Interacting Multiview Tracker

Ju Hong Yoon, Ming-Hsuan Yang, Senior Member, IEEE, and Kuk-Jin Yoon

Abstract—A robust algorithm is proposed for tracking a target object in dynamic conditions including motion blurs, illumination changes, pose variations, and occlusions. To cope with these challenging factors, multiple trackers based on different feature representations are integrated within a probabilistic framework. Each view of the proposed multiview (multi-channel) feature learning algorithm is concerned with one particular feature representation of a target object from which a tracker is developed with different levels of reliability. With the multiple trackers, the proposed algorithm exploits tracker interaction and selection for robust tracking performance. In the tracker interaction, a transition probability matrix is used to estimate dependencies between trackers. Multiple trackers communicate with each other by sharing information of sample distributions. The tracker selection process determines the most reliable tracker with the highest probability. To account for object appearance changes, the transition probability matrix and tracker probability are updated in a recursive Bayesian framework by reflecting the tracker reliability measured by a robust tracker likelihood function that learns to account for both transient and stable appearance changes. Experimental results on benchmark datasets demonstrate that the proposed interacting multiview algorithm performs robustly and favorably against state-of-the-art methods in terms of several quantitative metrics.

Index Terms—Object tracking, multiview representations, transition probability matrix, tracker interaction, multiple features

1 INTRODUCTION

*T*ISUAL tracking is a fundamental problem in computer vision, which finds a wide range of applications. For practical applications, it is essential for tracking algorithms to account for large appearance changes caused by illumination, pose variations, occlusions, and motion blurs [34] as shown in Fig. 1. To cope with large appearance changes, numerous methods based on multiple features have been proposed for robust visual tracking where different types of features are used complementarily for different scenarios. However, although significant progress has been made in the past decade, it remains a difficult problem to exploit and integrate multiple features for robust visual tracking. The most essential task is how to combine features adaptively to account for appearance changes. Here, it should be noted that each feature has different characteristics against appearance changes. For instance, representations based on histogram of oriented gradients (HOG) [7] are robust to pose variations, and appearance models based on Haar-like features [11] are effective to deal with occlusion.

In this paper, we propose a novel visual tracking algorithm that exploits and integrates multiple feature representations by considering their distinct characteristics to better account for appearance changes for robust tracking.

Features with different and complementary representation strengths are exploited, and multiple feature representations are used by trackers to describe object appearance. Each view (channel) of the multiview (multi-channel) feature learning framework is concerned with one particular representation of a target object [32]. Since each feature is defined in a different space, the likelihood probabilities by multiple trackers are computed at different scales. Consequently, the posterior distribution of each tracker is different even though the object state is defined in the same state space, as illustrated in Fig. 2. Hence, the scale difference should be taken into account when these posterior probabilities are used together for object state estimation. Nevertheless, it is difficult to assign the weights or to project different features to the same space. In this work, instead of combining multiple posterior distributions in mixture form directly, we select the most reliable tracker at each instance. In addition, to prevent unreliable trackers from drifts, the trackers are designed to share their sample distribution information via interaction. Consequently, unreliable trackers receive more reliable samples from reliable ones.

The main components of the proposed algorithm are shown in Fig. 3. At its core, a multiview feature representation [32] of a target object is proposed to account for appearance variations. Each tracker is developed based on one view (representation) of the target object. In addition, these trackers actively interact with each other to provide essential information of samples for effective visual tracking. To integrate multiple trackers for robust visual tracking, we propose the *tracker selection* and *tracker interaction* modules within a Bayesian framework. The tracker selection process determines the most reliable one in terms of tracker probabilities. The trackers share information of sample distributions through interaction based on a transition probability matrix (TPM) and a resampling method to remove unreliable samples. Through this interaction, the visual drifting problem can be alleviated. In the proposed algorithm, we

J. H. Yoon was with the School of Information and Communications, Gwangju Institute of Science and Technology, and he is currently with the Multimedia IP Center, Korea Electronics and Technology Institute, Seongnam-si, Gyeonggido, Republic of Korea. E-mail: jhyoon@keti.re.kr.

M.-H. Yang is with in the School of Engineering, University of California, Merced, USA. E-mail: mhyang@ucmerced.edu.

K.-J. Yoon is with the School of Information and Communications, Gwangju Institute of Science and Technology, Gwangju, Republic of Korea. E-mail: kjyoon@gist.ac.kr.

Manuscript received 27 May 2014; revised 15 Feb. 2015; accepted 4 May 2015. Date of publication 26 Aug. 2015; date of current version 8 Apr. 2016. Recommended for acceptance by S. Avidan.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TPAMI.2015.2473862

^{0162-8828 © 2015} IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See http://www.ieee.org/publications_standards/publications/rights/index.html for more information.



Fig. 1. Tracking results from videos with low contrast, drastic lighting changes, and pose variations (best viewed on high-resolution displays). The proposed algorithm (IMT) performs favorably against three top-ranked trackers (i.e., Struck [13], SCM [37], and ASLA [16]) from a recent benchmark study [31]. Quantitative results are presented in Table 3 and Fig. 7.

approximate the posterior distribution of each tracker by a set of samples. The interaction between trackers is implemented by two operations: retaining its own samples and receiving samples from other trackers. The objective of the transition probability matrix is to determine the number of samples for the aforementioned operations of each tracker.

In addition, to account for object appearance changes, we compute the tracker reliability and update the transition probability matrix to integrate trackers. The update of the transition probability matrix is formulated in a recursive Bayesian framework with a tracker likelihood function (TLF) measuring each tracker reliability at each frame. The reliability of each tracker is used in the tracker interaction and selection processes. Both abrupt and stable appearance changes are considered in the tracker likelihood function. Abrupt appearance changes are modelled by multiple feature representations. On the other hand, stable appearance chances are described by a set of representative templates.



Fig. 2. As trackers are constructed using different features, corresponding posterior distributions ($p(\mathbf{x}|\mathbf{z})$) are of different scales. σ_u denotes the standard deviation of u.

The contributions of the proposed interacting multiview tracking algorithm are as follows. First, we propose a novel tracking algorithm that integrates multiple trackers constructed by different feature representations via selection and interaction. Second, a robust likelihood function is proposed to measure tracker reliability, which is of great importance for robust tracking. Third, a novel tracker interaction scheme is proposed by using the transition probability matrix with a resampling technique. Experimental results on large-scale benchmark datasets show that the proposed tracking algorithm performs favorably against stateof-the-art methods.

Preliminary results of this work were presented in [35]. In this paper, we provide more detailed descriptions and analysis of the proposed interacting multiview tracking algorithm with full derivation and detailed implementation. We compare the proposed algorithm with 10 top performing trackers on 51 benchmark sequences from [31]. In



Fig. 3. Components of the proposed tracking algorithm.

addition, the three most related methods (CVT [22], MCS [4], and FCT [15]) are compared with detailed analysis. Furthermore, additional analysis is presented to demonstrate the effectiveness of the proposed interacting algorithm.

2 RELATED WORK AND PROBLEM CONTEXT

Numerous tracking methods have been proposed using multiple features over the past decade. In this section, we discuss the approaches that are closely related to our work, where appearance models are constructed based on different features. The tracking algorithms that use multiple features can be categorized as a single tracker with multiple observations [6], [30], [36], cascade trackers [10], [26], and parallel trackers [3], [4], [20], [22].

2.1 Multiple Observations

Assuming that features are conditionally independent, multiple observations are combined in product form for visual tracking [6], [30], [36]. However, the reliability of each observation model (based on one different feature) in these approaches is not estimated for combination, which is of crucial importance as each feature is effective for describing certain appearance changes (e.g., pose, illumination, and blur). In contrast, the reliability of each tracker in this work is measured by the tracker likelihood function and reflected in the tracker integration process.

2.2 Cascade Trackers

In [10], a visual tracking method based on a coupled hidden Markov model to combine particle filters and visual cues is proposed. The approach in [26] sequentially estimates object states using the Kalman and particle filters with multiple features including rectangular shape, discriminative cues between the foreground and background, color distribution, and object contour. The state predictions from Kalman filter based on rectangular shape are passed to the other particle filters for sequential processing. These estimated states are combined in a Bayesian filter to determine the object location in each frame. In [26], the adopted features are dependent and the sequential state predictions from early stages are forwarded to the next stage for processing and integration. Thus, it is difficult to add new trackers using other features for different tasks. In the proposed algorithm, all trackers operate in parallel and interact with others, thereby facilitating the addition of other trackers when necessary.

2.3 Parallel Trackers

In [22] and [4], two trackers with different features are combined and target locations are estimated by fusing tracking outputs [22] or selecting the most reliable one [4] based on covariance matrices of posterior distributions. However, a covariance matrix is not effective for measuring the reliability of a tracker when each posterior distribution is computed using observation models with different features (See Fig. 2). Different from [22] and [4], the proposed algorithm selects the most reliable tracker via the proposed tracker likelihood function rather than covariance matrices. The tracker likelihood function is designed to deal with both abrupt and stable appearance changes. Furthermore, the proposed method provides a more general framework that accommodates more than two feature representations. In [20], multiple trackers constructed from four observation models (based on hue, saturation, intensity, and edge features) and two motion models are used to account for appearance and motion changes. While all trackers operate in parallel, the interaction among trackers is based on heuristics as uniform sampling is carried out with a threshold computed by a normalized likelihood ratio. In contrast, the proposed interaction scheme utilizes the transition probability matrix which represents probabilistic dependencies between trackers. Since the transition probability matrix is recursively updated by measuring the reliability of each tracker, unreliable trackers become more dependent on reliable ones to draw samples. As a result, the drifting problem with unreliable trackers is alleviated.

3 ALGORITHMIC OVERVIEW

In the proposed algorithm, multiple interacting trackers based on different feature representations are used as shown in Fig. 3. The reliability of trackers as well as their inter-dependencies are taken into account, and in turn so are the drawn samples from an individual tracker. First, each tracker estimates the object state independently, and then the reliability of each estimated object state is measured by the robust tracker likelihood function. These likelihoods are used to update the tracker probabilities to select the most reliable one. In addition, the result from the most reliable tracker is used to update the object appearance in the representation update. To compute the current dependencies of each tracker on other trackers, the transition probability matrix is also updated by using the likelihoods from the TLF. By using the TPM, the tracker interaction makes unreliable trackers to depend more on the reliable ones to prevent the unreliable trackers from drifting. These interacted trackers are used to estimate the object state for the next frame.

