# Visual Tracking via Dynamic Graph Learning

Chenglong Li, Liang Lin, Wangmeng Zuo, Jin Tang, and Ming-Hsuan Yang

**Abstract**—Existing visual tracking methods usually localize a target object with a bounding box, in which the performance of the foreground object trackers or detectors is often affected by the inclusion of background clutter. To handle this problem, we learn a patch-based graph representation for visual tracking. The tracked object is modeled by with a graph by taking a set of non-overlapping image patches as nodes, in which the weight of each node indicates how likely it belongs to the foreground and edges are weighted for indicating the appearance compatibility of two neighboring nodes. This graph is dynamically learned and applied in object tracking and model updating. During the tracking process, the proposed algorithm performs three main steps in each frame. First, the graph is initialized by assigning binary weights of some image patches to indicate the object and background patches according to the predicted bounding box. Second, the graph is optimized to refine the patch weights by using a novel alternating direction method of multipliers. Third, the object feature representation is updated by imposing the weights of patches on the extracted image features. The object location is predicted by maximizing the classification score in the structured support vector machine. Extensive experiments show that the proposed tracking algorithm performs well against the state-of-the-art methods on large-scale benchmark datasets.

◆

## 1 INTRODUCTION

V ISUAL tracking is a fundamental and active research topic in computer vision due to its wide range of applications such as activity analysis, visual surveillance and self-driving systems. Despite significant progress has been made in recent years, it remains a challenging issue, partly due to the difficulty of constructing robust object representation to cope with various factors including camera motion, partial occlusion, background clutter and illumination change.

Numerous visual tracking methods recently adopt the tracking-by-detection paradigm, i.e., separating the foreground object from the background over time using a classifier. These methods usually localize the object using a bounding box, and draw positive (negative) samples from inside (outside) of the bounding box for the classifier update. Since the ground-truth object labeling is only available in the initial frame, incrementally updating the object classifier in subsequent frames often result in model drift due to inclusion of outlier samples.

Significant efforts have been made to alleviate the effects of outlier samples in visual tracking [1], [2], [3], [4], [5], [6], [7], [8], [9]. Several methods in [1], [3], [4] update the object classifiers by considering the distances of samples with respect to the bounding box center, e.g., the samples close to the center receiving higher weights. Some other methods [10], [11] segment foreground objects from the background during the tracking process

to exclude background clutter. However, these methods are limited in dealing with cluttered backgrounds (e.g., unreliable segmented object masks). To improve the robustness, Kim et al. [7] define an image patch based 8-neighbor graph to represent the tracked object, in which if two nodes are connected by an edge if they are are 8-neighbors and the edge weight is computed based on low-level feature distance. There are two main issues with this approach: i) it only considers the spatial neighbors and do not capture the intrinsic relationship between patches; ii) it uses low-level feature which are less effective in the presence of clutter and noise.

To handle these issues, we learn a robust object representation for visual tracking. Given one bounding box of the target object, we partition it into non-overlapping local patches, which are described by color and gradient histograms. Instead of using static structures in existing methods [12], [7], we learn a dynamic graph with patches as nodes (i.e., adaptive structure and node weights for each frame) for representing the target object, where the weight of each node describes how likely it belongs to the target object, and the edge weight indicates the appearance compatibility of two neighboring patches. Existing methods usually perform two steps for node weight computation, i.e., first constructing the graph with a static structure and low-level features, and then computing node weights based on some semi-supervised methods [12], [7]. In this work, we propose a novel representation model to jointly learn the graph that infers the graph structure, edge weights and node weights.

With the advances of compressed sensing [13], numerous methods exploiting the relationship of data representations have been proposed [14], [15], [16], [17]. The representations are generally utilized to define [14], [15], [16] or learn [17] the affinity matrix of a graph. Motivated by these methods, we represent each patch descriptor as a linear combination of other patch descriptors, and

C. Li and J. Tang are with the School of Computer Science and Technology, Anhui University, Hefei, 230601, China. C. Li is also with the Center for Research on Intelligent Perception and Computing (CRIPAC), National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA), Beijing, 100190, China. (e-mail: lcl1314@foxmail.com; jtang99029@foxmail.com)
L. Lin is with the School of Data and Computer Science, Sun Yat-Sen University, Guangzhou, 510006, China. (e-mail: linliang@ieee.org)
W. Zuo is with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin, 150001, China. (e-mail: cswmzuo@gmail.com)
M.-H. Yang is with the School of Engineering, University of California, Merced, CA 95344 USA. (e-mail: mhyang@ucmerced.edu)

develop a model to jointly optimize the graph structure, edge weights and node weights while suppressing the effects of noise from clutter and low-level features.

It is worth mentioning that our model has the following three distinctive properties: 1) it is capable of collaboratively optimize the graph structure, edge weights and node weights according to the underlying intrinsic relationship, which provides a flexible solution for visual tracking and other vision problems such as saliency detection [12] and semi-supervised object segmentation [18]; 2) it is effective to suppress the effects of noise from pixels and low-level features in computing the affinity matrix of the graph; 3) it is generic, and can incorporate other constraints (e.g., low-rank and sparse constraints) to further improve the robustness of graph learning.

To improve the tracking efficiency, we develop an alternating direction method of multipliers (ADMM) algorithm to seek the solution of the proposed model. In particular, the alternating direction method [19] is used to linearize the quadratic penalty term while avoiding an auxiliary variable and some matrix inversions such that each subproblem can be efficiently solved with a closed-form solution. We construct the robust object representations by combining patch features with the optimized weights, and then apply the structured support vector machine (SVM) [3] for object tracking and model update.

In each frame, the proposed algorithm is carried out with several steps. First, the graph is initialized with binary weights to according to the ground truth (first frame) or the predicted bounding box (subsequent frames). Second, the graph is optimized by a linearized ADMM algorithm. Third, the object feature representation is updated by imposing the patch weights on the extracted image features. The object location is finally predicted by adopting the structured SVM.

We make three major contributions for visual tracking and related applications in this work:

- We propose an effective approach to alleviate the effects of background clutter in visual tracking. Extensive experiments show that the proposed method outperforms most state-of-the-art trackers on four benchmark datasets.
- We present a novel representation model to learn a dynamic graph according to the intrinsic relationship among image patches. The proposed model is jointly optimizes the graph structure, edge weights and node weights while suppressing the effects of patch noise and/or corruption. It also provides a general solution for visual tracking and other vision problems such as saliency detection [12], [20] and interactive object segmentation [18], [21].
- We develop an ADMM algorithm to efficiently solve the associated optimization problem. Empirically, the proposed optimization algorithm exhibits stable convergence behavior on real image data.

This paper provides a more complete understanding of the early results [22], with more background, insights, analysis, and evaluation. In particular, our approach advances the early work in several aspects. First, we utilize the data representations to learn more meaningful graph affinity, instead of directly using data representations. Second, we generalize the graph learning algorithm to incorporating different constraints (or priors), such as the low rank, sparse and spatial smoothness constraints. We further discuss the merits for graph learning, and instantiate the sparse constraints into our framework. Third, scale estimation is considered in this work to improve visual tracking Finally, we carry out extensive experiments on large-scale benchmark datasets to demonstrate the effectiveness of the proposed algorithm, including quantitative comparisons with the state-of-the-art trackers and ablation studies.

## 2 LITERATURE REVIEW

A plethora of visual methods have been proposed in the literature [23], [24] and we discuss the most related work in this section. We discuss the advances of visual tracking in two aspects: constructing robust appearance models to alleviate the effect of background clutter, and learning data affinity for model construction.

### 2.1 Appearance Models

Various tracking methods have been proposed to improve the robustness to nuisance factors including label ambiguity, background clutter, corruption and occlusion. Grabner et al. [25] present a tracking approach that adapts to drastic appearance changes and limits the drifting problem. The knowledge from labeled data is used to construct static prior for online classifier while unlabeled samples are explored in a principled manner during tracking. Babenko et al. [26] use a bag of multiple samples, instead of a single sample, to update the classifier reliably. To avoid the label ambiguity, Hare et al. [3] exploit structured samples instead of binary-labeled samples when training the classifier in the structured SVM framework [27].

To alleviate the effects of background clutter, one representative approach is to assign weights to different pixels or patches in the bounding box. Comaniciu et al. [1] develop the kernel-based method to assign smaller weights to boundary pixels for histogram matching. In [4] He et al. also assume that pixels far from a box center should be less important. These methods do not perform well when a target shape cannot be well described by a rectangle or occluded. Some methods [10], [11] integrate segmentation results in visual tracking to alleviate the effects of background. These algorithms, however, reply heavily on the quality of segmentation results. Kim et al. [7] develop a random walk restart algorithm on a 8-neighbor graph to compute patch weights within the target object bounding box. Nevertheless, the constructed graph does not capture the relationship between patches well.

## 2.2 Data Affinity

In vision and learning problems, we often have a set of data $\mathbf{X} = \{\mathbf{x}_1, \ldots, \mathbf{x}_n\} \in \mathbb{R}^{d \times n}$ drawn from a union of $c$ subspaces $\{\mathbb{S}_{s=1}^c\}$, where $d$ is the feature dimension and $n$ is the number of data vectors. To characterize the relation between the data in $\mathbf{X}$, the key is to construct an effective affinity matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, in which $\mathbf{A}_{ij}$ reflects the similarity between data points $\mathbf{X}_i$ and $\mathbf{X}_j$. While computing Euclidean distances on the raw data is the most intuitive way to construct the data affinity matrix, such metric usually does not reveal the global subspace structure of data well.

With the advances of compressed sensing [13], significant efforts have been made to exploit the relationship of data representations [14], [15], [16], [17] where the general formulation is described by:

$$\min_{\mathbf{Z}, \mathbf{E}} \alpha\Theta(\mathbf{Z}) + \beta\Phi(\mathbf{E})$$
$$\text{s.t.} \quad \mathbf{X} = \mathbf{X}\mathbf{Z} + \mathbf{E}, \tag{1}$$

where $\mathbf{Z} \in \mathbb{R}^{n \times n}$ and $\mathbf{E} \in \mathbb{R}^{d \times n}$ denote the representation matrix and the residual matrix, respectively. In (1), $\Theta(\mathbf{Z})$ is the regularizer on $\mathbf{Z}$. $\Phi(\mathbf{E})$ is the model of $\mathbf{E}$, which can be with different forms depending on data characteristics; and $\alpha$ as well as $\beta$ are the weight parameters.

Numerous methods have been developed to extract compact information from image data including sparse representation (SR) [14]

$$\Theta(\mathbf{Z}) = ||\mathbf{Z}||_0, \tag{2}$$

and low-rank representation (LRR) [15],

$$\Theta(\mathbf{Z}) = \text{rank}(\mathbf{Z}). \tag{3}$$

Different from traditional methods, SR schemes can be used to exploit higher order relationships among more data points effectively, and hence provide more compact and discriminative models [16]. The main drawback of SR methods is that data is processed individually without taking the existence of inherent global structure into account. On the other hand, low-rank representation models use low rank constraints on data representations to capture the global structure of the whole data. It has been shown that, under mild conditions, LRR methods can preserve the membership of samples that belong to the same subspace well. Recently, Zhuang et al. [16] harnesses both sparsity and low-rankness of data to learn more informative representations.

In general, after solving the problem (1), the representation is used to define the affinity matrix of an undirected graph with $\mathbf{a}_{ij} = \frac{\mathbf{z}_{ij} + \mathbf{z}_{ji}}{2}$ for $\mathbf{x}_i$ and $\mathbf{x}_j$. However, the metric implies the affinity is already not the same as the original definition. This is because the affinity defines an approximation to the pairwise distances between data samples while the representation is the reconstruction coefficients of one sample from others. As such, Guo [17] proposes a method to simultaneously learn data representations and affinity matrix. Experimental results on the synthetic and real datasets demonstrate the effectiveness of learning model representation and affinity matrix jointly.

