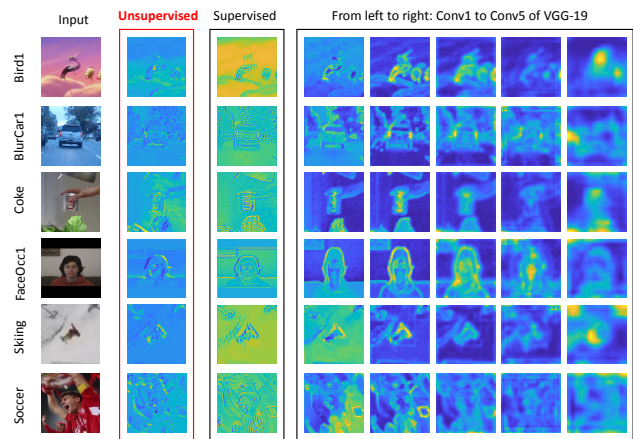
| | ECOhc | LUDT+ only ILSVRC | LUDT+ more data | LUDT+ more data + HOG |
|---|---|---|---|---|
| DP (%) | 85.4 | 84.3 | 85.5 | 85.8 |
| AUC (%) | 64.1 | 63.9 | 64.9 | 65.7 |
| Speed (FPS) | 60 | 55 | 55 | 42 |

ment) from the videos in OTB-2015 with only the ground-truth bounding box in the first frame available. Using these limited samples, we fine-tune our network by 100 iterations using the forward-backward pipeline, which takes about 6 minutes. As our learning method is unsupervised, we can utilize the frames from test sequences to adapt our tracker. Table 5 shows that performance is further improved by using this strategy. Our offline unsupervised training learns general feature representation, which can be transferred to an interested domain (e.g., OTB videos) using few-shot domain adaptation. This domain adaptation is similar to that in MD-Net [46], while our network parameters are initially offline learned in an unsupervised manner.

**Additional Internet Videos**. We also utilize more unlabeled videos to train our network. These videos are from the Ox-UvA dataset [57], where there are 337 videos in total. The OxUvA dataset is a subset of Youtube-BB [49] collected on YouTube. By adding these videos during training, our tracker improves the original one by 0.7% DP and 1.2% AUC as shown in Table 5. By leveraging another large-scale LaSOT dataset [16] where there are 1200 videos collected on the Internet, the tracking performance is further improved. It indicates that unlabeled data advances the unsupervised training. As our framework is fully unsupervised, it has the potential to take advantage of the in-the-wild videos on the Internet to boost the performance.

### 4.2.7 Empirical Features Embedding

As shown in Table 6, we train unsupervised LUDT+ using more unlabeled video sequences (both ILSVRC and LaSOT), which outperforms ECOhc [9] leveraging hand-crafted features including HOG [8] and ColorName [68]. In addition, we can combine the learned CNN features and empirical features to generate a more discriminative representation. We add the HOG feature to LUDT+ during tracking and evaluate its performance. Table 6 shows that this combination achieves a 65.7% AUC on OTB-2015. Moreover, embedding the HOG feature helps LUDT+ to outperform most state-of-the-art real-time trackers as shown in Table 7. Besides feature embedding and adaptive model update, there are still many improvements from [65,17,44,39] available to benefit our tracker. However, adding more addi-



**Fig. 8** Visualization of feature representations. First column: input image patches. Second and third columns: feature representations of our unsupervised LUDT and fully-supervised LUDT. The rest columns: feature maps from VGG-19 (from left to right: Conv1-2, Conv2-2, Conv3-4, Conv4-4, and Conv5-4). The feature map is visualized by averaging all the channels. Best viewed in color and zoom in.

tional mechanisms is out of the scope of this work. Following SiamFC and DCFNet, we currently use LUDT/LUDT+ trackers which are *only* trained on the ILSVRC dataset for *fair comparison* in the following evaluations.

### 4.3 Visualization of Unsupervised Representation

After learning the unsupervised Siamese tracking network, we visualize the network response to see how it differs from the same network trained using supervised learning. Fig. 8 shows the visualization performance. The first column shows the input frames. The network responses from unsupervised learning and supervised learning are shown in the second and third columns, respectively. The remaining columns show the feature responses from the off-the-shelf deep model VGG-19 [52]. We observe that the responses from the unsupervised learning and supervised learning are similar with minor differences. Specifically, the boundary responses from supervised learning are higher than those of unsupervised learning. This is because of the strong supervisions brought by the ground-truth labels. The network has learned to differentiate the target and background distractors according to labels, which increases the network attention around the object boundaries. In comparison, our unsupervised learning does not employ this process for attention enhancement, while still focusing on the center region of the object responses. From the viewpoint of Siamese network, both unsupervised and supervised feature representations focus on the target appearances, which facilitate the template matching through the correlation operation. Compared with the empirical features (e.g., HOG), we will show in the following that our unsupervised feature representa-

**Table 7** Evaluations with fully-supervised baseline (left) and state-of-the-art trackers (right) on the popular OTB-2015 benchmark [70]. The evaluation metrics are DP and AUC scores. Our unsupervised LUDT tracker performs favorably against popular baseline methods (left), while our LUDT+ tracker achieves comparable results with the recent state-of-the-art supervised trackers (right).

| Trackers | SiamFC | DCFNet | CFNet | LUDT | EAST | HP | SA-Siam | SiamRPN | RASNet | SACF | Siam-tri | RT-MDNet | MemTrack | StructSiam | LUDT+ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| DP (%) | 77.1 | - | 74.8 | 76.9 | - | 79.6 | 86.5 | 85.1 | - | 83.9 | 78.1 | 88.5 | 82.0 | 85.1 | 84.3 |
| AUC (%) | 58.2 | 58.0 | 56.8 | 60.2 | 62.9 | 60.1 | 65.7 | 63.7 | 64.2 | 63.3 | 59.2 | 65.0 | 62.6 | 62.1 | 63.9 |
| FPS | 86 | 70 | 65 | 70 | 25 | 159 | 69 | 160 | 83 | 23 | 86 | 50 | 50 | 45 | 55 |

tions achieve higher accuracy compared with hand-crafted features.