4 STATE ESTIMATION BY TRACKERS

The goal of visual tracking is to estimate an object state given the observations $\mathbf{z}_{1:t} = {\mathbf{z}_1, \dots, \mathbf{z}_t}$ up to time *t*. In this work, the object state is defined as $\mathbf{x}_t = [u_t, v_t, \theta_t, s_t, \alpha_t, \phi_t]^{\top}$ where (u_t, v_t) , θ_t , s_t , α_t , and ϕ_t denote the position, rotation angle, scale, aspect ratio, and skew direction, respectively, to account for affine motion. To robustly handle different kinds of appearance changes, we exploit multiple features for observation models of multiple trackers. Let $m_t \in$ $\{1, \ldots, M\}$ denote the index of *M* trackers constructed from M different features. For simplicity, we denote the *i*th tracker index as $m_t^i \triangleq \langle m_t = i \rangle$. We propose algorithms for the interaction and selection of M trackers. The tracker selection process determines the most reliable tracker at each frame. On the other hand, the drifting problem for the other M-1 trackers is alleviated via tracker interaction. Different from the method based on multiple models [5] where several motion predictions are used for feature point tracking, we exploit a number of representations in the proposed algorithm. Furthermore, we propose a novel tracker interaction approach using a particle filter.

The reliability of the *i*th tracker is represented by the tracker probability $P\{m_t^i | \mathbf{z}_{1:t}\}$. The posterior distribution of object state \mathbf{x}_t by the *i*th tracker is computed by

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}, m_t^i) = \frac{p(\mathbf{z}_k | \mathbf{x}_t, m_t^i) p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, m_t^i)}{\int p(\mathbf{z}_t | \mathbf{x}_t, m_t^i) p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, m_t^i) d\mathbf{x}_t}, \qquad (1)$$

where $p(\mathbf{z}_t | \mathbf{x}_t, m_t^i)$ is the observation model of the *i*th tracker and $p(\mathbf{x}_t | \mathbf{z}_{1:t-1}, m_t^i)$ is a sample distribution by the *i*th tracker given the observations up to time *t*-1 computed via interaction.

4.1 Tracker Interaction

The predicted distribution is computed with mixing probabilities $P\{m_{t-1}^j|m_t^i,\mathbf{z}_{1:t-1}\}$ by

$$p(\mathbf{x}_{t}|\mathbf{z}_{1:t-1}, m_{t}^{i}) = \int p(\mathbf{x}_{t}|\mathbf{x}_{t-1}, m_{t}^{i}) \tilde{p}(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t}^{i}) d\mathbf{x}_{t-1}, \text{and}$$

$$\tilde{p}(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t}^{i}) = \sum_{j=1}^{M} p(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t-1}^{j}) P\{m_{t-1}^{j}|m_{t}^{i}, \mathbf{z}_{1:t-1}\},$$
(2)

where $p(\mathbf{x}_t | \mathbf{x}_{t-1}, m_t^i)$ is a motion model and $\tilde{p}(\mathbf{x}_{t-1} | \mathbf{z}_{1:t-1}, m_t^i)$ is an interacted prior distribution. The mixing probability is computed by

$$P\{m_{t-1}^{j}|m_{t}^{i}, \mathbf{z}_{1:t-1}\} = \frac{P\{m_{t}^{i}|m_{t-1}^{j}, \mathbf{z}_{1:t-1}\}P\{m_{t-1}^{j}|\mathbf{z}_{1:t-1}\}}{\sum_{l=1}^{M} P\{m_{t}^{i}|m_{t-1}^{l}, \mathbf{z}_{1:t-1}\}P\{m_{t-1}^{l}|\mathbf{z}_{1:t-1}\}}.$$
(3)

Note that both the tracker probability and interaction probability are defined by the discrete probability $P\{\cdot\}$ as the number of the trackers is finite, and they satisfy

$$\sum_{i} P\{m_t^i | \mathbf{z}_{1:t}\} = 1, \quad \sum_{j} P\{m_{t-1}^j | m_t^i, \mathbf{z}_{1:t-1}\} = 1.$$

Motion smoothness is a constraint often considered in feature point tracking [5] and, thus, model probabilities $P\{m_{t-1}^{j}|\mathbf{z}_{1:t-1}\}$ at time *t*-1 are useful. However, in visual tracking, it is not effective to use previous model (tracker) probabilities to compute an interacted prior distribution, as occlusion, abrupt pose variations, or significant motion blurs scan cause abrupt appearance changes. Thus, we assume that all tracker probabilities are equal in the interaction scheme and then approximate the mixing probability in (3) by

$$P\{m_{t-1}^{j}|m_{t}^{i}, \mathbf{z}_{1:t-1}\} \approx P\{m_{t}^{i}|m_{t-1}^{j}, \mathbf{z}_{1:t-1}\},$$
(4)

where $P\{m_t^i | m_{t-1}^j, \mathbf{z}_{1:t-1}\}$ is an interaction probability.

4.2 Tracker Selection

We obtain the tracking result $\hat{\mathbf{x}}_t$ by selecting the most reliable tracker that has the highest tracker probability by

$$\hat{\mathbf{x}}_t = \arg\max_{\mathbf{x}_t} p(\mathbf{x}_t | \mathbf{z}_{1:t}, \hat{m}_t),$$

$$\hat{m}_t = \arg\max_{m_t^i} P\{m_t^i | \mathbf{z}_{1:t}\}, \ i = 1, \dots, M.$$
(5)

From (2), (4), and (5), both the tracker and interaction probabilities are utilized to estimate the object state and integrate

multiple trackers. In addition, both tracker and interaction probabilities are updated.

5 ONLINE UPDATE

IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE, VOL. 38, NO. 5, MAY 2016

In contrast to existing methods based on multiple trackers [4], [22], we estimate not only object states but also the tracker and interaction probabilities for efficient and effective integration. Since different features are effective in accounting for certain appearance changes, multiple representations are used to construct trackers. In addition, the reliability of each tracker varies since each one is designed in a different feature space. To achieve robust integration, we consider the reliability of each trackers in the interaction and selection processes.

For notation simplicity, we denote the tracker likelihood function of the *i*th tracker as

$$p(\mathbf{z}_t | m_t^i, \mathbf{z}_{1:t-1}) \triangleq \Lambda_t^i.$$
(6)

Similarly, the notations of the tracker and interaction probabilities are denoted by

$$P\{m_{t}^{i}|\mathbf{z}_{1:t}\} \triangleq T_{t}^{i},$$

$$P\{m_{t}^{i}|m_{t-1}^{j}, \mathbf{z}_{1:t-1}\} \triangleq \bar{\omega}_{t}^{j,i}.$$
(7)

These notations are used in the following sections for updating the tracker and interaction probabilities based on TLF.

5.1 Tracker Probability Update

The tracker probability is updated as

$$P\{m_{t}^{i}|\mathbf{z}_{1:t}\} = \frac{p(\mathbf{z}_{t}|m_{t}^{i}, \mathbf{z}_{1:t-1})}{p(\mathbf{z}_{t}|\mathbf{z}_{1:t-1})} P\{m_{t}^{i}|\mathbf{z}_{1:t-1}\}$$
$$= \frac{p(\mathbf{z}_{k}|m_{t}^{i}, \mathbf{z}_{1:t-1})}{p(\mathbf{z}_{t}|\mathbf{z}_{1:t-1})} \times$$
$$\sum_{j=1}^{M} P\{m_{t}^{i}|m_{t-1}^{j}, \mathbf{z}_{1:t-1}\} P\{m_{t-1}^{j}|\mathbf{z}_{1:t-1}\},$$
(8)

where the total probability $p(\mathbf{z}_t | \mathbf{z}_{1:t-1})$ is expressed by

$$p(\mathbf{z}_{t}|\mathbf{z}_{1:t-1}) = \sum_{i=1}^{M} p(\mathbf{z}_{t}|m_{t}^{i}, \mathbf{z}_{1:t-1}) \times \sum_{j=1}^{M} P\{m_{t}^{i}|m_{t-1}^{j}, \mathbf{z}_{1:t-1}\} P\{m_{t-1}^{j}|\mathbf{z}_{1:t-1}\}.$$
(9)

Based on (8) with the notations in (6) and (7), the sequential tracker probability update is described by

$$T_t^i = \frac{\Lambda_t^i \sum_{l=1}^M \bar{\omega}_{l-1}^{l,i} T_{l-1}^l}{\sum_{j=1}^M \Lambda_t^j \sum_{l=1}^M \bar{\omega}_{l-1}^{l,j} T_{l-1}^l}.$$
 (10)

5.2 Transition Probability Matrix Update

Fig. 4 shows the graphical model of the proposed algorithm based on multiple interacting trackers. A set of interaction probabilities is expressed in a transition probability matrix Ω that describes how trackers affect each other as

$$\mathbf{\Omega} = \begin{bmatrix} \omega^{j,i} \end{bmatrix}_{M \times M} = \begin{bmatrix} \omega^{1,1} & \cdots & \omega^{1,M} \\ \vdots & \ddots & \vdots \\ \omega^{M,1} & \cdots & \omega^{M,M} \end{bmatrix},$$
(11)



Fig. 4. Graphical model: Hidden variable (object state \mathbf{x}_t , a selected tracker index m_t , TPM Ω_t) and observation (observed image \mathbf{z}_t). 1) The TPM is updated using the current observation. 2) The tracker selection is conducted by updating the tracker probability based on the current observation and the TPM. 3) Each object state is estimated based on current observation, tracker selection, tracker interaction, and TPM.

where $\mathbf{\Omega}$ is an unknown matrix from some given prior distributions. The estimated $\overline{\mathbf{\Omega}}_t$ is computed by the minimum mean squared error based on its posterior distribution,

$$\bar{\mathbf{\Omega}}_{t} = \left[\bar{\omega}_{t}^{j,i}\right]_{M \times M} \triangleq E[\mathbf{\Omega} | \mathbf{z}_{1:t}] = \int \mathbf{\Omega} p(\mathbf{\Omega} | \mathbf{z}_{1:t}) d\mathbf{\Omega}.$$
(12)

The goal is to estimate the posterior distribution of the TPM within the Bayesian framework [17],

$$p(\mathbf{\Omega}|\mathbf{z}_{1:t}) = \frac{p(\mathbf{z}_t|\mathbf{\Omega}, \mathbf{z}_{1:t-1})}{p_{\mathbf{\Omega}}(\mathbf{z}_t|\mathbf{z}_{1:t-1})} p(\mathbf{\Omega}|\mathbf{z}_{1:t-1}),$$
(13)

where the TPM observation model $p(\mathbf{z}_t|\Omega, \mathbf{z}_{1:t-1})$ is derived in (16) by approximating the unknown Ω with $\overline{\Omega}_{t-1}$, and $\overline{\Omega}_{t-1}$ is the best estimate of the unknown Ω at time *t*-1 [17]. Thus, the TLF with the unknown Ω is equal to the TLF in (6) and the tracker probability with the unknown Ω is equal to the tracker probability in (7) as follows:

$$p(\mathbf{z}_{t}|m_{t}^{i}, \mathbf{\Omega}, \mathbf{z}_{1:t-1}) \approx p(\mathbf{z}_{t}|m_{t}^{i}, \mathbf{z}_{1:t-1}) = \Lambda_{t}^{i},$$

$$P\{m_{t-1}^{i}|\mathbf{\Omega}, \mathbf{z}_{1:t-1}\} \approx P\{m_{t-1}^{i}|\mathbf{z}_{1:t-1}\} = T_{t-1}^{i}.$$
(14)

