

# Parallel Tracking and Verifying

Heng Fan and Haibin Ling

**Abstract**—Being intensively studied, visual object tracking has witnessed great advances in either speed (e.g., with correlation filters) or accuracy (e.g., with deep features). Real-time and high accuracy tracking algorithms, however, remain scarce. In this paper we study the problem from a new perspective and present a novel *parallel tracking and verifying* (PTAV) framework, by taking advantage of the ubiquity of multi-thread techniques and borrowing ideas from the success of *parallel tracking and mapping* in visual SLAM. The proposed PTAV framework is typically composed of two components, a (base) tracker  $\mathcal{T}$  and a verifier  $\mathcal{V}$ , working in parallel on two separate threads. The tracker  $\mathcal{T}$  aims to provide a super real-time tracking inference and is expected to perform well most of the time; by contrast, the verifier  $\mathcal{V}$  validates the tracking results and corrects  $\mathcal{T}$  when needed. The key innovation is that,  $\mathcal{V}$  does not work on every frame but only upon the requests from  $\mathcal{T}$ ; on the other end,  $\mathcal{T}$  may adjust the tracking according to the feedback from  $\mathcal{V}$ . With such collaboration, PTAV enjoys both the high efficiency provided by  $\mathcal{T}$  and the strong discriminative power by  $\mathcal{V}$ . Meanwhile, to adapt  $\mathcal{V}$  to object appearance changes over time, we maintain a dynamic target template pool for adaptive verification, resulting in further performance improvements. In our extensive experiments on popular benchmarks including OTB2015, TC128, UAV20L and VOT2016, PTAV achieves the best tracking accuracy among all real-time trackers, and in fact even outperforms many deep learning based algorithms. Moreover, as a general framework, PTAV is very flexible with great potentials for future improvement and generalization.

**Index Terms**—Visual tracking, deep learning, correlation filter, verification, multi-thread, parallel tracking and verifying (PTAV).

## 1 INTRODUCTION

### 1.1 Background

As one of the most important components in computer vision, visual tracking has a long list of applications such as robotics, intelligent vehicles, visual surveillance, human-computer interaction and so forth [1]–[4]. Given an initial state (usually a bounding box) of a tracking target in the first frame, visual tracking aims at estimating the unknown states (e.g., position and scale) of the target object in subsequent consecutive frames. Although significant progresses have been made in recent decades, robust object tracking still remains challenging due to large appearance variations caused by many factors such as object occlusion, deformation, rotation, illumination variations, scale changes, motion blur and so on.

Recently, inspired by the success of deep convolutional neural networks (CNNs) [5] in image recognition (e.g., [6]), an emerging trend toward improving tracking accuracy is to utilize robust deep features for object appearance representation (e.g., [7]–[16]). Despite significant improvements obtained in accuracy, these algorithms often suffer from high computational burden due to either extracting expensive deep features (e.g., [7], [8], [13]–[16]) or online network fine-tuning (e.g., [9]–[12]), and hardly meet the real-time requirement (see Figure 1 for illustration).

Along a somewhat orthogonal direction, researchers have been proposing efficient visual trackers (e.g., [17]–[24]), notably represented by the series of trackers based on correlation filters. To achieve efficient computation, the correlation filter-based trackers usually represent object appearance with simple hand-crafted features such as raw pixels, HoG [25] and color names [26]. While easily running at real-time, these trackers usually perform less robustly compared to deep learning-based approaches (see again Figure 1).

Despite aforementioned progresses in either speed or accuracy, real-time high quality tracking algorithms remain scarce. A natural way is to seek a trade-off between speed and accuracy (e.g., [21],

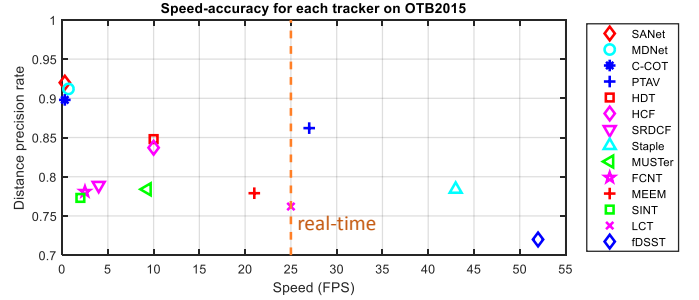


Fig. 1. Speed-accuracy plot of state-of-the-art trackers on OTB2015 [29]. For better illustration, only those trackers with accuracy higher than 0.7 are reported. Compared with high precision deep learning-based trackers (e.g., MDNet, SANet and C-COT) whose speeds are around 1 FPS, our PTAV runs in real-time without serious accuracy degradation. On the other hand, compared with other real-time trackers (e.g., Staple, LCT and fDSST), PTAV achieves a much higher accuracy. Moreover, PTAV even outperforms some deep learning-based trackers in both accuracy and speed (e.g., HCF, HDT, SINT and FCNT).

[27], [28]). In this paper we work toward this goal, but from a novel perspective described as the following.

### 1.2 Motivation

Our key idea is to decompose the original tracking task into two parallel but collaborative ones, one for fast tracking and the other for accurate verification. We are mainly inspired by the following observations or related works:

**Motivation 1.** When tracking a target from visual input, most of the time the target object moves smoothly and its appearance changes slowly or remains the same. Simple but efficient algorithms usually work fine for such easy cases. By contrast, hard cases (e.g., drastic object appearance variations) happen only occasionally, though they can cause serious consequences if not addressed properly. These hard cases usually require computationally expensive processes or analysis, such as the verification in our

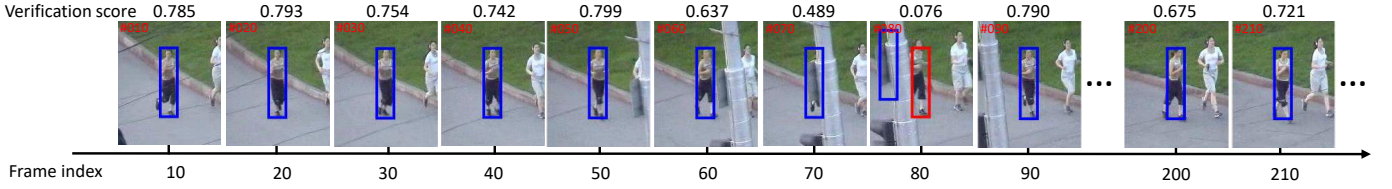


Fig. 2. Illustration of verifying scores on a typical sequence. Verifier validates tracking results every 10 frames. Most of the time the tracking results are reliable (showing in blue). Occasionally, e.g., frame #080, the verifier finds the original tracking result (showing in blue) unreliable and the tracker is corrected and resumes tracking based on detection result provided by verifier (showing in red).

approach. Intuitively, verifications are needed only occasionally instead of for every frame. Figure 2 shows a typical example with both cases.

**Motivation 2.** The ubiquity of multi-thread computing has already benefited computer vision systems, with notably in visual SLAM (*simultaneous localization and mapping*). By splitting tracking and mapping into two parallel threads, PTAM (*parallel tracking and mapping*) [30] provides one of the most popular SLAM frameworks with many important extensions (e.g., ORB-SLAM [31]). A key inspiration in PTAM is that mapping is not needed for every frame. Nor does verifying in our task.

**Motivation 3.** Last but not least, recent advances in either fast or accurate tracking algorithms provide promising building blocks and highly encourage us to seek a balanced system for real-time high accuracy visual tracking.

### 1.3 Contribution

With the motivations listed above, we propose to build real-time high accuracy trackers in a novel framework named *parallel tracking and verifying* (PTAV). PTAV typically consists of two components: a fast tracker<sup>1</sup> denoted by  $\mathcal{T}$  and an accurate verifier denoted by  $\mathcal{V}$ . The two components work in parallel on two separate threads while collaborating with each other. The tracker  $\mathcal{T}$  aims at providing a super real-time tracking inference and is expected to perform well most of the time, e.g., most frames in Figure 2. By contrast, the verifier  $\mathcal{V}$  checks the tracking results and corrects  $\mathcal{T}$  when needed, e.g., at frame #080 in Figure 2.

The key idea is, while  $\mathcal{T}$  needs to run on every frame,  $\mathcal{V}$  does not. As a general framework, PTAV allows the coordination between the tracker and the verifier:  $\mathcal{V}$  checks the tracking results provided by  $\mathcal{T}$  and sends feedback to  $\mathcal{V}$ ; and  $\mathcal{V}$  adjusts itself according to the feedback when necessary. By running  $\mathcal{T}$  and  $\mathcal{V}$  in parallel, PTAV inherits both the high efficiency of  $\mathcal{T}$  and the strong discriminative power of  $\mathcal{V}$ .

Implementing a PTAV algorithm requires three parts: a base tracker for  $\mathcal{T}$ , a base verifier for  $\mathcal{V}$ , and the coordination between them. For  $\mathcal{T}$ , we choose the Staple algorithm [21], which is correlation filter-based and runs efficiently by itself. For  $\mathcal{V}$ , we choose the Siamese network [32] for verification similar to that in [14]. For coordination,  $\mathcal{T}$  sends results to  $\mathcal{V}$  at an adaptive frequency that allows enough time for verification. On the verifier side, when an unreliable result is found,  $\mathcal{V}$  performs detection and sends the detected result to  $\mathcal{T}$  for correction. For  $\mathcal{V}$  to handle object appearance changes over time, we utilize  $k$ -means clustering to maintain a dynamic target template pool for adaptive verification, resulting in further improvements of PTAV in both accuracy and speed.

1. For conciseness, in the rest of this paper, we refer the *fast tracker* as a *tracker*, whenever no confusion caused.

The proposed PTAV algorithm is evaluated thoroughly on several popular benchmarks including OTB2015 [29], TC128 [33], UAV20L [34] and VOT2016 [35]. In these experiments, PTAV achieves the best tracking accuracy among all real-time trackers, and in fact performs even better than many deep learning-based solutions.

In summary, our first main contribution is the novel parallel tracking and verifying framework. With the framework, we make the second contribution by implementing a tracking solution that combines correlation kernel-based tracking and deep learning-based verification. Then, our solution demonstrates very promising results on thorough experiments in comparison with state-of-the-arts. Moreover, it is worth noting that PTAV is a very flexible framework and our implementation may not be optimal. We believe there are great rooms for future improvement and generalization.

This paper is an extended version of a preliminary conference publication [36]. The main new contributions or differences include: (1) a more robust base tracker (i.e., Staple) for  $\mathcal{T}$  in implementing PTAV, which brings clear performance improvement, (2) a dynamic target template pool for adaptive verification against target appearance variations, (3) various ablation studies on  $\mathcal{V}$  and  $\mathcal{T}$  to analyze PTAV, including two base verifiers (VGGNet [37] and AlexNet [6]) for  $\mathcal{V}$  and three base trackers (KCF [17], fDSST [19] and Staple [21]) for  $\mathcal{T}$ , and (4) more thorough experimental validation and analysis involving more state-of-the-art tracking algorithms and benchmarks.

The rest of this paper starts with an extensive overview of related work in Section 2. Then, Section 3 elaborates the proposed PTAV framework and implementation details. Section 4 shows experimental results, including comparisons with state-of-the-arts and ablation studies, followed by conclusion in Section 5.

## 2 RELATED WORK

**Visual tracking algorithms.** Visual tracking has been extensively studied and it is beyond our scope to review all previous studies. Instead, in the following we sample some representative works and discuss those closely related to ours. Some comprehensive reviews on object tracking can be found in [1]–[4].

The focus of this paper is on *model-free* single object tracking, for which existing algorithms are often categorized as either discriminative or generative. Discriminative algorithms usually treat tracking as a classification problem that distinguishes the target from changing background. Babenko *et al.* [38] apply multiple instant learning (MIL) to learn a classifier based on an adaptive appearance model for object tracking. Zhang *et al.* [18] propose a compressive sensing tracker by projecting high-dimensional features to low-dimensional compressed subspace. Grabner *et al.* [39] propose an online tracking algorithm via semi-supervised

boosting which treats samples from the first frame as labeled and other samples as unlabeled. Hare *et al.* [40] propose to leverage a kernelized structured output support vector machine (SVM) for robust visual tracking by mitigating the effect of wrong labeling samples.

