

Online multi-object tracking via robust collaborative model and sample selection



Mohamed A. Naiel^a, M. Omair Ahmad^{a,*}, M.N.S. Swamy^a, Jongwoo Lim^b,
Ming-Hsuan Yang^c

^a Department of Electrical and Computer Engineering, Concordia University, Montreal, QC H3G 1M8, Canada

^b Department of Computer Science and Engineering, Hanyang University, Seoul 113-791, Republic of Korea

^c School of Engineering, University of California, Merced, CA 95344 USA

ARTICLE INFO

Article history:

Received 5 September 2015

Revised 24 April 2016

Accepted 17 July 2016

Available online 21 August 2016

Keywords:

Multi-object tracking

Particle filter

Collaborative model

Sample selection

Sparse representation

ABSTRACT

The past decade has witnessed significant progress in object detection and tracking in videos. In this paper, we present a collaborative model between a pre-trained object detector and a number of single-object online trackers within the particle filtering framework. For each frame, we construct an association between detections and trackers, and treat each detected image region as a key sample, for online update, if it is associated to a tracker. We present a motion model that incorporates the associated detections with object dynamics. Furthermore, we propose an effective sample selection scheme to update the appearance model of each tracker. We use discriminative and generative appearance models for the likelihood function and data association, respectively. Experimental results show that the proposed scheme generally outperforms state-of-the-art methods.

© 2016 Elsevier Inc. All rights reserved.

1. Introduction

Multi-object tracking (MOT) is one challenging vision problem with numerous applications in automatic visual surveillance, behavior analysis, and intelligent transportation systems, to name a few. In the past decade, more attention has been paid on detecting and tracking one or more objects in videos. Recent advancement in object detection facilitates collaboration between the detection and tracking modules for multi-object tracking (Breitenstein et al., 2009).

Robust multi-object tracking entails solving many challenging problems such as occlusion, appearance variation, and illumination change. A pre-trained object detector robust to appearance variation of one specific class is often used as a critical module of most multi-object tracking methods. Specifically, one detector encodes the generic pattern information about a certain object class (e.g., cars, pedestrians and faces), and one tracker models the appearance of the specific target to maintain the target identity in an image sequence. However, an object detector is likely to generate false positives and negatives, thereby affecting the performance of a tracker in terms of data association and online model update.

In multi-object tracking, offline methods based on global optimization of all object trajectories usually perform better than online counterparts (Andriyenko and Schindler, 2011; Andriyenko et al., 2012; Brendel et al., 2011; Butt and Collins, 2012; Izadinia et al., 2012; Leal-Taixé et al., 2011; Shitrit et al., 2011; Wu et al., 2012; Zamir et al., 2012), and an experimental evaluation of recent methods can be found in Leal-Taixé et al. (2015). For instance, Brendel et al. proposed the maximum-weight independent set of a graph for data association (Brendel et al., 2011), and Zamir et al. used the generalized minimum clique graph to solve the data association (Zamir et al., 2012). In Butt and Collins (2012), the data association problem is solved by using a sliding window of three frames to generate short tracklets, and in case of inconsistencies, the algorithm uses larger tracklet optimization. The minimum-cost network flow is then used to optimize the overall object trajectories. For real-time applications, online methods (Breitenstein et al., 2009; Okuma et al., 2004; Shu et al., 2012; Wu et al., 2008) have been developed within the tracking-by-detection framework where data association between detections and trackers are carried out in an online manner.

Table 1 summarizes the multi-object tracking methods that are most related to this work. Online multi-object tracking can be carried out by using joint state-space model for multi-targets (Duffner and Odoñez, 2013; Eiselein et al., 2012; Jin and Mokhtarian, 2007; Maggio et al., 2008; Okuma et al., 2004; Vermaak et al., 2003).

* Corresponding author.

E-mail addresses: m_naiel@ece.concordia.ca (M.A. Naiel), omair@ece.concordia.ca (M.O. Ahmad), swamy@ece.concordia.ca (M.N.S. Swamy), jlim@hanyang.ac.kr (J. Lim), mhyang@ucmerced.edu (M.-H. Yang).

Table 1

Representative online multi-object tracking algorithms. AM: appearance model, J/I: joint/independent, MU: model update, PF: particle filter, KF: Kalman filter, CVM: constant velocity motion model, MCMC: Markov Chain Monte Carlo, CH: color histogram, LBP: local binary patterns, BOW: bag of words, DCD: detection confidence density, SS: sample selection, PGF: probabilistic gating function, $q(\cdot)$: proposal distribution, SGM: sparsity-based generative model, PGM: 2DPCA-based generative model, SDC: sparsity-based discriminative classifier.

Algorithm	Search and proposal distribution	J/I	Sample descriptor	Data association	MU	Likelihood function
Okuma et al. (2004)	Mixture particle filters, $q(\cdot)$ new observation and propagated particles	J	HSV CH	NA	No	Bhattacharyya similarity
Breitenstein et al. (2009)	PF with CVM	I	RGI, LBP	Boosted classifier, PGF and position	Yes	Distance between each particle and the associated detection, DCD, and Boosted classifier
Yang et al. (2009)	Bayesian filtering	I	RGB, shape, BOW	CVM, position and scale	Yes	Joint likelihood of AM features
Shu et al. (2012)	Detector based or KF	I	CH, LBP	SVM classifier, position, size	Yes + SS	KF (if no associated detection to the tracker)
Zhang et al. (2012)	Mean shift tracker or KF	I	CH, shift vector	Size, search area, tracker re-detection	Yes	Combination of Mean-shift and KF
Schumann et al. (2013)	PF with random walk or CVM	I	RGB CH	Overlap ratio	Yes	Detector confidence
Duffner and Odobez (2013)	MCMC with random walk, $q(\cdot)$ new detections and sampled particles	J	HSV CH	Overlap ratio and position, tracker re-detection	Yes	The product of the visible individual targets likelihoods
Proposed method	PF with detection based CVM, $q(\cdot)$ new associated detections and propagated particles	I	Grayscale	Overlap ratio, SGM and PGM, tracker re-detection	Yes + SS	SDC, different weights for newly created and propagated particles

For instance, a mixture particle filter has been proposed (Okuma et al., 2004) to compute the posterior probability via the collaboration between an object detector and the proposal distribution of the particle filter. However, the joint state-space tracking methods require high computational complexity. The probability hypothesis density filter (Mahler, 2003) has been incorporated in visual multi-target tracking (Maggio et al., 2007; Maggio et al., 2008) since the time complexity is linear with respect to the number of targets. However, it does not maintain the target identity, and consequently, requires an online clustering method to detect the peaks of the particle weights and applies data association to each cluster.

Numerous online multi-object tracking methods deal with each tracker independently (Breitenstein et al., 2009; Schumann et al., 2013; Shu et al., 2012; Yang et al., 2009; Zhang et al., 2012). In Breitenstein et al. (2009), a method based on a particle filter and two human detectors with different features was developed, where the observation model depends on the associated detection, the detector confidence density and the likelihood of appearance. In addition, Shu et al. (2012) introduced a part-based pedestrian detector for online multi-person tracking. This method combines the detection results with the Kalman filter, where data association is performed every frame, and the filter is used when occlusion occurs. Recently, Zhang et al. (2012) used the mean-shift trackers and the Kalman filter for multi-person tracking, where trackers are either weakly or strongly trained. We note that these methods are likely to have low recall as the detector and tracker are not integrated within the same framework.

The degeneracy problem of particle filters (Gordon et al., 1993) has been addressed in several methods (Huang and Djuric, 2004; Jinxia et al., 2012; Rui and Chen, 2001; Santhoshkumar et al., 2013) with more effective proposal distributions and re-sampling steps. Rui and Chen (2001) used the unscented Kalman filter for generating the proposal distribution, and Han et al. (2011) used a genetic algorithm to increase the diversity of the particles. Recently, the Metropolis Hastings algorithm has been used to sample particles from associated detections in the tracking-by-detection framework (Santhoshkumar et al., 2013). We note that the above-mentioned methods do not exploit the collaboration between detectors and trackers (Han et al., 2011; Rui and Chen, 2001), or do not consider the effect of false positive detections on the trackers (Santhoshkumar et al., 2013).

An adaptive appearance model is one of the important factors for effective object tracking as it accounts for appearance change (Salti et al., 2012; Wu et al., 2013). In Okuma et al. (2004), the appearance model is fixed during the tracking process and thus, may result in tracking failure. On the other hand, the trackers are updated with positive samples (Zhang et al., 2012) straightforwardly without differentiating whether they contain noise or not. As multiple objects are likely to be occluded, it is necessary to analyze the samples and reduce the likelihood of including noisy samples for model update. Recently, the appearance models (Shu et al., 2012) have been updated by the detected non-occluded object parts rather than the holistic samples.