Our unsupervised representation is compared with the off-the-shelf deep model VGG-19. Note that the VGG model is trained on an image classification task with supervised learning. We show the feature maps from different layers (i.e., Conv1-2, Conv2-2, Conv3-4, Conv4-4, and Conv5-4) of the VGG-19. From Fig. 8, we observe that our unsupervised feature representations share similarities with the low-level features (i.e., the first two layers) of VGG, which typically represents spatial details. It has been well studied in HCF [41] and C-COT [13] that merely using the first or second layer of the VGG model for DCF tracking contains limitations. However, our unsupervised representation better suits the tracking scenario since we jointly combine the feature representation learning with the DCF formulation in an end-to-end fashion. In the deeper layers of VGG-19 such as Conv4-4 and Conv5-4, the feature representation gradually loses spatial details but increases semantics, which can be combined with the low-level features to further boost the tracking performance [41, 13]. The semantic representation capability is obtained by distinguishing different object categories (i.e., image classification), while our unsupervised learning process lacks such image labels. In the future, we will investigate how to learn rich multiple-level representations for visual tracking under an unsupervised manner.

### 4.4 Comparison with State-of-the-art Methods

**OTB-2013 Dataset.** The OTB-2013 dataset [69] contains 50 challenging videos. On the OTB-2013 dataset, we evaluate our LUDT and LUDT+ trackers with state-of-the-art real-time trackers including ACT [4], ACFN [7], CFNet [58], SiamFC [1], SCT [6], CSR-DCF [39], DSST [10], and KCF [19] using precision and success plots.

As illustrated in Fig. 9, our unsupervised LUDT tracker outperforms CFNet and SiamFC in both distance precision and AUC score. It is worth mentioning that both LUDT and CFNet have similar network capabilities (network depth), leverage the same training data, and are not equipped with additional online improvements. Even though our approach is free of ground-truth supervision, it still achieves very competitive tracking accuracy. Our improved version, LUDT+, performs favorably against recent state-



**Fig. 9** Precision and success plots on the OTB-2013 dataset [69] for recent real-time trackers. The legend in each tracker shows the precision at 20 pixels of precision plot and AUC of success plot.
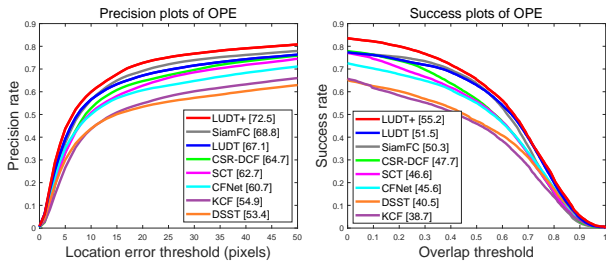


**Fig. 10** Precision and success plots on the OTB-2015 dataset [70] for recent real-time trackers. The legend in each tracker shows the precision at 20 pixels of precision plot and AUC of success plot.

of-the-art real-time trackers such as ACT and ACFN. Besides, our LUDT and LUDT+ trackers also exhibit a real-time speed of about 70 FPS and 55 FPS, respectively.

**OTB-2015 Dataset.** The OTB-2015 dataset [70] contains 100 challenging videos. On the OTB-2015 dataset [70], we evaluate LUDT and LUDT+ trackers with state-of-the-art real-time algorithms as that in OTB-2013. In Table 7, we further compare our methods with more state-of-the-art real-time trackers such as StructSiam [75], MemTrack [71], RT-MDNet [22], Siam-tri [14], SACF [74], RASNet [65], SiamRPN [32], SA-Siam [18], HP [15], and EAST [20].

From Fig. 10 and Table 7, we observe that our unsupervised LUDT tracker is comparable with supervised baseline methods (e.g., SiamFC and CFNet). On the OTB-2015 dataset, SiamFC achieves 77.1% DP and 58.2% AUC, while LUDT exhibits 76.9% DP and 60.2% AUC. Compared with CFNet, LUDT outperforms by 2.1% DP and 3.4% AUC. The DSST algorithm is a traditional DCF based tracker with an accurate target scale estimation. LUDT significantly outperforms it by 8.0% DP and 8.4% AUC, which illustrates that our unsupervised feature representation is more
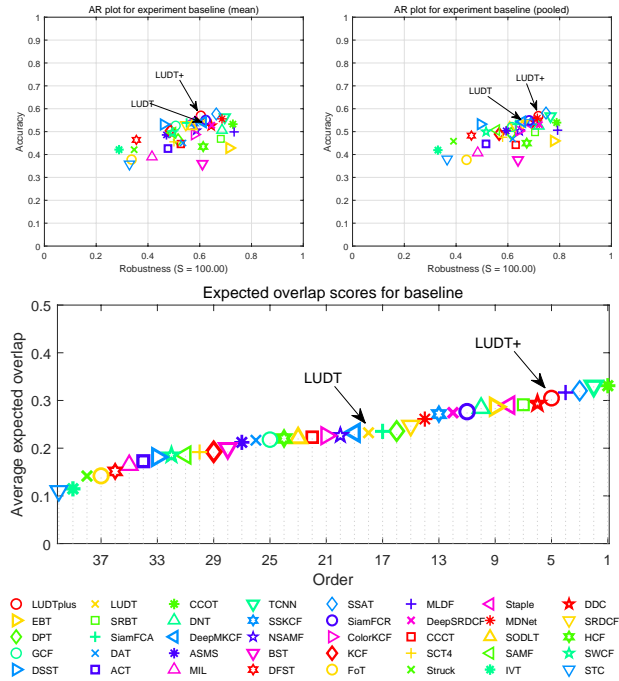
**Fig. 11** Precision and success plots on the Temple-Color dataset [35] for recent real-time trackers. The legend in each tracker shows the precision at 20 pixels of precision plot and AUC of success plot.