By contrast, generative algorithms usually formulate tracking as searching for regions most similar to the target. To this end, numerous object appearance modeling approaches have been proposed. In [41], Ross *et al.* propose an online incremental subspace learning method to adapt appearance changes for object tracking. Kwon *et al.* [42] present a modified particle filtering framework for tracking by combining multiple observation and motion models to handle large appearance and motion variations. Mei and Ling [43] model object appearance with sparse representation and propose the  $\ell_1$ -tracker, which is later improved via an accelerated proximal gradient algorithm [44]. In [45], Zhang *et al.* propose to incorporate target structure into sparse representation for robustness.

**Deep learning-based tracking.** Motivated by the power of deep features in visual recognition (e.g., [6], [37]), some trackers utilize deep features for object appearance modeling, and achieve excellent performance, though typically at the cost of low running speed. Wang *et al.* [11] introduce a stacked denoising autoencoder to learn generic image features for visual tracking. In [15], Wang *et al.* present a fully convolutional neural network tracking (FCNT) algorithm by transferring pre-trained CNN features to improve tracking accuracy. Ma *et al.* [7] replace HoG [25] with discriminative convolutional features for correlation filter tracking, resulting in remarkable performance gains. A similar idea is present by Qi *et al.* [8] by adaptively merging convolutional features from different layers. Hong *et al.* [16] propose to learn discriminative saliency map for online object tracking using CNNs. To address the problem of lack of training samples, Wang *et al.* [10] employ intermediate features in networks to learn a robust ensemble tracker. In [9], Nam *et al.* propose to impose multiple domain branches on a light architecture of CNNs to learn generic feature for tracking target, and then introduce an online tracking algorithm by updating network weights in each frame. In [12], recurrent neural networks (RNNs) are introduced to capture internal structure of a tracking target and the generated tracking algorithm achieves promising results on several tracking benchmarks. Though these approaches have achieved very impressive results, the heavy computation burden severely restricts their practical applications.

**Correlation filter-based tracking.** Recently, correlation filter has drawn increasing attention in visual tracking owing partly to its high computation efficiency. Bolme *et al.* [20] propose to use correlation filter for tracking through learning the minimum output sum of squared error (MOSSE). Benefitting from the high computation efficiency of correlation filter, this approach runs amazingly at hundreds of frames per second. Henriques *et al.* [22] incorporate kernel space into correlation filter and propose a circulant structure with kernel (CSK) method for visual object tracking, and later [17] extends CSK to the well-known kernelized correlation filters (KCF) tracker by substituting raw pixel intensities with HoG [25] for appearance representation. To deal with the scale issue, Danelljan *et al.* [19] suggest an extra scale filter into correlation filter tracking to adaptively estimate target scale. In [46], Li *et al.* adopt a similar strategy to address the problem of scale changes. Later, more efforts have been made to improve correlation filter tracking. Danelljan *et al.* [23] investigate color attributes to improve correlation filter tracking. To reduce the

risk of model drift, Ma *et al.* [27] introduce an auxiliary detector to re-locate the tracking target when sensing drift. In [47], Liu *et al.* propose a part-based correlation filter tracker by decomposing tracking target into fragments, which is robust to resist occlusion. Mueller *et al.* [48] propose to explicitly incorporate context into correlation filter to improve its discriminability. To alleviate boundary effect, the work in [49] presents spatially regularized correlation filters for tracking. Bertinetto *et al.* [21] introduce a complementary tracker by combining correlation filters and color histograms. In [50], Valmadre *et al.* propose correlation filter networks for tracking by enjoying end-to-end representation.

**Verification in tracking.** The idea of verification is not new for tracking. A notable example is the tracking-learning-detection (TLD) algorithm [24], in which tracking results are validated *per frame* to decide how learning/detection shall progress. Similar ideas have been used in later works. In [51], Hua *et al.* propose an occlusion and motion reasoning-based long-term tracker. An occlusion detector is applied in *each frame* to prevent tracking model from drift. Ma *et al.* [27] apply an additional object detector in correlation filter tracking for the same end. Unlike in previous studies, the verification in PTAV runs only on sampled frames. This mechanism allows PTAV to use strong verification algorithms without worrying much about running time efficiency. In fact, we utilize the Siamese network [32] that is designed for verification tasks.

Interestingly, tracking by itself can be also formulated as a verification problem that finds the best candidate similar to the tracking target [14], [28]. Bertinetto *et al.* [28] propose a fully-convolutional Siamese network for visual tracking by searching the tracking target within a local region. Tao *et al.* [14] formulate tracking problem as a task of object matching in each frame and develop a matching function based on the Siamese network. Despite obtaining excellent performance, the application of such trackers is limited by the heavy computation for extracting deep features in each frame. Compared with these studies, our solution treats verification only as a way to validate and correct the *fast tracker*, and does not run verification per frame.

**Ensemble tracking.** To achieve robustness in tracking, a natural solution is to leverage multiple different components to determine tracking result. Kwon *et al.* [42] combine multiple observation and motion models to handle large appearance changes in tracking. The TLD tracker [24] uses both tracker and detector for long-term tracking, and a similar idea is adopted in [27]. In [21], Bertinetto *et al.* propose a complementary tracking approach based on correlation filters and color histograms. Hong *et al.* [52] present a multi-store tracker (MUSTer) which consists of short- and long-term memory stores to process target appearance. The final tracking result is jointly determined by short- and long-term memories. Different than in these studies, our PTAV is comprised of two components, which run on two parallel threads asynchronously while collaborating with each other.

Though other ensemble approaches (e.g., [21], [27], [42], [52]) can be implemented using multiple threads as well, the proposed PTAV fundamentally differs from them. In these existing algorithms, different components are simultaneously used to determine tracking result in each frame, thus their multi-thread implementations are *synchronous*. By contrast, in PTAV,  $\mathcal{T}$  and  $\mathcal{V}$  function in different ways and run independent except for necessary interactions, and thus the multi-thread implementation is *asynchronous*.



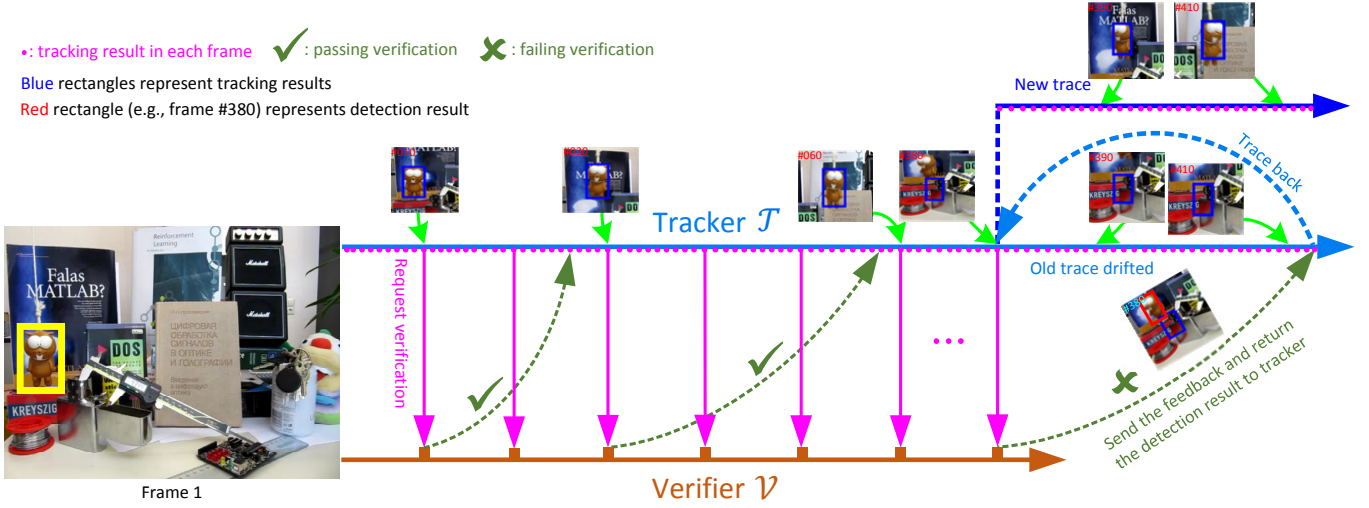


Fig. 3. Illustration of the PTAV framework in which tracking and verifying are processed asynchronously in two parallel threads.

### 3 PARALLEL TRACKING AND VERIFYING (PTAV)

#### 3.1 Framework

A typical PTAV consists of two components: a (fast) tracker  $\mathcal{T}$  and a (reliable) verifier  $\mathcal{V}$ . The two components work together toward real-time and high accuracy tracking.

- **The tracker  $\mathcal{T}$**  is responsible of the “real-time” requirement of PTAV, and needs to locate the target in each frame. Meanwhile,  $\mathcal{T}$  sends verification request to  $\mathcal{V}$  from time to time (though not every frame), and responds to feedback from  $\mathcal{V}$  by adjusting tracking or updating models. To avoid heavy computation,  $\mathcal{T}$  maintains a buffer of tracking information (e.g., intermediate status) in recent frames to facilitate fast tracing back when needed.
- **The verifier  $\mathcal{V}$**  is employed to pursue the “high accuracy” requirement of PTAV. Up on receiving a request from  $\mathcal{T}$ ,  $\mathcal{V}$  tries the best to first validate the tracking result (e.g., comparing it with the template), and then provide feedback to  $\mathcal{T}$ . To adapt  $\mathcal{V}$  to object appearance variations over time, the tracking target template is *not* fixed. Instead,  $\mathcal{V}$  collects a number of reliable tracking results, and then use  $k$ -means to cluster these results to obtain a target template pool for subsequent verification.

In PTAV,  $\mathcal{T}$  and  $\mathcal{V}$  run in parallel on two different threads with necessary interactions, as illustrated in Figure 3. The tracker  $\mathcal{T}$  and verifier  $\mathcal{V}$  are initialized in the first frame. After that,  $\mathcal{T}$  starts to process each arriving frame and generates the result (pink solid dot in Figure 3). In the meantime,  $\mathcal{V}$  validates the tracking result every several frames. Because tracking is much faster than verifying,  $\mathcal{T}$  and  $\mathcal{V}$  work asynchronously. Such mechanism allows PTAV to tolerate temporary tracking drift (e.g., at frame 380 in Figure 3), which will be corrected later by  $\mathcal{V}$ . When  $\mathcal{V}$  finds a tracking result unreliable, it searches the correct answer from a local region and sends it to  $\mathcal{T}$ . Upon the receipt of such feedback,  $\mathcal{T}$  stops current tracking job and traces back to resume tracking with the correction provided by  $\mathcal{V}$ .

It is worth noting that PTAV provides a very flexible framework, and some important designing choices are following. (1) The base algorithms for  $\mathcal{T}$  and  $\mathcal{V}$  may depend on specific applications and available computational resources. In addition, in

---

#### Algorithm 1: Parallel Tracking and Verifying (PTAV)

---

- 1 Initialize the tracking thread for tracker  $\mathcal{T}$ ;
  - 2 Initialize the verifying thread for verifier  $\mathcal{V}$ ;
  - 3 Initialize *current\_frame* as the second frame;
  - 4 Run  $\mathcal{T}$  (Alg. 2) and  $\mathcal{V}$  (Alg. 3) till the end of tracking;
- 

---

#### Algorithm 2: Tracking Thread $\mathcal{T}$

---

```

1 while current_frame is valid do
2   if received a message from  $\mathcal{V}$  then
3     if verification passed then
4       Update tracking model (optional);
5     else
6       Correct tracking;
7       Trace back and reset current_frame;
8   end
9 end
10 Tracking on the current_frame;
11 if time for verification according to  $N_{\text{int}}$  then
12   Send the current result to  $\mathcal{V}$  to verify;
13 end
14 current_frame  $\leftarrow$  next frame;
15 end

```

---



---

#### Algorithm 3: Verifying Thread $\mathcal{V}$

---

```

1 while not ended do
2   if received request from  $\mathcal{T}$  then
3     Verifying the tracking result;
4     Collect tracking result and perform  $k$ -means
       clustering if needed;
5   if verification failed then
6     Provide correction information in  $s$ ;
7     Adjust  $N_{\text{int}}$  if needed;
8   end
9   Send verification result  $s$  to  $\mathcal{T}$ ;
10 end
11 end

```

---

practice one may use more than one verifiers or even base trackers. (2) The response of  $\mathcal{T}$  to the feedback from  $\mathcal{V}$ , either positive or negative, can be largely designed to adjust to specific requests. (3) The correction of unreliable tracking results can be implemented in various ways, and it can even be conducted purely by  $\mathcal{T}$  (i.e., including target detection). (4)  $\mathcal{T}$  has numerous methods to use pre-computed and archived information for speeding up. Algorithms 1-3 summarize the general PTAV framework.