In this paper, we propose an online multi-object tracking scheme by using a robust collaborative model for interaction between a number of single-object trackers with sparse representation-based discriminative classifiers (Wright et al., 2009; Zhong et al., 2012), and a pre-trained object detector in the particle filter framework, where every target is tracked independently to avoid the high computational complexity of the joint probability with increasing number of targets. A novel sample selection scheme is used to update each tracker by using key samples with high confidence from the trajectory of an object, where the key sample represents the association between the tracker and a detection at time, t . In addition, we present a data association method with partial occlusion handling by using diverse generative models composed of sparsity-based generative model (Zhong et al., 2012), and two-dimensional principal component analysis (2DPCA) (Yang et al., 2004) generative model. Finally, we introduce a 2DPCA generative model to re-identify lost targets. Experimental results on benchmark datasets demonstrate that the proposed scheme generally outperforms state-of-the-art methods.

2. Overview of the proposed scheme

The proposed multi-object tracking scheme consists of three main components: a pre-trained object detector, a data association module and a number of single-object trackers. Fig. 1 shows the block diagram of the proposed scheme, wherein only one single-object tracker is shown. The object detector is applied on every frame and supports the data association module with a set of detections \mathcal{D}^t at time t . The object tracker adopts a hybrid motion

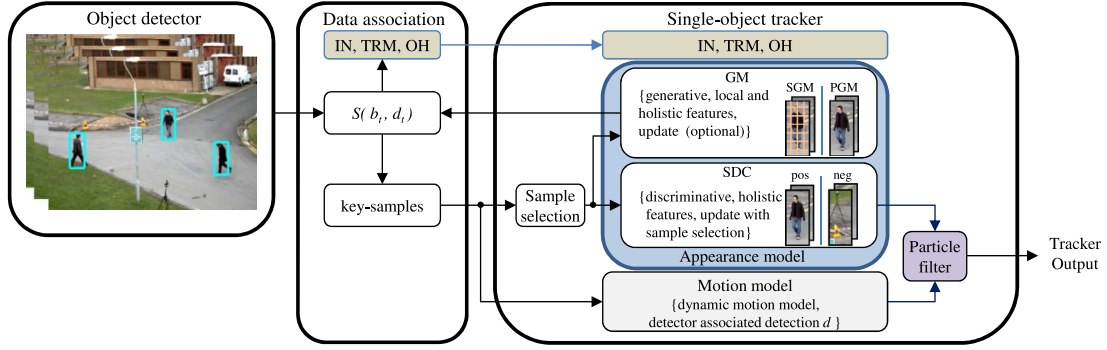


Fig. 1. Block diagram of the proposed multi-object tracking scheme, where IN, TRM, OH, pos, and neg denote initialization, termination, on-hold, positive, and negative, respectively (see text for details).

model, and a particle filter with a robust collaborative model is used to estimate the target location. The appearance model consists of a sparsity-based discriminative classifier (SDC) with holistic features, a sparsity-based generative model (SGM) with local features, and a 2DPCA-based generative model (PGM) with holistic features. The SDC is used to compute each sample confidence score of the particle filter, while the SGM and PGM are used to solve the data association problem. Each tracker also contains a sample selection scheme to update the appearance model with high confidence key samples. Finally, the data association module is used to construct the similarity matrix S to match detections, $d_t \in \mathcal{D}^t$, with existing trackers, $b_t \in \mathcal{B}_t^t$, at time t . Furthermore, it determines initialization, termination and on-hold states of the trackers, and supports the tracker with key samples from the target trajectory.

In this paper, we used the fast pedestrian detector (FPD) (Dollár et al., 2010) for multi-person tracking. In Section 4, we used other pre-trained detectors, such as the on-road vehicle detector proposed in Naiel et al. (2014), and the method in Dollár et al. (2014) to measure the tracking performance on several detection conditions and different types of objects.

3. Tracking scheme

Each object tracker is based on the particle filter tracking framework that uses the sparse representations and 2DPCA as the appearance model. We incorporate two measurements from the detector and tracker into the particle filter, and propose a novel collaborative model that directly affects the likelihood function to obtain the posterior estimate of the target location. We construct the appearance model of the target by using discriminative and generative appearance models, for the likelihood function and the data association. In the following, we use a gate function \mathcal{I}_{b_t} to represent the state of the tracker b_t when associated to the detection d_t at time t . The gate function is defined as

$$\mathcal{I}_{b_t} = \begin{cases} 1, & \text{if } b_t \text{ is associated with } d_t \text{ at time } t \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

3.1. Particle filter

In the Bayesian tracking framework, the posterior at time t is approximated by a weighted sample set $\{\mathbf{x}_t^i, \mathbf{w}_t^i\}_{i=1}^{N_s}$, where \mathbf{w}_t^i is the weight of particle, \mathbf{x}_t^i , and N_s is the total number of particles. The state \mathbf{x} consists of translation (x, y) , average velocity (v_x, v_y) , scale \hat{s} , rotation angle θ , aspect ratio η , and skew direction ϕ .

The measurement model of the proposed particle filter consists of two types. The first measurement is available every time t from the propagated particles $\mathbf{z}_{1:t}$. The second measurement is from the newly created particles that are available at time t when

a detection window, d_t , is associated to a tracker, b_t (i.e., $\mathcal{I}_{b_t} = 1$). Assume that at time t , the tracker b_t is associated to a detection d_t , then we sample candidate particles from the importance density, $q(\mathbf{x}_t^i | \mathbf{x}_{1:t-1}^i, \mathbf{z}_{1:t}, d_t)$. The posterior probability of the candidate location given the available measurements, $p(\mathbf{x}_t | \mathbf{z}_{1:t}, d_t)$, is

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}, d_t) \approx \sum_{i=1}^{N_s} \mathbf{w}_t^i \delta(\mathbf{x}_{1:t} - \mathbf{x}_{1:t}^i) \quad (2)$$

where

$$\mathbf{w}_t^i \propto \mathbf{w}_{t-1}^i \frac{p(\mathbf{z}_{1:t}, d_t | \mathbf{x}_t^i) p(\mathbf{x}_t^i | \mathbf{x}_{t-1}^i)}{q(\mathbf{x}_t^i | \mathbf{x}_{1:t-1}^i, \mathbf{z}_{1:t}, d_t)} \quad (3)$$

and $p(\mathbf{x}_t^i | \mathbf{x}_{t-1}^i)$ is the transition probability. In the proposed method, the particles are resampled every time t , and we have $\mathbf{w}_{t-1}^i = 1/N_s, \forall i$, and then we ignore \mathbf{w}_{t-1}^i term. In the current frame, since the propagated particles sampled at time t corresponding to the tracker position in the previous frame and the particles sampled at time t from the associated detection are independent, the particle weights are computed by

$$\mathbf{w}_t^i \propto \frac{p(\mathbf{z}_{1:t} | \mathbf{x}_t^i) p(d_t | \mathbf{x}_t^i) p(\mathbf{x}_t^i | \mathbf{x}_{t-1}^i)}{q(\mathbf{x}_t^i | \mathbf{x}_{1:t-1}^i, \mathbf{z}_{1:t}, d_t)} \quad (4)$$

where $p(d_t | \mathbf{x}_t^i)$ is the likelihood of detection d_t given the candidate location \mathbf{x}_t^i . To determine the term $p(d_t | \mathbf{x}_t^i)$ the confidence value of the object detector is required at every candidate location from each tracker, which is computationally expensive. Thus, we simplify this detection likelihood term using the Bayes rule as

$$p(d_t | \mathbf{x}_t^i) = p(\mathbf{x}_t^i | d_t) p(d_t) / p(\mathbf{x}_t^i) \propto p(\mathbf{x}_t^i | d_t) \quad (5)$$

where $p(\mathbf{x}_t^i | d_t)$ represents the probability of the candidate location given that the tracker is associated to a detection, d_t . Let the proposal distribution be defined as

$$q(\mathbf{x}_t^i | \mathbf{x}_{1:t-1}^i, \mathbf{z}_{1:t}, d_t) \propto p(\mathbf{x}_t^i | d_t) p(\mathbf{x}_t^i | \mathbf{x}_{t-1}^i) \quad (6)$$

By substituting into (4) with the simplification from (5) and (6), the simplified particle weights with Markovian model can be computed by

$$\mathbf{w}_t^i \propto p(\mathbf{z}_t | \mathbf{x}_t^i) \quad (7)$$

By normalizing the particle weights, the resulting state estimate is represented as a weighted average of the candidate locations. This makes the proposed scheme more robust to noisy detection results compared to maximum a posteriori methods.

When there is no detection associated to a tracker (i.e., $\mathcal{I}_{b_t} = 0$) the proposed particle filter can be simplified to the bootstrap particle filter (Gordon et al., 1993). In the bootstrap particle filter, the measurement model consists of the tracker measurements $\mathbf{z}_{1:t}$ and the importance density at time t can be defined as $q(\mathbf{x}_t^i | \mathbf{x}_{1:t-1}^i, \mathbf{z}_{1:t}) \propto p(\mathbf{x}_t^i | \mathbf{x}_{t-1}^i)$. It can be shown that the particle weights can be represented by (7).

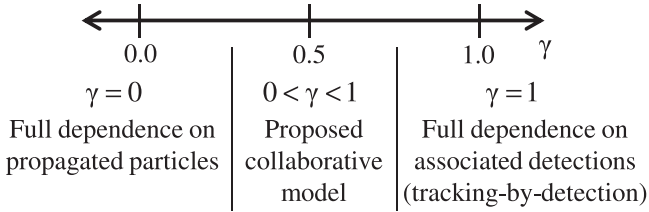


Fig. 2. Effect of changing the collaborative factor γ .