robust than empirical features (e.g., HOG). With a better DCF formulation and more advanced online update strategies [9], our LUDT+ tracker achieves comparable performance with the recent ACFN and ACT trackers. In Fig. 10 and Table 7, we do not compare with some remarkable non-realtime trackers. For example, MDNet [46] and ECO [9] can yield 67.8% and 69.4% AUC on the OTB-2015, but they are far from achieving a real-time speed.

Table 7 compares more recent supervised trackers. These latest approaches are mainly based on the Siamese network, which improve the baseline SiamFC method using various sophisticated techniques. Most trackers in Table 7 are trained using ILSVRC including LUDT+. However, it is worth mentioning that some algorithms (e.g., SA-Siam and RT-MDNet) adopt pre-trained CNN models (e.g., AlexNet [28] and VGG-M [3]) for network initialization. SiamRPN additionally uses more labeled training videos from the Youtube-BB dataset [49]. Compared with them, LUDT+ does not require data labels or off-the-shelf deep models, while still achieving comparable performance and efficiency.

**Temple-Color Dataset.** Temple-Color [35] is a more challenging benchmark with 128 color videos. In this dataset, we compare our trackers with some baselines and state-of-the-art trackers as on the OTB-2015 benchmark. Compared with the DCF trackers with empirical features (e.g., HOG feature), our tracker with unsupervised deep features exhibits a significant performance improvement as shown in Fig. 11. Specifically, SiamFC which is learned with full supervision achieves an AUC score of 50.3%, while LUDT exhibits a 51.3% AUC score. Compared with another representative supervised method CFNet, LUDT exceeds its performance by 6.4% DP and 4.7% AUC. Furthermore, our LUDT+ tracker performs favorably against existing state-of-the-art trackers.

**VOT2016 Dataset.** We report the evaluation results on the VOT2016 benchmark [25], which contains 60 videos selected from more than 300 videos. Different from the OTB dataset, the VOT toolkit will reinitialize when the tracker fails. The expected average overlap (EAO) is the final metric



**Fig. 12** Top: Accuracy-Robustness (AR) ranking plots generated by sequence mean (left) and sequence pooling (right) on the VOT2016 dataset [25]. Trackers closer to the upper right corner perform better. Bottom: Expected Average Overlap (EAO) graph with trackers ranked from right to left evaluated on VOT2016.

**Table 8** Comparison with state-of-the-art and baseline trackers on the VOT2016 benchmark [25]. The evaluation metrics include Accuracy, Failures (over 60 sequences), and Expected Average Overlap (EAO). The up arrows indicate that higher values are better for the corresponding metric and vice versa.
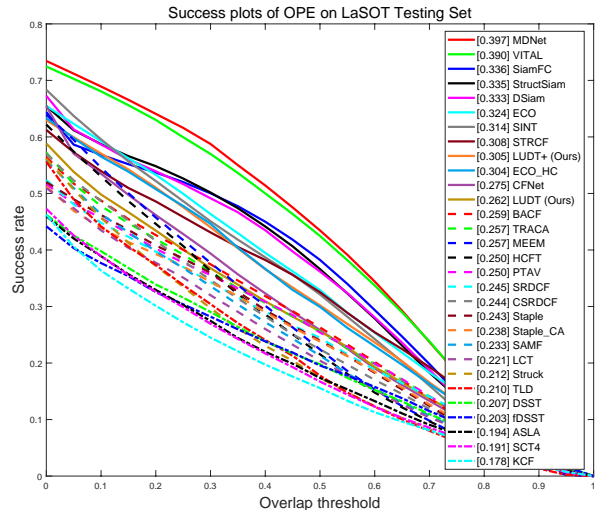
| Trackers | Accuracy (↑) | Failures (↓) | EAO (↑) | FPS (↑) |
|---|---|---|---|---|
| ECO | 0.54 | - | 0.374 | 6 |
| VITAL | - | - | 0.323 | 1 |
| DSLT | - | - | 0.332 | 6 |
| RTINet | 0.57 | - | 0.298 | 9 |
| C-COT | 0.52 | 51 | 0.331 | 0.3 |
| pyMDNet | - | - | 0.304 | 2 |
| HCF | 0.45 | 85 | 0.220 | 12 |
| ACT | - | - | 0.275 | 30 |
| SA-Siam | 0.53 | - | 0.291 | 50 |
| SiamRPN | 0.56 | - | 0.344 | 160 |
| SACF | - | - | 0.275 | 23 |
| StructSiam | - | - | 0.264 | 45 |
| MemTrack | 0.53 | - | 0.273 | 50 |
| SiamFC | 0.53 | 99 | 0.235 | 86 |
| SCT4 | 0.48 | 117 | 0.188 | 40 |
| DSST | 0.53 | 151 | 0.181 | 25 |
| KCF | 0.49 | 122 | 0.192 | 170 |
| LUDT (Ours) | 0.54 | 100 | 0.231 | 70 |
| LUDT+ (Ours) | 0.54 | 62 | 0.309 | 55 |

**Table 9** Comparison with state-of-the-art and baseline trackers on the VOT2017/2018 benchmark [25]. The evaluation metrics include Accuracy, Failures (over 60 sequences), and Expected Average Overlap (EAO). The up arrows indicate that higher values are better for the corresponding metric and vice versa.