## 3 PATCH-BASED GRAPH LEARNING

Given one bounding box that encloses the target object, we partition it into non-overlapping patches and assign each one with a weight that reflects the importance in describing the target object to alleviate the effects of background clutter. We concatenate these weighted patch descriptors into a feature vector and use the Struck [3] method for object tracking. In this section, we first describe a sparse low-rank model based on local patches, and then an efficient ADMM algorithm to compute the weights.

### 3.1 Formulation

Each bounding box of the target object is partitioned into $n$ non-overlapping patches, and a set of low-level appearance features are extracted and combined into one single $d$-dimensional feature vector $\mathbf{x}_i$ for characterizing the $i$-th patch. Using these patches as graph nodes, each bounding box can be represented with a graph, in which the weight of each node describes how likely it belongs to the target object and the edge weight between two neighboring patches indicates appearance compatibility.

For visual tracking, some patches in a target bounding box may belong to background due to irregular shape, scale variation and partial occlusion of the target object, as shown in Figure 1. Thus, we assign a weight for each graph node to alleviate the effects of background pixels on object tracking and model update. On the other hand, instead of constructing spatially ordered graphs [12], [7], the edges are dynamically learned for capturing the intrinsic relationship of data. In this work, we propose a novel graph learning approach to infer the edges and node weights jointly which performs well against the state-of-the-art alternatives for visual tracking.

All the feature vectors of $n$ patches in one bounding box form the data matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$. Each patch descriptor can be represented as a linear combination of remaining patch descriptors, and the representation of all patch vectors can then be formulated by $\mathbf{X} = \mathbf{X}\mathbf{Z}$, where $\mathbf{Z} \in \mathbb{R}^{n \times n}$ is the representation coefficient matrix. Since the patch feature matrix often contains noise, the representation can be obtained by solving the objective function (1).

The optimal representation coefficient matrix in (1) is often utilized to define the affinity matrix of an undirected graph in the way of $\frac{|\mathbf{z}_{ij}| + |\mathbf{z}_{ji}|}{2}$ for the feature vector $\mathbf{x}_i$ and $\mathbf{x}_j$. As $\mathbf{z}_{ij}$ and $\mathbf{z}_{ji}$ are the reconstruction coefficients, this encoded information is not the same as the original definition, which defines an approximation to the pairwise distances between $\mathbf{x}_i$ and $\mathbf{x}_j$ [17]. Therefore, we also learn the affinity matrix by assuming that the patch features should have larger probabilities to be

(a) Partial occlusion      (b) Deformation      (c) Scale reduction      (d) Scale increment
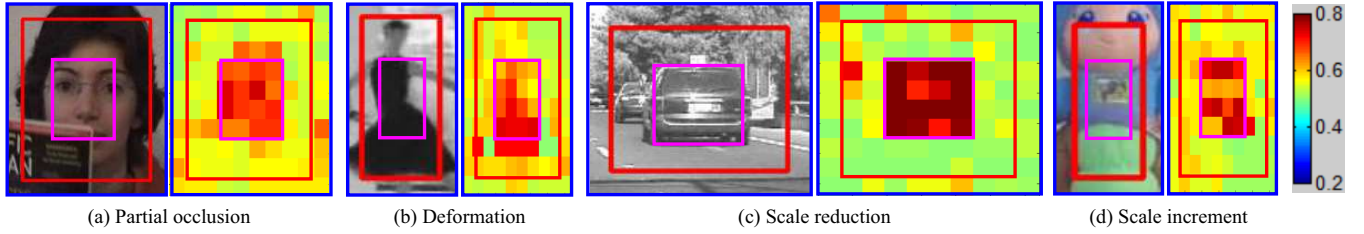
Fig. 1. Illustration of the original, shrunk and expanded bounding boxes on 4 video frames with different challenges, which are represented by the red, pink and blue colors, respectively. The optimized patch weights are also shown for clarity, in which the hotter color indicates the larger weight. One can see that the optimized patch weights are beneficial to suppressing the effects of background clutter.

in the same cluster if their representations have smaller distance, and impose the following constraints,

$$\min_{\mathbf{A1}=\mathbf{1}, \mathbf{A} \geq 0} \sum_{i,j=1}^{n} ||\mathbf{z}_i - \mathbf{z}_j||_F^2 \mathbf{a}_{ij}, \qquad (4)$$

where $\mathbf{A} \in \mathbb{R}^{n \times n}$ is the desired affinity matrix, whose element $\mathbf{a}_{ij}$ reflects the probability of the patch features $\mathbf{x}_i$ and $\mathbf{x}_j$ from the same cluster based on the distance between their representations $\mathbf{z}_i$ and $\mathbf{z}_j$. The constraints $\mathbf{A1} = \mathbf{1}$ and $\mathbf{A} \geq 0$ guarantee the probability property of each column of $\mathbf{A}$. With some simple algebra, we integrate these constraints into (1), and have

$$\min_{\mathbf{Z}, \mathbf{E}, \mathbf{A}} \alpha \Theta(\mathbf{Z}) + \beta \Phi(\mathbf{E}) + \gamma \operatorname{tr}(\mathbf{Z} \mathbf{L_A} \mathbf{Z}^\top) + \frac{\lambda}{2} ||\mathbf{A}||_F^2 \qquad (5)$$
$$\text{s.t.} \quad \mathbf{X} = \mathbf{XZ} + \mathbf{E}, \ \mathbf{A1} = \mathbf{1}, \ \mathbf{A} \geq 0,$$

where $\mathbf{L_A} = \mathbf{D_A} - \mathbf{A}$ is the Laplacian matrix of $\mathbf{A}$, and $\mathbf{D_A}$ is the degree matrix of $\mathbf{A}$, a diagonal matrix whose the $i$-th diagonal element is $\sum_j \mathbf{a}_{ij}$. In (5), $\gamma$ and $\lambda$ are weight parameters. In addition, the last term is used to avoid overfitting. Note that minimizing the term $\operatorname{tr}(\mathbf{Z} \mathbf{L_A} \mathbf{Z}^\top)$ could exclude the trivial solution $\mathbf{Z} = \mathbf{I}$, where $\mathbf{I}$ indicates the identity matrix. The trivial solution $\mathbf{E} = \mathbf{0}$ is also not achieved as it means the data are clean, which is an "ideal" case, and does not exist in real-world applications.

To alleviate the effects of background clutter, we assign a weight $\mathbf{w}_i$ for each patch $i$ using a semi-supervised formulation. Let $\mathbf{r} = \{\mathbf{r}_1, \mathbf{r}_2, \ldots, \mathbf{r}_n\}^\top$ be an initial weight vector, in which $\mathbf{r}_i = 1$ if $\mathbf{r}_i$ is a target object patch, and $\mathbf{r}_i = 0$ indicating a background patch. In this work, $\mathbf{r}$ is computed by the initial ground truth (for first frame) or the previous tracking result (for subsequent frames). For $i$-th patch, if it belongs to the shrunk region of the bounding box then $\mathbf{r}_i$ is 1, and if it belongs to the expanded region of the bounding box then $\mathbf{r}_i$ is 0. Figure 1 shows the one example how the weights are assigned. Although using a simple initialization strategy, we demonstrate empirically this scheme performs well empirically, and show the robustness to clutter and noise in Figure 2.

The remaining patches are non-determined, and are diffused by other patches. To this end, we define an
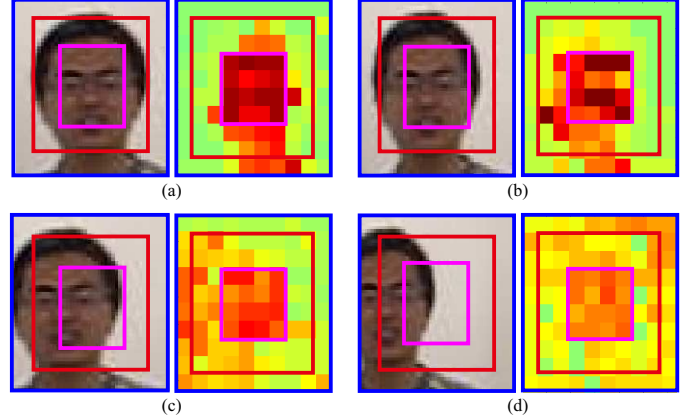


Fig. 2. Illustration of clutter and noise effects of initial seeds on a video frame in the sequence *boy* [23], please refer to Figure 1 for the detailed descriptions. One can see that the optimized patch weights are robust to clutter and noise, but bad if most of initial seeds are clutter and noise.

indicator vector $\Gamma$ that $\Gamma_i = 1$ indicates the $i$-th patch is foreground or background patch, and $\Gamma_i = 0$ denotes the $i$-th patch is non-determined patch. We integrate the patch weights into (5), and obtain

$$\min_{\mathbf{Z}, \mathbf{E}, \mathbf{A}, \mathbf{w}} \alpha \Theta(\mathbf{Z}) + \beta \Phi(\mathbf{E}) + \gamma \operatorname{tr}(\mathbf{Z} \mathbf{L_A} \mathbf{Z}^\top)$$
$$+ \lambda_1 \sum_{i,j} \mathbf{a}_{ij} (\mathbf{w}_i - \mathbf{w}_j)^2 + \frac{\lambda_2}{2} ||\Gamma \circ (\mathbf{w} - \mathbf{r})||^2 \qquad (6)$$
$$+ \frac{\lambda_3}{2} ||\mathbf{A}||_F^2 + \frac{\lambda_4}{2} ||\mathbf{w}||^2$$
$$\text{s.t.} \quad \mathbf{X} = \mathbf{XZ} + \mathbf{E}, \ \mathbf{A1} = \mathbf{1}, \ \mathbf{A} \geq 0, \ \mathbf{w} \geq 0,$$

where $\circ$ indicates the element-wise product. $\lambda_1$, $\lambda_2$, $\lambda_3$ and $\lambda_4$ are weight parameters. The third and fourth terms are the smoothness and fitting constraints. Since the indicator vector $\Gamma$ removes fitness constraint of non-determined patch weights, we introduce the last term to avoid overfitting. Specifically, the smoothness term of $\mathbf{w}$ constrains that $\mathbf{w}_i$ and $\mathbf{w}_j$ are similar to each other when $\mathbf{a}_{ij}$ is non-zero, and the fitting term of $\mathbf{w}$ controls that its elements are close to 0 or 1. However, the fitting constraint is partial, and we thus introduce $||\mathbf{w}||^2$ to

avoid its element amplitude too large.

## 3.2 Discussion

As discussed in Section 2, the regularizer $\Theta(\mathbf{Z})$ is usually based on sparse or low-rank priors, e.g., sparse representation (SR) and low-rank representation (LRR). The SR methods exploit higher order relationships among more data points and hence is more discriminative [14], [28]. The LRR approaches employ low rank constraints on data representations to capture the global structure of data points, and thus is robust to noise and corruption [16], [15]. However, the LRR methods require singular value decomposition (SVD) operations at each iteration, which is computationally demanding. Therefore, we impose the sparse constraints (i.e., $\ell_1$-norm, a convex surrogate for $\ell_0$-norm) on $\mathbf{Z}$ in this work for computational efficiency.