### 3.2 PTAV Implementation

#### 3.2.1 Tracking

We choose the Staple tracker [21] for  $\mathcal{T}$  in PTAV. The main idea of Staple is to combine two complementary cues, i.e., template and histogram, for tracking. To such end, given an image patch  $\mathbf{z}$ , a linear combination of tracking scores from template and histogram is proposed

$$y(\mathbf{z}) = (1 - \alpha)y_{\text{templ}}(\mathbf{z}) + \alpha y_{\text{hist}}(\mathbf{z}) \quad (1)$$

where  $\alpha$  denotes a trade-off parameter, and  $y_{\text{templ}}(\mathbf{z})$  and  $y_{\text{hist}}(\mathbf{z})$  represent the tracking responses based on template and on histogram information, respectively.

The tracking response on template is derived by learning the optimal correlation filter model  $\mathbf{w}$ , which is efficiently solved in frequency domain through the fast Fourier transformation (FFT). At time  $t$ , the FFT of the filter responses is first calculated using  $\mathbf{w}$  and an inverse FFT is then conducted to derive the final response  $y_{\text{templ}}(\mathbf{z})$ . The model  $\mathbf{w}$  is online updated in each frame.

The tracking response on histogram is based on a learned color statistic model  $\mathbf{h}$ , which is robust in resisting deformation. At time  $t$ ,  $\mathbf{h}$  is utilized to calculate  $y_{\text{hist}}(\mathbf{z})$ , and then dynamically updated. To adapt the tracker to scale changes, a scale filter is adopted to estimate the target scale. More details about the Staple tracker can be found in [21].

To efficiently leverage Staple as  $\mathcal{T}$  in PTAV, in addition to the original Staple algorithm,  $\mathcal{T}$  stores all intermediate results (e.g.,  $\mathbf{w}$  and  $\mathbf{h}$ ) for each frame after sending out last verification request. Let  $\mathcal{W} = \{\mathbf{w}_{\xi-\Delta}, \dots, \mathbf{w}_{\xi}\}$  and  $\mathcal{H} = \{\mathbf{h}_{\xi-\Delta}, \dots, \mathbf{h}_{\xi}\}$  represent the collections of  $\mathbf{w}$  and  $\mathbf{h}$ , where  $\xi$  is the index of the last frame processed by  $\mathcal{T}$ , and  $\Delta$  denotes a fixed size for temporal sliding window to store tracking models. These intermediate results in  $\mathcal{W}$  and  $\mathcal{H}$  allow  $\mathcal{T}$  for fast tracing back.

In particular, when  $\mathcal{V}$  detects unreliable tracking result in frame  $k$  while  $\mathcal{T}$  starts working on frame  $j$  ( $j > k$ ), a feedback consisting of correct target position and frame index information is sent to  $\mathcal{T}$ . Once receiving this feedback from  $\mathcal{V}$ ,  $\mathcal{T}$  first stops processing frame  $j$  and then utilizes the archived target position and tracking model (i.e.,  $\mathbf{w}_{k-1}$  and  $\mathbf{h}_{k-1}$ ) retrieved from  $\mathcal{W}$  and  $\mathcal{H}$  to resume subsequent tracking from frame  $k$ . Meanwhile, useless intermediate results in  $\mathcal{W}$  (i.e.,  $\mathbf{w}_k$  to  $\mathbf{w}_{j-1}$ ) and  $\mathcal{H}$  (i.e.,  $\mathbf{h}_k$  to  $\mathbf{h}_{j-1}$ ) will be discarded.

Note that we do not assume the correctness of  $\mathbf{w}_{k-1}$  and  $\mathbf{h}_{k-1}$  in the above strategy. In fact, one way is to trace backward from  $k-1$  to locate a reliable frame to resume tracking, at additional expense of more verification operations. In practice, however, we found that  $\mathbf{w}_{k-1}$  and  $\mathbf{h}_{k-1}$  typically provide sufficient initial guess for frame  $k$  and rely on the detection part to correct the incorrect tracking result. More details are given the following sections on verifying and detection.

To validate the tracking result,  $\mathcal{T}$  sends the verification results every  $N_{\text{int}}$  frames, where  $N_{\text{int}}$  denotes the dynamically adjustable verification interval as described later.

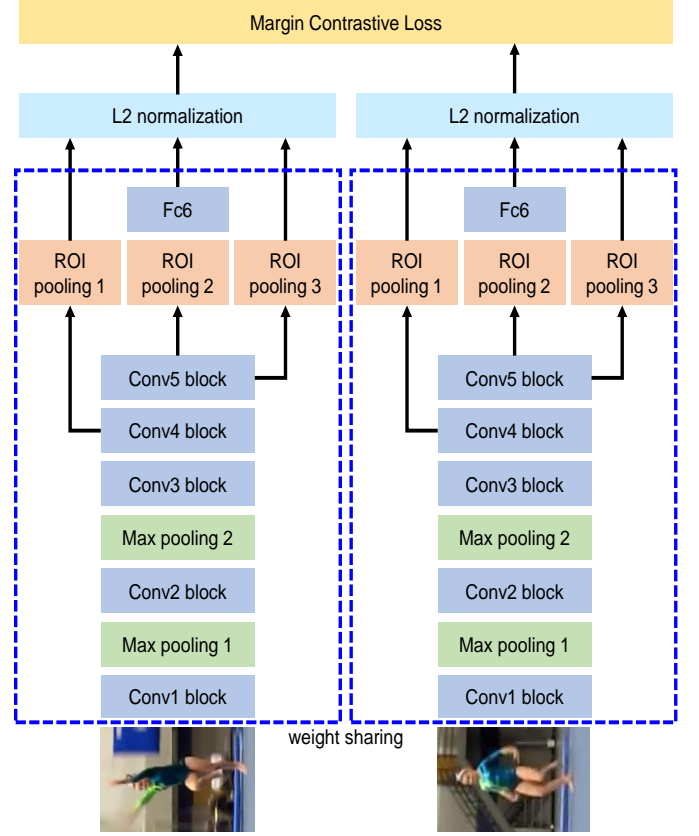


Fig. 4. Illustration of the architecture of the Siamese network for verifier.

#### 3.2.2 Verifying

The goal of verifying is to measure the similarity between a given sample and the target object. Inspired by [32], we use the Siamese network to develop verifier  $\mathcal{V}$  (similar to [14]) in PTAV, as depicted in Figure 4. The Siamese network contains two branches of CNNs, and processes two inputs separately. In this work, we borrow the architecture from VGGNet [37] for CNNs, but with an additional region of interest (RoI) pooling layer [53]. This is because, for detection,  $\mathcal{V}$  needs to process multiple regions in an image, from them the candidate most similar to the target object is selected as the final result. For efficiency, RoI pooling is used for simultaneously processing a set of regions.

In the Siamese network, the two CNN branches are connected with a single contrastive loss layer

$$\mathcal{L}(\mathbf{x}_i, \mathbf{x}_j, r_{ij}) = \frac{1}{2} r_{ij} D^2 + \frac{1}{2} (1 - r_{ij}) \max(0, \varepsilon - D^2) \quad (2)$$

where  $D = \|\psi(\mathbf{x}_i) - \psi(\mathbf{x}_j)\|_2$  is the Euclidean metric in which  $\psi(\cdot)$  represents feature transformation via the Siamese network,  $r_{ij} \in \{0, 1\}$  indicates that whether  $\mathbf{x}_i$  and  $\mathbf{x}_j$  are the same object or not, and  $\varepsilon$  represents the minimum distance margin.

Once training is finished<sup>2</sup>, one can use the learned verifying function  $\nu$  to compute verification score for each tracking result  $\mathbf{x}'$  via

$$\nu(\mathbf{x}_{\text{obj}}, \mathbf{x}') = \psi(\mathbf{x}_{\text{obj}})^T \psi(\mathbf{x}') \quad (3)$$

where  $\mathbf{x}_{\text{obj}}$  represents a fixed target template in the first frame. This strategy, as used in our preliminary work [36], may meet

2. In this work we adopt the same strategy as in [14] to train the verifier. We refer readers to [14] for detailed training process.

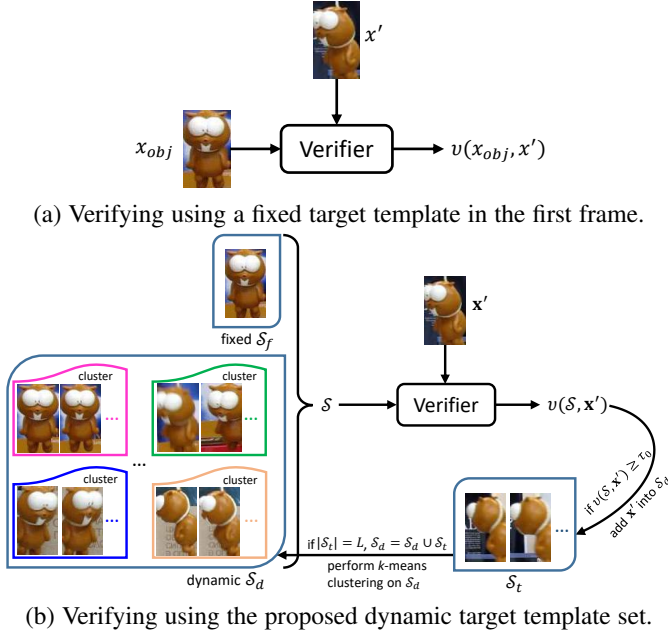


Fig. 5. Illustration of different verifying strategies in [36] (see image (a)) and this paper (see image (b)).

problems when the target object undergoes large appearance variations or deformations. As a result, using a fixed target template for verification may be unreliable for distant subsequent tracking results.

To alleviate this issue, we propose to employ a dynamic target template set  $\mathcal{S}$  for adaptive verification using  $k$ -means clustering. More specifically,  $\mathcal{S}$  is comprised of two components  $\mathcal{S}_f$  and  $\mathcal{S}_d$ . The  $\mathcal{S}_f = \{\mathbf{x}_{obj}\}$  contains only the target template  $\mathbf{x}_{obj}$  in the first frame and is fixed during tracking. The set  $\mathcal{S}_d$  is initially empty. During tracking, it is dynamically updated by collecting tracking results with high verification scores, as described later.