| Trackers | Accuracy (↑) | Failures (↓) | EAO (↑) | FPS (↑) |
|---|---|---|---|---|
| ECO | 0.48 | 59 | 0.280 | 6 |
| C-COT | 0.49 | 68 | 0.267 | 0.3 |
| SA-Siam | 0.50 | - | 0.236 | 50 |
| SiamRPN | - | - | 0.243 | 160 |
| SiamFC | 0.50 | 125 | 0.188 | 86 |
| Staple | 0.53 | 147 | 0.169 | 70 |
| TRACA | 0.42 | 183 | 0.137 | 100 |
| SRDCF | 0.49 | 208 | 0.119 | 5 |
| DSST | 0.40 | 310 | 0.079 | 25 |
| KCF | 0.45 | 165 | 0.135 | 170 |
| LUDT (Ours) | 0.46 | 149 | 0.154 | 70 |
| LUDT+ (Ours) | 0.49 | 88 | 0.230 | 55 |

for tracker ranking [27]. In Fig. 12, we show the accuracy-robustness (AR) plot and EAO ranking plot on VOT2016 with some participant trackers. The VOT2016 champion C-COT uses the pre-trained VGG-M model for feature extraction while not achieving real-time performance. The proposed LUDT+ method performs slightly worse than C-COT but runs much faster. It is worth mentioning that our real-time LUDT+ tracker even performs favorably against remarkable non-realtime deep trackers such as MDNet. Our LUDT tracker, without bells and whistles, surpasses classic DCF trackers such as DSST and KCF by a considerable margin and is comparable with some DCF methods with an off-the-shelf deep model (e.g., DeepMKCF and HCF).

In Table 8, we include more state-of-the-art trackers including VITAL [54], DSLT [37], RTINet [72], ACT [4], SA-Siam [18], SiamRPN [32], SACF [74], StructSiam [75], and MemTrack [71] on the VOT2016 benchmark. Compared with the baseline SiamFC, our LUDT tracker yields favorable results. Compared with fully-supervised trackers, LUDT+ overall exhibits competitive performance as well as efficiency.

**VOT2017/2018 Dataset.** The VOT2017 [26] and VOT2018 [24] are the same benchmark with more challenging videos compared with those in the VOT2016 dataset. In Table 9, we present the Accuracy, Failures, and EAO of the state-of-the-art trackers on VOT2017/VOT2018. The proposed LUDT tracker is still superior to the standard DCF trackers using hand-crafted features such as DSST and KCF. Our LUDT+ yields an EAO score of 0.230, which is comparable with the advanced Siamese trackers such as SA-Siam and SiamRPN that take advantage of additional backbone networks or training data.

**LaSOT Dataset.** We further evaluate our unsupervised approach on the large-scale LaSOT testing dataset [16] with
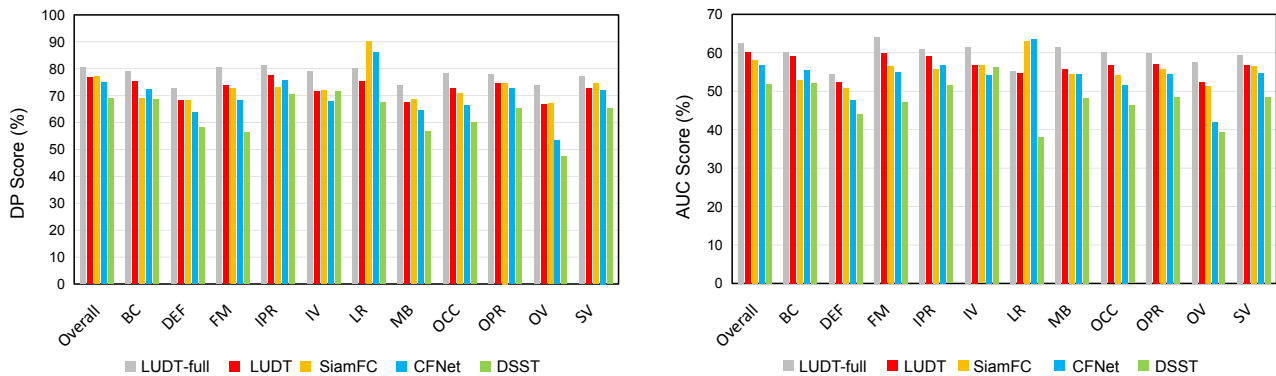