With these approximations for $p(\mathbf{z}_t|\mathbf{\Omega}, \mathbf{z}_{1:t-1})$ in (16), the total probability $p_{\mathbf{\Omega}}(\mathbf{z}_t|\mathbf{z}_{1:t-1})$ is also approximated as described in (17). Based on (16) and (17), the sequential update of the TPM posterior distribution in (13) is expressed by

$$p(\mathbf{\Omega}|\mathbf{z}_{1:t}) \approx \frac{\mathbf{T}_{t-1}^{\top} \mathbf{\Omega} \mathbf{\Lambda}_t}{\mathbf{T}_{t-1}^{\top} \bar{\mathbf{\Omega}}_{t-1} \mathbf{\Lambda}_t} p(\mathbf{\Omega}|\mathbf{z}_{1:t-1}),$$
(15)

where $\mathbf{\Lambda}_t = [\mathbf{\Lambda}_t^1, \dots, \mathbf{\Lambda}_t^M]^\top$ and $\mathbf{T}_{t-1} = [T_{t-1}^1, \dots, T_{t-1}^M]^\top$.

$$p(\mathbf{z}_{t}|\mathbf{\Omega}, \mathbf{z}_{1:t-1}) = \sum_{i=1}^{M} p(\mathbf{z}_{t}|m_{t}^{i}, \mathbf{\Omega}, \mathbf{z}_{1:t-1}) P\{m_{t}^{i}|\mathbf{\Omega}, \mathbf{z}_{1:t-1}\}$$

$$= \sum_{i=1}^{M} \underbrace{p(\mathbf{z}_{t}|m_{t}^{i}, \mathbf{\Omega}, \mathbf{z}_{1:t-1})}_{\approx \Lambda_{t}^{i}}$$

$$\sum_{j=1}^{M} \underbrace{P\{m_{t}^{i}|m_{t-1}^{j}, \mathbf{\Omega}, \mathbf{z}_{1:t-1}\}}_{\triangleq \omega^{j,i}} P\{m_{t-1}^{j}|\mathbf{\Omega}, \mathbf{z}_{1:t-1}\}}_{\approx T_{t-1}^{j}}$$

$$\approx \sum_{i=1}^{M} \Lambda_{t}^{i} \sum_{j=1}^{M} \omega^{j,i} T_{t-1}^{j} = \Lambda_{t}^{\top} \mathbf{\Omega}^{\top} \mathbf{T}_{t-1} = \mathbf{T}_{t-1}^{\top} \mathbf{\Omega} \Lambda_{t},$$
(16)

where

$$\mathbf{\Lambda}_t = [\mathbf{\Lambda}_t^1, \dots, \mathbf{\Lambda}_t^M]^\top, \quad \mathbf{T}_{t-1} = [T_{t-1}^1, \dots, T_{t-1}^M]^\top.$$

$$p_{\mathbf{\Omega}}(\mathbf{z}_{t}|\mathbf{z}_{1:t-1}) = \int \underbrace{p(\mathbf{z}_{t}|\mathbf{\Omega}, \mathbf{z}_{1:t-1})}_{\approx \mathbf{T}_{t-1}^{\top} \mathbf{\Omega} \mathbf{\Lambda}_{t} \operatorname{in}(16)} p(\mathbf{\Omega}|\mathbf{z}_{1:t-1}) d\mathbf{\Omega}$$

$$\approx \mathbf{T}_{t-1}^{\top} \left(\int \mathbf{\Omega} p(\mathbf{\Omega}|\mathbf{z}_{1:t-1}) d\mathbf{\Omega} \right) \mathbf{\Lambda}_{t} = T_{t-1}^{\top} \bar{\mathbf{\Omega}}_{t-1} \mathbf{\Lambda}_{t}.$$
(17)

5.2.1 Sample Approximation

For practical implementations, the TPM posterior distribution in (15) is approximated by first or second order numerical integration methods as they are shown to be more robust and accurate than other approaches [17]. In numerical integration, since the prior information of probabilistic interactions between trackers is not usually given, the interaction probabilities are defined on a finite grid. Thus, the TPM prior distribution is approximated by a set of samples $\{ \mathbf{\Omega}^{q} | q = 1, \ldots, N_{\mathbf{\Omega}} \}$ with corresponding weights $\{ p(\mathbf{\Omega}^{q} | \mathbf{z}_{1:t-1}) | q = 1, \ldots, N_{\mathbf{\Omega}} \}$, and the TPM posterior distribution in (15) is described by

$$p(\mathbf{\Omega}^{q}|\mathbf{z}_{1:t}) = \frac{\mathbf{T}_{t-1}^{\top}\mathbf{\Omega}^{q}\mathbf{\Lambda}_{t}}{\mathbf{T}_{t-1}^{\top}\bar{\mathbf{\Omega}}_{t-1}\mathbf{\Lambda}_{t}}p(\mathbf{\Omega}^{q}|\mathbf{z}_{1:t-1}).$$
 (18)

We obtain the updated TPM $\bar{\mathbf{\Omega}}_t$ at time *t* as

$$\bar{\mathbf{\Omega}}_{t} = \frac{1}{C} \sum_{q=1}^{N_{\mathbf{\Omega}}} \mathbf{\Omega}^{q} p(\mathbf{\Omega}^{q} | \mathbf{z}_{1:t}),$$
(19)

where $C = \sum_{q=1}^{N_{\Omega}} p(\Omega^{q} | \mathbf{z}_{1:t})$ is a normalization term and each TPM sample is expressed by $\Omega^{q} = \left[\omega_{q}^{j,i}\right]_{M \times M}$. The interaction probabilities are chosen as $0 \le \omega_{q}^{j,i} \le 1$ and satisfy the condition $\sum_{i=1}^{M} \omega_{q}^{j,i} = 1$.

6 ROBUST TRACKER LIKELIHOOD FUNCTION

The reliability of each tracker is used to update the tracker probability and TPM within the Bayesian framework. The tracker likelihood function computes the reliability of each one by measuring the tracking results individually. The estimated object state from the *i*th tracker at time t is

$$\hat{\mathbf{x}}_t^i = \arg \max_{\mathbf{x}_t} \ p(\mathbf{x}_t | \mathbf{z}_{1:t}, m_t^i). \tag{20}$$

Since $\hat{\mathbf{x}}_t^i$ is obtained from the *i*th tracker, the accuracy of $\hat{\mathbf{x}}_t^i$ is considered as the reliability of the *i*th tracker. Hence, the TLF is expressed by

$$p(\mathbf{z}_t | m_t^i, \mathbf{z}_{1:t-1}) = p_{\text{TLF}}(\mathbf{z}_t | \hat{\mathbf{x}}_t^i).$$
(21)

For measuring the tracker reliability, we use instantaneous and reconstruction features to account for transient and stable appearance changes. These two representations are assumed to be independent and all M features ($\mathbf{f}^k, k = 1, \dots, M$) are used for computing the TLF to measure the reliability of each tracker. Thus, the TLF is formulated by

$$p_{\text{TLF}}(\mathbf{z}_t | \hat{\mathbf{x}}_t^i) \approx p_{\text{I}}(\mathbf{z}_t | \hat{\mathbf{x}}_t^i) p_{\text{R}}(\mathbf{z}_t | \hat{x}_t^i)$$
$$= \prod_{k=1}^M p(\mathbf{z}_t | \hat{\mathbf{x}}_t^i, \mathbf{f}_{\text{I},t}^k) p(\mathbf{z}_t | \hat{\mathbf{x}}_t^i, \mathbf{f}_{\text{R},t}^k), \qquad (22)$$

where *k* is the feature index, $p_{I}(\mathbf{z}_{t}|\hat{\mathbf{x}}_{t}^{i})$ is the TLF based on the instantaneous appearance model (IAM), and $p_{R}(\mathbf{z}_{t}|\hat{\mathbf{x}}_{t}^{i})$ is the TLF based on the reconstruction appearance model (RAM). The instantaneous object appearance $\bar{\mathbf{f}}_{I,t}^{k}$ is obtained from a set of recent observations $\mathbf{f}_{I,t}^{k}$. The reconstructed object appearance $\bar{\mathbf{f}}_{R,t}^{i,k}$ is computed from the stable appearance $\mathbf{f}_{R,t}^{k}$ using the *k*th feature and the tracking result $\mathbf{z}_{t}^{i,k}$ from the *i*th tracker. Each TLF is computed by

$$p(\mathbf{z}_t | \hat{\mathbf{x}}_t^i, \mathbf{f}_{\mathbf{I},t}^k) = \exp(-\rho \| \bar{\mathbf{f}}_{\mathbf{I},t}^k - \mathbf{z}_t^{i,k} \|^2),$$
(23)

$$p(\mathbf{z}_t | \hat{\mathbf{x}}_t^i, \mathbf{f}_{\mathrm{R},t}^k) = \exp(-\rho \| \overline{\mathbf{f}}_{\mathrm{R},t}^{i,k} - \mathbf{z}_t^{i,k} \|^2),$$
(24)

where ρ is a control parameter and

$$\mathbf{z}_t^{i,k} = \frac{\operatorname{Vec}(F^k(I(\hat{\mathbf{x}}_t^i)))}{\|\operatorname{Vec}(F^k(I(\hat{\mathbf{x}}_t^i)))\|},$$
(25)

where $Vec(\cdot)$ represents vectorization, $I(\mathbf{x}_t)$ denotes an image region based on a state vector \mathbf{x}_t , $F^k(\cdot)$ denotes the *k*th feature extraction, and $\mathbf{z}_t^{i,k} \in \mathfrak{R}^{d^k}$ where d^k is the dimension of the *k*th feature. The IAM and RAM are computed as follows.