3.1.1. Motion model

In the proposed method, we adopt a hybrid motion model based on the first-order Markov chain and the associated detection. The new candidate state \mathbf{x}_t^d at time t is provided to the motion model if a detection is successfully associated to the tracker (i.e., $\mathcal{I}_{b_t} = 1$) and the initial velocity is set to be the average velocity of the tracker particles. The candidate state at time t , \mathbf{x}_t , relates to the set of propagated particles X^{b_t} and the set of associated detection X^{b_t, d_t} , by

$$\mathbf{x}_t = \begin{cases} F\mathbf{x}_{t-1} + Q & \text{if } \mathbf{x}_t \in X^{b_t} \\ \mathbf{x}_t^d + P & \text{if } \mathbf{x}_t \in X^{b_t, d_t} \end{cases} \quad (8)$$

where Q and P are the Gaussian noise vectors, $N_s = N_s^P + N_s^\Gamma$, and N_s^P and N_s^Γ are the cardinality of X^{b_t} and X^{b_t, d_t} , respectively. In the above equation, F denotes the transition matrix of size 8×8 , which is defined as

$$F = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \quad (9)$$

3.1.2. Robust collaborative model

The object detector applies computationally expensive space-scale search to the entire image to localize specific class of objects, and proposes candidate locations that have high probability of existence. To exploit high confidence associated detections, we incorporate a set of new particles, X^{b_t, d_t} , in the likelihood function to allow the object detector to guide the trackers. Let $H(\mathbf{x}_t^i)$ denote SDC tracker confidence score of candidate \mathbf{x}_t^i . The likelihood of the measurement, \mathbf{z}_t , can be computed by

$$p(\mathbf{z}_t | \mathbf{x}_t^i) = \pi^i H(\mathbf{x}_t^i) \quad (10)$$

where

$$\pi^i = \begin{cases} 1 - \gamma & \text{if } \mathcal{I}_{b_t} = 1, \mathbf{x}_t^i \in X^{b_t} \\ \gamma & \text{if } \mathcal{I}_{b_t} = 1, \mathbf{x}_t^i \in X^{b_t, d_t} \\ 1 & \text{otherwise, i.e., } \mathcal{I}_{b_t} = 0 \end{cases} \quad (11)$$

and $\gamma \in [0, 1]$ is the collaborative factor. In (10), the particles from the associated detections and previously propagated particles are weighted differently. Fig. 2 shows the effect of changing the collaborative factor value. Fig. 3(a) and (b) shows an example of particle weights for the detector particles and the propagated particles using $\gamma = 0.54$. If $\mathcal{I}_{b_t} = 1$ and $\gamma > 0.5$, the weight π^i allows the detector to guide the tracker by giving more weights to the newly associated particles than the propagated particles. However, a detector may have false positives, and thus, the tracker should not depend completely on the detector. From our experiments, we find that the proposed scheme with the value of γ between 0.5 and 0.85 performs best. If the detector suffers from missing detections (i.e., $\mathcal{I}_{b_t} = 0$), the likelihood function in (10) will only depend

on the previously propagated particles $\mathbf{x}_t^i \in X^{b_t}$, which represent the bootstrap particle filter (Gordon et al., 1993). Our collaborative model is based on the hybrid motion model that incorporates associated detections with object dynamics. In contrast, the motion model adopted in Breitenstein et al. (2009) depends only on propagated particles, and the likelihood function depends on tracker appearance model and the detector confidence density. The collaborative model in Okuma et al. (2004) only exists in the proposal distribution and the likelihood is without weighting collaborative factor.

3.1.3. Resampling

In each frame, the set of candidate particles $\{\mathbf{x}_t^i, \mathbf{w}_t^i\}_{i=1}^{N_s}$ are resampled to avoid the degeneracy problem. The resampling process also allows the detector to guide the tracker effectively. As each tracker resamples particles based on particle weights computed from the proposed collaborative model (10), the propagated particles with low weights are replaced with newly created particles from the associated detections.

3.2. Appearance model

In the proposed method, the SGM and SDC are used in a way different from that in Zhong et al. (2012). First, we do not use the collaboration between SGM and SDC (Zhong et al., 2012), instead we use SGM with PGM to compute the similarity matrix of the data association module for occlusion handling (23), and the modified SDC model is used to compute the likelihood of the particle filter (10). The number of particles in the filter is usually larger than the number of detections and trackers at time t , and the computational complexity of SDC is lower than SGM. Therefore, the resulting tracker is more efficient. Second, our SDC uses the down-sampled grayscale image without the feature selection method used in Zhong et al. (2012). Third, our SDC confidence measure depends on the sparsity concentration index (Wright et al., 2009). Finally, we propose the key sample selection scheme to update the appearance models with high confidence samples.

3.2.1. Sparsity-based discriminative classifier

We construct a discriminative sparse appearance model to compute the confidence score as used in (10). The initial training samples are collected in a similar way to Zhong et al. (2012), where each SDC tracker is initialized using N_p positive samples drawn from the object center with a small variation from the center of the detection state \mathbf{x}_t^d , and N_n negative samples are taken from the annular region surrounding the target center without overlap with a detection window d_t . Next, each sample is normalized to a canonical size of $(m \times n)$, and vectorized to be one column of the matrix $A \in \mathbb{R}^{r \times N^t}$, where $r = mn$ and $N^t = N_p + N_n + N_{p,u}^t + N_{n,u}^t$, such that $N_{p,u}^t$ and $N_{n,u}^t$ denote the buffer size of the selected key samples up to time t . Let the measurement corresponding to the candidate location \mathbf{x}_t^i be denoted by $\mathbf{z}_t^i \in \mathbb{R}^r$. We obtain the sparse coefficients α^i for the i th candidate by solving the following optimization problem,

$$\min_{\alpha^i} \|\mathbf{z}_t^i - A\alpha^i\|_2^2 + \lambda_{SDC} \|\alpha^i\|_1 \quad (12)$$

We compute the classifier confidence score by

$$H(\mathbf{x}_t^i) = \exp\left(-\frac{(\varepsilon_+^i - \varepsilon_-^i)}{\sigma}\right) \Omega(\alpha^i) \quad (13)$$

where $\varepsilon_+^i = \|\mathbf{z}_t^i - A_+ \alpha_+^i\|_2^2$ is the reconstruction error of the candidate \mathbf{z}_t^i with respect to the template set of the positive class A_+ , and the sparse coefficient vector of the i th candidate that corresponds to the positive class, α_+^i . Similarly, $\varepsilon_-^i = \|\mathbf{z}_t^i - A_- \alpha_-^i\|_2^2$ is the reconstruction error of the same candidate \mathbf{z}_t^i with respect to

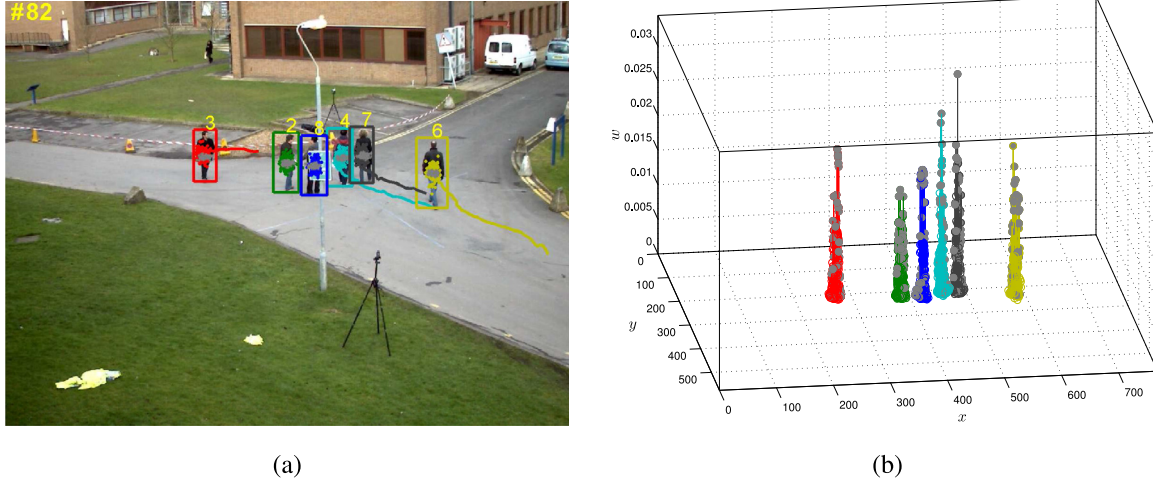


Fig. 3. Effect of the proposed collaborative model on the tracker particles. (a) Illustrates the candidate particles proposed by the object detector (masked as gray) and propagated particles (colored). (b) Particles weights for new (masked as gray) and propagated particles (colored).

the template set of the negative class A_- , and the corresponding sparse coefficient vector α_-^t . The parameter σ adjusts the confidence measure, and $\Omega(\alpha^i)$ represents the sparsity concentration index (SCI) (Wright et al., 2009) defined as

$$\Omega(\alpha^i) = \frac{J \cdot \max_j \|\delta'_j(\alpha^i)\|_1 / \|\alpha^i\|_1 - 1}{J - 1} \in [0, 1] \quad (14)$$

where δ'_j is a function that selects the coefficients corresponding to the j th class and suppresses the rest, and J is the number of classes ($J = 2$ in this work). The SCI checks the validity of a candidate such that it can be represented by a linear combination of the training samples in one class. When the sparse coefficients concentrate in a certain class, the SCI value is high. This index allows each tracker to assign high weights to candidates resembling the positive training samples, and rejects others related to other targets or background structures.