With the dynamic set  $\mathcal{S}$ , we can compute the verification score for each tracking result  $\mathbf{x}'$  as follows

$$\nu(\mathcal{S}, \mathbf{x}') = \omega_o \psi(\mathbf{x}_{obj})^T \psi(\mathbf{x}') + \omega_c \sum_{i=1}^{N_C} \sum_{\mathbf{x}_j \in C_i} \psi(\mathbf{x}_j)^T \psi(\mathbf{x}') \quad (4)$$

where  $\omega_o$  denotes the weight for  $\mathcal{S}_f$ ,  $\omega_c$  represents the weight for each cluster  $C_i$  obtained by performing  $k$ -means clustering on  $\mathcal{S}_d$ , and  $N_C = |\mathcal{S}_d|/L$  is the number of clusters ( $L$  is roughly a predefined size of each cluster, and  $|\mathcal{S}_d|$  denotes the size of  $\mathcal{S}_d$ ). The weights  $\omega_o$  and  $\omega_c$  are calculated as

$$\omega_o = \frac{\exp(0.5)}{\exp(0.5) + N_C \times \exp(0.5/N_C)} \quad (5)$$

$$\omega_c = \frac{1}{N_C} (1 - \omega_o) \quad (6)$$

The set  $\mathcal{S}_d$  is updated as follows. For each tracking result  $\mathbf{x}'$ , we use Equ. 4 to calculate its verification score  $\nu(\mathcal{S}, \mathbf{x}')$ . If  $\nu(\mathcal{S}, \mathbf{x}')$  is greater than a predefined threshold  $\tau_0$ , we treat  $\mathbf{x}'$  as a reliable target template and add it into a temporal set  $\mathcal{S}_t$ . This process is repeated until the number of elements in  $\mathcal{S}_t$  is equal to  $L$ . We then move all elements in  $\mathcal{S}_t$  to  $\mathcal{S}_d$  and thus leave  $\mathcal{S}_t$  empty. If the number of elements in  $\mathcal{S}_d$  is greater than  $L \times N_{C_{max}}$ , where  $N_{C_{max}}$  denotes the maximum number of clusters,

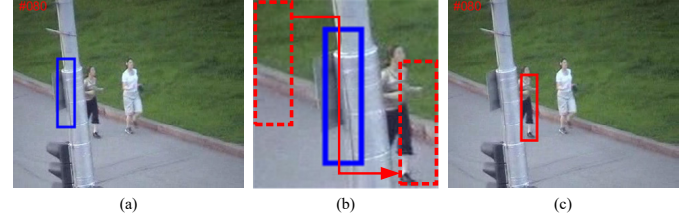


Fig. 6. Verification-based detection. When an unreliable tracking result is found (showing in blue in (a)), the verifier  $\mathcal{V}$  searches/detects the target in a local region (shown in (b)). The dashed red rectangles in (b) represent object candidates generated by sliding window. The red rectangle in (c) is the detection result.

the oldest  $L$  elements will be removed from  $\mathcal{S}_d$ . Afterwards,  $k$ -means clustering [54] is applied on  $\mathcal{S}_d$  to obtain new clusters

$$\{C_i\}_{i=1}^{N_C} = k\text{-means}(\mathcal{S}_d, N_C) \quad (7)$$

Note that when performing  $k$ -means clustering, the elements in  $\mathcal{S}_d$  are represented with HoG features [25] for the sake of efficiency. After obtaining new clusters, we employ Equ. 5 and 6 to calculate weights  $\omega_o$  and  $\omega_c$ . Figure 5 illustrates the process of adaptive verification.

With the dynamic target template set  $\mathcal{S}$ ,  $\mathcal{V}$  can make smarter decisions than when only a fixed template is used (i.e., Equ. 3), and hence reduces the number of unnecessary verifications to speed up the entire system. Besides, now that verification is more precise, the verification-based detection (see Section 3.2.3) is improved as well.

### 3.2.3 Verification-based detection

Given a tracking result from  $\mathcal{T}$ , we use Equ. 4 to compute its verification score. If the verification score is lower than a predefined threshold  $\tau_1$ ,  $\mathcal{V}$  will treat it as a tracking failure. In this case,  $\mathcal{V}$  needs to detect the target, again using the Siamese network. Unlike for verification, detection requires to verify multiple image patches from a local region<sup>3</sup> and finds the best one. Thanks to the RoI pooling layer, these candidates can be simultaneously processed in just one pass, resulting in significant reduction in computation. Let  $\{\mathbf{c}_i\}_{i=1}^N$  denote the candidate set generated by sliding window, and the detection result  $\hat{\mathbf{c}}$  is determined by

$$\hat{\mathbf{c}} = \underset{\mathbf{c}_i}{\operatorname{argmax}} \nu(\mathcal{S}, \mathbf{c}_i), \quad i = 1, 2, \dots, N \quad (8)$$

where  $\nu(\mathcal{S}, \mathbf{c}_i)$  returns the verification score between the target template set  $\mathcal{S}$  and candidate  $\mathbf{c}_i$ .

After obtaining the detection result  $\hat{\mathbf{c}}$ , we determine whether or not to take it to be an alternative for tracking result according to its verification score. If  $\nu(\mathcal{S}, \hat{\mathbf{c}})$  is less than a predefined threshold  $\tau_2$ ,  $\hat{\mathbf{c}}$  is considered to be unreliable, and we do not replace tracking result with  $\hat{\mathbf{c}}$ . Instead, we decrease the verifying interval  $N_{int}$  to 1, and enlarge the local searching region for target detection. Until detection result  $\hat{\mathbf{c}}$  passes verification (i.e.,  $\nu(\mathcal{S}, \hat{\mathbf{c}}) \geq \tau_2$ ), we then restore  $N_{int}$  and the size of local searching region to initial settings. Figure 6 describes the detection process.

3. The local region is a square of size  $\gamma(w^2 + h^2)^{\frac{1}{2}}$  centered at the location of the tracking result in this validation frame, where  $w$  and  $h$  denote the width and height of the tracking result, and  $\gamma$  controls the scale and is dynamically adjusted based on detection result.

### 3.3 Implementation Details

Our PTAV is implemented in C++ and its verifier uses Caffe [55] on a single NVIDIA GTX TITAN Z GPU with 6GB memory. The merging factor  $\alpha$  in Eq. (1) is set to 0.3. Other parameters for tracking remain the same as in [21]. The Siamese network for verification is initialized with the VGGNet [37] and trained based on the approach in [14]. The clustering interval  $L$  is empirically set to 5 and the maximum number of clusters  $N_{C_{\max}}$  to 10. The verification interval  $N_{\text{int}}$  is initially set to 10. The thresholds  $\tau_0$ ,  $\tau_1$  and  $\tau_2$  are set to 0.6, 0.33 and 0.53, respectively. The parameter  $\gamma$  is initialized to 1.5, and is adaptively adjusted based on the detection result. If the detection result with  $\gamma = 1.5$  is not reliable, the verifier will increase  $\gamma$  for a larger searching region. Meanwhile, the verification interval  $N_{\text{int}}$  is decreased to 1. When the new detection result becomes faithful,  $\gamma$  and  $N_{\text{int}}$  are then restored to 1.5 and 10.

The source code of our implementation, as well as tracking results, are made publicly available at <http://www.dabi.temple.edu/~hbling/code/PTAV/ptav.htm>.

## 4 EXPERIMENTS

### 4.1 Experiment on OTB2015

**Dataset and evaluation settings.** The OTB2015 benchmark [29] contains 100 fully annotated challenging video sequences. These sequences are labeled based on 11 attributes, including deformation (DEF), occlusion (OCC), scale variation (SV), illumination variations (IV), motion blur (MB), fast motion (FM), background clutter (BC), out-of-view (OV), low resolution (LR), in-plane rotation (IPR) and out-of-plane rotation (OPR).

Following the protocol in [29], we use three metrics, *distance precision rate* (DPR), *overlap success rate* (OSR) and *center location error* (CLE), to evaluate different tracking algorithms. DPR demonstrates the percentage of frames whose estimated average center location errors are within the given threshold distance (e.g., 20 pixels) to groundtruth. OSR shows the percentage of successful frames at the threshold ranging from 0 to 1, and can be defined as the overlap score more than a fixed value (e.g., 0.5), where the overlap ratio is defined as  $\text{score} = \text{area}(R_{GT} \cap R_T) / \text{area}(R_{GT} \cup R_T)$  with the groundtruth  $R_{GT}$  and tracking result  $R_T$ . CLE represents the Euclid distance between centers of tracking result and groundtruth.

#### 4.1.1 Overall performance

We evaluate PTAV on OTB2015 [29] and compare it with twelve state-of-the-art trackers from three typical categories: (i) deep feature-based tracking algorithms, including SINT [14], HCF [7], SiamFC [28], HDT [8] and CFNet [50]; (ii) correlation filter based trackers, including fDSST [19], LCT [27], KCF [17] and Staple [21]; and (iii) other representative tracking methods, including TLD [24], MEEM [56] and Struck [40]. We also note that there are other state-of-the-art tracking algorithms such as MDNet [9], SANet [12] and C-COT [13] (see Figure 1). However, the speeds of these trackers are around 1 frames per second (fps). Since this work is focused on real-time object tracking, we compare PTAV with trackers whose speeds are no less than 10 fps, except for SINT [14] since it can be viewed as the baseline for tracking by verification. In particular, for SINT [14], we use its tracking results without optical flow because no optical flow part is provided from the released source code. Another baseline for PTAV is the Staple

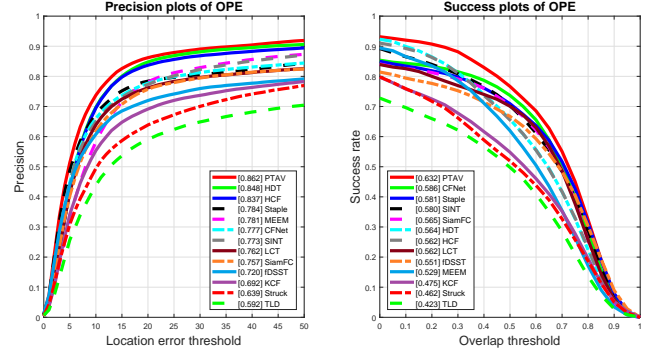


Fig. 7. Comparison with pseudo real-time trackers on OTB2015 [29] using distance precision rate (DPR) and overlap success rate (OSR).

TABLE 1

Comparisons with pseudo real-time tracking methods on OTB2015 [29] in distance precision rate (DPR%) at a threshold of 20 pixels, overlap success rate (OSR%) at an overlap threshold of 0.5, center location error (CLE) in pixels and speed (fps). The best two results are highlighted in red and blue fonts, respectively

	Algorithms	DPR	OSR	CLE	Speed
(i)	PTAV (Ours)	<b>86.2</b>	<b>77.9</b>	<b>18.9</b>	27
	HCF [7]	83.7	65.6	22.8	10
	HDT [8]	<b>84.8</b>	64.9	<b>20.6</b>	10
	SINT [14]	77.3	70.3	26.3	2
	SiamFC [28]	75.7	70.9	37.1	<b>58</b>
	CFNet [50]	77.7	73.2	35.2	43
(ii)	Staple [21]	78.4	70.9	31.9	43
	LCT [27]	76.2	70.1	67.1	25
	fDSST [19]	72.0	67.6	51.1	51
	KCF [17]	69.2	54.8	45.0	<b>243</b>
(iii)	MEEM [56]	78.1	62.2	27.7	21
	TLD [24]	59.2	48.3	35.0	20
	Struck [40]	63.9	51.6	47.1	10

tracker [21], which provides the (fast) tracking part of PTAV. It is worth noting that other tracking algorithms may also be used for the tracking part in PTAV.