6.1 Transient Object Appearance

The short-term object appearance changes are model by a set of recent object observations $\mathbf{f}_{\mathbf{I},t}^k = [\mathbf{f}_{\mathbf{I},t-l}^k, \dots, \mathbf{f}_{\mathbf{I},t-1}^k]$. The instantaneous appearance model $\mathbf{\bar{f}}_{\mathbf{I},t}^k$ is obtained by averaging the recent *L* appearances as

$$\bar{\mathbf{f}}_{\mathrm{I},t}^{k} = \frac{1}{L} \sum_{l=1}^{L} \mathbf{f}_{\mathrm{I},t-l}^{k}.$$
 (26)

6.2 Stable Object Appearance

The long-term object appearance $\mathbf{z}_t^{i,k}$ can be represented by a linear combination of stable features $\mathbf{f}_{\mathrm{R},t}^k$ that are r representative features,

$$\mathbf{z}_{t}^{i,k} \approx \mathbf{f}_{\mathrm{R},t}^{k} \boldsymbol{\alpha}_{t}^{i,k} = \mathbf{f}_{1,t}^{k} \boldsymbol{\alpha}_{1,t}^{i,k} + \mathbf{f}_{2,t}^{k} \boldsymbol{\alpha}_{2,t}^{i,k} + \ldots + \mathbf{f}_{r,t}^{k} \boldsymbol{\alpha}_{r,t}^{i,k}, \quad (27)$$

where $\mathbf{f}_{\mathrm{R},t}^{k} = [\mathbf{f}_{1,t}^{k}, \dots, \mathbf{f}_{r,t}^{k}] \in \mathfrak{R}^{d^{k} \times r}$, $\boldsymbol{\alpha}_{t}^{i,k} = [\boldsymbol{\alpha}_{1,t}^{i,k}, \dots, \boldsymbol{\alpha}_{r,t}^{i,k}]^{\top} \in \mathfrak{R}^{r}$ is an coefficient vector. By including the noise vector $\boldsymbol{\epsilon}^{i,k}$, we have

$$\mathbf{z}_{t}^{i,k} = \mathbf{f}_{\mathrm{R},t}^{k} \boldsymbol{\alpha}_{t}^{i,k} + \boldsymbol{\epsilon}^{i,k} = \begin{bmatrix} \mathbf{f}_{\mathrm{R},t}^{k} & \mathbf{I}^{k} \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}_{t}^{i,k} \\ \boldsymbol{\beta}_{t}^{i,k} \end{bmatrix}.$$
 (28)

We use a set of non-target (trivial) templates from a d^k -dimensional identity matrix $\mathbf{I}^k \in \mathfrak{R}^{d^k \times d^k}$ [24] with a non-target coefficient vector $\boldsymbol{\beta}_t^{i,k} = [\beta_{1,t}^{i,k}, \beta_{2,t}^{i,k}, \dots, \beta_{d^k,t}^{i,k}]^\top \in \mathfrak{R}^{d^k}$. If the observation contains little noise, then the non-target coefficient vector has only a few nonzero coefficients in $\boldsymbol{\beta}_t^{i,k}$.

In the proposed tracking algorithm, we obtain M tracking results at each frame, $\{\hat{\mathbf{x}}_t^i | i = 1, ..., M\}$. Based on the result of the *i*th tracker, the candidate image region

represented by the *k*th feature is denoted as $\mathbf{z}_{t}^{i,k}$ in (25). The reconstructed appearance for $\mathbf{z}_{t}^{i,k}$ is denoted as $\mathbf{f}_{\mathrm{R},t}^{k} \boldsymbol{\alpha}_{t}^{i,k}$. We obtain the coefficient vector $\boldsymbol{\alpha}_{t}^{i,k}$ by using ℓ_{1} sparse coding as it is robust to wide range of image corruptions, especially occlusions, [19], [24]. The coefficient vector $\mathbf{c}_{t}^{i,k}$ is computed by

$$\min_{\mathbf{c}_t^{i,k}} \|\mathbf{c}_t^{i,k}\|_1, \quad \text{s.t.} \quad \|\mathbf{z}_t^{i,k} - \mathbf{D}_t^k \mathbf{c}_t^{i,k}\|_2^2 \le \lambda,$$
(29)

where $\lambda = 0.01$, and

$$\mathbf{D}_{t}^{k} = [\mathbf{f}_{\mathrm{R},t}^{k}, \mathbf{I}^{k}], \quad \mathbf{c}_{t}^{i,k} = [(\boldsymbol{\alpha}_{t}^{i,k})^{\top}, (\boldsymbol{\beta}_{t}^{i,k})^{\top}]^{\top}.$$
 (30)

The reconstructed object appearance $\mathbf{\bar{f}}_{\mathrm{R},t}^{i,k}$ for $\mathbf{z}_{t}^{i,k}$ is computed as $\mathbf{\bar{f}}_{\mathrm{R},t}^{i,k} = \mathbf{f}_{\mathrm{R},t}^{i,k} \alpha_{t}^{i,k}$.

7 REPRESENTATION UPDATE

In this section, we present the update mechanisms for transient and stable object appearance as well as observation models for trackers based on M feature representations $\{\hat{\mathbf{f}}_t^k = \mathbf{z}_t^{\hat{m}_t,k} | k = 1, ..., M\}$ where $\mathbf{z}_t^{i,k}$ is from (25) and \hat{m}_t is the index of the selected tracker in (5).

7.1 Transient Features

We use transient features to account for abrupt appearance changes of a target object. Each transient feature consists of the recently estimated observation as $\mathbf{f}_{\mathrm{I},t+1}^{k} = \left[\mathbf{f}_{\mathrm{I},t-\theta}^{k}, \ldots, \mathbf{f}_{\mathrm{I},t}^{k}\right]$, where $\mathbf{f}_{\mathrm{I},t}^{k} = \hat{\mathbf{f}}_{t}^{k}$ and θ is a variable that determines the duration.

7.2 Stable Features

Each stable feature $\mathbf{f}_{\mathrm{R},t}^k$ is updated based on whether it can be sparsely represented by the current templates. Similar to [25], each feature is updated by analyzing the non-zero elements in the non-target coefficient vector $\boldsymbol{\beta}_t^{i,k}$. When occlusion occurs, a target object cannot be sparsely represented by the target template set. Consequently, numerous nonzero coefficients correspond to the non-target templates, and noise is measured by $\boldsymbol{\beta}_t^{\hat{m}_t,k} \in \Re^{d^k}$ in (29) where \hat{m}_t is the index of the selected tracker. We count non-zero elements in $\boldsymbol{\beta}_t^{\hat{m}_t,k}$ and compute a noise ratio R_{noise}^k as $R_{\text{noise}}^k = B^k/d^k$ where B^k is the number of non-zero elements in $\boldsymbol{\beta}_t^{\hat{m}_t,k}$. If the noise ratio R_{noise}^k is smaller than a threshold, one feature $\mathbf{f}_{i,t}^k \in \mathbf{f}_{\mathrm{R},t}^k$ with the lowest value is replaced by the feature of the estimated observation $\hat{\mathbf{f}}_t^k$.

7.3 Observation Model

In this work, the observation model for each tracker (i.e., $p(\mathbf{z}_t | \mathbf{x}_t, m_t^i)$ in (1)) is based on the incremental subspace model [27] for its computational efficiency over ℓ_1 sparse coding. For online tracking, it is known that error accumulation is inevitable when an appearance model is updated with new observations [12], [28]. Note that not every observation model is updated at every frame. For the selected tracker of a given frame, the corresponding appearance model is not updated since it describes the target object



Fig. 5. Representation update examples. The transient and stable features are shown in the red and blue boxes, respectively. The learned principal components are shown in the green boxes. The yellow circles demonstrate the updated stable features at different frames.

well. On the other hand, the observation models of all the other trackers are updated with the new observation.

Examples of representation updates (i.e., transient and stable features as well as observation model discussed in Section 7) are shown in Fig. 5. To show difference of each representation, we only show the intensity features for comparisons. In the *Coke* sequence, partial occlusions with illumination changes occur frequently. As introduced in Section 7, the transient features better account for frequent appearance changes of the object in such cases while the stable features are rarely updated. The principal components of the object appearance from an observation model are shown in green boxes. These principal components are incrementally updated in each observation model to account for appearance changes.

8 INTERACTING MULTIVIEW TRACKER

The main components of the proposed interacting multiview tracker (IMT) are described in Fig. 3 and Algorithm 1. We present the algorithmic details in this section.

8.1 Estimated Object States of Multiple Trackers

We use a particle filter for state prediction. The prior distribution of each tracker $p(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t-1}^i)$ in (2) is approximated by a set of N samples as

$$p(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t-1}^{i}) \approx \sum_{q=1}^{N} s_{q,t-1}^{i} \delta(\mathbf{x}_{q,t-1}^{i} - \mathbf{x}_{t-1}), \quad (31)$$

where $\delta(\cdot)$ is a delta function centered at sample $\mathbf{x}_{q,t-1}^{i}$, and $s_{q,t-1}^{i}$ is a sample weight.

8.1.1 Interacted Prior via Tracker Interaction

At each frame, multiple trackers interact with each other by mixing their posterior distributions described in (2) based on the TPM. The interaction is efficiently carried out via the proposed interaction method by Algorithm 2, i.e.,

$$\begin{bmatrix} \tilde{\mathcal{X}}_{t-1}^1, \dots, \tilde{\mathcal{X}}_{t-1}^M \end{bmatrix} = \text{Tracker_Interaction} \begin{bmatrix} \bar{\mathbf{\Omega}}_{t-1}, \mathcal{X}_{t-1}^1, \dots, \mathcal{X}_{t-1}^M \end{bmatrix},$$
(32)

where $\mathcal{X}_{t-1}^i = {\{\mathbf{x}_{q,t-1}^i, s_{q,t-1}^i\}_{q=1}^N}$ is the sample approximation of the prior distribution of the *i*th tracker and $\tilde{\mathcal{X}}_{t-1}^i$ is the interacted prior distribution. The tracker interaction approach in this work is similar in spirit to [1], [4] where the posterior

distribution of the unreliable tracker is replaced by the most reliable one. In addition, the reliability of the tracker is measured by exploring the covariance of the posterior distribution at each frame. However, the proposed interaction method forces trackers to interact with each other via the TPM. Hence, not all samples are transfered to other trackers.