**Fig. 13** Success plots on the LaSOT testing set [16]. The legend in each tracker shows the AUC of the success plot. Best viewed in color and zoom in.
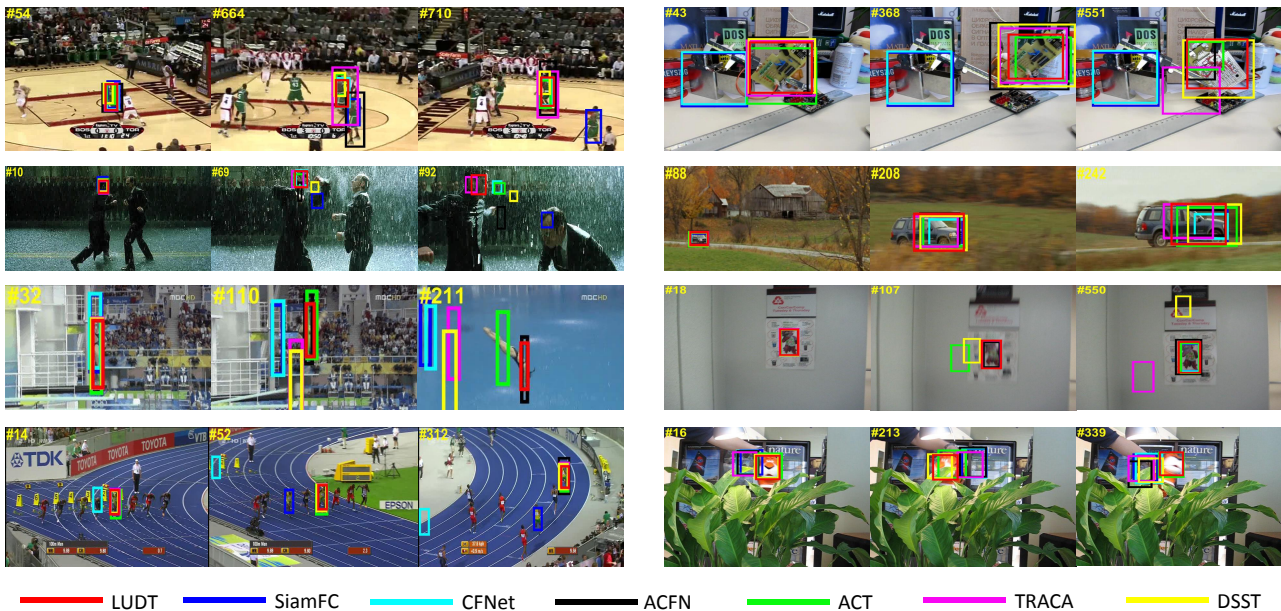
**Table 10** Comparison with state-of-the-art and baseline trackers on the TrackingNet benchmark [45]. The evaluation metrics include Precision, Normalized Precision, and Success (AUC score).

| Trackers | Precision | Norm.Prec | Success |
|---|---|---|---|
| MDNet | 0.565 | 0.705 | 0.606 |
| CFNet | 0.533 | 0.654 | 0.578 |
| SiamFC | 0.533 | 0.663 | 0.571 |
| ECO | 0.492 | 0.618 | 0.554 |
| ECOhc | 0.476 | 0.608 | 0.541 |
| CSRDCF | 0.480 | 0.622 | 0.534 |
| Staple_CA | 0.468 | 0.605 | 0.529 |
| Staple | 0.470 | 0.603 | 0.528 |
| BACF | 0.461 | 0.580 | 0.523 |
| SRDCF | 0.455 | 0.573 | 0.521 |
| SAMF | 0.477 | 0.598 | 0.504 |
| ASLA | 0.406 | 0.536 | 0.478 |
| SAMF_AT | 0.447 | 0.560 | 0.472 |
| DLSSVM | 0.418 | 0.562 | 0.470 |
| DSST | 0.460 | 0.588 | 0.464 |
| MEEM | 0.386 | 0.545 | 0.460 |
| Struck | 0.402 | 0.539 | 0.456 |
| DCF | 0.419 | 0.548 | 0.448 |
| KCF | 0.419 | 0.546 | 0.447 |
| CSK | 0.368 | 0.503 | 0.429 |
| TLD | 0.336 | 0.460 | 0.417 |
| TLD | 0.292 | 0.438 | 0.400 |
| MOSSE | 0.326 | 0.442 | 0.388 |
| LUDT (Ours) | 0.469 | 0.593 | 0.543 |
| LUDT+ (Ours) | 0.495 | 0.633 | 0.563 |

280 videos. The videos in LaSOT are more challenging with an average length of about 2500 frames. As shown in Fig. 13, our LUDT tracker still outperforms hand-crafted feature based DCF trackers such as BACF [17], CSR-DCF [38], DSST [10], and SCT4 [6]. Furthermore, the proposed LUDT+ approach achieves an AUC score of 30.5%, which is even comparable with some state-of-the-art deep DCF track-

**Fig. 14** Attribute-based evaluation on the OTB-2015 dataset [70]. The 11 attributes are background clutter (BC), deformation (DEF), fast motion (FM), in-plane rotation (IPR), illumination variation (IV), low resolution (LR), motion blur (MB), occlusion (OCC), out-of-plane rotation (OPR), out-of-view (OV), and scale varition (SV), respectively.



**Fig. 15** Qualitative evaluation of our proposed LUDT and other trackers including SiamFC [1], CFNet [58], ACFN [7], ACT [4], TRACA [5], and DSST [10] on 8 challenging videos from OTB-2015. From left to right and top to down are *Basketball*, *Board*, *Matrix*, *CarScale*, *Diving*, *BlurOwl*, *Bolt*, and *Tiger1*, respectively. Best viewed in color.

ers including ECO (32.4%) [9], STRCF (30.8%) [33], and TRACA (25.7%) [5] that leverage off-the-shelf deep models as feature extractors.

**TrackingNet Dataset.** The recently released large-scale TrackingNet dataset [45] contains more than 30K videos with more than 14 million dense bounding box annotations. The videos are collected on the Internet (YouTube), providing large-scale high-quality data for assessing trackers in the wild. We test our LUDT and LUDT+ on the testing set with 511 videos. Following [45], we adopt three metrics including Precision, Normalized Precision, and Success (AUC) for performance evaluation. In Table 10, we exhibit the results of our methods and all the evaluated trackers on this benchmark. On this dataset, our LUDT achieves an AUC score
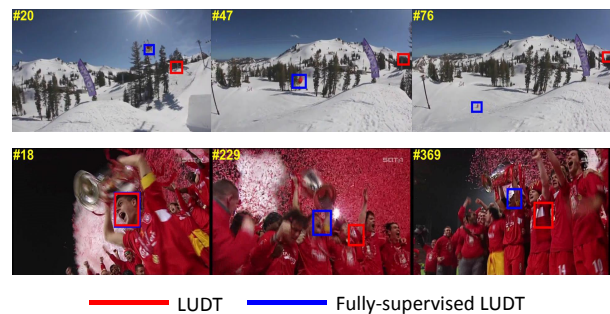
of 54.3%, which obviously outperforms other hand-crafted feature based DCF trackers such as ECOhc, CSR-DCF, and BACF by 0.2%, 0.9%, and 2.0%, respectively. Note that the above DCF trackers are improved versions with additional regularization terms, while ours merely utilizes a standard DCF formulation. Our superior performance illustrates the representational power of our unsupervised features. Besides, it is worth mentioning that our LUDT, on this large-scale benchmark, is even comparable with the state-of-the-art ECO, which leverages both hand-crafted and off-the-shelf deep features. Without labeled data for model training, our improved LUDT+ achieves better performance and slightly outperforms ECO by 0.9% in terms of AUC.