We report the results in one-pass evaluation (OPE) using DPR and OSR as shown in Figure 7. Overall, PTAV performs favorably against other tracking algorithms. In addition, we present quantitative comparison of DPR at 20 pixels, OSR at 0.5, center location error (CLE) in pixels and tracking speed (fps) in Table 1. It demonstrates that PTAV outperforms other trackers in all three metrics. Among the trackers under comparison, HCF [7] utilizes deep hierarchical features to represent object appearance and obtains the DPR of 83.7% and OSR of 65.6%. Likewise, HDT [8] exploits all layers in VGGNet [37] for tracking and achieves the DPR of 84.8% and OSR of 64.8%. Compared to these two deep feature-based approaches, our tracker achieves better performance with DPR of 86.2% and OSR of 77.9%. Besides, owing to the adoption of parallel framework, PTAV (27 fps) is more than twice faster than the HCF [7] (10 fps) and HDT [8] (10 fps). Compared with SINT [14], which uses similar Siamese network for tracking, PTAV improves DPR from 77.3% to 86.2% and OSR from 70.3% to 77.6%. In addition, PTAV runs at real-time while SINT [14] needs large improvement in speed. Compared with the baseline Staple [21], PTAV achieves significant improvements on DPR (by 7.8%) and OSR (by 7.0%). Compared to representative MEEM [56] with DPR of 78.1% and OSR of 62.2%, PTAV obtains



TABLE 2

Average DPR (%) in terms of individual attributes on OTB2015 [29]. The best two results are highlighted in red and blue fonts, respectively

Attributes	PTAV	HDT [8]	HCF [7]	Staple [21]	MEEM [56]	CFNet [50]	SINT [14]	SiamFC [28]	LCT [27]	fDSST [19]	KCF [17]	Struck [40]	TLD [24]
IV	<b>84.7</b>	<b>82.0</b>	81.7	79.1	74.0	75.7	80.9	73.5	74.6	72.8	70.8	54.9	55.9
OPR	<b>83.5</b>	80.8	<b>81.0</b>	74.2	79.8	75.3	79.4	74.5	75.0	66.4	67.5	59.9	57.1
SV	<b>82.5</b>	<b>81.1</b>	80.2	73.1	74.0	74.8	74.2	74.3	68.6	66.9	63.9	60.4	56.4
OCC	<b>81.4</b>	<b>77.4</b>	76.7	72.6	74.1	71.3	73.1	69.6	68.2	62.6	62.2	53.3	52.4
DEF	<b>81.2</b>	<b>82.1</b>	79.1	74.8	75.4	66.9	75.0	67.6	68.9	59.9	61.7	52.7	48.4
MB	<b>80.5</b>	79.4	<b>79.7</b>	72.6	72.1	76.1	72.8	69.8	67.3	68.4	61.7	59.4	53.6
FM	78.1	<b>80.6</b>	<b>79.7</b>	70.3	73.4	74.1	72.5	73.0	67.5	69.3	62.8	62.0	54.8
IPR	82.8	<b>84.4</b>	<b>85.4</b>	77.0	79.3	80.3	81.1	74.8	78.2	72.5	69.3	63.4	60.3
OV	<b>79.0</b>	66.3	67.7	66.1	68.3	65.0	72.5	67.8	59.2	57.7	49.8	49.1	45.2
BC	<b>87.6</b>	<b>84.7</b>	<b>84.7</b>	77.0	75.1	73.7	75.1	69.4	74.0	78.4	71.6	57.3	46.1
LR	<b>84.0</b>	76.6	78.7	60.9	60.5	<b>86.1</b>	78.8	83.4	49.0	61.7	54.5	62.8	55.2
Overall	<b>86.2</b>	<b>84.8</b>	83.7	78.4	78.1	77.7	77.3	75.7	76.2	72.0	69.2	63.9	59.2

TABLE 3

Average OSR (%) in terms of individual attributes on OTB2015 [29]. The best two results are highlighted in red and blue fonts, respectively.

Attributes	PTAV	HDT [8]	HCF [7]	Staple [21]	MEEM [56]	CFNet [50]	SINT [14]	SiamFC [28]	LCT [27]	fDSST [19]	KCF [17]	Struck [40]	TLD [24]
IV	<b>64.2</b>	53.5	54.0	59.8	51.7	57.4	<b>61.8</b>	54.9	56.6	55.6	47.4	42.0	41.4
OPR	<b>60.4</b>	53.6	53.7	53.8	52.8	55.3	<b>58.6</b>	54.4	54.1	50.1	45.4	42.7	39.0
SV	<b>59.1</b>	48.9	48.8	52.9	47.3	55.5	<b>55.8</b>	55.5	49.2	51.0	39.9	40.7	38.8
OCC	<b>60.6</b>	52.8	52.5	54.8	50.3	53.6	<b>55.8</b>	52.3	50.7	47.8	43.8	39.3	36.3
DEF	<b>59.9</b>	54.3	53.0	55.4	48.9	49.2	<b>55.5</b>	49.0	49.9	46.1	43.6	38.3	34.1
MB	<b>61.2</b>	56.3	57.3	55.8	54.3	<b>59.3</b>	57.4	55.5	53.2	54.8	45.6	46.1	42.6
FM	<b>58.3</b>	55.4	55.5	54.1	52.8	<b>57.0</b>	55.7	56.4	52.7	55.4	45.5	46.1	41.8
IPR	<b>59.0</b>	55.5	55.9	55.2	52.8	<b>59.0</b>	<b>58.5</b>	55.7	55.7	54.5	46.5	45.2	42.5
OV	<b>56.9</b>	47.2	47.4	48.1	48.4	48.0	<b>55.9</b>	51.1	45.2	45.7	39.3	37.8	33.5
BC	<b>64.1</b>	58.0	<b>58.7</b>	57.4	52.1	54.5	56.7	50.4	55.3	58.5	49.8	44.2	35.2
LR	54.6	42.0	42.4	41.1	35.5	<b>61.9</b>	53.9	<b>60.4</b>	33.0	44.6	30.6	34.7	37.2
Overall	<b>63.2</b>	56.4	56.2	58.1	52.9	<b>58.6</b>	58.0	56.5	56.2	55.1	47.5	46.2	42.3

performance gains by DPR of 8.1% and OSR of 15.7%.

#### 4.1.2 Attribute-based evaluation

We further analyze the performance of PTAV under the eleven different attributes on OTB2015. In Tables 2 and 3, we summarize the evaluation results in terms of DPR and OSR.

For DPR, PTAV achieves the best results under 7 out of 11 attributes including IV (84.7%), OPR (83.5%), SV (82.5%), OCC (81.4%), MB (80.5%), OV (79.0%) and BC (87.6%). For sequences with deformation, HDT [8] performs the best with average DPR of 82.1% owing to the use of richer deep features. Our tracker uses Staple [21] as the tracking part, which leverages color information to handle object deformation. Accompanied by an accurate verifier, PTAV achieves competitive performance and ranks the second in the case of deformation with average DPR of 81.2%. Furthermore, compared to the baseline Staple [21], we obtain large gain on the average DPR by 6.4% under deformation. For low resolution sequences, CFNet [50] obtains the best result with average DPR of 86.1% by taking advantage of deep features. PTAV ranks the second with competitive average DPR of 84.0%.

For sequences with fast motion and in-plane rotation, the two deep feature-based trackers HDT [8] and HCF [7] perform better than ours, because in these two situations deep features are more efficacious than hand-crafted features to represent object appearance. In PTAV, we utilize simple HoG [25] features and RGB histograms to model object appearance in tracking, which are sensitive to in-plane rotation and fast motion (we can see that from the performance of Staple [21]). Nevertheless, with the help of useful feedbacks from a robust and accurate verifier, PTAV still achieves competitive performance with average DPRs of 78.1% and 82.8% under these two challenges, respectively.

For OSR, on the other hand, PTAV achieves the best results under 10 of 11 attributes including IV (64.2%), OPR (60.4%), SV (59.1%), OCC (60.6%), DEF (59.9%), MB (61.2%), FM (58.3%), IPR (59.0%), OV (56.9%) and BC (64.1%). Low resolution (LR) is the only attribute for which PTAV does not rank the best, while CFNet [50] and SiamFC [28] obtain better results than PTAV. Specifically, PTAV achieves an OSR of 54.6%, higher than all other trackers including its two baselines.

#### 4.1.3 Qualitative evaluation

To further analyze and demonstrate the performance of PTAV, we conduct rich qualitative evaluation described as following.

**Occlusions.** Figure 8(a) demonstrates the sampled tracking results on sequences *Box*, *Lemming*, *Girl2* and *Jogging-1*, all involving heavy target occlusions. In *Box*, the target undergoes not only occlusions but also background clutters (e.g., #476). In *Lemming*, the tracking target suffers from occlusions and motion blur (e.g., #312 and #489). In *Girl2*, the target is fully occluded by the background (e.g., #114). In sequence *Jogging-1*, the tracking target is heavily occluded with significant deformation (e.g., #72). From Figure 8(a) we can see that PTAV handles well the occlusion in these sequences. Though the tracking part may lose the target temporally because of occlusions, it can quickly be corrected by the verifier and resumes tracking. Compared to deep feature based trackers (HCF [7], HDT [8] and CFNet [50]), which lose the tracking target when occlusions happen (e.g., #342 in *Girl2* and #489 in *Lemming*), PTAV performs more robustly. Since correlation filter is sensitive to occlusion and no re-detection module is adopted, KCF [17], fDSST [19] and Staple [21] lose the tracking target in all sequences. LCT [27] applies an additional



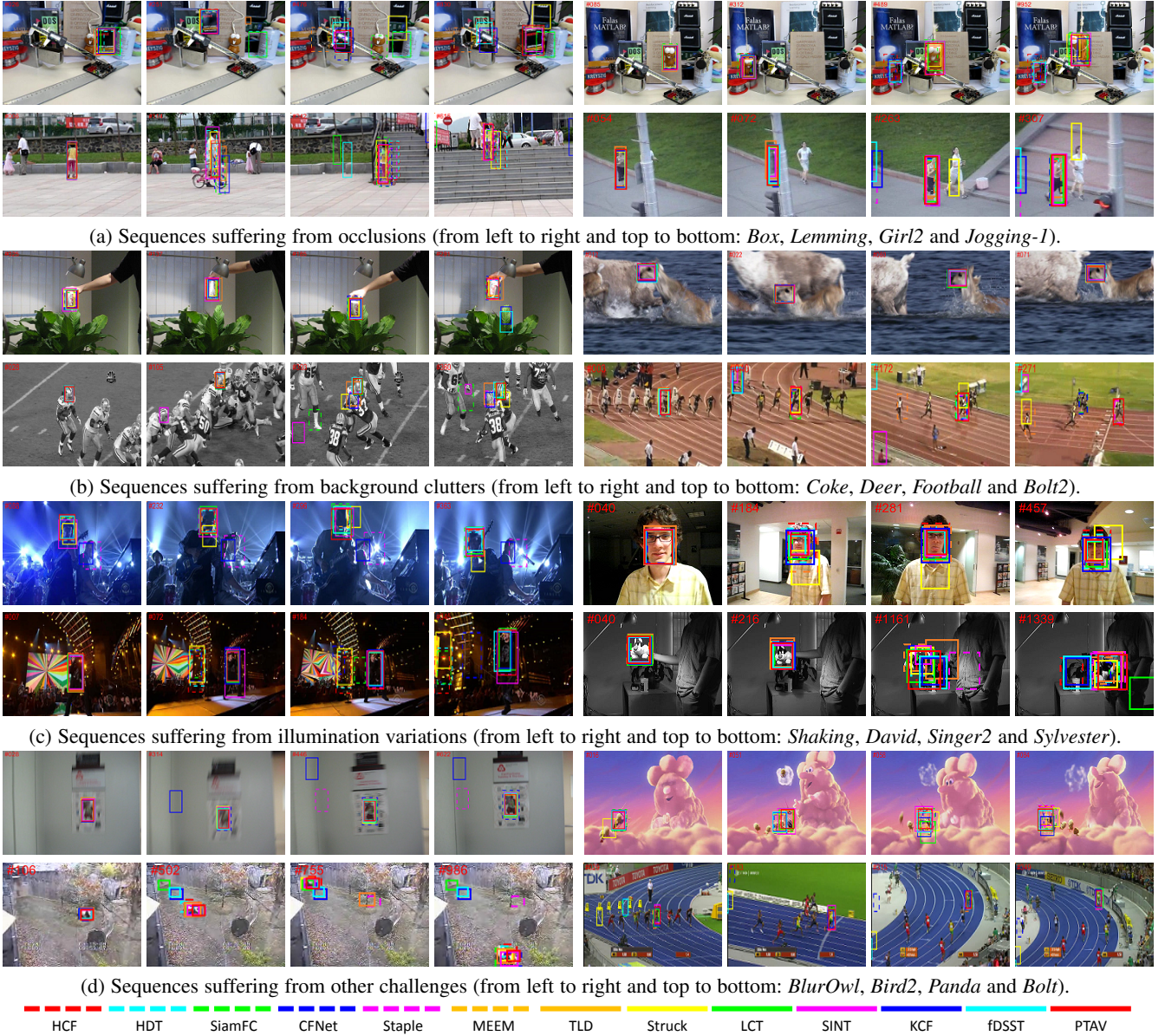


Fig. 8. Qualitative evaluation of the proposed algorithm and other twelve state-of-the-art trackers on sixteen challenging sequences.

detector to re-localize tracking target when it recovers from occlusion. Nevertheless, it still drifts to the background when occlusions are accompanied with background clutters (e.g., #476 in *Box*). MEEM [56] stores multiple memories of target during tracking, and re-detects the target using old memories when it recovers from occlusion (e.g., *Lemming*, *Girl2* and *Jogging-1*). However, it meets problems in presence of background clutters (e.g., #476 and #830 in *Box*). SiamFC [28] and SINT [14] deal with occlusion by searching in a local region when the target recovers from occlusions. In addition, no model update is adopted in SiamFC [28] and SINT [14], avoiding introducing background into the tracking model. Other approaches such as TLD [24] and Struck [40] does not deal well with occlusion and drift to background (e.g., #342 in *Girl2*).