In (6), the model $\Phi(\mathbf{E})$ can be in different forms based on the characteristic of data. For visual tracking, as some image patches are corrupted (e.g., occluded by background or other objects), we employ $\ell_{2,1}$-norm (a convex surrogate for $\ell_{2,0}$-norm) on $\mathbf{E}$. Putting the data terms and prior together, we have:

$$
\begin{aligned}
\min_{\mathbf{Z},\mathbf{E},\mathbf{A},\mathbf{w}} \quad & \alpha||\mathbf{Z}||_1 + \beta||\mathbf{E}||_{2,1} + \gamma \ \mathrm{tr}(\mathbf{Z}\mathbf{L_A}\mathbf{Z}^\top) \\
& + \lambda_1 \sum_{i,j} \mathbf{a}_{ij}(\mathbf{w}_i - \mathbf{w}_j)^2 + \frac{\lambda_2}{2}||\Gamma \circ (\mathbf{w} - \mathbf{r})||^2 \\
& + \frac{\lambda_3}{2}||\mathbf{A}||_F^2 + \frac{\lambda_4}{2}||\mathbf{w}||^2 \\
\text{s.t.} \quad & \mathbf{X} = \mathbf{X}\mathbf{Z} + \mathbf{E}, \ \mathbf{A1} = \mathbf{1}, \ \mathbf{A} \geq 0, \ \mathbf{w} \geq 0,
\end{aligned}
\tag{7}
$$

where $||\cdot||_1$ and $||\cdot||_{2,1}$ denote $\ell_1$-norm and $\ell_{2,1}$-norm of a matrix, respectively. (7) is reasonable as two patches should prefer to be sparsely represented by same set of patches if they are similar. In particular, for optimizing $\mathbf{Z}$, we exploit higher order relationship among patches by minimizing $l_1$ norm on $\mathbf{Z}$, and also penalty inconsistency between $\mathbf{z}_i$ and $\mathbf{z}_j$ when patch features $\mathbf{x}_i$ and $\mathbf{x}_j$ are similar (i.e., large $\mathbf{a}_{ij}$) by minimizing $\mathrm{tr}(\mathbf{Z}\mathbf{L_A}\mathbf{Z}^\top)$.

It is worth noting that although $\mathbf{A}$ is a non-symmetrical affinity matrix, as shown in next section, the solutions of the variables that rely on $\mathbf{A}$ (i.e., $\mathbf{Q}$ and $\mathbf{w}$) are based on a symmetrical affinity matrix, i.e., $(\mathbf{A} + \mathbf{A}^\top)/2$.

## 3.3 Optimization

Although (7) is not jointly convex on $\mathbf{Z}$, $\mathbf{E}$, $\mathbf{A}$ and $\mathbf{w}$, but it is convex with respect to each of them when others are fixed. The ADMM (Alternating Direction Method of Multipliers) algorithm [19] has shown to be an efficient and effective solver of problems similar to (7). To apply ADMM for the above problem, we need to make the objective function separable. Therefore, we introduce an

auxiliary variable $\mathbf{Q} \in \mathbb{R}^{n \times n}$ to replace $\mathbf{Z}$ in (7):

$$
\begin{aligned}
\min_{\mathbf{Z},\mathbf{E},\mathbf{A},\mathbf{w},\mathbf{Q}} \quad & \alpha||\mathbf{Z}||_1 + \beta||\mathbf{E}||_{2,1} + \gamma \ \mathrm{tr}(\mathbf{Q}\mathbf{L_A}\mathbf{Q}^\top) \\
& + \lambda_1 \sum_{i,j} \mathbf{a}_{ij}(\mathbf{w}_i - \mathbf{w}_j)^2 + \frac{\lambda_2}{2}||\Gamma \circ (\mathbf{w} - \mathbf{r})||^2 \\
& + \frac{\lambda_3}{2}||\mathbf{A}||_F^2 + \frac{\lambda_4}{2}||\mathbf{w}||^2 \\
\text{s.t.} \quad & \mathbf{X} = \mathbf{X}\mathbf{Z} + \mathbf{E}, \ \mathbf{Z} = \mathbf{Q}, \mathbf{A1} = \mathbf{1}, \mathbf{A} \geq 0, \ \mathbf{w} \geq 0.
\end{aligned}
\tag{8}
$$

The augmented Lagrangian function of (8) is

$$
\begin{aligned}
& \mathbb{L}_{\{\mathbf{A1}=\mathbf{1},\mathbf{A}\geq 0,\mathbf{w}\geq 0\}}(\mathbf{Z},\mathbf{Q},\mathbf{E},\mathbf{A},\mathbf{w}) \\
& = \alpha||\mathbf{Z}||_1 + \beta||\mathbf{E}||_{2,1} + \gamma \ \mathrm{tr}(\mathbf{Q}\mathbf{L_A}\mathbf{Q}^\top) \\
& + \lambda_1 \sum_{i,j} \mathbf{a}_{ij}(\mathbf{w}_i - \mathbf{w}_j)^2 + \frac{\lambda_2}{2}||\Gamma \circ (\mathbf{w} - \mathbf{r})||^2 \\
& + \frac{\lambda_3}{2}||\mathbf{A}||_F^2 + \frac{\lambda_4}{2}||\mathbf{w}||^2 - \frac{1}{2\mu}(||\mathbf{Y}_1||_F^2 + ||\mathbf{Y}_2||_F^2) \\
& + f(\mathbf{Z},\mathbf{Q},\mathbf{E},\mathbf{w},\mathbf{Y}_1,\mathbf{Y}_2,\mu),
\end{aligned}
\tag{9}
$$

where $\mu > 0$ is the penalty parameter, and $f(\mathbf{Z},\mathbf{Q},\mathbf{E},\mathbf{w},\mathbf{Y}_1,\mathbf{Y}_2,\mu) = \frac{\mu}{2}(||\mathbf{X} - \mathbf{X}\mathbf{Z} - \mathbf{E} + \mathbf{Y}_1/\mu||_F^2 + ||\mathbf{Z} - \mathbf{Q} + \mathbf{Y}_2/\mu||_F^2)$. In the above equation, $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are the Lagrangian multipliers. The ADMM alternatively updates one variable by minimizing $\mathbb{L}$ with fixing other variables. In addition to the Lagrangian multipliers, there are 5 variables, including $\mathbf{Z}$, $\mathbf{Q}$, $\mathbf{E}$, $\mathbf{A}$ and $\mathbf{w}$, to be solved. The solutions of these subproblems are discussed below.

### 3.3.1 Solving $\mathbf{Z}$

With other variables in (9) fixed, the $\mathbf{Z}$-subproblem can be written as:

$$
\min_{\mathbf{Z}} \quad \alpha||\mathbf{Z}||_1 + f(\mathbf{Z},\mathbf{Q}^k,\mathbf{E}^k,\mathbf{w}^k,\mathbf{Y}_1^k,\mathbf{Y}_2^k,\mu^k).
\tag{10}
$$

To avoid using an auxiliary variable and matrix inversions, we use the linearized ADMM method [19] to minimize the $\mathbf{Z}$-subproblem of (9). The quadratic term $f$ is replaced by its first order approximation at the previous iteration and adding a proximal term. Thus, $\mathbf{Z}^{k+1}$ can be updated by:

$$
\arg\min_{\mathbf{Z}} \quad ||\mathbf{Z}||_1 + \frac{\eta\mu^k}{2\alpha}||\mathbf{Z} - \mathbf{Z}^k||_F^2 + \langle \nabla_{\mathbf{Z}}f^k, \mathbf{Z} - \mathbf{Z}^k\rangle,
\tag{11}
$$

where $f^k$ is the shorthand of $f(\mathbf{Z}^k,\mathbf{Q}^k,\mathbf{E}^k,\mathbf{w}^k,\mathbf{Y}_1^k,\mathbf{Y}_2^k,\mu^k)$. In (11), $\nabla_{\mathbf{Z}}f$ is the partial differential of $f$ with respect to $\mathbf{Z}$, and $\eta = ||\mathbf{X}||_F^2$. With some manipulation, we have: $\nabla_{\mathbf{Z}}f = -\mu(\mathbf{X}^\top(\mathbf{X} - \mathbf{X}\mathbf{Z} - \mathbf{E} + \mathbf{Y}_1/\mu) - (\mathbf{Z} - \mathbf{Q} + \mathbf{Y}_3/\mu))$.

Generally, the solution of $\mathbf{Z}^{k+1}$ is obtained by the soft-threshold (or shrinkage) method [29]:

$$
\mathbf{Z}^{k+1} = \mathbb{S}_{\frac{\alpha}{\eta\mu^k}}(\mathbf{P}^k),
\tag{12}
$$

where $\mathbf{P}^k = \mathbf{Z}^k - \tau^k\nabla_{\mathbf{Z}}f^k \in \mathbb{R}^{n \times n}$, and $\mathbb{S}_{\frac{\alpha}{\eta\mu^k}}(\mathbf{P}^k)$ is the soft-threshold operator on $\mathbf{P}^k$ with parameter $\frac{\alpha}{\eta\mu^k}$.

### 3.3.2 Solving $\mathbf{Q}$

By fixing other variables in (9), the $\mathbf{Q}$-subproblem can be formulated as:

$$\min_{\mathbf{Q}} \quad \gamma \, \mathrm{tr}(\mathbf{Q}\mathbf{L_A}\mathbf{Q}^\top) + ||\mathbf{Z} - \mathbf{Q} + \mathbf{Y}_2/\mu||_F^2. \tag{13}$$

To compute $\mathbf{Q}$, we take the derivative of $\mathbb{L}$ with respect to $\mathbf{Q}$, and set it to be 0. With some manipulation, we have:

$$\mathbf{Q}^{k+1} = (\mathbf{Z}^{k+1} + \mathbf{Y}_2^k/\mu^k)(\mathbf{I} + \gamma(\mathbf{L}_{\mathbf{A}^k} + \mathbf{L}_{\mathbf{A}^k}^\top))^{-1}, \tag{14}$$

where $\mathbf{I}$ is the identity matrix.

### 3.3.3 Solving $\mathbf{E}$

The $\mathbf{E}$-subproblem can be formulated as follows when other variables in (9) are fixed:

$$\mathbf{E}^{k+1} = \arg\min_{\mathbf{E}} \quad ||\mathbf{E}||_{2,1} + \frac{\mu^k}{2\beta}||\mathbf{X} - \mathbf{X}\mathbf{Z}^{k+1} - \mathbf{E} + \frac{\mathbf{Y}_1^k}{\mu^k}||_F^2. \tag{15}$$

which is computed by the $\ell_{2,1}$ minimization method [15]:

$$\mathbf{E}^{k+1} = \mathcal{S}_{\frac{\beta}{\mu^k}}(\mathbf{X}\mathbf{Z}^{k+1} - \mathbf{X} - \frac{\mathbf{Y}_1^k}{\mu^k}), \tag{16}$$

where $\mathcal{S}_{\frac{\beta}{\mu^k}}(\cdot)$ is the $\ell_{2,1}$ minimization operator with parameter $\frac{\beta}{\mu^k}$.

### 3.3.4 Solving $\mathbf{A}$

When other variables in (9) are fixed, the $\mathbf{A}$-subproblem can be formulated as:

$$\min_{\mathbf{A}\mathbf{1}=\mathbf{1},\mathbf{A}\geq 0} \quad \gamma \, \mathrm{tr}(\mathbf{Q}^{k+1}\mathbf{L_A}\mathbf{Q}^{(k+1)T})$$
$$+ \lambda_1 \sum_{i,j} \mathbf{a}_{ij}(\mathbf{w}_i^k - \mathbf{w}_j^k)^2 + \frac{\lambda_3}{2}||\mathbf{A}||_F^2. \tag{17}$$

We separate (17) into a set of independent problems, and each $\mathbf{a}_i$ can be computed efficiently with a closed-form solution (please see the appendix of [17] for details) as:

$$\mathbf{a}_i^{k+1} = (\frac{1 + \sum_{j=1}^{\xi} \hat{\mathbf{u}}_{ij}^{\mathbf{Q}^{k+1}}}{\xi}\mathbf{1} - \mathbf{u}_i^{\mathbf{Q}^{k+1}})_+, \tag{18}$$

where $\mathbf{u}_i^{\mathbf{Q}^{k+1}} \in \mathbb{R}^{n\times 1}$ is a vector whose $j$-th element is $\mathbf{u}_{ij}^{\mathbf{Q}^{k+1}} = \frac{\frac{\gamma}{2}||\mathbf{q}_i^{k+1} - \mathbf{q}_j^{k+1}||_F^2 + \lambda_1(\mathbf{w}_i^k - \mathbf{w}_j^k)^2}{\lambda_3}$. Notice that the parameter $\xi \in \{1,\ldots,n\}$ is introduced to control the number of nearest neighbors of $\mathbf{q}_i$ (or $\mathbf{x}_i$) that could have chance to connect edges with $\mathbf{q}_i$ (or $\mathbf{x}_i$).