**Attribute Analysis.** The videos on OTB-2015 [70] are annotated with 11 different attributes, namely: background clutter (BC), deformation (DEF), out-of-plane rotation (OPR), scale variation (SV), occlusion (OCC), illumination variation (IV), motion blur (MB), in-plane rotation (IPR), out-of-view (OV), fast motion (FM), and low resolution (LR). On the OTB-2015 benchmark, we further analyze the performances over different challenges in Fig. 14. On the majority of challenging scenes, our LUDT tracker outperforms the popular SiamFC [1] and CFNet [58] trackers. However, our performance advantage in the DP metric is less obvious than that in the AUC metric. Compared with the fully-supervised LUDT tracker, the main performance gaps are from illumination variation (IV), occlusion (OCC), and fast motion (FM) attributes. Unsupervised learning can be further improved on these attributes.

**Qualitative Evaluation.** We evaluate LUDT with supervised trackers (e.g., ACT, ACFN, SiamFC, TRACA, and CFNet) and a baseline DCF tracker (DSST) on eight challenging videos, as shown in Fig. 15. On *Matrix* and *Tiger1* videos, the targets undergo partial occlusion and background clutter, while on *BlurOwl*, the target is extremely blurry due to the drastic camera shaking. In these videos, DSST based on empirical features fails to cope with the challenging factors while LUDT is able to handle. This illustrates the robustness of our unsupervised feature representation, which achieves favorable performance compared to the empirical features. The SiamFC and CFNet trackers tend to drift when the target and distractors are similar (e.g., *Bolt* and *Basketball* sequences) while LUDT is able to handle these challenging scenes because of the discriminating capability of DCF and its online model update mechanism. Without online improvements, LUDT is still able to track the target accurately, especially on the challenging *Board* and *Diving* videos. It is worth mentioning that such a robust tracker is learned from raw videos under an unsupervised manner.

### 4.5 Limitations

Fig. 16 shows the limitations of our unsupervised learning. First, compared with the fully supervised learning, our tracker trained via unsupervised learning tends to drift when an occlusion or a drastic appearance change occurs (e.g., the targets in *Skiing* and *Soccer* sequences). The semantic representations brought by ground-truth annotations are missing. Second, our unsupervised learning involves both forward and backward trackings. The computational load during the training phase is a potential drawback although the learning process is offline.



**Fig. 16** Failure cases of our LUDT tracker. The top and bottom videos are *Skiing* and *Soccer*, respectively. Compared to its fully-supervised version, our unsupervised method is not robust enough when the target undergoes drastic appearance change and occlusion.

## 5 Concluding Remarks

In this paper, we present how to train a visual tracker using unlabeled videos in the wild, which is rarely investigated in visual tracking. By designing an unsupervised Siamese correlation filter network, we verify the feasibility and effectiveness of our forward-backward based unsupervised training pipeline. To further facilitate the unsupervised training, we extend our framework to consider multiple frames and employ a cost-sensitive loss. Extensive experiments exhibit that the proposed unsupervised tracker, without bells and whistles, performs as a solid baseline and achieves comparable results with the classic fully-supervised trackers. Equipped with additional online improvements such as a sophisticated update scheme, our LUDT+ tracker performs favorably against the state-of-the-art tracking algorithms. Furthermore, we provide a deep analysis of our unsupervised representation by feature visualization and extensive ablation studies. Our unsupervised framework shows a promising potential in visual tracking, such as utilizing more unlabeled data or weakly labeled data to further improve the tracking accuracy.

## References

1. Bertinetto, L., Valmadre, J., Henriques, J.F., Vedaldi, A., Torr, P.H.: Fully-convolutional siamese networks for object tracking. In: Proceedings of the European Conference on Computer Vision Workshops (ECCV Workshop) (2016)
2. Bolme, D.S., Beveridge, J.R., Draper, B.A., Lui, Y.M.: Visual object tracking using adaptive correlation filters. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2010)

3. Chatfield, K., Simonyan, K., Vedaldi, A., Zisserman, A.: Return of the devil in the details: Delving deep into convolutional nets. In: British Machine Vision Conference (BMVC) (2014)

4. Chen, B., Wang, D., Li, P., Wang, S., Lu, H.: Real-time 'actor-critic' tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

5. Choi, J., Jin Chang, H., Fischer, T., Yun, S., Lee, K., Jeong, J., Demiris, Y., Young Choi, J.: Context-aware deep feature compression for high-speed visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

6. Choi, J., Jin Chang, H., Jeong, J., Demiris, Y., Young Choi, J.: Visual tracking using attention-modulated disintegration and integration. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

7. Choi, J., Jin Chang, H., Yun, S., Fischer, T., Demiris, Y., Young Choi, J.: Attentional correlation filter network for adaptive visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

8. Dalal, N., Triggs, B.: Histograms of oriented gradients for human detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2005)