The SDC tracker is updated every R_u frames using the selected key samples, K_u^t (Section 3.3). At each key sample location, we collect positive and negative samples as part of the initialization process. To leverage between computational load and memory requirement, we set the maximum number of positive and negative samples. If the number of positive, $N_{p,u}^t$ or negative, $N_{n,u}^t$ samples exceeds the limit, we replace the old samples (other than those collected in the first frame) with the new selected key samples.

3.2.2. Sparsity-based generative model

We use a sparsity-based generative model to measure similarity in the data association module. Fig. 4 illustrates the block diagram of the proposed SGM in the training and test modes. The training template consists of M local patches, $\{y_i\}_{i=1}^M$ and each patch of size $\hat{m} \times \hat{n}$. These M patches are vectorized¹ and quantized into N_k centroids using the k -means algorithm to construct the dictionary $D \in \mathbb{R}^{\hat{f} \times N_k}$ ($\hat{f} = \hat{m}\hat{n}$). For the i th patch, y_i , the sparse-coefficients, $\beta_i \in \mathbb{R}^{N_k \times 1}$, is computed by

$$\min_{\beta_i} \|y_i - D\beta_i\|_2^2 + \lambda_{SGM} \|\beta_i\|_1 \quad (15)$$

The adopted SGM is concerned with representing the appearance of the positive class of the tracker by using the sparse coefficients of M local patches of the object and candidate location c , where

¹ The vectorization function is defined as $\text{Mat2Vec}: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^r$, where $r = mn$ is the dimension of the vector, and $(m \times n)$ is the order of the input matrix. The inverse of the vectorization function is defined as $\text{Vec2Mat}: \mathbb{R}^r \rightarrow \mathbb{R}^{m \times n}$.

each location is represented by a sparse histogram feature vector $\rho = [\beta_1, \beta_2, \dots, \beta_M]^T$, and $\rho^c = [\beta_1^c, \beta_2^c, \dots, \beta_M^c]^T$, corresponding to the initial object and the candidate location, respectively. To handle occlusion, the patch reconstruction error, $\{\varepsilon_i = \|y_i^c - D\beta_i^c\|_2^2\}_{i=1}^M$, is used to suppress the coefficients of occluded patches. Let ψ_i be the non-occlusion indicator for the i th patch and is computed by

$$\psi_i = \begin{cases} 1_{N_k,1} & \text{if } \varepsilon_i < \varepsilon_0 \\ 0_{N_k,1} & \text{otherwise} \end{cases} \quad (16)$$

where $1_{N_k,1}$ and $0_{N_k,1}$ denote the vector of size N_k of ones and zeros. The final histogram can be represented by $\varphi = \psi \odot \rho$, and $\varphi^c = \psi \odot \rho^c$, corresponding to the training template, and the candidate location, where \odot denotes the element-wise multiplication. By taking the spatial representation into consideration, the resulting histogram, φ can handle occlusion effectively. Fig. 5 illustrates the effect of the partial occlusion handling scheme. If the reconstruction error is greater than the threshold, ε_0 , then the non-occlusion indicator, ψ , suppresses these patches. The generative model similarity, $G_{SGM}(b_t, c)$, between the candidate φ_c and the model φ is measured by using the intersection kernel.

As in Zhong et al. (2012), the dictionary, D , is fixed during the tracking process, while the sparse histogram of the initial template, $\rho_{initial}$, is updated every update rate, R_u . The sparse histogram is updated by

$$\rho_{new} = \mu \rho_{initial} + (1 - \mu) \rho_K \quad (17)$$

where μ is the learning rate, and ρ_K represents the sparse histogram corresponding to the selected key sample from the set K_u^t that provides the maximum similarity to the training templates (see Section 3.3 for the sample selection scheme). This conservative update scheme by using the confidence key samples and maintaining the initial template provide effective tracking.

3.2.3. 2DPCA-based generative model

In addition to part-based SGM, we use a holistic generative model based on the 2DPCA scheme (Yang et al., 2004), referred to as PGM, to solve the data association problem. The reason being that a combination of PGM and SGM increases the tracking performance (see Section 4). For each tracker b_t , we use N positive samples, $\{Y_j\}_{j=1}^N$ each of size $m \times n$, where samples are taken from the positive class of the initial target location, or selected key samples, K_u^t . Each j th sample Y_j is projected by the orthonormal matrix $V \in \mathbb{R}^{n \times r_1}$, $r_1 \leq n$ and form $F^j = Y_j V$, of size $m \times r_1$. The image

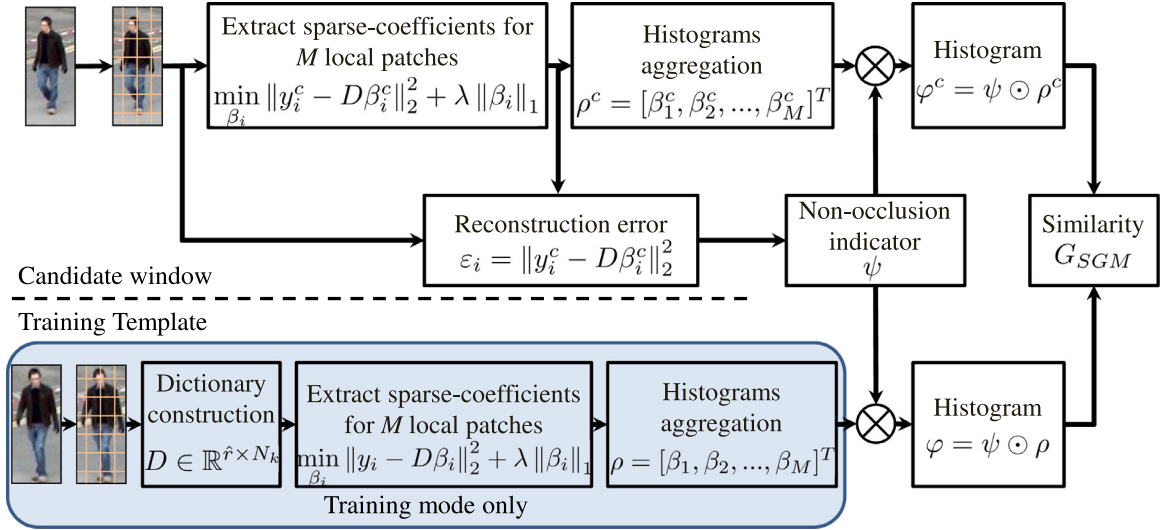


Fig. 4. Block diagram of the sparsity-based generative model.

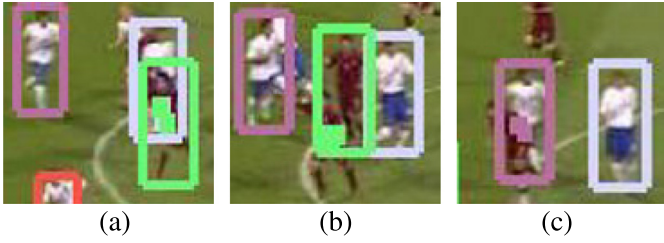


Fig. 5. Sample results for SGM partial occlusion handling scheme, where the marked patches with the same tracker color are the patches at which SGM reconstruction error is greater than the SGM error threshold.

covariance matrix Γ is defined by

$$\Gamma = \frac{1}{N} \sum_{j=1}^N (Y_j - \bar{Y})^\top (Y_j - \bar{Y}) \quad (18)$$

where \bar{Y} is the average image of all training samples, and Γ is the nonnegative definite matrix. The objective of 2DPCA is to find the optimal orthonormal matrix, V_{opt} , that maximizes the total scatter in the learned subspace. The total scatter criterion $J(V)$ is defined by

$$J(V) = V^\top \Gamma V \quad (19)$$

The optimal projection matrix V_{opt} is composed of the r_1 eigenvectors of matrix Γ corresponding to the first r_1 largest eigenvalues, where the vectors are stacked together in matrix V of size $n \times r_1$. We extract features of the j th training example, Y_j , through projecting on matrix V , as $F^j = Y_j V$, and then we vectorize the resulting feature matrix and have the feature vector f^j of size $(1 \times mr_1)$.