**Background clutters.** Background clutters are prone to result in tracking drift and even failures because the tracker may mix the tracking target with cluttering background. Figure 8(b) displays the sampled experimental results on sequences *Coke*, *Deer*, *Foot-*

*ball* and *Bolt2*. In *Coke*, the object appearance is visually similar to the background, and occlusion and illumination variation occur as well. In *Deer*, *Football* and *Bolt2*, there exist similar distracters involved with other challenges including occlusions (e.g., *Football*), deformations (e.g., *Bolt2*) and motion blur (e.g., *Deer*). From Figure 8(b), we can observe that SiamFC [28], SINT [14], KCF [17], Struck [40] and Staple [21] drift to background in *Football* (e.g., #360). CFNet [50] and fDSST [19] localize well the target in *Football*, but still fail in *Coke* (e.g., #291) where background clutters are accompanied with occlusions and illumination change. MEEM [56] handles nicely occlusion and motion blur, but fails in presence of deformations in *Bolt2* (e.g., #40). LCT [27] and TLD [24] are robust against background clutters since they utilize an extra object detector to seek for the target after it moves away from the cluttering region. However, they lose the tracking target when heavy deformation happens in *Bolt2* (e.g., #271). HCF [7] and HDT [8] perform robustly in these sequences because of the powerful deep features. Likewise, our tracker is able to deal with

these situations owing to the verifier. Besides, the cooperation between tracking and verifying allows PTAV to run in real-time.

**Illumination variations.** Illumination variation often causes drift problem. Figure 8(c) shows sampled results of sequences *Shaking*, *David*, *Singer2* and *Sylvester*. In *Shaking* and *Sylvester*, the target suffers from not only illumination variations but also background clutters and rotations. In *David* and *Singer2*, illumination variations are accompanied by rotations and scale changes. We can see from Figure 8(c) that deep feature-based trackers SiamFC [28], CFNet [50], HCF [7] and HDT [8] lose the target in *Singer2* (e.g., #345) where background clutters happen. MEEM [56] can handle rotations but still fails in *Singer2* (e.g., #345). Staple [21] and KCF [17] are sensitive to rotation and fail in *Shaking* (e.g., #363) and *Sylvester* (e.g., #1339). LCT works well in *Shaking*, *David* and *Singer2*, yet have problems when heavy rotation exists in *Sylvester* (e.g., #1339). Our tracker performs well on these sequences. Though its tracking part may drift to background due to rotation (e.g., #1161 in *Sylvester*), this situation is found and immediately corrected by its verifier (e.g., #1339 in *Sylvester*).

**Other challenges.** Figure 8(d) demonstrates results of sequences *BlurOwl*, *Bird2*, *Panda* and *Bolt2*, which contain other challenges including fast motion, motion blur, rotation, scale change, deformation and so forth. In *BlurOwl*, the camera moves quickly, causing serious motion blur (e.g., #622). KCF [17] and Staple [21] lose the target, while PTAV well localizes the tracking target thanks to its verifier which corrects tracker. In *Bolt2*, TLD [24], Struck [40], CFNet [50] and fDSST [19] drift to background because of deformations. On *Panda*, LCT [27], fDSST [19], Staple [21] and KCF [17] lose the tracking target owing to scale changes and rotations. By contrast, MEEM [56], HCF [7], HDT [8], SiamFC [28], SINT [14] and our PTAV performs favorably because of powerful feature representation.

## 4.2 Experiment on TC128

The TC128 benchmark [33] is comprised of 128 fully annotated challenging color sequences. On TC128 [33], PTAV runs at 24 *fps* and is with eleven state-of-the-art trackers including MEEM [56], HCF [7], HDT [8], Staple [21], SiamFC [28], SRDCF [49], DeepSRDCF [57], fDSST [19], KCF [17], LCT [27] and Struck [40]. Following the protocol in [33], we report evaluation results in OSR and DPR as shown in Figure 9.

Among the eleven compared trackers, DeepSRDCF [57] extends SRDCF [49] by replacing hand-crafted features with convolutional features and obtains the best performance with DPR of 74.0% and OSR of 53.6%. By contrast, PTAV improves the state-of-the-art methods on DPR to 77.2% and OSR to 56.3%, obtaining the gains of 3.2% and 2.7%, respectively. In comparison with SiamFC [28] with DPR of 69.6% and OSR of 49.7%, PTAV achieves improvements of 7.2% and 6.6% on DPR and OSR, respectively. Compared with the other baseline, Staple [21], which obtains a DPR of 66.7% and an OSR of 49.7%, PTAV achieves significant improvements as well, showing clearly the benefits of introducing a verifier. For more detailed analysis, we show the average DPR for five trackers on different attributes in Figure 10. PTAV can well handle various challenging factors and outperform the other four trackers in nine out of eleven attributes.

## 4.3 Experiment on UAV20L

The recently proposed UAV20L dataset [34] consists of 20 fully annotated sequences, with length ranging from 1,717 to 5,527

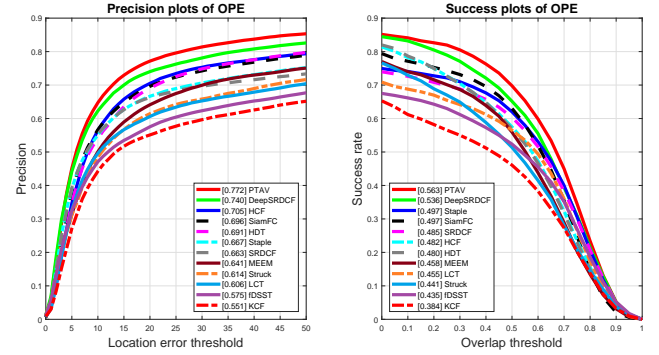


Fig. 9. Comparison with eleven state-of-the-art trackers on TC128 [33] using distance precision rate and overlap success rate.

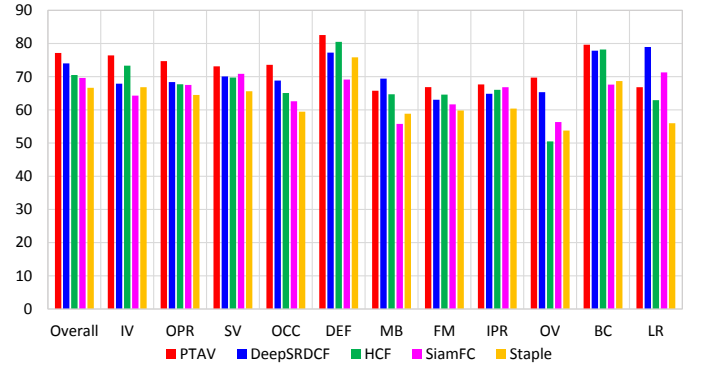


Fig. 10. Average DPR (%) in term of individual attributes on TC128 [33].

frames. These videos are challenging because the tracking target suffers from appearance changes caused by various factors. The proposed PTAF tracker runs at 30 *fps*, and compared with ten state-of-the-art trackers including SiamFC [28], MUSTer [52], SRDCF [49], HCF [7], MEEM [56], SAMF [46], Struck [40], fDSST [19], LCT [27] and KCF [17].

Following [34], we report evaluation results in Figure 11. PTAV achieves the best performance in both DPR (73.2%) and OSR (50.4%), outperforming other approaches by large margins (6.2% and 10.1% compared with the second best in DPR and OSR, respectively). Furthermore, we analyze the performance of PTAV on twelve individual attributes provided with UAV20L [34], including scale variation (SV), aspect ratio change (ARC), low resolution (LR), fast motion (FM), full occlusion (FOC), partial occlusion (POC), out-of-view (OV), background clutter (BC), illumination variation (IV), viewpoint change (VC), camera motion (CM) and similar object (SOB). Figure 12 displays the average DPR for five trackers on different attributes, and PTAV achieves the best results on each attribute.

## 4.4 Experiment on VOT2016

Finally, we test PTAV on the VOT2016 challenge [35], which contains 60 challenging sequences. VOT2016 aims at evaluating short-term tracking performance and thus a tracker is re-initialized whenever failure happens. In other words, a tracker is reset if its tracking results are found unreliable. Nevertheless, this protocol is not directly applicable to our tracker since PTAV automatically detects failures by itself and rolls back to resume tracking.



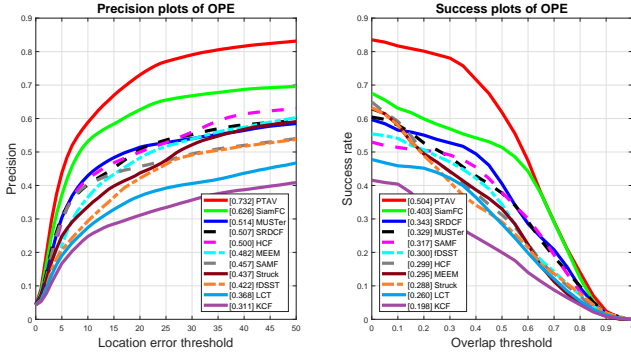


Fig. 11. Comparison with ten state-of-the-art trackers on UAV20L [34] using distance precision rate and overlap success rate.

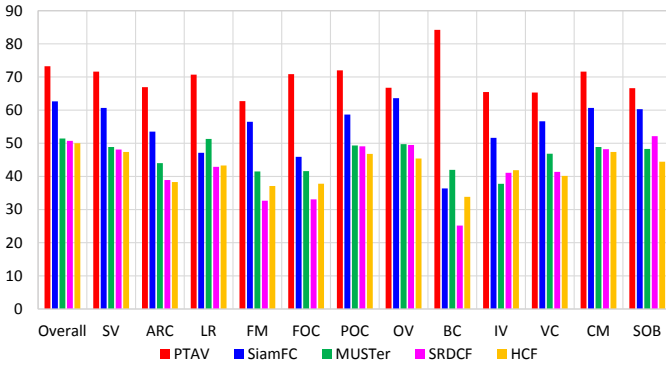


Fig. 12. Average DPR (%) in term of attributes on UAV20L [34].

TABLE 4

Comparisons with state-of-the-art tracking methods on VOT2016 [35] in terms of expected average overlap (EAO%), accuracy (%), robustness (%) and no-reset average overlap (AO%). The best two results are highlighted in red and blue fonts, respectively.

Algorithms	EAO	Accuracy	Robustness	AO
PTAV (Ours)	31.2	<b>56.1</b>	27.9	43.2
C-COT [13]	<b>33.1</b>	53.9	<b>23.8</b>	46.9
TCNN [58]	<b>32.5</b>	55.4	26.8	<b>48.5</b>
SSAT [35]	32.1	<b>57.7</b>	29.1	<b>51.5</b>
MLDF [35]	31.1	49.0	<b>23.3</b>	42.8
Staple [21]	29.5	54.4	37.8	38.8
DDC [35]	29.3	54.1	34.5	39.1
EBT [59]	29.1	46.5	25.2	37.0
SRBT [35]	29.0	49.6	35.0	33.3
Staple+ [35]	28.6	55.7	36.8	39.2
DNT [60]	27.8	51.4	32.9	42.7

To follow the above evaluation protocol, we modify PTAV by running it multiple rounds with different starting frames. In particular, in each round, we run PTAV at current starting frame without resetting. For the first round, the first frame in the input sequence is used as the start frame. Afterward, we compare the tracking results with groundtruth to find the first failure using the VOT2016 protocol, and then we re-initialize PTAV from the failure frame for the next round. We repeat this process until no failure is detected. On VOT2016, PTAV runs at 25 *fps*.