9. Danelljan, M., Bhat, G., Shahbaz Khan, F., Felsberg, M.: Eco: Efficient convolution operators for tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

10. Danelljan, M., Häger, G., Khan, F., Felsberg, M.: Accurate scale estimation for robust visual tracking. In: British Machine Vision Conference (BMVC) (2014)

11. Danelljan, M., Häger, G., Khan, F.S., Felsberg, M.: Adaptive decontamination of the training set: A unified formulation for discriminative visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

12. Danelljan, M., Hager, G., Shahbaz Khan, F., Felsberg, M.: Learning spatially regularized correlation filters for visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2015)

13. Danelljan, M., Robinson, A., Khan, F.S., Felsberg, M.: Beyond correlation filters: Learning continuous convolution operators for visual tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2016)

14. Dong, X., Shen, J.: Triplet loss in siamese network for object tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

15. Dong, X., Shen, J., Wang, W., Liu, Y., Shao, L., Porikli, F.: Hyperparameter optimization for tracking with continuous deep q-learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

16. Fan, H., Lin, L., Yang, F., Chu, P., Deng, G., Yu, S., Bai, H., Xu, Y., Liao, C., Ling, H.: Lasot: A high-quality benchmark for large-scale single object tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

17. Galoogahi, H.K., Fagg, A., Lucey, S.: Learning background-aware correlation filters for visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2017)

18. He, A., Luo, C., Tian, X., Zeng, W.: A twofold siamese network for real-time object tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

19. Henriques, J.F., Caseiro, R., Martins, P., Batista, J.: High-speed tracking with kernelized correlation filters. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **37**(3), 583–596 (2015)

20. Huang, C., Lucey, S., Ramanan, D.: Learning policies for adaptive tracking with deep feature cascades. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2017)

21. Huang, D., Luo, L., Chen, Z., Wen, M., Zhang, C.: Applying detection proposals to visual tracking for scale and aspect ratio adaptability. International Journal of Computer Vision (IJCV) **122**(3), 524–541 (2017)

22. Jung, I., Son, J., Baek, M., Han, B.: Real-time mdnet. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

23. Kalal, Z., Mikolajczyk, K., Matas, J.: Tracking-learning-detection. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **34**(7), 1409–1422 (2012)

24. Kristan, M., Leonardis, A., Matas, J., Felsberg, M., Pflugfelder, R., Cehovin Zajc, L., et al.: The sixth visual object tracking vot2018 challenge results. In: Proceedings of the European Conference on Computer Vision Workshops (ECCV Workshop) (2018)

25. Kristan, M., Matas, J., Leonardis, A., Felsberg, M., Cehovin, L., Fernández, G., Vojir, T., Hager, et al.: The visual object tracking vot2016 challenge results. In: Proceedings of the European Conference on Computer Vision Workshops (ECCV Workshop) (2016)

26. Kristan, M., Matas, J., Leonardis, A., Felsberg, M., Cehovin, L., Fernández, G., Vojir, T., Hager, et al.: The visual object tracking vot2017 challenge results. In: Proceedings of the IEEE International Conference on Computer Vision Workshops (ICCV Workshop) (2017)

27. Kristan, M., Matas, J., Leonardis, A., Vojíř, T., Pflugfelder, R., Fernandez, G., Nebehay, G., Porikli, F., Čehovin, L.: A novel performance evaluation methodology for single-target trackers. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **38**(11), 2137–2155 (2016)

28. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: Advances in Neural Information Processing Systems (NeurIPS) (2012)

29. Le, Q.V., Ranzato, M., Monga, R., Devin, M., Chen, K., Corrado, G.S., Dean, J., Ng, A.Y.: Building high-level features using large scale unsupervised learning. arXiv preprint arXiv:1112.6209 (2011)

30. Lee, D.Y., Sim, J.Y., Kim, C.S.: Multihypothesis trajectory analysis for robust visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015)

31. Lee, H.Y., Huang, J.B., Singh, M., Yang, M.H.: Unsupervised representation learning by sorting sequences. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2017)

32. Li, B., Yan, J., Wu, W., Zhu, Z., Hu, X.: High performance visual tracking with siamese region proposal network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

33. Li, F., Tian, C., Zuo, W., Zhang, L., Yang, M.H.: Learning spatial-temporal regularized correlation filters for visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

34. Li, F., Yao, Y., Li, P., Zhang, D., Zuo, W., Yang, M.H.: Integrating boundary and center correlation filters for visual tracking with aspect ratio variation. In: Proceedings of the IEEE International Conference on Computer Vision Workshops (ICCV Workshop) (2017)

35. Liang, P., Blasch, E., Ling, H.: Encoding color information for visual tracking: algorithms and benchmark. IEEE Transactions on Image Processing (TIP) **24**(12), 5630–5644 (2015)

36. Liu, S., Zhang, T., Cao, X., Xu, C.: Structural correlation filter for robust visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

37. Lu, X., Ma, C., Ni, B., Yang, X., Reid, I., Yang, M.H.: Deep regression tracking with shrinkage loss. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

38. LukeźIǎź, A., Vojíř, T., Čehovin Zajc, L., Matas, J., Kristan, M.: Discriminative correlation filter tracker with channel and spatial reliability. International Journal of Computer Vision (IJCV) **126**(7), 671–688 (2018)

39. Lukezic, A., Vojir, T., Cehovin Zajc, L., Matas, J., Kristan, M.: Discriminative correlation filter with channel and spatial reliability. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

40. Luo, W., Sun, P., Zhong, F., Liu, W., Zhang, T., Wang, Y.: End-to-end active object tracking and its real-world deployment via reinforcement learning. IEEE transactions on pattern analysis and machine intelligence (TPAMI) **42**(6), 1317–1332 (2019)