For each candidate location, we project the candidate sample, Y^c , using the matrix V , and vectorize the resulting matrix to obtain the test feature vector f^c of size $1 \times mr_1$. The nearest neighbor classifier is used to infer the index of the j th training example, \hat{j} closest to the test vector f^c

$$\hat{j} \leftarrow \underset{j \in \{1, 2, \dots, N\}}{\operatorname{argmin}} \|f^c - f^j\|_2 \quad (20)$$

where $\|\cdot\|_2$ denotes the l_2 -norm. The reconstruction error between the test image and the training examples is $\varepsilon_{PGM} = \|a_{\hat{j}} - a_c\|_2$, where $a_{\hat{j}} = \operatorname{Mat2Vec}(F^{\hat{j}} V^\top)$ and $a_c = \operatorname{Mat2Vec}(F^c V^\top)$. The similarity between the test and training features is computed by

$$G_{PGM} = \exp(-\varepsilon_{PGM} / \hat{\sigma}^2) \quad (21)$$

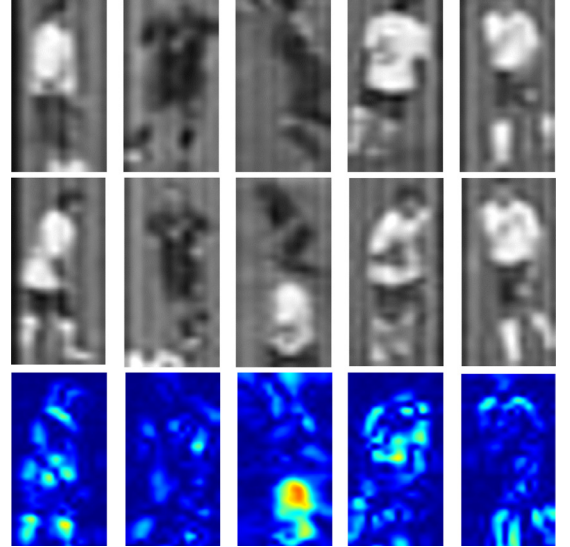


Fig. 6. (Top) Reconstructed nearest neighbor training samples by PGM. (Middle) Reconstructed patches at candidate locations. (Bottom) Absolute reconstruction error, where the pixel with brighter color means high error value.

Fig. 6 shows a sample intermediate output from the proposed PGM scheme. The PGM is able to retrieve the closest training patches in 2DPCA feature subspace, which provides accurate similarity measures in (21).

Similar to SDC tracker, PGM is updated every R_u frames, by using the initial positive and the selected key samples at time t , where $N = N_p + N_{p,u}^t$. To update 2DPCA feature space, we used a batch learning technique. In this scheme, we update the optimal projection matrix, V_{opt} , and extract the feature vectors, $\{f^j\}_{j=1}^N$. While the incremental 2DPCA learning has been used in Wang et al. (2007), we find that batch learning performs more efficiently than the incremental learning scheme, since we replace some samples every update rate with newly selected key samples.

3.3. Sample selection

We propose a sample selection scheme to learn and adapt the appearance model for each tracker by using the samples with high confidence from the object trajectory, in a way similar to

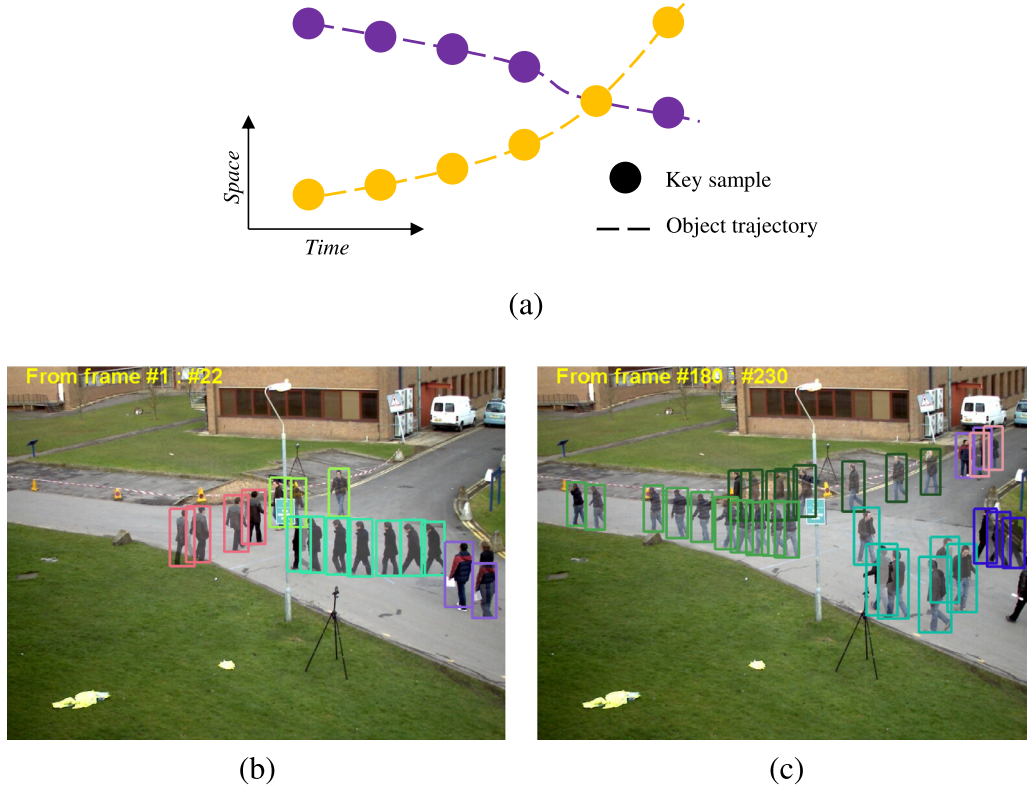


Fig. 7. (a) Key samples in the object trajectories and occlusion issues that should be handled, (b and c) Examples for key samples selected from object trajectories, using a sequence from the *PETS09-S2L1* dataset.

existing methods (Breitenstein et al., 2009; Kuo et al., 2010; Shu et al., 2012). Examples for key sample locations in the object trajectory are shown in Fig. 7, where two scenarios for the key samples are selected from the tracker history. The sample selection scheme alleviates the problem of including occluded samples for more effective model update and thus, reduces the drifting problem. The proposed sample selection scheme is based on the following criteria:

1. We measure the goodness of the key samples. A good key sample is one at which the tracker b_t does not intersect with other trackers or nearby detections except the associated detection d_t . We denote the set of good key samples at time t by K_g^t .
2. We use the online trained SDC tracker to measure the similarity between the current appearance model of the tracker, b_t , and the i th good key sample $K_{g,i}^t \in K_g^t$ by

$$S_{DC}(b_t, K_{g,i}^t) = \exp(-(\varepsilon_+^i - \varepsilon_-^i)/\sigma^2) \quad (22)$$

where $\varepsilon_+^i = \|\mathbf{z}_t^i - A_+ \alpha_+^i\|_2^2$, $\varepsilon_-^i = \|\mathbf{z}_t^i - A_- \alpha_-^i\|_2^2$, and α_+^i and α_-^i are computed by using (12).

3. If $S_{DC}(b_t, K_{g,i}^t) > s_0 \geq 0$, where s_0 is the SDC similarity threshold, then this key sample is selected for the model update. The final set of selected key samples, K_u^t , which have high similarity with the SDC tracker, are used to update the tracker appearance model (Section 3.2). Note that when $s_0 = 0$, all the samples are selected.

3.4. Data association

The similarity matrix S for data association measures the relation between a tracker $b_t \in \mathcal{B}_t^t$ and a detection $d_t \in \mathcal{D}^t$ by

$$S(b_t, d_t) = G(b_t, d_t)O(b_t, d_t) \quad (23)$$

where $G(b_t, d_t) = G_{SGM}(b_t, d_t) + G_{PGM}(b_t, d_t)$ considers the appearance similarity between the tracker b_t and detection d_t , and $O(b_t, d_t)$ represents the overlap ratio between the tracker and the detection to suppress confusing detections, where the overlap ratio is based on the PASCAL VOC criterion (Everingham et al., 2010).

The association is computed online by using the Hungarian algorithm to match a tracker to a detection in a way similar to existing methods (Breitenstein et al., 2009; Shu et al., 2012). The proposed data association scheme iteratively finds the maximum in the matrix S , and associates the tracker b_t to a detection d_t if $S(b_t, d_t)$ is larger than a threshold s_1 . The row and the column corresponding to $S(b_t, d_t)$ are removed. As the object detector is likely to miss some objects, using the similarity threshold, s_1 , can alleviate the tracker to be updated with confusing nearby detections. Furthermore, we select a number of key samples to update the appearance model (Sections 3.3 and 3.2). We initialize new trackers with non-associated detection windows if the maximum overlap with other existing trackers is less than o_1 to avoid creating multiple trackers for the same target.

Re-detection module

A pre-trained object detector usually suffers from false positives and negatives, thereby causing trackers to drift. On the other hand, a tracker does not perform well in the presence of heavy occlusion or background clutters. To handle these challenging cases, we introduce the inactive or on-hold states before tracker termination in case the tracker misses a high number of detections. Let the set of trackers on-hold be denoted as \mathcal{B}_h^t . When the tracker does not estimate the target location at an inactive state, we adopt the PGM (Section 3.2) to measure the similarity between the tracker on-hold $b_t \in \mathcal{B}_h^t$ and the new candidate location. When the tracker is in the inactive state b_t^h , it still can be reinitialized after checking the similarity with the new un-associated detection, d_t^u by $S_h(b_t^h, d_t^u) = G_{PGM}(b_t^h, d_t^u)$ (where G_{PGM} is computed by (21)).