PTAV is compared with top ten trackers in the VOT2016 challenge, including C-COT [13], TCNN [58], SSAT [35], MLDF [35], Staple [21], DDC [35], EBT [59], SRBT [35], Staple+ [35] and DNT [60]. Table 4 demonstrates comparison results

in VOT2016. It shows that C-COT [13] and TCNN [58] achieve the best results with EAOs of 33.1% and 32.5%, respectively. C-COT [13] utilizes deep features to model object appearance and TCNN [58] proposes tree-structured CNNs for tracking with on-line update. Despite obtaining superior performances, their speeds are around 1 and 2 *fps*. By contrast, PTAV achieves competitive result (EAO of 31.2%), while running in real-time.

## 4.5 Ablation Study

### 4.5.1 Different trackers for $\mathcal{T}$

In PTAV,  $\mathcal{T}$  is required to be efficient and accurate most of the time. To demonstrate the effects of different  $\mathcal{T}$ , we compare three different base tracking algorithms including Staple [21] (the choice in this paper), fDSST [19] and KCF [17]. Among these trackers, KCF [17] runs the most efficiently while least accurately in short time. Compared with KCF [17] and fDSST [19], Staple [21] performs more robustly since it utilizes color information for tracking, which results in its relative inefficiency. The comparison results on OTB2015 [29], TC128 [33] and UAV20L [34] are shown in Table 5.

From Table 5, we can see that PTAV with the Staple base tracker (PTAV<sub>Staple</sub>) performs better than those with fDSST (PTAV<sub>fDSST</sub>) and KCF (PTAV<sub>KCF</sub>). Though KCF runs the fastest among these trackers, it performs least accurately in short time, resulting in more requests for verifications and detections, and significantly increased computations. As shown in Table 5, the speeds of PTAV<sub>KCF</sub> on OTB2015 [29], TC128 [33] and UAV20L [34] are respectively 24, 19 and 20 *fps*, which are much slower than PTAV<sub>Staple</sub> (27, 23 and 30 *fps*, respectively) and PTAV<sub>fDSST</sub> (27, 24 and 26 *fps*, respectively).

In terms of tracking accuracy, on OTB2015 [29], PTAV<sub>fDSST</sub> achieves competitive performance (85.2% of DPR and 77.9 % of OSR) compared to PTAV<sub>Staple</sub> (86.2% of DPR and 77.9% of OSR). However, on the more challenging UAV20L [34], PTAV<sub>Staple</sub> significantly outperforms PTAV<sub>fDSST</sub> in both accuracy and efficiency. Specifically, PTAV<sub>Staple</sub> obtains a DPR of 73.2%, an OSR of 62.4% and speed of 30 *fps* while PTAV<sub>fDSST</sub> with DPR of 63.1%, OSR of 47.8% and speed of 26 *fps*. The main reason accounting for this is that the baseline Staple leverages color cues for tracking. In UAV20L [34], the tracking target frequently suffers from severe view changes, which are fatal to HoG features. Nevertheless, Staple is able to deal with view changes using color statistics, and thus performs better than fDSST in short periods. As a consequence, PTAV<sub>Staple</sub> performs more favorably than PTAV<sub>fDSST</sub> and requires less verifications and detections, further improving efficiency. Besides, we observe from Table 5 that all the three PTAV versions improve their baseline trackers by large margins.

### 4.5.2 Different verifiers for $\mathcal{V}$

The verifier  $\mathcal{V}$  plays a crucial role in PTAV by validating tracking results and correcting  $\mathcal{T}$  if needed. To guarantee the quality of verification,  $\mathcal{V}$  is required to be as accurate as possible. To such purpose, we adopt the Siamese networks [32] for  $\mathcal{V}$ . To study the effects of different  $\mathcal{V}$ , we compare two different alternatives: one is based on VGGNet [37] in previous experiments, and the other utilizes the much lighter AlexNet [6]<sup>4</sup>. Compared to the VGGNet-based  $\mathcal{V}$ , the AlexNet-based  $\mathcal{V}$  runs more efficiently but

4. For the verifier based on AlexNet [6], one just needs to replace and initialize the five convolutional blocks in Figure 4 with AlexNet.



TABLE 5

Comparisons of DPR (%), OSR (%), CLE in pixels and speed (fps) among different  $\mathcal{T}$  with VGGNet [37] based  $\mathcal{V}$  on three benchmarks.

	OTB2015 [29]				TC128 [33]				UAV20L [34]			
	DPR	OSR	CLE	Speed	DPR	OSR	CLE	Speed	DPR	OSR	CLE	Speed
PTAV <sub>Staple</sub>	86.2	77.9	18.9	27	77.2	70.0	32.1	23	73.2	62.4	56.1	30
PTAV <sub>fDSST</sub>	85.2	77.9	19.4	27	75.2	66.1	32.7	24	63.1	47.8	102.7	26
PTAV <sub>KCF</sub>	74.2	58.6	34.5	24	63.9	54.6	53.6	19	41.6	32.1	195.3	20
Staple [21]	78.4	70.9	31.9	43	66.7	62	57.5	43	48.5	44.3	223.1	49
fDSST [19]	72.0	67.6	51.1	51	57.5	52.4	82.1	51	42.2	34.4	256.8	52
KCF [17]	69.2	54.8	45.0	243	55.1	46.1	77.4	242	31.1	20.0	282.4	245

TABLE 6

Comparisons of DPR (%), OSR (%), CLE in pixels and speed (fps) between different  $\mathcal{V}$  with same tracker on three benchmarks.

		OTB2015 [29]				TC128 [33]				UAV20L [34]			
		DPR	OSR	CLE	Speed	DPR	OSR	CLE	Speed	DPR	OSR	CLE	Speed
Staple	PTAV <sub>VGGNet</sub>	86.2	77.9	18.9	27	77.2	70	32.1	23	73.2	62.4	56.1	30
	PTAV <sub>AlexNet</sub>	84.0	75.5	20.7	34	75.0	68.9	40.3	31	65.9	58.6	70.7	33
fDSST	PTAV <sub>VGGNet</sub>	85.2	77.9	19.4	27	75.2	66.1	32.7	24	63.1	47.8	102.7	26
	PTAV <sub>AlexNet</sub>	79.4	74.3	31.2	29	67.5	61.1	43.8	26	54.7	43.2	121.0	31
KCF	PTAV <sub>VGGNet</sub>	74.2	58.6	34.5	24	63.9	54.6	53.6	19	41.6	32.1	195.3	20
	PTAV <sub>AlexNet</sub>	72.3	58.0	37.1	31	62.4	52.6	54.9	22	39.6	27.6	216.1	27

TABLE 7

Comparisons of DPR (%), OSR (%), CLE in pixels and speed (fps) between fixed template and dynamic template set using  $\mathcal{V}$  based on Staple [21] and  $\mathcal{T}$  based on VGGNet [37] on three benchmarks.

	OTB2015 [29]				TC128 [33]				UAV20L [34]			
	DPR	OSR	CLE	Speed	DPR	OSR	CLE	Speed	DPR	OSR	CLE	Speed
dynamic templates	86.2	77.9	18.9	27.0	77.2	70.0	32.1	23.0	73.2	62.4	56.1	30.0
fixed template	85.6	77.1	20.4	25.0	76.3	68.9	32.3	22.0	72.6	61.6	62.8	28.0

TABLE 8

Comparisons of different  $N_{\text{int}}$  in DPR and speed on OTB2015 [29].

	$N_{\text{int}} = 5$	$N_{\text{int}} = 10$	$N_{\text{int}} = 15$
DPR (%)	86.3	86.2	84.6
Speed (fps)	25	27	30

TABLE 9

Comparisons of tracking speed (fps) on three benchmarks.

	OTB2015 [29]	TC128 [33]	UAV20L [34]	VOT2016 [35]
Single thread	15	14	17	15
Two threads	27	23	30	25

less accurately. Specifically, the speed of VGGNet based- $\mathcal{V}$  is 6 *fps* while its AlexNet counterpart runs at 17 *fps*. The comparison results of different verifiers with the same tracker on OTB2015, TC128 and UAV20L are reported in Table 6.

From Table 6, we observe that, when using the same base tracker, PTAV with VGGNet-based  $\mathcal{V}$  (PTAV<sub>VGGNet</sub>) outperforms that with AlexNet-based  $\mathcal{V}$  (PTAV<sub>AlexNet</sub>) on all three benchmarks. For example, by selecting Staple as the tracking part, PTAV<sub>VGGNet</sub> achieves DPRs of 86.2%, 77.2% and 73.2% on OTB2015, TC128 and UAV20L, respectively, obtaining improvements of respectively 2.2%, 2.2% and 7.3% compared to PTAV<sub>AlexNet</sub> with DPRs of 84.0%, 75.0% and 65.9%, respectively. In PTAV, an accurate  $\mathcal{V}$  is able to effectively reduce detections, which in turn decreases computation for verifications and leads to running time efficiency. Consequently, though AlexNet is computationally much more efficient than VGGNet, the speed of PTAV<sub>VGGNet</sub> is competitive to that of PTAV<sub>AlexNet</sub>, and it runs in real-time (27, 23 and 30 *fps* on OTB2015, TC128 and UAV20L, respectively), as shown in Table 6.

#### 4.5.3 Fixed template v.s. dynamic template set

To adapt  $\mathcal{V}$  to target appearance variation, we propose the dynamic template set for adaptive verification, which can take advantages

of confident tracking results to improve validation quality, leading to reduction of verifications and detections for efficiency. Table 7 shows the comparison results of PTAV using dynamic template set versus using a fixed template. From Table 7, we observe that using dynamic template set for verification improves PTAV in both accuracy and efficiency. In specific, the DPRs on three benchmarks are improved from 85.6%, 76.3% and 72.6% to 86.2%, 77.2% and 73.2%, respectively. The speeds are boosted from 25, 22 and 28 *fps* to 27, 23 and 30 *fps*, respectively.

#### 4.5.4 Different verification interval $N_{\text{int}}$

In PTAV, different verification interval  $N_{\text{int}}$  may affect both the accuracy and efficiency. A smaller  $N_{\text{int}}$  implies more frequent verification, which requires more computation and thus degrades the efficiency. A larger  $N_{\text{int}}$ , on the contrary, may cost less computation but may put PTAV at the risk when the target appearance changes quickly. If the tracker loses the tracking target, it may update vast backgrounds in its appearance model until next verification. Even if the verifier re-locates the target and offers a correct detection result, the tracker may still lose it due to heavy changes in the target appearance model. Table 8 demonstrates sampled results with three different  $N_{\text{int}}$  on OTB2015 [29].

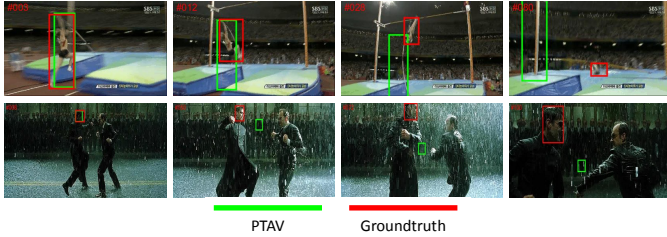


Fig. 13. Failure cases of *Jump* and *Matrix* for PTAV on OTB2015 [29].

Taking into account both accuracy and speed, we set  $N_{\text{int}}$  to 10 in our experiments.