### 3.3.5 Solving $\mathbf{w}$

By fixing other variables in (9), the $\mathbf{w}$-subproblem can be formulated as:

$$\min_{\mathbf{w}\geq 0} \quad \lambda_1 \sum_{i,j} \mathbf{a}_{ij}^{k+1}(\mathbf{w}_i - \mathbf{w}_j)^2 + \frac{\lambda_2}{2}||\Gamma \circ (\mathbf{w} - \mathbf{r})||^2$$
$$+ \frac{\lambda_4}{2}||\mathbf{w}||^2. \tag{19}$$

Similar to the solution for updating $\mathbf{Q}$, we take the derivative of $\mathbb{L}$ with respect to $\mathbf{w}$, and set it to be 0.

---

**Algorithm 1** Optimization Procedure to (8)

---

**Input:** The patch feature matrix $\mathbf{X}$ and the initial weight vector $\mathbf{r}$, the parameters $\alpha$, $\beta$, $\gamma$, $\lambda_1$, $\lambda_2$, $\lambda_3$ and $\lambda_4$; Set $\mathbf{Z}_0 = \mathbf{Q}_0 = \mathbf{A}_0 = \mathbf{Y}_{2,0} = \mathbf{0}$, $\mathbf{E}_0 = \mathbf{Y}_{1,0} = \mathbf{0}$, $\mathbf{w} = \mathbf{1}$, $\mu_0 = 0.1$, $\mu_{max} = 10^{10}$, $\rho = 1.1$ and $k = 0$.

**Output:** $\mathbf{Z}$, $\mathbf{Q}$, $\mathbf{E}$, $\mathbf{A}$ and $\mathbf{w}$.

1: **while** not converged **do**
2:    Update $\mathbf{Z}^{k+1}$ by (12);
3:    Update $\mathbf{Q}^{k+1}$ by (14);
4:    Update $\mathbf{E}^{k+1}$ by (16);
5:    **for** $i$ from 1 to $n$ **do**
6:       Update $\mathbf{a}_i^{k+1}$ by (18);
7:    **end for**
8:    Update $\mathbf{w}^{k+1}$ by (20);
9:    Update Lagrange multipliers;
10:   Update $\mu^{k+1}$ by $\mu^{k+1} = \min(\mu_{max}, \rho\mu^k)$;
11:   Update $k$ by $k = k + 1$.
12: **end while**

---

With some manipulation, the closed-form solution of this subproblem can be computed by:

$$\mathbf{w}^{k+1} = [(2\lambda_1(\mathbf{D}^{k+1} - \mathbf{A}^{k+1} - \mathbf{A}^{(k+1)T}) + \lambda_2\Gamma'$$
$$+ \lambda_4\mathbf{I})^{-1}(\lambda_2\Gamma \circ \mathbf{r})]_+, \tag{20}$$

where $\mathbf{D}$ is the degree matrix of $(\mathbf{A} + \mathbf{A}^\top)$ that $\mathbf{D} = \mathrm{diag}\{\mathbf{d}_{11}, \mathbf{d}_{22}, \ldots, \mathbf{d}_{nn}\}$, and $\mathbf{d}_{ii} = \sum_j (\mathbf{a}_{ij} + \mathbf{a}_{ji})$, and $\Gamma' = \mathrm{diag}\{\Gamma_1, \Gamma_2, \ldots, \Gamma_n\}$.

The procedure of solving (9) terminates when the maximum element changes of $\mathbf{Z}$, $\mathbf{Q}$, $\mathbf{E}$, $\mathbf{A}$ and $\mathbf{w}$ between two consecutive iterations are less than a threshold (e.g., $10^{-6}$ in this work) or the maximum number of iterations reaches a pre-defined number (e.g., 100 in this work).

Algorithm 1 summarizes the optimization procedures. Since each subproblem of (9) is convex, the solution by the proposed algorithm satisfies the Nash equilibrium conditions [30].

## 4 STRUCTURED SVM TRACKING

In this section, we incorporate the optimized weights of patches into the conventional tracking-by-detection algorithm, Struck [3], for visual tracking. Although we use the Struck method in this work, the optimized patch weights can also be incorporated into other tracking-by-detection algorithms. The Struck method selects the optimal target bounding box $b_t^*$ in the $t$-th frame by maximizing a classification score:

$$b_t^* = \arg\max_b \langle \mathbf{h}_{t-1}, \mathbf{x}_{t,b} \rangle, \tag{21}$$

where $\mathbf{h}_{t-1}$ is the normal vector of a decision plane of the $(t-1)$-th frame, and $\mathbf{x}_{t,b} = [\mathbf{x}_{t,1}; \mathbf{x}_{t,2}; \ldots; \mathbf{x}_{t,n}]$ denotes the descriptor representing a bounding box $b$ in the $t$-th frame. Instead of using binary-labeled samples, the Struck method employs structured samples that consist of a target bounding box and nearby boxes in the same frame to alleviate the labeling ambiguity in training

the classifier. Specifically, it enforces that the confidence score of a target bounding box is larger than that of a nearby box by a margin determined by the overlap ratio between two boxes. As such, the Struck method can reduce adverse the labeling ambiguity problems.

For robust tracking, we decompose the problem of target state estimation into the two subproblems of translation estimation and scale estimation [31], [32], [33]. Motivated by Bayesian filtering algorithms [34], [35], we propose a simpler yet effective random strategy for target state refinement.

## 4.1 Translation Estimation

We incorporate the optimized patch weights into the Struck method, in which we improve the robustness to drastic appearance changes and unreliable tracking results of a target object. Given the bounding box of the target object in the previous frame $t - 1$, we first set a searching window in current frame $t$. For $i$-th candidate bounding box within the search window, we weight its patch feature descriptor $\mathbf{x}_{t,i}$ by the weight $\hat{\mathbf{w}}_{t-1,i} = 1/(1 + \exp(-\sigma \overline{\mathbf{w}}_{t-1,i}))$, and concatenate them into a vector as the feature representation:

$$\hat{\mathbf{x}} = [\hat{\mathbf{w}}_{t-1,1}\mathbf{x}_{t,1}; \hat{\mathbf{w}}_{t-1,2}\mathbf{x}_{t,2}; \dots; \hat{\mathbf{w}}_{t-1,n}\mathbf{x}_{t,n}], \quad (22)$$

where we normalize $\mathbf{w}$ as $\overline{\mathbf{w}}$ so that all elements of $\overline{\mathbf{w}}$ sum to 1, and the parameter $\sigma$ is fixed to a pre-defined number (e.g., 37 in this work). Herein, we use the Sigmoid function to map the normalized weights into the range of 0 to 1, which has a parameter $\sigma$ to control the steepness of normalized weights. The optimal bounding box $b_t^*$ can be selected to update the object location by maximizing the classification score:

$$\hat{b}_t = \arg\max_b \ (\omega \langle \mathbf{h}_{t-1}, \hat{\mathbf{x}}_{t,b} \rangle + (1 - \omega) \langle \mathbf{h}_0, \hat{\mathbf{x}}_{t,b} \rangle), \quad (23)$$

where $\mathbf{h}_0$ is learned in the initial frame, which can alleviate the issue of learning drastic appearance changes, and $\omega$ is a weight parameter.

## 4.2 Scale Estimation

Given the estimated location $\hat{b}_t$, we sample a set $\mathbb{B}_t$ of bounding boxes from the Gaussian distribution centered at $\hat{b}_t$, in which the elements of the covariance are the variations of the affine parameters, and its setting depends on motion variations of the target object. To simultaneously estimate scales and refine locations, we utilize four independent affine parameters to draw samples including the scale factor, aspect ratio and translation. For example, we empirically set to these parameters (scale factor, aspect ratio, and translation) to 0.05, 0.01, 1 and 1, respectively in this paper. As object translation is estimated before, we use 100 samples in this paper to compute scale while slightly adjusting translation for a trade-off between efficiency and accuracy. The bounding box is updated by the one with the highest score,

$$b_t^* = \arg\max_{b \in \mathbb{B}_t} \ (\omega \langle \mathbf{g}_{t-1}, \hat{\mathbf{x}}_{t,b} \rangle + (1 - \omega) \langle \mathbf{g}_0, \hat{\mathbf{x}}_{t,b} \rangle), \quad (24)$$

where $\mathbf{g}_{t-1}$ and $\mathbf{g}_0$ are classifiers trained in scale spaces at time $t - 1$ and 0, respectively.

To update the classifier $\mathbf{g}_t$, we use a similar method to translation estimation. Given the optimal estimate $b_t^*$, we extract bounding boxes $\mathbb{S}_t$ around $b_t^*$ at different scales and the corresponding feature representations for scale factors $\{0.50, 0.52, \dots, 0.98, 1.02, \dots, 1.48, 1.50\}$ excluding the positive sample with the scale factor 1 [36]. We then find the optimal $\mathbf{g}_t^*$ by

$$\mathbf{g}_t^* = \arg\min_{\mathbf{g}_t} \ \xi||\mathbf{g}_t||^2 + \sum_{b \in \mathbb{S}_t} \max\{0, \triangle(b, b_t^*) - \langle \mathbf{g}_t, \mathbf{x}_b - \mathbf{x}_{b_t^*} \rangle\}, \quad (25)$$

where $\triangle(b_t, b) = 1 - IoU(b_t, b)$, and $IoU$ indicates the Intersection-over-Union operation. To optimize (25), we use the stochastic variance reduced gradient scheme [37]. To reduce the sensitivity to noises of scale update, we carry out scale estimation at an interval of 3 frames.

## 4.3 Model Update

To alleviate the issues of model drift, we update the classifier and patch weights only when the confidence score of tracking result is larger than a threshold $\theta$. In this paper, the confidence score of tracking result in $t$-th frame is defined as the average similarity between the weighted descriptor of the tracked bounding box and the positive support vectors

$$\frac{1}{|\mathbb{P}_t|} \sum_{\mathbf{s} \in \mathbb{P}_t} \langle \mathbf{s}, \hat{\mathbf{x}}_{t, b_t^*} \rangle, \quad (26)$$

where $\mathbb{P}_t$ is the set of the positive support vectors at time $t$. Algorithm 2 shows the main steps of the proposed tracking method.

## 4.4 Discussion

It should be noted that the proposed tracking algorithm is significantly different from the recently proposed approaches that use sparse representation for object tracking [38], [39], [40], [34] in which reconstruction errors or representation coefficients are used to compute the confidence of candidates in the Bayesian filtering framework. While we employ the sparse representation to learn a dynamic graph for representing a target object, the node weights are used to suppress the effects of background clutter in the tracking-by-detection framework.

In addition, our approach is also significantly different from the SOWP [7] method in several aspects. First, the proposed algorithm learns a dynamic graph to represent a target object that better captures the intrinsic relationship among image patches. Second, our method optimizes the edge and node weights jointly while the SOWP method first computes the edge weights and then the node weights. Third, the proposed tracker considers the initial foreground and background clutter in a unified model, while the SOWP method requires two passes to compute the final patch weights, one for foreground and another for background.