Algorithm 1. Proposed Interacting Multiview Tracker (IMT)

- 1: (Initial Step)
- 2: at time t = 0
- 3: The initial states of multiple trackers are set to $\{\mathbf{x}_0^i = \mathbf{x}_0 | i = 1, \dots, M\}.$
- 4: The initial set of samples for the particle filter $\{\mathcal{X}_0^i = \{\mathbf{x}_{q,0}^i, s_{q,0}^i = \frac{1}{N}\}_{q=1}^N | i = 1, \dots, M\}.$
- 5: The initial TPM is given by $\bar{\Omega}_0 = \frac{1}{N_0} \sum_q \Omega^q$ >Section 9.2.
- 6: The initial tracker probability is set to

$$\{T_0^i = \frac{1}{M} | i = 1, \dots, M\}.$$

- 7: (Tracking Step)
- 8: for $t \ge 1$ do
- 9: **for** i = 1 : M **do** \triangleright i is a tracker index
- 10: 1) Compute the interacted prior distribution $\tilde{\mathcal{X}}_{t-1}^{i} = \{\tilde{\mathbf{x}}_{q,t-1}^{i}, \tilde{s}_{q,t-1}^{i}\}_{q=1}^{N}$ using $\{\mathcal{X}_{t-1}^{i}|i=1,\ldots,M\}$ with the TPM $\bar{\mathbf{\Omega}}_{t-1}$ and the tracker probability $\{T_{t-1}^{i}|i=1,\ldots,M\}$ using Algorithm 2.
- 11: 2) Predict state samples $\{\mathbf{x}_{q,t}^i, s_{q,t|t-1}^i\}_{q=1}^N$ using (34).
- 12: 3) Update state samples $\{\mathbf{x}_{a,t}^i, s_{a,t}^i\}_{a=1}^N$ using (37).
- 13: 4) Obtain the estimated state $\hat{\mathbf{x}}_{t}^{i}$ from the *i*th tracker using (39).
- 14: end for
- 15: 5) Compute the TLFs {Λⁱ_t|i = 1,..., M} using (22) and the set of *M* estimated object states { xⁱ_t|i = 1,..., M}.
- 16: 6) The tracker probability update with the TLFs
 - $\{\Lambda^i_t|i=1,\ldots,M\}$ using (10).
- 17: 7) The TPM update with the tracker probabilities $\{T_t^i | i = 1, ..., M\}$ and TLFs $\{\Lambda_t^i | i = 1, ..., M\}$ using (18) and (19).
- 18: 8) Compute the tracking result $\hat{\mathbf{x}}_t$ using (5).
- 19: 9) Update representations as described in Section 7.

The interacted prior distribution in (2) can be expressed by a sample representation as

$$\tilde{p}(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t-1}^{i}) \approx \sum_{q=1}^{N} \tilde{s}_{q,t-1}^{i} \delta(\tilde{\mathbf{x}}_{q,t-1}^{i} - \mathbf{x}_{t-1}).$$
(33)

By tracker interaction, we first remove the samples far from the selected tracking result $\hat{\mathbf{x}}_{t-1}$ based on a kernel. As described in Algorithm 2, a uniform kernel is defined in terms of position with respect to range R with standard deviations (q_u, q_v) along u and v image coordinates. In addition, H is a transformation matrix that returns position parameters as from a previous state by $[p_{u,t-1}, p_{v,t-1}]^{\top} =$ $H\hat{x}_{t-1}$. Second, multiple trackers interact with each other based on the TPM and a resampling technique [9]. The TPM contains information of how samples are transferred or retained. For instance, $N \times \bar{\omega}_{t-1}^{i,i}$ represents that the number of samples is retained in the *i*th tracker sample set after interaction, and $N \times \bar{\omega}_{t-1}^{j,i}$ represents that the number of samples from the *j*th tracker is transferred to the *i*th tracker. If the *i*-tracker is effective for some frames, then $\bar{\omega}_{t-1}^{i,i}$ becomes greater than $\bar{\omega}_{t-1}^{j,i}$ $(j \neq i)$ due to an update of the TPM. Hence, most samples of the *i*th are retained, and the *i*th tracker obtains a few samples from other trackers. Finally, we select samples according to the interaction probabilities, $\bar{\omega}_{t-1}^{j,i}$ of the TPM by resampling such that reliable samples with large weights in each tracker are retained.

Algorithm 2. Tracker Interaction: $[\tilde{X}_{t-1}^1, \dots, \tilde{X}_{t-1}^M] =$ **Tracker_Interaction** $[\bar{\mathbf{\Omega}}_{t-1}, \mathcal{X}_{t-1}^1, \dots, \mathcal{X}_{t-1}^M]$

1: Input

- 2: Given $\{X_{t-1}^i = \{\mathbf{x}_{q,t-1}^i, s_{q,t-1}^i\}_{q=1}^N | i = 1, \dots, M\}$ 3: \triangleright Sample representation of a posterior distribution of *i* the tracker

```
4:
```

```
5: for i = 1 : M do
```

```
for q = 1 : N do
6:
```

 $s_{q,t-1}^{*i} = s_{q,t-1}^{i} \operatorname{Kernel}(\mathbf{H}\mathbf{x}_{q,t-1}^{i} - \mathbf{H}\mathbf{x}_{t-1}, \mathbf{R})$ 7:

 $s_{q,t-1}^{*i} := s_{q,t-1}^{*i} / \sum_{q} s_{q,t-1}^{*i}, \ q = 1, \dots, N$ 9:

- 10: end for
- 11:

12: Given $\bar{\omega}_{t-1}^{j,i} \in \bar{\mathbf{\Omega}}_{t-1}$

- 13: for i = 1 : M do
- $\tilde{\mathcal{X}}_{t-1}^i = \phi$ 14:
- 15: for j = 1 : M do
- $$\begin{split} \boldsymbol{\tilde{\mathcal{X}}} &= \operatorname{Resampling}(\{\boldsymbol{\mathbf{x}}_{q,t-1}^{j}, \boldsymbol{s}_{q,t-1}^{*j}\}_{q=1}^{N}, \quad N \times \bar{\boldsymbol{\omega}}_{t-1}^{j,i}) \\ \boldsymbol{\tilde{\mathcal{X}}}_{t-1}^{i} &:= \boldsymbol{\tilde{\mathcal{X}}}_{t-1}^{i} \cup \boldsymbol{\mathcal{X}} \end{split}$$
 16:

⊳TPM

- 17:
- 18: end for
- 19: end for

20: 21: Output

22: $\{\tilde{\mathcal{X}}_{t-1}^i = \{\tilde{\mathbf{x}}_{q,t-1}^i, \tilde{s}_{q,t-1}^i = \frac{1}{N}\}_{q=1}^N | i = 1, \dots, M\}$

23: \triangleright Sample representation of an interacted prior of *i*th tracker

24:

25: Given parameters

26: $\mathbf{H} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}$ ⊳position conversion matrix 27: $\mathbf{R} = \sqrt{(2 \times q_v)^2 + (2 \times q_u)^2}$ ⊳kernel range

8.1.2 Sampling via Motion Models

We draw new state samples from the interacted prior distribution $\tilde{p}(\mathbf{x}_{t-1}|\mathbf{z}_{1:t-1}, m_{t-1}^i)$. In this work, we use the zero and first order motion models for state prediction $p(\mathbf{x}_t | \mathbf{x}_{t-1}, m_t^i)$. The zero-order motion is identical to the random walk motion, and the first-order motion utilizes the prior translation $\Delta \mathbf{x}_t = [\Delta u, \Delta v, 0, 0, 0, 0]^{\top}$ by computing the difference between estimated positions at time t - 1 and t - 2. Thus, samples are drawn based on

$$\mathbf{x}_{q,t}^{i} \sim p(\mathbf{x}_{t} | \mathbf{x}_{t-1}, m_{t}^{i}) = \begin{cases} \mathcal{N}(\mathbf{x}_{q,t-1}^{i}, \mathbf{Q}_{0}) & \text{if } \tau < 0.5 \\ \mathcal{N}(\mathbf{x}_{q,t-1}^{i} + \Delta \mathbf{x}_{t}, \mathbf{Q}_{1}) & \text{otherwise,} \end{cases}$$
(34)

where \mathbf{Q}_0 and \mathbf{Q}_1 denote the zero-and first-order motion covariances, respectively, as given in Section 9.1. We use a uniform random variable τ distributed within [0, 1] to select the motion model for drawing each sample. The set of the predicted samples is $\{\mathbf{x}_{q,t}^i, s_{q,t|t-1}^i\}_{q=1}^N$ where $s_{q,t|t-1}^i = \tilde{s}_{q,t-1}^i$.

8.1.3 Sample Update via Observation Models

An observation for the *i*th tracker is expressed by

$$\mathbf{z}_t^i = Vec(F^i(I(\mathbf{x}_t))) + \mathbf{v}_t^i, \quad i = 1, \dots, M,$$
(35)

where $I(\mathbf{x}_t)$ denotes an image template based on a state vector \mathbf{x}_t , $F^i(\cdot)$ represents the *i*th feature extraction, and \mathbf{v}_t^i is noise. In the incremental subspace based observation model [27], we compute the mean and principal eigenvectors with updates for the appearance model in each tracker. Based on the template mean \overline{O}^i and L principal eigenvectors \mathbf{g}^i_{l} $l = 1, \ldots, L$, the *i*th observation model based on the *i*th feature is given by

$$p(\mathbf{z}_{t}|\mathbf{x}_{t}, m_{t}^{i}) = \exp(-\rho_{T} \|\mathbf{z}_{t}^{i} - \sum_{l} \mathbf{c}_{l} \mathbf{g}_{l}^{i}\|^{2}),$$

$$\mathbf{c}_{l} = (\mathbf{g}_{l}^{i})^{\top} (\mathbf{z}_{t}^{i} - \bar{O}^{i}), \quad l = 1, \dots, L,$$
(36)

where ρ_T is a control parameter and \mathbf{c}_l is the coefficient from the projection of the template onto each principal eigenvector (16 eigenvectors are used for each observation model).

We note that the TLF in (22) is not related to the observation model in (36). The TLF is only used to update the tracker probability and TPM because it is time-consuming to measure all particle samples if we use TLF instead of (36). For efficient implementation, we use (36) as an observation model to measure particle samples of a single tracker as it can be computed efficiently to adapt object appearance changes. Based on (35) and (36), the weight of each sample is updated by

$$s_{q,t}^{i} = \frac{p(\mathbf{z}_{t}|\mathbf{x}_{q,t}^{i}, m_{t}^{i})s_{q,t|t-1}^{i}}{\sum_{q=1}^{N} p(\mathbf{z}_{t}|\mathbf{x}_{q,t}^{i}, m_{t}^{i})s_{q,t|t-1}^{i}}.$$
(37)

With the samples and weights in (34) and (37), we obtain the sample representation of the posterior distribution $p(\mathbf{x}_t)$ $\mathbf{z}_{1:t}, m_t^i$ in (1) as

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}, m_t^i) \approx \sum_{q=1}^N s_{q,t}^i \delta(\mathbf{x}_{q,t}^i - \mathbf{x}_t),$$
(38)

which is described by a set of samples with weights $\{\mathbf{x}_{q,t}^{i}, s_{q,t}^{i}\}_{q=1}^{N}$.