#### 4.5.5 Two threads v.s. one

In PTAV, the tracker  $\mathcal{T}$  does not rely on the verifier  $\mathcal{V}$  most of the time, and two separate threads process tracking and verifying in parallel for efficiency. Consequently,  $\mathcal{T}$  does not have to wait for the feedback from  $\mathcal{V}$  to process next frame, and it traces back and resumes tracking only when receiving a correction feedback from  $\mathcal{V}$ . Owing to storing intermediate results,  $\mathcal{T}$  is able to quickly trace back without little extra computation. Table 9 shows the comparison of speed between using two threads and using only a single thread. It shows that using two threads in parallel clearly improves the efficiency of the system.

#### 4.6 Failure Cases

With the collaboration between  $\mathcal{T}$  and  $\mathcal{V}$ , PTAV usually performs well in dealing with various challenging situations; however, there exist scenarios in which PTAV may fail. Shown in Figure 13(a) on *Jump*, though  $\mathcal{V}$  can accurately detect unreliable tracking results of  $\mathcal{T}$ , it does not provide correct feedbacks for subsequent tracking due to heavy deformation. On *Matrix* shown in Figure 13(b), the target undergoes severe illumination variation, occlusion, rotation and background cluttering, causing difficulties for  $\mathcal{T}$  to localize the target. Even though  $\mathcal{V}$  corrects  $\mathcal{T}$  when validating unreliable tracking results,  $\mathcal{T}$  still drifts to background quickly. In certain cases where the target suffers from heavy appearance changes,  $\mathcal{V}$  cannot provide effective feedbacks to  $\mathcal{T}$ , resulting in irrecoverable drift and failures.

### 5 CONCLUSION

In this paper, we propose a new visual tracking framework, *parallel tracking and verifying* (PTAV), which decomposes object tracking into two sub-tasks, fast tracking and reliable verifying. We show that, by carefully distributing the two tasks into two parallel threads and allowing them to work together, PTAV can achieve the best known tracking accuracy among all real-time tracking algorithms. Furthermore, to adapt the verifier to object appearance variations, we propose using dynamic target templates for adaptive verification, resulting in further improvements in both accuracy and efficiency. The encouraging results are demonstrated in extensive experiments on four popular benchmarks. Moreover, PTAV is a flexible framework with great rooms for improvement and generalization, and thus is expected to inspire the designing of more efficient tracking algorithms in the future.

### REFERENCES

- [1] A. W. Smeulders, D. M. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah, "Visual tracking: An experimental survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 36, no. 7, pp. 1442–1468, 2014.
- [2] A. Yilmaz, O. Javed, and M. Shah, "Object tracking: A survey," *ACM Computing Surveys*, vol. 38, no. 4, p. 13, 2006.
- [3] X. Li, W. Hu, C. Shen, Z. Zhang, A. Dick, and A. V. D. Hengel, "A survey of appearance models in visual object tracking," *ACM Transactions on Intelligent Systems and Technology*, vol. 4, no. 4, pp. 1–58, 2013.
- [4] M. Kristan, J. Matas, A. Leonardis, T. Vojnir, R. Pflugfelder, G. Fernandez, G. Nebehay, F. Porikli, and L. vCehovin, "A novel performance evaluation methodology for single-target trackers," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 11, pp. 2137–2155, 2016.
- [5] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Proceedings of Advances in Neural Information Processing Systems*, 2012.
- [7] C. Ma, J.-B. Huang, X. Yang, and M.-H. Yang, "Hierarchical convolutional features for visual tracking," in *Proceedings of IEEE International Conference on Computer Vision*, 2015.
- [8] Y. Qi, S. Zhang, L. Qin, H. Yao, Q. Huang, J. Lim, and M.-H. Yang, "Hedged deep tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [9] H. Nam and B. Han, "Learning multi-domain convolutional neural networks for visual tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [10] L. Wang, W. Ouyang, X. Wang, and H. Lu, "Stct: Sequentially training convolutional networks for visual tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [11] N. Wang and D.-Y. Yeung, "Learning a deep compact image representation for visual tracking," in *Proceedings of Advances in Neural Information Processing Systems*, 2013.
- [12] H. Fan and H. Ling, "Sanet: Structure-aware network for visual tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 2017.
- [13] M. Danelljan, A. Robinson, F. S. Khan, and M. Felsberg, "Beyond correlation filters: Learning continuous convolution operators for visual tracking," in *Proceedings of European Conference on Computer Vision*, 2016.
- [14] R. Tao, E. Gavves, and A. W. Smeulders, "Siamese instance search for tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [15] L. Wang, W. Ouyang, X. Wang, and H. Lu, "Visual tracking with fully convolutional networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015.
- [16] S. Hong, T. You, S. Kwak, and B. Han, "Online tracking by learning discriminative saliency map with convolutional neural network," in *Proceedings of International Conference on Machine Learning*, 2015.
- [17] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-speed tracking with kernelized correlation filters," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 3, pp. 583–596, 2015.
- [18] K. Zhang, L. Zhang, and M.-H. Yang, "Real-time compressive tracking," in *Proceedings of European Conference on Computer Vision*, 2012.
- [19] M. Danelljan, G. Häger, F. S. Khan, and M. Felsberg, "Discriminative scale space tracking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 8, pp. 1561–1575, 2017.
- [20] D. S. Bolme, J. R. Beveridge, B. A. Draper, and Y. M. Lui, "Visual object tracking using adaptive correlation filters," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [21] L. Bertinetto, J. Valmadre, S. Golodetz, O. Miksik, and P. H. Torr, "Staple: Complementary learners for real-time tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [22] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "Exploiting the circulant structure of tracking-by-detection with kernels," in *Proceedings of European conference on computer vision*, 2012.
- [23] M. Danelljan, F. Shahbaz Khan, M. Felsberg, and J. Van de Weijer, "Adaptive color attributes for real-time visual tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2014.
- [24] Z. Kalal, K. Mikolajczyk, and J. Matas, "Tracking-learning-detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 7, pp. 1409–1422, 2012.

- [25] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2005.
- [26] J. Van De Weijer, C. Schmid, J. Verbeek, and D. Larlus, "Learning color names for real-world applications," *IEEE Transactions on Image Processing*, vol. 18, no. 7, pp. 1512–1523, 2009.
- [27] C. Ma, X. Yang, C. Zhang, and M.-H. Yang, "Long-term correlation tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [28] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr, "Fully-convolutional siamese networks for object tracking," in *Proceedings of European Conference on Computer Vision Workshop*, 2016.
- [29] Y. Wu, J. Lim, and M.-H. Yang, "Object tracking benchmark," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1834–1848, 2015.
- [30] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *Proceedings of IEEE International Symposium on Mixed and Augmented Reality*, 2007.
- [31] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "ORB-SLAM: a versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.
- [32] S. Chopra, R. Hadsell, and Y. LeCun, "Learning a similarity metric discriminatively, with application to face verification," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2005.
- [33] P. Liang, E. Blasch, and H. Ling, "Encoding color information for visual tracking: Algorithms and benchmark," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 5630–5644, 2015.
- [34] M. Mueller, N. Smith, and B. Ghanem, "A benchmark and simulator for uav tracking," in *Proceedings of European Conference on Computer Vision*, 2016.
- [35] M. Kristan and et al, "The visual object tracking vot2016 challenge results," in *Proceedings of European Conference on Computer Vision Workshop*, 2016.
- [36] H. Fan and H. Ling, "Parallel tracking and verifying: A framework for real-time and high accuracy visual tracking," in *Proceedings of IEEE International Conference on Computer Vision*, 2017.
- [37] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proceedings of International Conference on Learning Representations*, 2014.
- [38] B. Babenko, M.-H. Yang, and S. Belongie, "Robust object tracking with online multiple instance learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 8, pp. 1619–1632, 2011.
- [39] H. Grabner, C. Leistner, and H. Bischof, "Semi-supervised on-line boosting for robust tracking," in *Proceedings of European Conference on Computer Vision*, 2008.
- [40] S. Hare, S. Golodetz, A. Saffari, V. Vineet, M.-M. Cheng, S. L. Hicks, and P. H. Torr, "Struck: Structured output tracking with kernels," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 10, pp. 2096–2109, 2016.
- [41] D. A. Ross, J. Lim, R.-S. Lin, and M.-H. Yang, "Incremental learning for robust visual tracking," *International Journal of Computer Vision*, vol. 77, no. 1-3, pp. 125–141, 2008.
- [42] J. Kwon and K. M. Lee, "Visual tracking decomposition," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2010.
- [43] X. Mei and H. Ling, "Robust visual tracking using  $\ell_1$  minimization," in *Proceedings of IEEE International Conference on Computer Vision*, 2009.
- [44] C. Bao, Y. Wu, H. Ling, and H. Ji, "Real time robust l1 tracker using accelerated proximal gradient approach," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
- [45] T. Zhang, S. Liu, C. Xu, S. Yan, B. Ghanem, N. Ahuja, and M.-H. Yang, "Structural sparse tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [46] Y. Li and J. Zhu, "A scale adaptive kernel correlation filter tracker with feature integration," in *Proceedings of European Conference on Computer Vision Workshop*, 2014.
- [47] T. Liu, G. Wang, and Q. Yang, "Real-time part-based visual tracking via adaptive correlation filters," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [48] M. Mueller, N. Smith, and B. Ghanem, "Context-aware correlation filter tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [49] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg, "Learning spatially regularized correlation filters for visual tracking," in *Proceedings of IEEE International Conference on Computer Vision*, 2015.
- [50] J. Valmadre, L. Bertinetto, J. F. Henriques, A. Vedaldi, and P. H. Torr, "End-to-end representation learning for correlation filter based tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [51] Y. Hua, K. Alahari, and C. Schmid, "Occlusion and motion reasoning for long-term tracking," in *Proceedings of European Conference on Computer Vision*, 2014.
- [52] Z. Hong, Z. Chen, C. Wang, X. Mei, D. Prokhorov, and D. Tao, "Multi-store tracker (MUSTER): A cognitive psychology inspired approach to object tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [53] R. Girshick, "Fast R-CNN," in *Proceedings of IEEE International Conference on Computer Vision*, 2015.
- [54] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient k-means clustering algorithm: Analysis and implementation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 881–892, 2002.
- [55] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of ACM Multimedia Conference*, 2014.
- [56] J. Zhang, S. Ma, and S. Sclaroff, "Meem: robust tracking via multiple experts using entropy minimization," in *Proceedings of European Conference on Computer Vision*, 2014.
- [57] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg, "Convolutional features for correlation filter based visual tracking," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2015.
- [58] H. Nam, M. Baek, and B. Han, "Modeling and propagating cnns in a tree structure for visual tracking," *arXiv*, 2016.
- [59] G. Zhu, F. Porikli, and H. Li, "Beyond local search: Tracking objects everywhere with instance-specific proposals," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [60] Z. Chi, H. Li, H. Lu, and M.-H. Yang, "Dual deep network for visual tracking," *IEEE Transactions on Image Processing*, vol. 26, no. 4, pp. 2005–2015, 2017.



**Heng Fan** received his B.E. degree in College of Science, Huazhong Agricultural University (HZAU), Wuhan, China, in 2013. He is currently a Ph.D. student in the Department of Computer and Information Science, Temple University, Philadelphia, USA. His research interests include computer vision, pattern recognition and machine learning.



**Haibin Ling** received the BS and MS degrees from Peking University, China, in 1997 and 2000, respectively, and the PhD degree from the University of Maryland College Park in 2006. From 2000 to 2001, he was an assistant researcher at Microsoft Research Asia. From 2006 to 2007, he worked as a postdoctoral scientist at the University of California Los Angeles. After that, he joined Siemens Corporate Research as a research scientist. Since fall 2008, he has been with Temple University where he is now an Associate Professor.

Ling's research interests include computer vision, augmented reality, medical image analysis, and human computer interaction. He received the Best Student Paper Award at the ACM Symposium on User Interface Software and Technology (UIST) in 2003, and the NSF CAREER Award in 2014. He serves as associate editors for IEEE Transactions on Pattern Analysis and Machine Intelligence, Pattern Recognition, and Computer Vision and Image Understanding, and has served as area chairs for CVPR 2014 and CVPR 2016.