---

**Algorithm 2** Proposed Object Tracking Algorithm

---

**Input:** Input video sequence,
    initial target bounding box $b_0$.
**Output:** Estimated object state $b_t^*$.

 1: // Initialization
 2: Compute $\hat{\mathbf{w}}_0$, $\mathbf{h}_0$ and $\mathbf{g}_0$ according to $b_0$;
 3: **repeat**
 4:     Set the searching window in $t$-th frame according to $b_{t-1}^*$ and extract features $\mathbf{x}_t$;
 5:     // Feature construction
 6:     Construct feature representation $\hat{\mathbf{x}}_t$ using $\mathbf{x}_t$ and $\hat{\mathbf{w}}_{t-1}$;
 7:     // Translation estimation
 8:     Estimate object location $\hat{b}_t$ by (23);
 9:     // Scale estimation
10:     Estimate final object state $b_t^*$ by (24);
11:     // Model and weight update
12:     **if** $\frac{1}{|\mathbb{P}_t|}\sum_{\mathbf{s}\in\mathbb{P}_t}\langle\mathbf{s},\hat{\mathbf{x}}_{t,b_t^*}\rangle < \theta$ **then**
13:         // Weight computation
14:         Run Algorithm 1 for computing the patch weights $\hat{\mathbf{w}}_t$ according to $b_t^*$;
15:         Update $\mathbf{h}_t$ and $\mathbf{g}_t$;
16:     **end if**
17: **until** *End of video sequence*.

---

# 5 PERFORMANCE EVALUATION

The proposed tracker based weighted patch-based graph (WPG) representation is implemented in C++. All experiments are carried out on a machine with an Intel i7 4.0 GHz CPU and 32 GB RAM. We test runtime of WPG on the OTB100 dataset [23], and scale each frame such that the minimum side length of a bounding box is 32 pixels for efficiency. The proposed algorithm is able to track a target object at 5 frames per second where the optimization method converges within 50 iterations. We use the benchmark datasets and protocols [23], [24], [41] to evaluate the proposed approach. In addition, we evaluate several variants of the proposed method to demonstrate the contribution of main modules.

## 5.1 Experimental Setup

### 5.1.1 Parameters

For fair comparisons, we fix all parameters and other settings on all datasets in our experiments. We partition all bounding box into 64 non-overlapping patches as a trade-off between accuracy and efficiency, and extract color and gradient histograms for each patch, where the dimension of gradients and each color channel is set to be 8. We evaluate different number of patches from $\{36, 49, 64, 81, 100\}$, and empirically determine that the proposed method performs best with 64 patches as a trade-off between accuracy and complexity. Note that we fix patch number as square to adapt patch size to the size of object bounding box, which makes patches have consistent shape with target object. Otherwise, it is hard to find a unified partition method for all sequences.
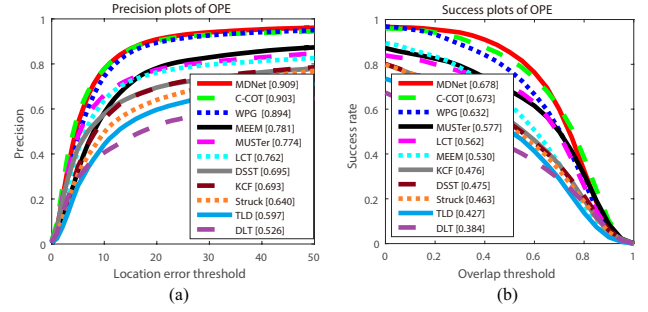


Fig. 3. Precision and success plots of OPE (one-pass evaluation) [23] of the proposed tracker against other state-of-the-art trackers on OTB100. The representative score of PR is presented in the legend.

To improve efficiency, each frame is scaled such that the minimum side length of a bounding box is 32 pixels. A bounding box is described by with and height of $W$ and $H$ pixels. The side length of a search window is initially set to a small range, $0.8\sqrt{WH}$, to reduce false positives, and then set to a large range, $\sqrt{WH}$, to handle abrupt motions if the center distance of the object box between two consecutive frames is above a predefined threshold (e.g., 5 pixels in this work).

For seed selection, we shrink and expand the tracked bounding box $(lx, ly, W, H)$ as $(lx + 0.2W, ly + 0.2H, 0.6W, 0.6H)$ and $(lx - W', ly - H', W + 2W', H + 2H')$, respectively, where $(lx, ly)$ denotes the upper left coordinate of the tracked bounding box, and $W'$ and $H'$ indicate the patch width and height, respectively [7]. In the proposed model (6), there involves several parameters, which are set as follows. On one hand, similar to [17], we simplify the settings as $\alpha = \beta$ and $\lambda_3 = \lambda_4$. Following [17], we set $\{\alpha, \beta, \gamma, \lambda_3, \lambda_4, \xi\} = \{0.1, 0.1, 10, 1, 1, 6\}$. Although $\gamma$ is 2 orders of magnitude higher than $\alpha$ and $\beta$, we find that these terms can balance well by outputting each term after optimization. On the other hand, $\lambda_1$ and $\lambda_2$ are to control the balance of smoothness and fitness of $\mathbf{w}$. According to the setting of similar models [12], [22], we set $\{\lambda_1, \lambda_2\} = \{5, 0.5\}$. For the Struck method, we empirically set $\{\omega, \theta\} = \{0.67, 0.25\}$ [7].

### 5.1.2 OTB100 Dataset

We evaluate the proposed tracking method on the OTB100 benchmark dataset [23]. The OTB100 dataset contains 100 image sequences with ground-truth object locations and attributes for performance analysis. We use precision rate (PR) and success rate (SR) with the threshold of 20 pixels for quantitative performance.

### 5.1.3 Temple Color Dataset

For further validating the effectiveness of the proposed approach, we also compare with other tracking approaches on the Temple Color dataset [42]. This database contains 128 challenging image sequences of human,

(a) Basketball

(b) Bolt2

(c) Diving

(d) Iornman

(e) Human3

(f) Box

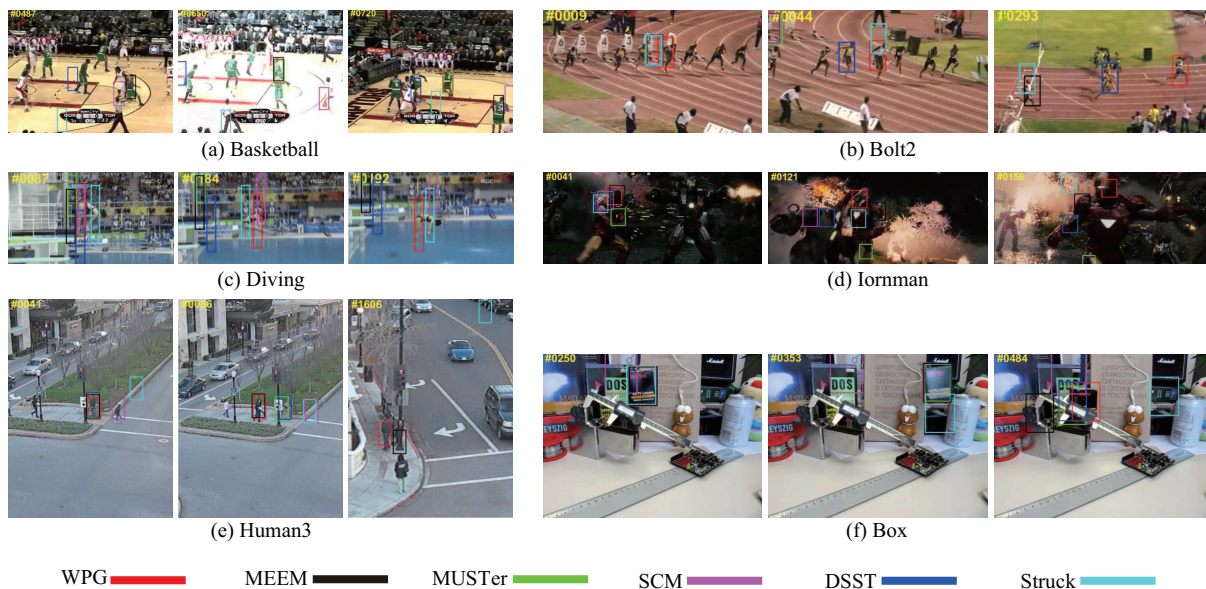WPG ▬▬▬  MEEM ▬▬▬  MUSTer ▬▬▬  SCM ▬▬▬  DSST ▬▬▬  Struck ▬▬▬

Fig. 4. Sample results of our method against Struck [3], MEEM [5], MUSTer [32], DSST [33] and SCM [39].
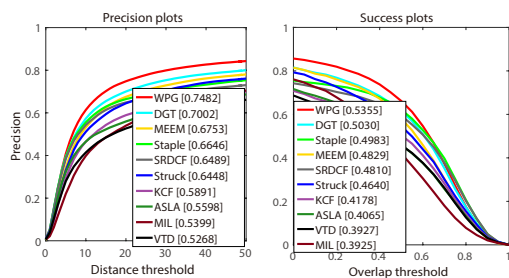


Fig. 5. PR and SR curves on the Temple Color dataset where ten trackers are shown here.

animals and rigid objects. In addition to tracking ground truth, each sequence in the dataset is also annotated by its challenging factors as defined in [23]. The evaluation metrics are also defined in [23].

### 5.1.4 NUS PRO Dataset

We also compare the proposed algorithm with other tracking approaches on the NUS-PRO dataset [24]. This dataset contains 365 challenging image sequences of pedestrians and rigid objects, mainly captured from moving cameras. Aside from target locations, each sequence is annotated with occlusion level for evaluation. We employ the threshold-response relationship (TRR) with three criteria for occlusion computation [24] on the entire dataset to evaluate the proposed tracking method.

### 5.1.5 VOT Challenge Dataset

For more comprehensive evaluation, we also run the proposed tracker on the VOT2014 challenge dataset [41], whose dataset contains more deformations and the aligned bounding boxes contain more noise. Accuracy

(ACC) and robustness (ROB) are used to assess the performance of a tracker. The accuracy computes the overlap ratio between an estimated bounding box and the ground truth. The robustness indicates the number of tracking failures, i.e., the number of frames in which the overlap ratios are zero.

### 5.2 Evaluation on the OTB100 Dataset

We first evaluate the proposed algorithm on the OTB100 dataset against tracking methods. Next we analyze the performance of evaluated methods based on attributes of image sequences.

#### 5.2.1 Tracking Methods Without Deep Features

We evaluate the proposed algorithm against the state-of-the-art tracking methods without using deep features, e.g., Struck [3], DSST [33], MEEM [5], MUSTer [32] and SOWP [7]. Figure 3 shows the OPE plots on the OTB100 dataset, and Figure 4 presents some qualitative results. Overall, the proposed algorithm performs favorably against the state-of-the-art methods, e.g., 9.1% over SOWP in the precision score and 5.5% over MUSTer in the success score. Figure 4 shows that the proposed approach effectively handles scenes with illumination variation (*Basketball* and *Ironman*), background clutter (*Diving*, *Ironman* and *Box*), deformation (*Basketball*, *Bolt2* and *Diving*) and partial occlusion (*Ironman*, *Box* and *Human3*).