The inactive tracker is reactivated if $S_h(b_t^h, d_t^u) > s_2$, where s_2 is a pre-defined threshold. During the inactive state, the proposed tracker can re-identify lost targets and discriminate among trackers using the 2DPCA feature space learned from selected key samples.

4. Experimental results

4.1. Datasets

We evaluate the tracking performance of the proposed algorithm using seven challenging sequences, namely, the *PETS09-S2L1*, *PETS09-S2L2* (Ferryman, 2009), *UCF Parking Lot (UCF-PL)* dataset (Shu et al., 2012), *Soccer* dataset (Wu et al., 2008), *Town Center* dataset (Benfold and Reid, 2011), and *Urban* as well as *Sunny* sequences from *LISA 2010* dataset (Sivaraman and Trivedi, 2010), and compare it with that of several state-of-the-art online MOT methods.

The *PETS09-S2L1* sequence consists of 799 frames of 768×576 pixels recorded at 7 frames per second with medium crowd density. The *PETS09-S2L2* sequence consists of 442 frames with the same resolution and frame rate as the *PETS09-S2L1* sequence, but it contains heavy crowd density and illumination changes. The target objects undergo scale changes, long-term occlusion, and with similar appearance. The ground truth (GT) data from Yang and Nevatia (2012), (<http://iris.usc.edu/people/yangbo/downloads.html>). Last retrieved March 14, 2016) and (<http://research.milanton.de/data.html>). Last retrieved March 14, 2016) are used for evaluating the tracking results on *PETS09-S2L1* and *PETS09-S2L2*, respectively. The *Soccer* sequence consists of 155 frames of 960×544 pixels recorded at 3 to 5 frames per second. The challenging factors of this sequence include heavy occlusion, sudden change of motion direction of players, high similarity among players of the same team, and scale changes. The GT data provided by Wu et al. (2008) are used for evaluation. On the *PETS09-S2L1*, *PETS09-S2L2* and *Soccer* sequences, the FPD detector Dollár et al. (2010) is used as the baseline detector for the proposed tracking scheme.

The *UCF-PL* dataset consists of 998 frames of 1920×1080 pixels recorded at 29 frames per second with medium crowd density, long-term occlusion, and targets of similar appearance. On this dataset, the detection results of the part-based pedestrian detector proposed in Shu et al. (2012) are used for evaluation based on the GT data provided by (<http://crcv.ucf.edu/data/ParkingLOT/index.php>). Last retrieved March 14, 2016).

The *Town Center* dataset consists of 4500 frames of 1080×1920 pixels recorded at 25 frames per second. The dataset contains medium crowd density, heavy occlusion, and scale changes. In Benfold and Reid (2011), two categories of GT annotations are provided based on the full body and head regions of pedestrians. On this dataset, the aggregated channel feature (ACF) detector proposed by Dollár et al. (2014) is used for performance evaluation. In the case of the full body of pedestrians, it has been observed that the ACF detector does not perform well on this sequence as the false positive rate is high. To alleviate this problem, the first 500 frames of this sequence are used to collect hard-negative samples related to the background clutters, and the ACF detector is re-trained using both the INRIA dataset (Dalal and Triggs, 2005) and hard-negative samples. In case of tracking multiple people based on the head regions, the positive training examples provided in Benfold and Reid (2011) and negative samples collected from the first 500 frames of this sequence are used to train the ACF detector.

The *Urban* and *Sunny* sequences from the *LISA 2010* dataset (Sivaraman and Trivedi, 2010) contain car images of 704×480 collected at 30 frames per second from a camera mounted on a moving vehicle. The *Urban* sequence (300 frames) was captured from an urban area with a low traffic density on a cloudy day, while the

Sunny sequence (300 frames) was captured from a highway with medium traffic density on a sunny day. The challenging factors of these sequences include the effect of camera vibration, illumination changes, and the targets' scale changes; the GT data are provided by Sivaraman and Trivedi (2010). The pre-trained vehicle detector proposed in Naiel et al. (2014) is used for evaluation on this dataset.

4.2. Qualitative results

In this section, we study the qualitative performance of the proposed tracking scheme using the datasets mentioned above. Figs. 8 and 9 show some of the tracking results and videos are available at <https://youtu.be/InAUnU596UE>.

PETS09-S2L1. Fig. 8(a) shows the sample tracking results of the proposed scheme on the *PETS09-S2L1* sequence. The proposed method performs well despite several short-term occlusions, scale and pose changes. Furthermore, it should be mentioned that the pre-trained FPD detector (Dollár et al., 2010) misses objects that are close to the camera or those located far from the camera.

PETS09-S2L2. Fig. 8(b) shows that non-occluded targets are tracked well although targets with long-term occlusions or located far from the camera are missed. Again, it should be mentioned that the FPD detector (Dollár et al., 2010) misses numerous detections in this sequence due to the high crowd density.

Soccer. This sequence contains soccer players with similar visual appearance and fast motion. The FPD detector (Dollár et al., 2010) is not trained to detect the soccer players at different poses. Nevertheless, the proposed scheme performs well with accurate short tracklets, as shown in Fig. 8(c).

UCF-PL. This sequence contains crowds of medium density, with occlusions. Fig. 8(d) shows some tracking results for the proposed scheme using the detector in Shu et al. (2012). Despite the challenges of the sequence, the proposed tracking scheme maintains long trajectories.

Town Center. The crowd density of this sequence is medium with a number of long-term occlusions. Fig. 8(e) and (f) shows sample tracking results corresponding to full body and head, respectively. While it is difficult to track the full human body due to heavy occlusions, or the head due to false positives, the proposed method performs well.

LISA 2010. Fig. 9(a) and (b) shows the sample results of our tracker using the detector in Naiel et al. (2014) on the *Urban* and *Sunny* sequences. The *Urban* sequence contains only one vehicle, but there is illumination change and the effect of camera vibrations. The *Sunny* sequence contains, on average, three non-occluded vehicles with different velocities. In spite of these challenges, the proposed scheme tracks the vehicles very well in both cases.

4.3. Quantitative results

We use the CLEAR MOT metrics (Bernardin and Stiefelwagen, 2008) including multiple object tracking accuracy (MOTA), multiple object tracking precision (MOTP), false negative rate (FNR), false positive rate (FPR), and identity switches (IDSW) for evaluating the performance of the proposed tracker. We use the overlap threshold of 0.5 for all experiments. For this study, we set the various parameters to be $N_s^p = 150$, $N_s^f = 100$, $N_p = N_{p,u}^c = 10$, $N_n = N_{n,u}^c = 20$, $R_u = 10$, $\lambda_{SDC} = 0.02$, $\lambda_{SGM} = 0.01$, $\hat{\sigma} = 10^4$, $\varepsilon_0 = 0.8$, $\mu = 0.6$, $\hat{\sigma} = 5 \times 10^6$, $s_0 = 1.0$, $s_1 = 2.5$, $s_2 = 0.7$, and $o_1 = 0.2$. For the multi-person tracking sequences, namely, *PETS09-S2L1*, *PETS09-S2L2*, *UCF-PL*, *Soccer*, and *Town Center (Body)*, we use $m = 32$, $n = 16$, $M = 84$,



Fig. 8. Sample tracking results for five sequences, the arrangement from top to bottom as (a) and (b) *PETS09-S2L1*, and *PETS09-S2L2*, respectively, (c) Soccer sequence, (d) *UCF-PL* sequence, (e) *Town Center* dataset (body), and (f) *Town Center* dataset (head).

$\hat{m} = \hat{n} = 6$ and $N_k = 50$. Further, for the multi-head tracking sequence, namely, *Town Center (Head)*, as well as the multi-vehicle tracking sequences, namely, *Urban* and *Sunny*, we use $m = n = 16$, $M = 16$, $\hat{m} = \hat{n} = 6$ and $N_k = 16$.

Effect of the collaborative factor. To measure the effect of the proposed collaborative model, we changed the value of the collab-

orative factor γ in the interval $[0,1]$ in increments of 0.2. Fig. 10 shows the performance of the proposed method with different values of γ for the *PETS09-S2L1* sequence. When $\gamma = 0$, the likelihood function of the particle filter is based completely on the propagated particles, and the proposed method does not perform well due to the degeneracy problem. When $\gamma = 1$, the likelihood function is based on the associated detections, and the tracker does not



Fig. 9. Sample tracking results for *LISA 2010* dataset, where (a) and (b) correspond to *Urban* and *Sunny* sequences, respectively.