8.1.4 Estimated Object States

From the updated posterior distributions, we obtain a set of M estimated states using the maximum a posterior estimates (i = 1, ..., M),

$$\hat{\mathbf{x}}_{t}^{i} = \mathbf{x}_{\hat{q},t}^{i}, \quad \hat{q} = \arg\max_{q}(\{s_{q,t}^{i}|q=1,\ldots,N\}).$$
 (39)

8.2 Tracker Selection and TPM Update

To select the most reliable tracker and update the TPM, we compute the reliability of trackers using the TLF $p_{\text{TLF}}(\mathbf{z}_t | \hat{\mathbf{x}}_t^i) = \Lambda_t^i$ in (22) and M estimated states, as well as $\{\hat{\mathbf{x}}_t^i | i = 1, ..., M\}$ from M multiple trackers. With the TLFs $\{\Lambda_t^i | i = 1, ..., M\}$, we update the tracker probability using (10) and obtain updated tracker probabilities $\{T_t^i | i = 1, ..., M\}$. By selecting the highest tracker probability, we obtain the tracking result $\hat{\mathbf{x}}_t$ as described in (5). The tracking result $\hat{\mathbf{x}}_t$ is then used for a representation update (See Section 7). After computing the set of the updated tracker probabilities $\{T_t^i | i = 1, ..., M\}$ and the set of TLFs $\{\Lambda_t^i | i = 1, ..., M\}$, we update the TPM $\overline{\mathbf{\Omega}}_t$ using (18) and (19).

9 EXPERIMENTS

We evaluate the proposed IMT algorithm with the stateof-the-art methods using several benchmark datasets (http://vision.ucsd.edu/~bbabenko/project_miltrack. [2] shtml) and [31] (http://visual-tracking.net), as well as our own sequences (i.e., Startrek and Starwars). In this work, we use three trackers with different feature representations based on HOG, intensity, and Haar-like features, which have been shown to be effective for handling occlusions, motion blurs, pose variations, and illumination changes. We discuss motion parameter settings in Section 9.1, the sampling scheme for the TPM in Section 9.2, and feature extraction in Section 9.3. We analyze the TPM and show how it is used by multiple trackers in Section 9.4. In Section 9.5, we demonstrate the effects of the proposed TLF, and in Section 9.6, we compare the proposed IMT algorithm with other tracking methods based on one single feature representation of HOG, intensity, and Haar-like features (denoted as SHOG, SI, and SHaar methods). We evaluate the proposed algorithm in Section 9.7 with methods based on multiple trackers or representations including the approaches with combination of visual trackers (CVT) [22], the multi-cue switching tracker (MCS) [4], and a single tracker with multiple observation models (SMO) similar to [29] where the tracker reliability is not measured. For fair comparisons, each single tracker of the proposed IMT algorithm and parameters are the same as those used in the CVT, MCS, SMO, SHOG, SI, and SHaar methods. Furthermore, in Section 9.8, we compare the IMT algorithm with state-of-the-art trackers including the MIL [2], TLD [18], VTD [20], VTS [21], Struck [13], ASLA [16], SCM [37], CXT [8], LSK [23], CSK [14], and KCF [15] methods.

FOR quantitative comparisons, we present the success rate rather than center location error as it is not fully reflected, especially after tracking drifting [31]. The code and datasets are available at https://cvl.gist.ac.kr/project/imt.html.

9.1 Motion Parameters

In this work, an object state is expressed by six parameters of the affine transformation [27] based on a diagonal covariance matrix $\mathbf{Q} = \text{diag}(q_u^2, q_v^2, q_\theta^2, q_s^2, q_\alpha^2, q_\phi^2)$ with the following variables q_u and q_v are standard deviations of position and q_θ , q_s , q_α , and q_ϕ are standard deviations of rotation angle, scale, aspect ratio, and skew, respectively. For all the experiments, we fix four parameters as $q_\theta = 0.02$, $q_s = 0.01$, $q_\alpha = 0$, $q_\phi = 0.001$. The translation standard deviation of the zero-order motion \mathbf{Q}_0 are fixed as $q_u = q_v = 6$. The translation standard deviation \mathbf{Q}_1 are fixed as $q_u = q_v = 3$. Since the SI, SHOG, SHaar, SMO, MCS, CVT,

Tracker 1	Tracker 2
$\boldsymbol{\omega}_{1}^{1} = [0.7, \ 0.15, \ 0.15]^{\top}$	$\boldsymbol{\omega}_1^2 = \begin{bmatrix} 0.15, \ 0.7, \ 0.15 \end{bmatrix}^{\top}$
$\boldsymbol{\omega}_2^1 = [0.6, \ 0.20, \ 0.20]^{\top}$	$\boldsymbol{\omega}_2^2 = [0.20, \ 0.6, \ 0.20]^{\top}$
$\boldsymbol{\omega}_3^1 = [0.5, \ 0.25, \ 0.25]^{\top}$	$\boldsymbol{\omega}_3^2 = [0.25, \ 0.5, \ 0.25]^{\top}$
$\boldsymbol{\omega}_{4}^{1} = [0.4, \ 0.30, \ 0.30]^{\top}$	$\boldsymbol{\omega}_4^2 = [0.30, \ 0.4, \ 0.30]^{\top}$
$\boldsymbol{\omega}_{5}^{1} = [0.3, \ 0.35, \ 0.35]^{\top}$	$\boldsymbol{\omega}_5^2 = [0.35, \ 0.3, \ 0.35]^{\top}$
$\boldsymbol{\omega}_{6}^{1} = \left[0.2, \ 0.40, \ 0.40 \right]^{\top}$	$\boldsymbol{\omega}_{6}^{2} = \begin{bmatrix} 0.40, \ 0.2, \ 0.40 \end{bmatrix}^{\top}$
Tracker 3	3
$\omega_1^3 = [0.15, 0.15]$	$[5, 0.7]^ op$
$\omega_2^3 = [0.20, 0.20]$	$[0, 0.6]^{ op}$
$\omega_3^3 = [0.25, 0.25]$	$[5, 0.5]^ op$
$\omega_4^3 = [0.30, 0.30]$), 0 .4] [⊤]
$\omega_5^3 = [0.35, 0.35]$	$[5, 0.3]^ op$
$m{\omega}_6^3 = [0.40, 0.40]$	$[0, 0.2]^{ op}$

and IMT methods are based on the same single tracker [27], we use the same parameter settings as mentioned above. We note that the results in [35] are based on optimized parameters for each sequence, whereas in this work the parameters are fixed for all experiments.

9.2 TPM Setting for Three Trackers

As discussed in Section 5.2, we approximate the TPM posterior distribution by a set of TPM samples $\{\hat{\mathbf{\Omega}}^q|q =$ $1, \ldots, N_{\Omega}$. To construct the TPM, we use the interaction probability basis defined on a finite grid in Table 1 where each vector represents the interaction probabilities describing how samples are retained and transferred. For instance, if we use 600 state samples for each tracker, the interaction probability basis ω_1^1 represents that $\omega_1^1 \times 600 = [420, 90,$ 90]^{\top} where the first tracker retains its own 420 samples and receives 90 samples from the second and 90 samples from the third trackers, respectively. In this work, we only set the maximum and minimum values for the diagonal entries of the TPM. The diagonal values are set to $\omega_s^{i,i} \in \{0.2, 0.3, ...\}$ 0.4, 0.5, 0.6, 0.7 to make each tracker retain, at most, 70 percent of its own samples and at least 20 percent of its own samples. The off-diagonal values of the TPM are set with a given diagonal value by $\omega_s^{j,i} = \frac{1-\omega_s^{i,i}}{2}$ (See Table 1). Using the interaction probability basis in Table 1, we obtain a TPM sample as

$$\mathbf{\Omega}^{q} = \left[\boldsymbol{\omega}_{s_{1}}^{1}, \boldsymbol{\omega}_{s_{2}}^{2}, \boldsymbol{\omega}_{s_{3}}^{3} \right], \quad s_{1}, s_{2}, s_{3} = 1, \dots, 6.$$

Consequently, 216 TPM samples { $\mathbf{\Omega}^{q}|q = 1, ..., 216$ } are generated by considering all combinations of the basis in Table 1. These TPM samples are fixed in all experiments. The initial TPM $\overline{\mathbf{\Omega}}_{0}$ is obtained by averaging all of TPM samples. Note that the TPM method is not sensitive to initial values as it is updated at each frame. To demonstrate this, we compare the performance of the IMT method with two different initial TPMs (TPM_{ave} and TPM_{naive}) as shown in Table 2, where TPM_{ave} is obtained by averaging all of TPM

TABLE 2						
Average Tracking Success Rate on 16 Benchmark						
Sequences in Table 3						

	IMT with $\mathrm{TPM}_{\mathrm{ave}}$	IMT with $\mathrm{TPM}_{\mathrm{naive}}$
average success rate	92	90

The IMTs with different initial TPM settings show similar performance.

samples as discussed above, and TPM_{naive} is a matrix whose elements are equally set to $\frac{1}{3}$.

9.3 Feature Extraction

In this work, the size of an image template is 32-by-32 pixels, from which a 1,024-dimensional intensity feature vector is formed. To generate HOG features, we use 36 blocks, each block has four cells within an image template, and the dimensions of HOG feature for each block is 36 (i.e., each HOG feature vector has 1,296 dimensions). The Haar-like features are generated with two horizontal and vertical edge filters within a 32-by-32 template to 1,760-dimensional vectors.

9.4 Analysis of TPM and Tracker Probability

We analyze how TPM is used among multiple trackers to account for different object appearance changes. In Fig. 6, the diagonal interaction probabilities $(\bar{\omega}_t^{i,i})$ of the TPM and tracker probabilities are shown according to object appearance changes over time. When the diagonal entry $\bar{\omega}_t^{i,i}$ decreases, then the off-diagonal entries $\bar{\omega}_t^{j,i}, j \neq i$ increases (as $\sum_{j=1}^M \bar{\omega}_t^{j,i} = 1$). The increase of the off-diagonal entries represent that the *i*th tracker becomes more dependent on other trackers. It also shows that when the *i*th tracker probability $\bar{\omega}_t^{i,i}$ of the TPM tends to increase. The increase of the diagonal entry represents that the *i*th tracker becomes less dependent on other trackers.