The excellent performance of WPG suggests that the proposed tracker is able to mitigate outlier effects by integrate local patch weights into feature representations, which brings biggest performance gain for achieving state-of-the-art tracking performance. In addition to it, the following components are also beneficial to promoting tracking performance. First, local patch representa-

TABLE 1
Attribute-based PR scores on OTB benchmark compared with recent trackers, where the best results of deep and non-deep trackers divided by dash line are in red and green colors, respectively.

| | MDNet | C-COT | HCF | DLT | SOWP | MEEM | MUSTer | KCF | LCT | DSST | Struck | WPG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IV | 0.911 | 0.878 | 0.817 | 0.522 | 0.777 | 0.740 | 0.782 | 0.708 | 0.746 | 0.723 | 0.545 | 0.873 |
| SV | 0.892 | 0.881 | 0.802 | 0.542 | 0.750 | 0.740 | 0.715 | 0.639 | 0.686 | 0.667 | 0.600 | 0.858 |
| OCC | 0.857 | 0.904 | 0.767 | 0.454 | 0.754 | 0.741 | 0.734 | 0.622 | 0.682 | 0.615 | 0.537 | 0.863 |
| DEF | 0.899 | 0.865 | 0.791 | 0.451 | 0.741 | 0.754 | 0.689 | 0.617 | 0.689 | 0.568 | 0.527 | 0.878 |
| MB | 0.866 | 0.899 | 0.797 | 0.427 | 0.710 | 0.722 | 0.699 | 0.617 | 0.673 | 0.636 | 0.594 | 0.817 |
| FM | 0.885 | 0.883 | 0.797 | 0.426 | 0.719 | 0.735 | 0.691 | 0.628 | 0.675 | 0.602 | 0.626 | 0.824 |
| IPR | 0.910 | 0.877 | 0.854 | 0.471 | 0.828 | 0.794 | 0.773 | 0.693 | 0.782 | 0.724 | 0.637 | 0.877 |
| OPR | 0.900 | 0.899 | 0.810 | 0.517 | 0.790 | 0.798 | 0.748 | 0.675 | 0.750 | 0.675 | 0.593 | 0.882 |
| OV | 0.825 | 0.895 | 0.677 | 0.558 | 0.633 | 0.685 | 0.591 | 0.498 | 0.558 | 0.487 | 0.503 | 0.802 |
| BC | 0.925 | 0.882 | 0.847 | 0.509 | 0.781 | 0.752 | 0.786 | 0.716 | 0.740 | 0.708 | 0.566 | 0.885 |
| LR | 0.942 | 0.975 | 0.787 | 0.615 | 0.713 | 0.605 | 0.677 | 0.545 | 0.490 | 0.595 | 0.674 | 0.948 |
| All | 0.909 | 0.903 | 0.837 | 0.526 | 0.803 | 0.781 | 0.774 | 0.692 | 0.762 | 0.695 | 0.640 | 0.894 |

TABLE 2
Attribute-based SR scores on OTB benchmark compared with recent trackers, where the best results of deep and non-deep trackers divided by dash line are in red and green colors, respectively.

| | MDNet | C-COT | HCF | DLT | SOWP | MEEM | MUSTer | KCF | LCT | DSST | Struck | WPG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| IV | 0.684 | 0.674 | 0.540 | 0.408 | 0.554 | 0.517 | 0.600 | 0.474 | 0.566 | 0.489 | 0.422 | 0.632 |
| SV | 0.658 | 0.654 | 0.488 | 0.399 | 0.478 | 0.474 | 0.518 | 0.399 | 0.492 | 0.413 | 0.404 | 0.579 |
| OCC | 0.646 | 0.674 | 0.525 | 0.335 | 0.528 | 0.504 | 0.554 | 0.438 | /0.507 | 0.426 | 0.394 | 0.619 |
| DEF | 0.649 | 0.614 | 0.530 | 0.295 | 0.527 | 0.489 | 0.524 | 0.436 | 0.499 | 0.412 | 0.383 | 0.605 |
| MB | 0.679 | 0.706 | 0.573 | 0.353 | 0.557 | 0.545 | 0.557 | 0.456 | 0.532 | 0.465 | 0.468 | 0.622 |
| FM | 0.675 | 0.676 | 0.555 | 0.345 | 0.542 | 0.529 | 0.539 | 0.455 | 0.527 | 0.440 | 0.470 | 0.603 |
| IPR | 0.655 | 0.627 | 0.559 | 0.348 | 0.567 | 0.529 | 0.551 | 0.465 | 0.557 | 0.485 | 0.453 | 0.606 |
| OPR | 0.661 | 0.652 | 0.537 | 0.376 | 0.549 | 0.528 | 0.541 | 0.454 | 0.541 | 0.453 | 0.424 | 0.608 |
| OV | 0.627 | 0.648 | 0.474 | 0.384 | 0.497 | 0.488 | 0.469 | 0.393 | 0.452 | 0.374 | 0.384 | 0.586 |
| BC | 0.676 | 0.652 | 0.587 | 0.553 | 0.575 | 0.523 | 0.579 | 0.498 | 0.481 | 0.373 | 0.438 | 0.652 |
| LR | 0.631 | 0.629 | 0.424 | 0.422 | 0.416 | 0.355 | 0.477 | 0.306 | 0.330 | 0.311 | 0.313 | 0.575 |
| All | 0.678 | 0.673 | 0.562 | 0.384 | 0.560 | 0.530 | 0.577 | 0.475 | 0.562 | 0.475 | 0.463 | 0.632 |

tions are robust to object deformation and partial occlusion. Second, the classification and update schemes are used to avoid model contagious by drastic appearance changes and unreliable tracking results of a target object. Finally, the scale handling strategy is employed to adapt to scale variations and also refine object translation.

### 5.2.2 Tracking Methods Based on Deep Features

We evaluate the proposed algorithm against the state-of-the-art tracking methods using deep features including DLT [43], HCF [44], C-COT [45] and MDNet [46]. Figure 3, Table 1 and Table 2 show the evaluation results. Overall, the proposed tracker performs well against the DLT and HCF methods in all aspects. The proposed tracker performs equally well against the C-COT and MDNet schemes in terms of precision and slightly worse in terms of success rate. Furthermore, the proposed algorithm differs from the C-COT and MDNet methods in several aspects.

- The proposed tracking method does not require laborious pre-training or a large training set. In addition, it does not need to save a large pre-trained deep model. We initialize the proposed model using the ground truth bounding box in the first frame, and update it in subsequent frames.
- It is easy to implement as each subproblem of the proposed model has a closed-form solution.

- It performs more robustly than the MDNet and C-COT methods in some situations. In particular, it outperforms the C-COT method on sequences with background clutters in terms of precision and success rate, which suggests the effectiveness of our approach in suppressing the background effects during tracking.

### 5.2.3 Attribute-based Evaluation

We present the precision plots with 11 different attributes in Table 1 and Table 2. The attributes include background clutter (BC), deformation (DEF), fast motion (FM), illumination variation (IV), in-plane rotation (IPR), low resolution (LR), motion blur (MB), occlusion (OCC), out-of-plane rotation (OPR), out of view (OV) and scale variation (SV).

The comparison plots show that our tracker significantly outperforms other non-DL-based tracking methods, and achieves comparable performance with DL-based ones on the attribute-based subsets (e.g., BC and DEF), which validates the effectiveness of introducing the optimized weights in the object representation that suppresses background clutter and noises. The performance of our tracker against others on OCC and OV suggests that the adopted classification and update schemes can re-track objects in case of tracking failure, e.g., totally occlusion and re-entering the field of view,
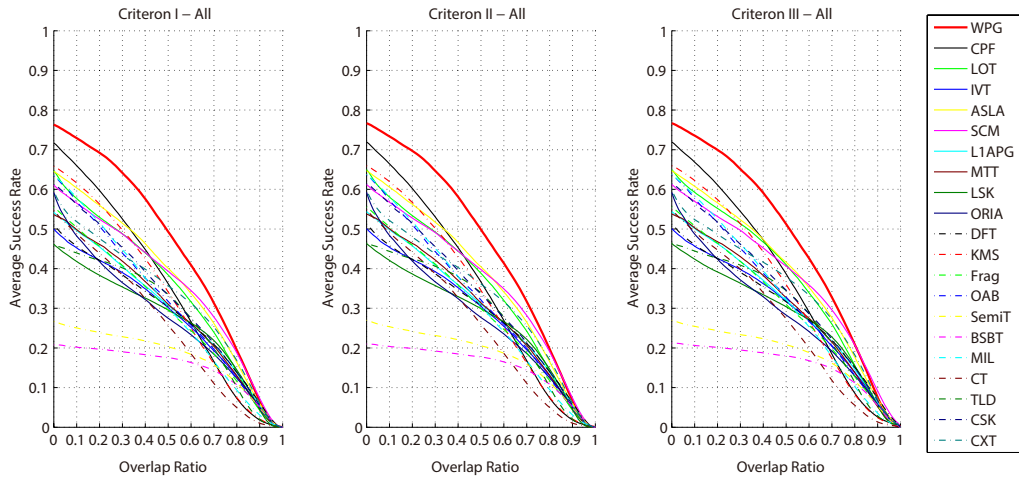
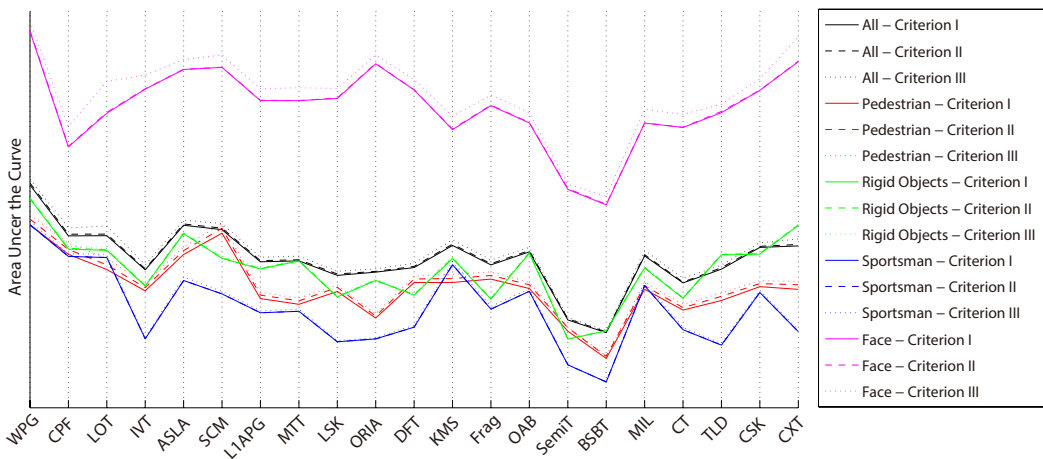Fig. 6. TRR curves on NUS-PRO, where twenty trackers are shown here.



Fig. 7. AUC plots of TRR curves with different object categories on NUS-PRO, where twenty trackers are shown here.

and alleviate incorrect update of noisy samples. The even worse performance of our tracker against others on FM and LR suggests the weakness of our used features (color and gradient) in representing the target object and search strategy, and we will address these issues in future work.

In particular, we compare our WPG with the SOWP method [7] that is most related to us as follows. For the PR score, WPG outperforms the SOWP method significantly, especially on the sequences with deformation, out of view, background clutter and low resolution. It demonstrates advances of WPG over SOWP in learning robust object feature representations under background inclusion and less information, and also in re-tracking objects after they back to view. For the SR score, WPG also excels SOWP with a large margin, especially on the challenges of scale variation, occlusion and low resolution, which verify the effectiveness of scale handling, background suppression and reliability highlighting in WPG, while SOWP does not handle scale variations and is also limited by its weight computation scheme.

## 5.3 Evaluation on the Temple Color Dataset

We evaluate the proposed algorithm on the Temple Color dataset [42]. Figure 5 shows the evaluation results against 9 state-of-the-art tracking approaches, including DGT [22], Staple [47], MEEM [5], SRDCF [48], Struck [49], KCF [50], ASLA [51], MIL [26], and VTD [52]. Overall, the proposed algorithm performs favorably against the other trackers, e.g., DGT (Our previous version) (PR/SR: 4.8%/3.2%), Staple (8.4%/3.7%) and SRDCF (10.0%/5.4%).