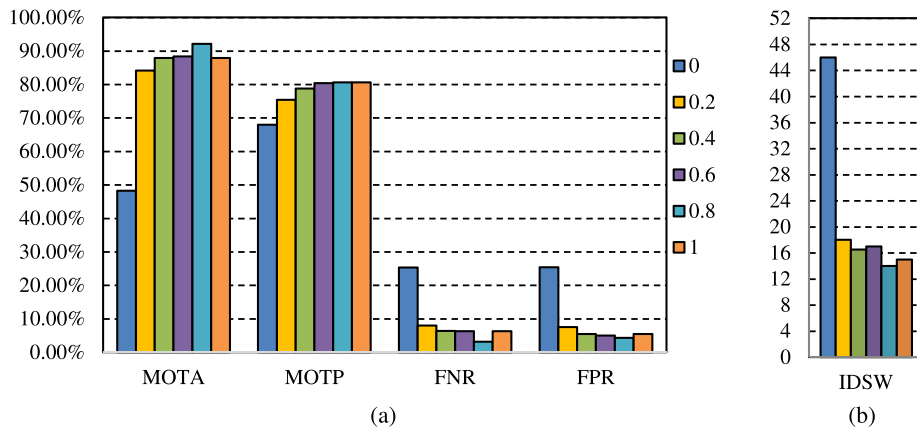


Fig. 10. Performance of the proposed method on the *PETS09-S2L1* sequence for different values of the collaborative factor γ .

perform well due to false positives and missed detections. The proposed method performs best for this sequence when $\gamma = 0.8$, as can be seen from Fig. 10. It is worth noting that for high tracking performance, the value of γ should be adjusted according to the detector used. For detectors with high precision and recall (the ones used in the *PETS09-S2L1*, *UCF-PL*, *Town Center (Head)*, *Urban* and *Sunny* sequences), the proposed tracker provides a high MOTA value when γ is in the interval of $[0.65, 0.85]$. On the other hand, when the detector has low precision and recall (the ones used in the case of *PETS09-S2L2*, *Soccer* and *Town Center (Body)* sequences), the proposed tracker provides a high MOTA value when γ is in the interval of $[0.5, 0.6]$.

Number of key samples. We analyze the effect of the number of key samples retained on MOTA using the *PETS09-S2L1* sequence. The appearance model (SDC, SGM, and PGM) is updated online at an update rate R_u of 10. Fig. 11 shows the performance of the proposed tracker when the number of key samples retained is varied. We choose the number of retained key samples to be 20 at which the highest MOTA performance is exhibited, as seen from Fig. 11.

Key sample selection. To demonstrate the strength of the proposed sample selection scheme, we examine the performance of the proposed tracking scheme by varying the SDC similarity threshold, s_0 , from 0 to 1.5 in increments of 0.1. Fig. 12 shows the performance of the proposed scheme at different SDC tracker similarity threshold values. When $s_0 = 1$, the proposed tracker exhibits the best performance in terms of MOTA. If $0 \leq s_0 < 1$, the performance is not as good in view of the fact that only a few or none of the key samples are rejected, and hence, occluded samples are

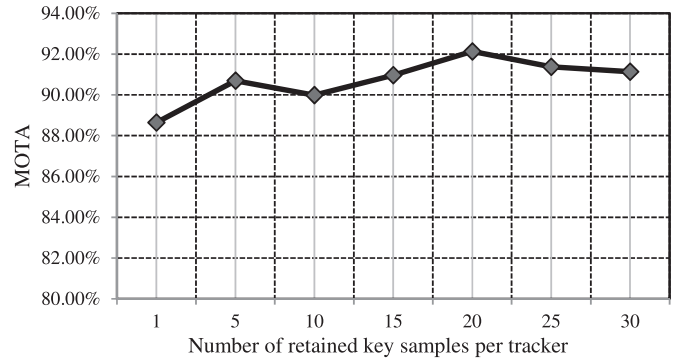


Fig. 11. MOTA vs. number of retained key samples for the proposed tracker on the *PETS09-S2L1* sequence.

likely to be selected. When $s_0 > 1.2$, the proposed tracker performs worse than that at $s_0 = 1$, since a large number of key samples are rejected. As such, we choose $s_0 = 1.0$ for all the experiments.

Effect of tracker re-detection. We analyze the effect of using the re-detection module on MOT tracking. Fig. 13 shows that the proposed method with tracker re-detection scheme achieves slightly lower FNR and FPR than that obtained without using the tracker re-detection scheme, while maintaining approximately the same performance in terms of MOTA and MOTP values. The tracker re-detection scheme aims to reduce the number of identity switches and maintains long trajectories, without reducing the tracking performance.

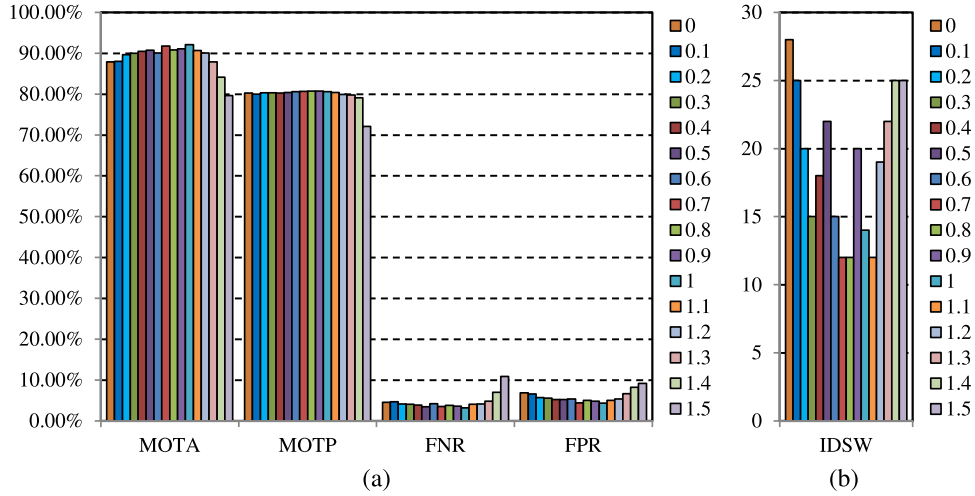


Fig. 12. Performance of the proposed tracking scheme with respect to the SDC similarity threshold, s_0 , using the *PETS09-S2L1* sequence.

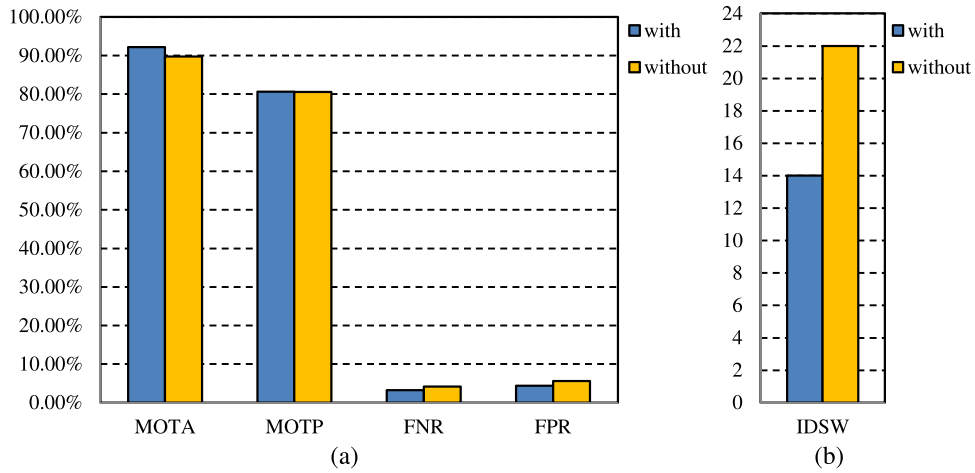


Fig. 13. Performance of the proposed method with and without tracker re-detection on the *PETS09-S2L1* sequence.

Generative appearance models. We study the tracking performance of the proposed method by using several types of generative models to solve the data association problem in (23). These generative models are (1) SGM, as outlined in Section 3.2.2, which is based on local patch features (by substituting in (23) by $G = G_{SGM}$); (2) 2DPCA generative model, as proposed in Section 3.2.3, which is based on holistic features (by substituting in (23) by $G = G_{PCGM}$); (3) combination of SGM and 2DPCA generative models as mentioned in Section 3.4; (4) principal component analysis (PCA)² generative model (instead of using the 2DPCA generative model); and (5) combination of SGM and PCA generative models.

The main differences between 2DPCA versus PCA are as follows. The covariance matrix in the case of 2DPCA can be computed directly from the image samples in 2D matrices rather than 1D vectors as in the case of PCA (Yang et al., 2004; Zhang and Zhou, 2005). The complexity for computing the covariance matrix using a 2DPCA-based appearance model is $\mathcal{O}(mn^2N)$, whereas the corresponding complexity using a PCA-based appearance model is $\mathcal{O}(m^2n^2N)$, when a set of N image samples, each of size $m \times n$ pixels, is used. Further, it may be pointed out that 2DPCA encodes

² The function `pcaApply` from toolbox Dollár 2015 has been used to calculate the PCA.

the relationship among neighboring rows in a given set of image samples (Zhang and Zhou, 2005). Such a relationship should have a positive effect on the tracking performance.

Table 2 shows the results on the seven sequences. Overall, the proposed scheme with SGM in conjunction with 2DPCA performs better than that by using SGM with PCA. In most sequences, the method of using SGM with 2DPCA or SGM with PCA performs better than that using only SGM. On a machine with 2.9 GHz CPU, the average tracking time per frame (over all the seven sequences without counting the time for object detection) for the proposed tracker with SGM and 2DPCA is 2.88 s whereas the corresponding time in the case of SGM and PCA is 2.90 s. Hence, this improvement in the performance of the proposed tracker is achieved without loss in speed.