41. Ma, C., Huang, J.B., Yang, X., Yang, M.H.: Hierarchical convolutional features for visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2015)

42. Ma, C., Huang, J.B., Yang, X., Yang, M.H.: Adaptive correlation filters with long-term and short-term memory for object tracking. International Journal of Computer Vision (IJCV) **126**(8), 771–796 (2018)

43. Meister, S., Hur, J., Roth, S.: Unflow: Unsupervised learning of optical flow with a bidirectional census loss. In: AAAI Conference on Artificial Intelligence (AAAI) (2018)

44. Mueller, M., Smith, N., Ghanem, B.: Context-aware correlation filter tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

45. Müller, M., Bibi, A., Giancola, S., Al-Subaihi, S., Ghanem, B.: Trackingnet: A large-scale dataset and benchmark for object tracking in the wild. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

46. Nam, H., Han, B.: Learning multi-domain convolutional neural networks for visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

47. Olshausen, B.A., Field, D.J.: Sparse coding with an overcomplete basis set: A strategy employed by v1? Vision research **37**(23), 3311–3325 (1997)

48. Pu, S., Song, Y., Ma, C., Zhang, H., Yang, M.H.: Deep attentive tracking via reciprocative learning. In: Advances in Neural Information Processing Systems (NeurIPS) (2018)

49. 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: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

50. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **39**(6), 1137–1149 (2016)

51. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. International Journal of Computer Vision (IJCV) **115**(3), 211–252 (2015)

52. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)

53. Song, Y., Ma, C., Gong, L., Zhang, J., Lau, R., Yang, M.H.: Crest: Convolutional residual learning for visual tracking. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2017)

54. Song, Y., Ma, C., Wu, X., Gong, L., Bao, L., Zuo, W., Shen, C., Lau, R.W., Yang, M.H.: Vital: Visual tracking via adversarial learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

55. Sui, Y., Zhang, Z., Wang, G., Tang, Y., Zhang, L.: Exploiting the anisotropy of correlation filter learning for visual tracking. International Journal of Computer Vision (IJCV) pp. 1–22 (2019)

56. Tomasi, C., Kanade, T.: Detection and tracking of point features (1991)

57. Valmadre, J., Bertinetto, L., Henriques, J.F., Tao, R., Vedaldi, A., Smeulders, A., Torr, P., Gavves, E.: Long-term tracking in the wild: A benchmark. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

58. Valmadre, J., Bertinetto, L., Henriques, J.F., Vedaldi, A., Torr, P.H.: End-to-end representation learning for correlation filter based tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

59. Vondrick, C., Pirsiavash, H., Torralba, A.: Anticipating visual representations from unlabeled video. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

60. Vondrick, C., Shrivastava, A., Fathi, A., Guadarrama, S., Murphy, K.: Tracking emerges by colorizing videos. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

61. Wang, N., Song, Y., Ma, C., Zhou, W., Liu, W., Li, H.: Unsupervised deep tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

62. Wang, N., Yeung, D.Y.: Learning a deep compact image representation for visual tracking. In: Advances in Neural Information Processing Systems (NeurIPS) (2013)

63. Wang, N., Zhou, W., Tian, Q., Hong, R., Wang, M., Li, H.: Multi-cue correlation filters for robust visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

64. Wang, Q., Gao, J., Xing, J., Zhang, M., Hu, W.: Dcfnet: Discriminant correlation filters network for visual tracking. arXiv preprint arXiv:1704.04057 (2017)

65. Wang, Q., Teng, Z., Xing, J., Gao, J., Hu, W., Maybank, S.: Learning attentions: Residual attentional siamese network for high performance online visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

66. Wang, X., Gupta, A.: Unsupervised learning of visual representations using videos. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2015)

67. Wang, X., Jabri, A., Efros, A.A.: Learning correspondence from the cycle-consistency of time. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

68. Weijer, J.V.D., Schmid, C., Verbeek, J., Larlus, D.: Learning color names for real-world applications. IEEE Transactions on Image Processing (TIP) **18**(7), 1512–1523 (2009)

69. Wu, Y., Lim, J., Yang, M.H.: Online object tracking: A benchmark. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2013)

70. Wu, Y., Lim, J., Yang, M.H.: Object tracking benchmark. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) **37**(9), 1834–1848 (2015)

71. Yang, T., Chan, A.B.: Learning dynamic memory networks for object tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

72. Yao, Y., Wu, X., Zhang, L., Shan, S., Zuo, W.: Joint representation and truncated inference learning for correlation filter based tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

73. Yin, Z., Shi, J.: Geonet: Unsupervised learning of dense depth, optical flow and camera pose. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)

74. Zhang, M., Wang, Q., Xing, J., Gao, J., Peng, P., Hu, W., Maybank, S.: Visual tracking via spatially aligned correlation filters network. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

75. Zhang, Y., Wang, L., Qi, J., Wang, D., Feng, M., Lu, H.: Structured siamese network for real-time visual tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)

76. Zhipeng, Z., Houwen, P., Qiang, W.: Deeper and wider siamese networks for real-time visual tracking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2019)

77. Zhou, T., Brown, M., Snavely, N., Lowe, D.G.: Unsupervised learning of depth and ego-motion from video. In: Proceedings

of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)

78. Zhou, T., Krahenbuhl, P., Aubry, M., Huang, Q., Efros, A.A.: Learning dense correspondence via 3d-guided cycle consistency. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)

79. Zhou, X., Zhu, M., Daniilidis, K.: Multi-image matching via fast alternating minimization. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2015)

80. Zhu, Z., Wang, Q., Li, B., Wu, W., Yan, J., Hu, W.: Distractor-aware siamese networks for visual object tracking. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)