In the Startrek sequence (See Fig. 6a), both object and background appearances are drastically changed due to abrupt illumination variations. In such scenarios, the tracker based on intensity features is not reliable; hence, its tracker probability is usually low, and likewise, its interaction probability is consistently low. In the David sequence, the tracker based on HOG features is more robust than others when large pose variations occur, which can be explained by that face contour is more effective for tracking in such scenarios (See Fig. 6b). On the other hand, trackers based on all the other features perform well when moderate appearance changes occur. In the Lemming sequence (See Fig. 6c), when the target object undergoes partial occlusions, the interaction probability for the tracker with Haar-like feature increases and its tracker probability is greater than that of other trackers. When the motion blurs suddenly occur, the interaction probability for the tracker based on HOG features increases and the interaction probabilities of other trackers decrease. Similarly, the tracker probability of the tracker based on HOG features is greater than that of other trackers as the shape of the object is consistent. The



Fig. 6. Changes of interaction probabilities on the diagonal of the TPM and tracker probabilities. Each color line represents one type of trackers. Each color box represents one type of appearance changes. The results are obtained by running the IMT 10 times.

tracker based on Haar-like features adaptively learns the appearance changes. As a result, its interaction and tracker probabilities increase after a few frames.



Fig. 7. The area under curve (AUC) of each success plot [31]. OPE: Running the trackers throughout each sequence with initializations of the ground truth positions. TRE: Running the trackers with initialization from the ground truth position at different frames. SRE: Running the trackers with initialization from the different bounding boxes at the first frame. In all evaluation metrics, the IMT performs well against the other state-of-the-art methods.

9.5 Analysis of TLF

To show the effectiveness of the combination of the instantaneous and reconstruction appearance models in the TLF (See Section 6), we evaluate the tracking results using three combinations. The first one is the IMT-all, which uses both IAM and RAM together, as proposed in this work; the second one is the IMT-IAM which uses only the instantaneous appearance model; and the third one is the IMT-RAM, which utilizes only the reconstruction appearance model. As shown in Table 3, the IMT-all achieves more robust and consistent performance than the other two alternatives.

9.6 Comparison with Single-Feature Trackers

Table 3 shows the results of three trackers based on one single feature (i.e., SI, SHOG, and SHaar). These trackers are the same as the single tracker used in the IMT, and their observation models are described in (36). Overall, the proposed multiview tracking algorithm performs better than these trackers with a single feature. In addition, the trackers based on multiple features (i.e., SMC, MCS, and CVT) perform better than the SI, SHOG, and SHaar methods. These results demonstrate the merits of using multiple features for robust object tracking.

9.7 Comparison with Most Related Trackers

The SMO, CVT [22], and MCS [4] methods are related to the proposed method, but the integration approach of multiple features are different, as discussed in Section 2. As shown in Table 3, the proposed IMT algorithm performs favorably against these tracking algorithms.

The SMO tracker exploits multiple observation models in a particle filter framework. However, it does not perform well, as all observation models contribute equally to the estimation of object states without considering their reliability. Hence, posterior distributions and tracking performance may be affected by one tracker with an unreliable observation model.

The CVT method fuses tracking results from multiple trackers with their reliability weight where each one is determined solely by the covariance information of its posterior distribution. As discussed in Section 2 and shown in Fig. 2, the covariance-based approach may not achieve reliable results as the covariance of each posterior distribution does not accurately represent tracker reliability because each one is constructed from a different feature space (i.e., no calibration of tracking results). In addition, similar to SMO, it does not consider the reliability information in the interaction step, which has the interaction scheme in computing the likelihood.

TABLE 3 Success Rate Using *the Same Default Parameters*

	IMT	IMT	IMT	SI	SHOG	SHaar	SMO	MCS	CVT	Struck	ASLA	SCM	KCF	MIL	TLD	VTD
	-RAM	-IAM	-All					[4]	[22]	[13]	[16]	[37]	[15]	[2]	[18]	[20]
Startrek	91	47	86	1	12	36	44	56	76	78	1	56	74	36	3	89
Starwars	86	83	90	2	75	13	20	19	79	40	85	68	92	45	1	40
David	98	100	99	34	99	34	99	62	100	67	97	95	75	62	96	68
Girl	97	80	98	28	87	85	73	73	81	100	74	99	84	68	46	98
Football	86	79	87	64	76	59	73	57	64	66	65	57	70	73	41	76
CAVIAR	100	99	100	49	44	89	100	100	100	41	97	100	38	38	19	41
Woman	98	93	100	16	9	100	92	97	67	100	100	100	100	16	31	15
Singer1	100	66	100	39	50	98	94	63	70	29	99	100	29	27	99	43
Sylv	68	81	77	45	72	44	45	63	75	92	74	88	81	54	92	80
Trellis	93	99	98	36	68	82	90	62	89	78	85	85	84	24	47	50
Deer	100	100	100	32	98	98	77	33	2	100	2	2	82	12	73	4
Jumping	96	92	95	21	28	7	70	17	10	79	16	12	28	47	84	11
Board	80	89	86	10	77	70	65	50	52	70	71	89	86	51	11	34
Lemming	72	68	85	23	52	17	46	39	38	80	69	30	44	83	4	52
Tiger1	90	87	96	10	50	47	35	42	43	84	83	52	69	62	45	85
Coke	75	59	75	3	44	48	68	58	57	78	69	69	69	32	48	7

The top and second best results are denoted by red and blue.

In contrast, the MCS method selects the most reliable tracker at each frame where the reliability is determined by the acceptance ratio using the covariance of the posterior and prior distributions. If the acceptance ratio is below the threshold (e.g., 0.2 in the experiments), the tracker is considered to be unreliable. In the sampling stage, the MCS method simply replaces the probability distribution of unreliable trackers by that of the most reliable tracker. However, the covariance information is not reliable, as discussed above. This sampling process is likely to cause tracking failure as it does consider all information of unreliable trackers, which can be incorrectly selected due to inaccurate covariance information.

Different from the MCS, CVT, and SMO methods, each tracker of the proposed IMT algorithm generates tracking results independently, and the most reliable one is selected using the TLF, which measures the tracker reliability robustly at each frame as shown in Fig. 6, by considering stability and effectiveness of feature representations (See also Fig. 5). In addition, the reliability information is effectively utilized in the tracker interaction process. Thus, the IMT algorithm performs favorably against these methods based on multiple trackers.

9.8 Comparison with State-the-of-Art Trackers

9.8.1 Benchmark Dataset

We compare the proposed IMT algorithm with 29 state-ofthe-art trackers using a large benchmark dataset [31] that contains 51 sequences. Three evaluation metrics are used to evaluate whether the tracking algorithms are sensitive to different initial settings. For the one-pass evaluation (OPE), we use a ground truth bounding box in the first frame for initialization. For the temporal robustness evaluation (TRE), we initialize each tracker with ground truth locations at different frames. For the spatial robustness evaluation (SRE), we use the perturbed ground truth locations in the first frames for experiments. The top 10 tracking algorithms are shown in Fig. 7 for presentation clarity. Fig. 7 shows that the IMT algorithm performs robustly and favorably against the top nine trackers using all the evaluation metrics (OPE, TRE, and SRE).

9.8.2 Startrek and Starwars

The target objects undergo drastic illumination changes and motion blurs in low resolution and contrast image sequences. As shown in Table 3, Figs. 8a and 8b, most of the trackers do not perform well. On the other hand, the IMT algorithm tracks the objects well in both sequences due to the use of tracker reliability to weigh less on the unreliable tracker (i.e., a tracker with intensity feature) and more on reliable trackers in the tracker integration scheme (via tracker selection and interaction), as shown in Fig. 6a.

9.8.3 David, Girl, and Football

The objects in these sequences undergo large pose variations with occlusions. The VTD method drifts away from the target objects when large appearance changes occur (e.g., #167 in Fig. 8c). When the target object is partially occluded by other similar objects (e.g., #441 in Fig. 8d and #297 in Fig. 8e), the VTD, MIL, and TLD methods do not perform well. Although the KCF tracks the object center location well, it cannot estimate the size of the objects. The IMT algorithm tracks the target objects reliably as different trackers are selected to handle different tracking scenarios, as shown in Fig. 6b.

9.8.4 Woman and CAVIAR

The objects in both sequences undergo heavy occlusions. In addition, the scale of the object in the *CAVIAR* sequence changes significantly, as shown in Fig. 8f. The Struck and TLD methods do not perform well when large-scale change occurs. Due to significant scale changes in the *CAVIAR* sequence, the KCF shows limited tracking performance. When heavy occlusions occur in the *Woman* sequence (#60 in Fig. 8g), the MIL and VTD methods start to drift away from the target object. On the other hand, the IMT algorithm tracks the target objects well by efficiently using Haar-like features, which are more robust for handling occlusion than other features, as shown in Table 3.

9.8.5 Singer1, Sylv, and Trellis

The objects in these sequence undergo large appearance changes due to illumination and pose variations. As shown in Figs. 8h, 8i, the MIL methods do not perform well. The VTD, ASLA, TLD, and Struck approaches do not track the object reliably when illumination and pose variations occur together (#248 and #398 in Fig. 8i). In addition, the Struck and VTD methods do not perform well when scale and large illumination changes occur simultaneously (#54 and #190 in Fig. 8h). The KCF does not deal with large-scale changes well, as shown in the *Singer1* sequence. Different from other tracking methods, the IMT algorithm tracks the object favorably by using complementary features for various appearance changes.

9.8.6 Jumping and Deer

The object appearances change significantly due to fast motion and blurs with noise in both sequences. Except for the IMT, Struck, and TLD methods, other trackers do not handle drastic motion blurs well, as shown in Table 3, Figs. 8j, and 8k. The IMT algorithm effectively uses shape features (HOG) to deal with motion blurs. Table 3 shows that better results are obtained by trackers based on SHOG features. Furthermore, by using stable features in the TLF, large noise caused by motion blurs is well handled by the IMT algorithm, especially in the *Jumping* sequence.

9.8.7 Tiger1, Coke, Board, and Lemming

The target objects in these sequences undergo various appearance changes including motion blurs, illumination changes, occlusions, and pose variations. When the target object undergoes motion blurs and illumination changes simultaneously in the *Coke* sequence (#190 and #216 in Fig. 8p), the ASLA, SCM, and KCF methods do not perform well. When frequent partial occlusions occur (e.g., #316 in Fig. 8o and #190 in Fig. 8p), the ASLA, TLD, KCF, and MIL methods drift away from the target objects. On the other hand, the TLD, VTD, and MIL methods fail to track the objects well (#68 and #249 in Fig. 8m and #383 and #709 in



Fig. 8. Experimental results of state-of-the-art tracking methods.

Fig. 8n) when motion blurs occur. The ASLA and Struck methods do not perform well when large pose changes occur (#540 in the *Board* sequence and #1128 in the *Lemming* video). In contrast, the IMT algorithm performs well, which can be attributed to adaptive use of HOG features to handle motion blurs and Haar-like features to deal with occlusions, as shown in Fig. 6c. As the IMT algorithm utilizes transient and stable features for tracker selection and interaction, it is more robust in dealing with large object appearance changes.