## 5.4 Evaluation on the NUS-PRO Dataset

We evaluate the proposed algorithm against the state-of-the-art trackers on the NUS-PRO [24] dataset.

### 5.4.1 Overall Performance

We present the evaluation results of our method against 20 conventional trackers on the NUS-PRO dataset [24] in Figure 6. Overall, the proposed tracker performs favorably against other trackers on the NUS-PRO dataset.
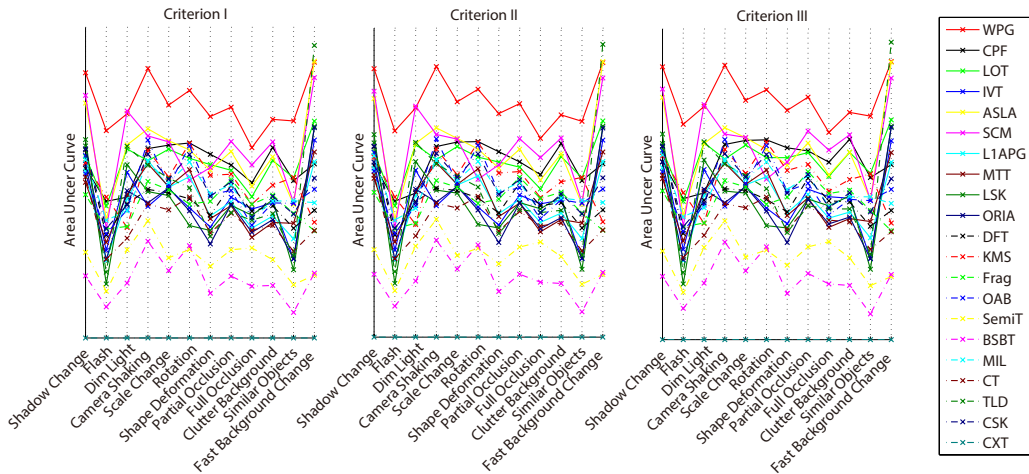
Fig. 8. AUC plots of TRR curves with different challenges on NUS-PRO, where twenty trackers are shown here.

TABLE 3
Comparison of WPG against the SOWP method [7] and the top three trackers on the VOT2014 challenge dataset [41]. 'ACC w/o' denotes the ACC score without the re-initialization step. The best performance is in boldface.

| | Baseline | | | Region noise | | |
|---|---|---|---|---|---|---|
| | ACC | ROB | ACC w/o | ACC | ROB | ACC w/o |
| DSST | **0.66** | 1.16 | 0.46 | **0.57** | 1.28 | 0.43 |
| SAMF | 0.61 | 1.28 | 0.50 | **0.57** | 1.43 | 0.43 |
| KCF | 0.62 | 1.32 | 0.39 | **0.57** | 1.51 | 0.36 |
| SOWP | 0.57 | 0.56 | 0.51 | 0.55 | 0.68 | 0.48 |
| WPG | 0.57 | **0.53** | **0.52** | 0.55 | **0.50** | **0.50** |

The results of the top 4 performing methods (CPF [53], ASLA [51], SCM [39] and LOT [54]) show that the combination of local feature representations and particle filter search models can achieve the state-of-the-art performance. Although adopting only the local feature representation, the proposed tracking algorithm performs well on the NUS-PRO dataset.

### 5.4.2 Category-based Evaluation

We present how the proposed tracker performs on 4 object categories in the NUS-PRO database. The AUC plots of TRR curves in Figure 7 show that the proposed method performs well in the *rigid object*, *sportsman* and *face* sequences, and comparably with the SCM scheme in the *pedestrian* sequences. The sportsman category is the most challenging among 4 object types in the NUS-PRO database, followed by the classes of *pedestrians*, *rigid objects* and *faces*.

### 5.4.3 Attribute-based Evaluation

We present the AUC plots of TRR curves of the evaluated tracking algorithms based on 12 attributes, including shadow change (SC), flash (FL), dim light (DL), camera shaking (CS), scale change (SC), rotation (RO), shape

deformation (SD), partial occlusion (PO), full occlusion (FO), clutter background (CB), similar objects (SO) and fast background change (FBC). The proposed tracker performs well on scenes with most attributes including SC, DL, CS, SC, RO, SD, PO, FO, CB and SO. The evaluation results are consistent with the findings on the OTB100 dataset except that the FL, DL, CS and FBC attributes are not reported on the OTB100 dataset and the proposed method performs slightly worse than the others on the scenes with the DL and FBC attributes. The performance on the sequences with the DL attribute may be explained by the adopted features (color and gradient) of the proposed method for representing target objects under low illumination conditions, which can be improved by using integrating more features [34]. On the other hand, the performance of the proposed algorithm on sequences with the FBC attribute can be explained by the search strategy, and can be further improved by using robust motion or search models to leverage more temporal and spatial information.

## 5.5 Evaluation on the VOT Challenge Dataset

Finally, we report the evaluation results of WPG against SOWP [7] and the top three trackers (i.e., DSST [33], SAMF [55] and KCF [50]) on the VOT2014 challenge dataset [41], as shown in Table 3. In Baseline evaluation, a tracker is initialized with a ground truth. In Region noise evaluation, a tracker inputs a perturbed ground truth.

From Table 3, we can see that WPG obtains low ACC scores, achieves the best ROB results in both evaluations. In the VOT challenge, a re-initialization step is triggered using a new ground truth when a tracker is detected as failure. Therefore, the compared trackers fail to track more frequently than WPG, and thus they obtain higher overlap ratios. To mitigate these effects of re-initialization, we remove re-initialization step in evaluations, and denote overlap ratios as ACC w/o.

TABLE 4
Performance of 4 variants of the proposed method
against the SOWP method [7].

|     | SOWP  | $\text{WPG}_A$ | $\text{WPG}_Z$ | $\text{WPG}_E$ | $\text{WPG}_W$ |
|-----|-------|------|------|------|------|
| PR  | 0.803 | 0.870 | 0.882 | 0.873 | 0.811 |
| SR  | 0.560 | 0.610 | 0.626 | 0.612 | 0.597 |

|     | WPG   | $\text{WPG}'_A$ | $\text{WPG}'_Z$ | $\text{WPG}'_E$ |
|-----|-------|------|------|------|
| PR  | 0.894 | 0.861 | 0.883 | 0.878 |
| SR  | 0.632 | 0.612 | 0.627 | 0.624 |

The results show that WPG yields the best ACC scores without the re-initialization.

It is worth noting that sequences of Region noise evaluation contain more clutter and noise, but WPG performs better on Region noise evaluation than on Baseline. It suggests that WPG can handle region noise more effectively than others.

### 5.6 Analysis

To demonstrate the effectiveness of the main components, we present empirical results using 4 variants of the proposed algorithm on the OTB100 dataset . These variants are: 1) $\text{WPG}_W$: We remove the patch weights in our tracking algorithm, 2) $\text{WPG}_A$: We remove the affinity learning and directly utilize the representation coefficients to diffuse patch weights. The objective function is:

$$
\begin{aligned}
\min_{\mathbf{Z},\mathbf{E},\mathbf{w}} \quad & \alpha||\mathbf{Z}||_0 + \beta||\mathbf{E}||_{2,0} + \lambda_1 \sum_{i,j} \mathbf{z}_{ij}(\mathbf{w}_i - \mathbf{w}_j)^2 \\
& + \frac{\lambda_2}{2}||\Gamma \circ (\mathbf{w} - \mathbf{r})||^2 + \frac{\lambda_4}{2}||\mathbf{w}||^2 \\
\text{s.t.} \quad & \mathbf{X} = \mathbf{XZ} + \mathbf{E}, \ \mathbf{w} \geq 0.
\end{aligned}
\tag{27}
$$

We also use the ADMM algorithm to solve (27). 3) $\text{WPG}_Z$: We remove the sparse constraints on $\mathbf{Z}$, but enforce minimizing $||\mathbf{Z}||_F^2$ to avoid the trivial solution. Thus, $\mathbf{Z}$ can be updated with the closed-form solution:

$$
\begin{aligned}
\mathbf{Z}^{k+1} = (\frac{\alpha + \mu^k}{\mu^k}\mathbf{I} + \mathbf{X}^\top\mathbf{X})^{-1}, \\
(\mathbf{X}^\top(\mathbf{X} - \mathbf{E}^k + \frac{\mathbf{Y}_1^k}{\mu^k}) + \mathbf{Q}^k - \frac{\mathbf{Y}_2^k}{\mu^k}).
\end{aligned}
\tag{28}
$$

4) $\text{WPG}_E$: We remove the sample-specific sparse constraints on $\mathbf{E}$ but enforce minimizing $||\mathbf{E}||_F^2$ to avoid the trivial solution. Thus, $\mathbf{E}$ can be updated with a closed-form solution:

$$
\mathbf{E}^{k+1} = \frac{\mu^k}{\mu^k + \beta}(\mathbf{X} - \mathbf{XZ}^{k+1} + \frac{\mathbf{Y}_1^k}{\mu^k}).
\tag{29}
$$

To rules out the implementation flaw or optimization differences, we set parameters $\gamma$, $\alpha$ and $\beta$ to a ridiculously low number (e.g., $10^{-10}$ in this work) to render contribution of each term, and denote them as $\text{WPG}'_A$, $\text{WPG}'_Z$ and $\text{WPG}'_E$, respectively.

Table 4 shows the evaluation results against the SOWP method [7]. The performance gains achieved by the proposed algorithm over the SOWP method demonstrate the significances of the main components. In particular, the results show that: 1) Introducing patch weights into the object representations helps suppress the effects of background clutters in visual tracking by comparing the performance of $\text{WPG}_W$ against the other schemes. WPG is a spatially reliability learning method, which has been proven to be an effective way to mitigate outlier effects, and thus bring big performance gains for achieving state-of-the-art tracking performance [56], [57]. 2) The WPG, $\text{WPG}_A$, $\text{WPG}_Z$ and $\text{WPG}_E$ methods perform well against the SOWP scheme, which suggests that the dynamic graph facilitates optimizing the patch weights by capturing the intrinsic relationship among image patches. Comparing with restriction of spatial neighbors in SOWP, variants of WPG are good at exploring long-range relationships among patches, and also mitigating noise effects of low-level features. Hence, the patch weights optimized by WPG variants are more accurate and robust. 3) The WPG algorithm performs better than the $\text{WPG}_A$, $\text{WPG}_Z$ and $\text{WPG}_E$ schemes, thereby justifying the effectiveness of learning graph affinity matrix $\mathbf{A}$, sparse constraints on $\mathbf{Z}$, and sample-specific sparse constraints on $\mathbf{E}$, respectively. First, sparse representation based graph [58], [28] could automatically select the most informative neighbors for each patch, and explore higher order relationships among patches, hence is more powerful and discriminative. Second, the learned graph could suppress corrupted and noisy image patches by modelling noise in sparse representation.

The performances of $\text{WPG}'_A$, $\text{WPG}'_Z$ and $\text{WPG}'_E$ against WPG further demonstrate above observations and conclusions.

## 6 CONCLUSIONS

In this paper, we propose an effective algorithm for visual tracking by suppressing the effects of background clutters. A patch-based graph is learned dynamically by capturing the intrinsic relationship among patches. To reduce the computational complexity, we develop an efficient algorithm for the proposed model by solving several convex subproblems. Finally, the optimized patch weights are incorporated into the structured SVM framework to carry out the tracking task. Extensive experimental results on three benchmark datasets demonstrate the effectiveness of the proposed algorithm over the state-of-the-art methods. Our future work will focus on: 1) learning the dynamic spatio-temporal graphs to explore more relations among image patches, 2) developing robust motion or search models for addressing fast object or background motions, and 3) replacing the hand-craft features with hierarchical appearance models for more effective representations.
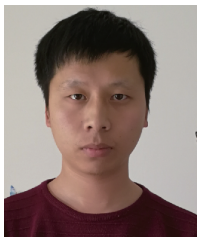
## REFERENCES

[1]  D. Comaniciu, V. Ramesh, and P. Meer, "Kernel-based object tracking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2003. 1, 2

[2] A. Adam, E. Rivlin, and I. Shimshoni, "Robust fragments-based tracking using the integral histogram," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2006. 1

[3] S. Hare, A. Saffari, and P. H. S. Torr, "Struck: Structured output tracking with kernels," in *Proceedings of IEEE International Conference on Computer Vision*, 2011. 1, 2, 3, 6, 9

[4] S. He, Q. Yang, R. Lau, J. Wang, and M.-H. Yang, "Visual tracking via locality sensitive histograms," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2013. 1, 2

[5] 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. 1, 9, 11

[6] T. Zhang, K. Jia, C. Xu, Y. Ma, and N. Ahuja, "Partial occlusion handling for visual tracking via robust part matching," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2014. 1

[7] H.-U. Kim, D.-Y. Lee, J.-Y. Sim, and C.-S. Kim, "Sowp: Spatially ordered and weighted patch descriptor for visual tracking," in *Proceedings of IEEE International Conference on Computer Vision*, 2015. 1, 2, 3, 7, 8, 9, 11, 12, 13

[8] Y. Li, J. Zhu, and S. C. Hoi, "Reliable patch trackers: Robust visual tracking by exploiting reliable patches," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 1

[9] R. Timofte, J. Kwon, and L. V. Gool, "Picaso: pixel correspondences and soft match selection for real-time tracking," *Computer Vision and Image Understanding*, vol. 153, pp. 151–162, 2016. 1

[10] S. Duffner and C. Garcia, "Pixeltrack: A fast adaptive algorithm for tracking non-rigid objects," in *Proceedings of IEEE International Conference on Computer Vision*, 2013. 1, 2

[11] F. Yang, H. Lu, and M.-H. Yang, "Robust superpixel tracking," *IEEE Transactions on Image Processing*, vol. 23, no. 4, pp. 1639–1651, 2014. 1, 2

[12] C. Yang, L. Zhang, H. Lu, X. Ruan, and M.-H. Yang, "Saliency detection via graph-based manifold ranking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2013. 1, 2, 3, 8

[13] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006. 1, 3

[14] E. Elhamifar and R. vidal, "Sparse subspace clustering," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2009. 1, 3, 5

[15] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust recovery of subspace structures by low-rank representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 1, pp. 171–184, 2013. 1, 3, 5, 6

[16] L. Zhuang, H. Gao, Z. Lin, Y. Ma, X. Zhang, and N. Yu, "Nonnegative low rank and sparse graph for semi-supervised learning," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 1, 3, 5

[17] X. Guo, "Robust subspace segmentation by simultaneously learning data representations and their affinity matrix," in *Proceedings of International Joint Conference on Artificial Intelligence*, 2015. 1, 3, 6, 8

[18] C. Li, L. Lin, W. Zuo, W. Wang, and J. Tang, "An approach to streaming video segmentation with sub-optimal low-rank decomposition," *IEEE Transactions on Image Processing*, vol. 25, no. 5, pp. 1947–1960, 2016. 2

[19] Z. Lin, R. Liu, and Z. Su, "Linearized alternating direction method with adaptive penalty for low rank representation," in *Proceedings of Annual Conference on Neural Information Processing Systems*, 2011. 2, 5

[20] K. Wang, L. Lin, J. Lu, C. Li, and K. Shi, "Pisa: Pixelwise image saliency by aggregating complementary appearance contrast measures with edge-preserving coherence," *IEEE Transactions on Image Processing*, vol. 24, no. 10, pp. 3019–3033, 2015. 2

[21] C. Li, L. Lin, W. Zuo, S. Yan, and J. Tang, "Sold: Sub-optimal low-rank decomposition for efficient video segmentation," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 2

[22] C. Li, L. Lin, W. Zuo, and J. Tang, "Learning patch-based dynamic graph for visual tracking," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, 2017. 2, 8, 11

[23] Y. Wu, J. Lim, and M.-H. Yang, "Object tracking benchmark," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015. 2, 4, 8, 9

[24] A. Li, M. Li, Y. Wu, M.-H. Yang, and S. Yan, "Nus-pro: A new visual tracking challenge," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 335–349, 2016. 2, 8, 9, 11

[25] H. Grabner, M. Grabner, and H. Bischof, "Semi-supervised on-line boosting for robust tracking," in *Proceedings of European Conference on Computer Vision*, 2008. 2

[26] 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. 7, pp. 1619–1632, 2011. 2, 11

[27] I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun, "Large margin methods for structured and interdependent output variables," *Journal of Machine Learning Research*, vol. 6, pp. 1453–1484, 2005. 2

[28] J. Wright, Y. Ma, J. Mairal, G. Sapiro, T. S. Huang, and S. Yan, "Sparse representation for computer vision and pattern recognition," *Proceedings of IEEE*, vol. 98, no. 6, pp. 1031–1044, 2010. 5, 13

[29] Z. Lin, A. Ganesh, J. Wright, M. Chen, L. Wu, and Y. Ma, "Fast convex optimization algorithms for exact recovery of a corrupted low-rank matrix," *UIUC Technical Report UILU-ENG-09-2214*, 2009. 5

[30] Y. Xu and W. Yin, "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion," *SIAM Journal on imaging sciences*, vol. 6, no. 3, pp. 1758–1789, 2013. 6

[31] 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. 7

[32] 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. 7, 9

[33] M. Danelljan, G. Hager, F. Khan, and M. Felsberg, "Accurate scale estimation for robust visual tracking," in *Proceedings of British Machine Vision Conference*, 2014. 7, 9, 12

[34] C. Li, H. Cheng, S. Hu, X. Liu, J. Tang, and L. Lin, "Learning collaborative sparse representation for grayscale-thermal tracking," *IEEE Transactions on Image Processing*, vol. 25, no. 12, pp. 5743–5756, 2016. 7, 12

[35] T. Zhang, C. Xu, and M.-H. Yang, "Multi-task correlation particle filter for robust object tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 7

[36] H.-U. Kim and C.-S. Kim, "Locator-checker-scaler object tracking using spatially ordered and weighted patch descriptor," *IEEE Transactions on Image Processing*, vol. 26, no. 8, pp. 3817–3830, 2017. 7

[37] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," in *Proceedings of Annual Conference on Neural Information Processing Systems*, 2013. 7

[38] X. Mei and H. Ling, "Robust visual tracking and vehicle classification via sparse representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 11, pp. 2259–2272, 2011. 7

[39] W. Zhong, H. Lu, and M.-H. Yang, "Robust object tracking via sparsity-based collaborative model," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 7, 9, 12

[40] C. Bao, Y. Wu, H. Ling, and H. Ji, "Real time robust $l_1$ tracking using accelerated proximal gradient approach," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 7

[41] M. Kristan, R. Pflugfelder, and et al., "The visual object tracking vot2014 challenge results," in *Proceedings of European Conference on Computer Vision Workshop*, 2014. 8, 9, 12

[42] P. Liang, E. Blasch, and H. Ling, "An approach to streaming video segmentation with sub-optimal low-rank decomposition," *IEEE Transactions on Image Processing*, vol. 24, no. 12, pp. 5630–5644, 2015. 8, 11

[43] N. Wang and D.-Y. Yeung, "Learning a deep compact image representation for visual tracking," in *Proceedings of Annual Conference on Neural Information Processing Systems*, 2013. 10

[44] 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. 10

[45] 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. 10

[46] 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

[47] 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. 11

[48] M. Danelljan, G. Hager, F. Shahbaz Khan, and M. Felsberg, "Learning spatially regularized correlation filters for visual tracking," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015. 11

[49] 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. 11

[50] 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*, 2015. 11, 12

[51] X. Jia, H. Lu, and M.-H. Yang, "Visual tracking via adaptive structural local sparse appearance model," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 11, 12

[52] J. Kwon and K. M. Lee, "Visual tracking decomposition," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2010. 11

[53] P. Perez, C. Hue, J. Vermaak, and M. Gangnet, "Color-based probabilistic tracking," in *Proceedings of European Conference on Computer Vision*, 2002. 11

[54] S. Oron, A. Bar-Hillel, D. Levi, and S. Avidan, "Locally orderless tracking," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 12

[55] L. Agapito, M. M. Bronstein, and C. Rother, "A scale adaptive kernel correlation filter tracker with feature integration," in *Proceedings of European Conference on Computer Vision Workshop*, 2014. 12

[56] A. Lukezic, T. Vojir, L. C. Zajc, J. Matas, and M. Kristan, "Discriminative correlation filter with channel and spatial reliability," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 13

[57] C. Sun, D. Wang, H. Lu, and M.-H. Yang, "Correlation tracking via joint discrimination and reliability learning," in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 13

[58] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 2, pp. 210–227, 2009. 13

**Liang Lin** received the BS and PhD degrees from the Beijing Institute of Technology (BIT), Beijing, China, in 1999 and 2008, respectively. He is currently a full Professor with the School of Advanced Computing, Sun Yat-Sen University, Shunde, China. From 2006 to 2007, he was a joint PhD student with the Department of Statistics, University of California, Los Angeles (UCLA), CA. He was a Post-Doctoral Research Fellow with the Center for Vision, Cognition, Learning, and Art of UCLA. He was supported by several promotive programs or funds for his work such as Program for New Century Excellent Talents of Ministry of Education (China) in 2012 and Guangdong Distinguished Young Scholar Fund in 2013. Dr. Lin was a recipient of Best Paper Runners-Up Award in ACM NPAR 2010, Google Faculty Award in 2012, and Best Student Paper Award in IEEE ICME 2014.



**Wangmeng Zuo** received the Ph.D. degree in computer application technology from the Harbin Institute of Technology, Harbin, China, in 2007. From August 2009 to February 2010, he was a Visiting Professor in Microsoft Research Asia. He is currently a Professor in the School of Computer Science and Technology, Harbin Institute of Technology. His current research interests include image modeling and blind restoration, discriminative learning, biometrics, and computer vision. Dr. Zuo is an Associate Editor of the IET Biometrics.



**Jin Tang** received the B.Eng. degree in the School of Automation and the Ph.D. degree in the School of Computer Science and Technology at Anhui University, Hefei, China, in 1999 and 2007, respectively. He is currently a Professor in the School of Computer Science and Technology at Anhui University.



**Chenglong Li** received the M.S. and Ph.D. degrees from the School of Computer Science and Technology, Anhui University, Hefei, China, in 2013 and 2016, respectively. From 2014 to 2015, he was a visiting student with the School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China. He is currently a lecturer at the School of Computer Science and Technology, Anhui University, and also a postdoctoral research fellow at the Center for Research on Intelligent Perception and Computing (CRIPAC), National Laboratory of Pattern Recognition (NLPR), Institute of Automation, Chinese Academy of Sciences (CASIA), China. He was a recipient of the ACM Hefei Doctoral Dissertation Award in 2016.
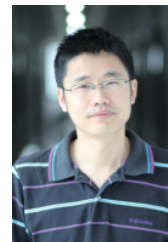


**Ming-Hsuan Yang** is a professor in Electrical Engineering and Computer Science at University of California, Merced. He received the PhD degree in computer science from the University of Illinois at Urbana-Champaign in 2000. Yang served as an associate editor of the IEEE Transactions on Pattern Analysis and Machine Intelligence from 2007 to 2011, and is an associate editor of the International Journal of Computer Vision, Image and Vision Computing and Journal of Artificial Intelligence Research. He received the NSF CAREER award in 2012 and the Google Faculty Award in 2009. He is a senior member of the IEEE and the ACM.