9.9 Run Time Performance

We implement the proposed and evaluated methods (i.e., IMT, MCS, and CVT) using MATLAB. For each method, we use 600 samples for every tracker. The most time-consuming part of the proposed IMT algorithm is to extract multiple features. As the MCS and CVT methods use the same features (HOG, Haar-like, and intensity), the run time performance is comparable to that of the IMT algorithm (0.8 seconds versus 1.4 seconds per frame). The run time of the

IMT is higher as it entails solving an ℓ_1 minimization problem for computing the TLF using (29), which can be further reduced by recent efficient ℓ_1 solvers [33].

10 CONCLUSIONS

In this paper, we propose a robust visual tracking algorithm that integrates multiple trackers based on different feature representations via tracker interaction and selection. The tracker interaction is carried out based on the transition probability matrix, which is designed to alleviate the drifting problems of less reliable tracking methods. The update of the transition probability matrix and tracker selection are computed based on the reliability of each tracker via the proposed tracker likelihood function. To better account for abrupt and gradual appearance changes, each likelihood function is formulated based on transient and stable features. The proposed tracking algorithm selects the best one among multiple trackers to account for object appearance changes. Experimental results on benchmark datasets demonstrate that the proposed tracking algorithm performs favorably against state-of-the-art methods.

ACKNOWLEDGMENTS

This work was supported by the ICT R&D program of MSIP/IITP [B0101-15-0552], an NRF grant funded by the MSIP (No. NRF-2015R1A2A1A01005455), the Giga KOREA Project [GK130100], and the Global Frontier Project (CISS-2011-0031868). M.-H. Yang was supported in part by US National Science Foundation CAREER Grant 1149783 and IIS Grant 1152576.

REFERENCES

- [1] S. Avidan, "Ensemble tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 2, pp. 261–271, Feb. 2007.
- [2] B. Babenko, M.-H. Yang, and S. Belongie, "Robust object tracking with online multiple instance learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 8, pp. 1619–1632, Aug. 2011.
- [3] V. Badrinarayanan, "Probabilistic graphical models for visual tracking of objects," PhD thesis, INRIA Bretagne-Atlantique and Technicolor Corporate Research Labs, Rennes, France, 2009.
- [4] V. Badrinarayanan, P. Perez, F. L. Clerc, and L. Oisel, "Probabilistic color and adaptive multi-feature tracking with dynamically switched priority between cues," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2007, pp. 1–8.
- [5] Y. Bar-Shalom, T. Kirubarajan, and X.-R. Li, *Estimation with Applications to Tracking and Navigation*. New York, NY, USA: Wiley, 2002.
- [6] P. A. Brasnett, L. Mihaylova, N. Canagarajah, and D. Bull, "Particle filtering with multiple cues for object tracking," in *Proc. SPIE's Annu. Symp.*, 2005, pp. 430–441.
 [7] N. Dalal and B. Triggs, "Histograms of oriented gradients for
- [7] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in Proc. IEEE Conf. Comput. Vis. Pattern Recog., 2005, pp. 886–893.
- [8] T. B. Dinh, N. Vo, and G. Medioni, "Context tracker: Exploring supporters and distracters in unconstrained environments," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2011, pp. 1177–1184.
- [9] A. Doucet, S. Godsill, and C. Andrieu, "On sequential monte carlo sampling methods for Bayesian filtering," *Statist. Comput.*, vol. 10, no. 3, pp. 197–208, 2000.
- [10] W. Du and J. Piater, "A probabilistic approach to integrating multiple cues in visual tracking," in *Proc. Eur. Conf. Comput. Vis.*, 2008, pp. 225–238.
- [11] H. Grabner and H. Bischof, "On-line boosting and vision," in Proc. IEEE Conf. Comput. Vis. Pattern Recog., 2006, pp. 260–267.
- [12] H. Grabner, C. Leistner, and H. Bischof, "Semi-supervised on-line boosting for robust tracking," in *Proc. Eur. Conf. Comput. Vis.*, 2008, pp. 234–247.
- [13] S. Hare, A. Saffari, and P. H. S. Torr, "Struck: Structured output tracking with kernels," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2011, pp. 263–270.
- [14] J. A. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "Exploiting the circulant structure of tracking-by-detection with kernels," in *Proc. Eur. Conf. Comput. Vis.*, 2012, pp. 702–715.
- [15] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-speed tracking with kernelized correlation filters," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 3, pp. 583–596, Mar. 2015.
- [16] X. Jia, H. Lu, and M.-H. Yang, "Visual tracking via adaptive structural local sparse appearance model," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2012, pp. 1822–1829.
- [17] V. P. Jilkov and X. R. Li, "Online Bayesian estimation of transition probabilities for Markovian jump systems," *IEEE Trans. Signal Process.*, vol. 52, no. 6, pp. 1620–1630, Jun. 2004.
 [18] Z. Kalal, J. Matas, and K. Mikolajczyk, "P-n learning: Bootstrap-
- [18] Z. Kalal, J. Matas, and K. Mikolajczyk, "P-n learning: Bootstrapping binary classifiers by structural constraints," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2010, pp. 49–56.
- [19] S.-J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky, "An interior-point method for large-scale l1-regularized logistic regression," J. Mach. Learn. Res., 2007, pp. 1519–1555.

- [20] J. Kwon and K. M. Lee, "Visual tracking decomposition," in Proc. IEEE Conf. Comput. Vis. Pattern Recog., 2010, pp. 1269–1276.
- [21] J. Kwon and K. M. Lee, "Tracking by sampling trackers," in Proc. IEEE Int. Conf. Comput. Vis., 2011, pp. 1195–1202.
- [22] I. Leichter, M. Lindenbaum, and Ê. Rivlin, "A general framework for combining visual trackers—the "black boxes" approach," Int. J. Comput. Vis., vol. 67, no. 3, pp. 343–363, 2006.
 [23] B. Liu, J. Huang, L. Yang, and C. Kulikowsk, "Robust tracking
- [23] B. Liu, J. Huang, L. Yang, and C. Kulikowsk, "Robust tracking using local sparse appearance model and k-selection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2011, pp. 1313–1320.
- [24] X. Mei and H. Ling, "Robust visual tracking and vehicle classification via sparse representation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 11, pp. 2259–2272, Nov. 2011.
 [25] X. Mei, H. Ling, Y. Wu, E. Blasch, and L. Bai, "Minimum error
- [25] X. Mei, H. Ling, Y. Ŵu, E. Blasch, and L. Bai, "Minimum error bounded efficient 1 tracker with occlusion detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2011, pp. 1257–1264.
- [26] F. Moreno-Noguer, A. Sanfeliu, and D. Samaras, "Dependent multiple cue integration for robust tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 4, pp. 670–685, Apr. 2008.
- [27] D. A. Ross, J. Lim, R.-S. Lin, and M.-H. Yang, "Incremental learning for robust visual tracking," Int. J. Comput. Vis., vol. 77, no. 1–3, pp. 125–141, 2008.
- [28] J. Santner, C. Leistner, A. Saffari, T. Pock, and H. Bischof, "Prost: Parallel robust online simple tracking," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2010, pp. 723–730.
- *put. Vis. Pattern Recog.*, 2010, pp. 723–730.
 [29] M. Spengler and B. Schiele, "Towards robust multi-cue integration for visual tracking," *Mach. Vis. Appl.*, vol. 14, no. 1, pp. 50–58, 2003.
- [30] H. Wang and D. Suter, "Efficient visual tracking by probabilistic fusion of multiple cues," in *Proc. Int. Conf. Pattern Recog.*, 2006, pp. 892–895.
 [31] Y. Wu, J. Lim, and M.-H. Yang, "Online object tracking: A bench-
- [31] Y. Wu, J. Lim, and M.-H. Yang, "Online object tracking: A benchmark," in Proc. IEEE Conf. Comput. Vis. Pattern Recog., 2013, pp. 2411–2418.
- pp. 2411–2418.
 [32] C. Xu, D. Tao, and C. Xu, "A survey on multi-view learning," *CoRR*, abs/1304.5634, 2013.
- [33] A. Y. Yang, A. Ganesh, Z. Zhou, S. Sastry, and Y. Ma, "A review of fast 11-minimization algorithms for robust face recognition," *CoRR*, abs/1007.3753, 2010.
- [34] A. Yilmaz, O. Javed, and M. Shah, "Object tracking: A survey," ACM Comput. Survey, vol. 38, no. 4, pp. 1–13, 2006.
- [35] J. H. Yoon, D. Y. Kim, and K.-J. Yoon, "Visual tracking via adaptive tracker selection with multiple features," in *Proc. Eur. Conf. Comput. Vis.*, 2012, vol. 7575, pp. 28–41.
- [36] E. Zelniker, T. M. Hospedales, S. Gong, and T. Xiang, "A unified Bayesian framework for adaptive visual tracking," in *Proc. Brit. Mach. Vis. Conf.*, 2009, pp. 1–18.
- [37] W. Zhong, H. Lu, and M.-H. Yang, "Robust object tracking via sparsity-based collaborative model," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, 2012, pp. 1838–1845.



Ju Hong Yoon received the BS degree in electrical and electronic engineering from Sungkyunkwan University in 2008 and the MS and PhD degrees from the Gwangju Institute of Science and Technology in 2009 and 2014, respectively. He is currently a senior researcher at the Korea Electronics Technology Institute. His current research includes multi-object tracking, stereo vision, filtering theory, etc.

YOON ET AL.: INTERACTING MULTIVIEW TRACKER



Ming-Hsuan Yang received the PhD degree in computer science from the University of Illinois at Urbana-Champaign in 2000. He is an associate professor in electrical engineering and computer science at the University of California, Merced. Prior to joining UC Merced in 2008, he was a senior research scientist at the Honda Research Institute working on vision problems related to humanoid robots. He 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 US National Science Foundation CAREER award in 2012, the Senate Award for Distinguished Early Career Research at UC Merced in 2011, and the Google Faculty Award in 2009. He is a senior member of the IEEE and the ACM.



Kuk-Jin Yoon received the BS, MS, and PhD degrees in electrical engineering and computer science from the Korea Advanced Institute of Science and Technology (KAIST) in 1998, 2000, 2006, respectively. He was a post-doctoral fellow in the PERCEPTION team in INRIA-Grenoble, France, for two years from 2006 and 2008 and joined the School of Information and Communications in the Gwangju Institute of Science and Technology (GIST), Korea, as an assistant professor in 2008. He is currently an associate pro-

fessor and the director of the Computer Vision Laboratory in GIST.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.