4.4. Performance comparison

In this section, we evaluate the performance of the proposed algorithm with two online MOT methods in Yoon et al.(2015); Bae and Yoon (2014) using the seven challenging sequences described in Section 4.1. Table 3 shows the performance of these two methods (using the original source code) along with that of the proposed tracker in terms of the various CLEAR MOT metrics. In addition, the performance of the proposed scheme is compared with the reported results of state-of-the-art online MOT methods

Table 2
Performance of the proposed scheme using different generative models.

Sequence	Generative model	MOTA	MOTP	FNR	FPR	IDSW
<i>PETS09-S2L1</i>	SGM	89.08%	79.89%	5.11%	5.42%	17
	PCA	89.86%	79.97%	<u>5</u>	4.76%	16
	SGM + PCA	<u>90</u>	<u>80</u>	5.34%	4.30%	13
	2DPCA	89.81%	79.82%	5.40%	4.41%	20
	Proposed	92.13%	80.62%	3.19%	<u>4</u>	<u>14</u>
<i>PETS09-S2L2</i>	SGM	36.43%	71.19%	39.38%	26.31%	263
	PCA	45.69%	<u>71</u>	<u>35</u>	20.25%	218
	SGM + PCA	44.35%	71.54%	36.04%	21.30%	237
	2DPCA	<u>46</u>	71.77%	36.59%	19.33%	<u>221</u>
	Proposed	46.88%	71.66%	34.92%	<u>19</u>	258
<i>Soccer</i>	SGM	67.36%	70.28%	<u>16</u>	14.49%	45
	PCA	70.33%	70.64%	18.85%	<u>10</u>	<u>38</u>
	SGM + PCA	70.21%	70.99%	18.20%	11.03%	36
	2DPCA	<u>71</u>	70.73%	17.00%	10.66%	49
	Proposed	73.54%	<u>70</u>	16.20%	9.45%	<u>38</u>
<i>UCF-PL</i>	SGM	82.30%	71.84%	10.77%	6.27%	<u>16</u>
	PCA	82.14%	<u>71</u>	<u>10</u>	6.56%	21
	SGM + PCA	<u>83</u>	71.81%	10.64%	5.40%	<u>16</u>
	2DPCA	81.89%	71.75%	11.47%	5.90%	18
	Proposed	85.02%	71.89%	8.70%	<u>5</u>	15
<i>Town Center (Body)</i>	SGM	69.41%	73.82%	17.08%	12.81%	444
	PCA	<u>70</u>	73.83%	18.18%	11.08%	351
	SGM + PCA	69.83%	73.89%	19.29%	10.37%	320
	2DPCA	71.24%	74.02%	<u>18</u>	<u>10</u>	<u>337</u>
	Proposed	70.16%	<u>73</u>	19.35%	9.95%	342
<i>Town Center (Head)</i>	SGM	70.32%	<u>68</u>	14.96%	14.48%	164
	PCA	<u>72</u>	68.71%	<u>14</u>	<u>13</u>	163
	SGM + PCA	69.37%	68.78%	15.62%	14.77%	166
	2DPCA	70.43%	68.82%	15.06%	14.29%	158
	Proposed	74.54%	69.15%	13.02%	12.21%	158
<i>LISA10 Urban</i>	SGM	100.00%	82.67%	0.00%	0.00%	0
	PCA	100.00%	82.68%	0.00%	0.00%	0
	SGM + PCA	100.00%	82.68%	0.00%	0.00%	0
	2DPCA	100.00%	82.68%	0.00%	0.00%	0
	Proposed	100.00%	82.68%	0.00%	0.00%	0
<i>LISA10 Sunny</i>	SGM	97.22%	78.28%	0.78%	1.98%	0
	PCA	97.22%	78.28%	0.78%	1.98%	0
	SGM + PCA	97.22%	78.28%	0.78%	1.98%	0
	2DPCA	97.22%	78.28%	0.78%	1.98%	0
	Proposed	97.22%	78.28%	0.78%	1.98%	0
Average	SGM	76.51%	74.60%	13.10%	10.22%	-
	PCA	78.45%	74.72%	<u>12</u>	8.53%	-
	SGM + PCA	78.05%	<u>74</u>	13.24%	8.64%	-
	2DPCA	<u>78</u>	74.73%	13.04%	<u>8</u>	-
	Proposed	79.94%	74.87%	12.02%	7.87%	-

Note: The best and the second best results on each dataset are shown in boldface and underscored, respectively. The proposed method is SGM + 2DPCA.

(Benfold and Reid, 2011; Breitenstein et al., 2011; Gomez et al., 2012; Poiesi et al., 2013; Shu et al., 2012; Zhang et al., 2012; Zhou et al., 2014) using the sequences considered in these papers.

On the *PETS09-S2L1* and *PETS09-S2L2* sequences, the proposed scheme provides the second highest MOTA values. It also offers the highest and second highest MOTP values on the *PETS09-S2L1* and *PETS09-S2L2* sequences, respectively. This can be attributed to the proposed update mechanism, and the inactive or on-hold states of the tracker.

For the *Soccer* sequence, the proposed scheme performs better than the methods in Yoon et al. (2015); Bae and Yoon (2014) despite fast camera motion and the presence of similar objects in the scenes. For the *UCF-PL* sequence, the MOTA value of the proposed method is higher than that of the methods in Yoon et al. (2015); Shu et al. (2012); Bae and Yoon (2014), using the same detector as in Shu et al. (2012). On the other hand, the MOTP value of the proposed technique is close to that of Shu et al. (2012). In addition, the proposed method has lower values for FNR and FPR than the methods in Yoon et al. (2015); Shu et al. (2012); Bae and Yoon (2014) do.

For the *Town Center* dataset, the proposed scheme is first evaluated to track the full body of pedestrians. In this case, the proposed scheme yields the second highest MOTP, FNR and FPR values compared to the methods in Yoon et al. (2015); Bae and Yoon (2014); Benfold and Reid (2011); Zhang et al. (2012); Shu et al. (2012). Next, the proposed scheme is evaluated on tracking the heads of pedestrians from the same dataset. The head regions in this sequence are less occluded than the full body, although the head detector has higher FPR than the full-body detector. As shown in Table 3, the proposed method performs well against other approaches (Yoon et al., 2015; Bae and Yoon, 2014; Poiesi et al., 2013; Benfold and Reid, 2011) in terms of MOTA. For the *Urban* and *Sunny* sequences from *LISA 2010* dataset, the proposed scheme provides a better performance than that provided by the methods in Yoon et al. (2015); Bae and Yoon (2014) for tracking multiple vehicles on-road.

We note that the proposed scheme uses grayscale images as features, whereas the methods in Breitenstein et al. (2011); Shu et al. (2012); Zhang et al. (2012) are based on the color or gradient information of the targets. In addition, the proposed scheme does not require the detector confidence density or a gate function

Table 3
Performance measures of CLEAR MOT metrics.

Sequence	Method	MOTA	MOTP	FNR	FPR	IDSW	
PETS09-S2L1	Proposed	<u>92</u>	80.62%	3.19%	4.33%	<u>14</u>	
	Yoon et al. (2015)*	66.64%	57.46%	17.99%	15.14%	34	
	Bae and Yoon (2014)*	89.94%	<u>79</u>	<u>4</u>	<u>4</u>	23	
	Zhang et al. (2012)	93.27%	68.17%	-	-	19	
	Zhou et al. (2014)	87.21%	58.47%	-	-	-	
	Breitenstein et al. (2011)	79.70%	56.30%	-	-	-	
PETS09-S2L2	Proposed	<u>46</u>	<u>71</u>	34.92%	<u>19</u>	258	
	Yoon et al. (2015)*	26.85%	47.99%	51.27%	28.86%	<u>218</u>	
	Bae and Yoon (2014)*	45.98%	71.77%	<u>35</u>	19.06%	325	
	Zhang et al. (2012)	66.72%	58.21%	-	-	215	
	Soccer	Proposed	73.54%	70.77%	16.20%	9.45%	38
	Yoon et al. (2015)*	29.99%	53.77%	52.89%	26.19%	10	
UCF-PL	Bae and Yoon (2014)*	<u>54</u>	<u>69</u>	<u>35</u>	<u>12</u>	<u>24</u>	
	Proposed	85.02%	71.89%	8.70%	5.65%	15	
	Yoon et al. (2015)*	29.50%	45.33%	38.04%	33.95%	15	
Town Center (Body)	Bae and Yoon (2014)*	<u>82</u>	<u>73</u>	<u>10</u>	<u>6</u>	15	
	Shu et al. (2012)	79.30%	74.10%	18.30%	8.70%	-	
	Proposed	70.16%	<u>73</u>	<u>19</u>	<u>9</u>	342	
	Yoon et al. (2015)*	62.93%	48.66%	20.00%	17.14%	330	
	Bae and Yoon (2014)*	79.07%	73.46%	11.19%	9.44%	307	
	Benfold and Reid (2011)	61.30%	80.30%	21.00%	18.00%	-	
Town Center (Head)	Zhang et al. (2012)	<u>73</u>	68.75%	-	-	421	
	Shu et al. (2012)	72.90%	71.30%	-	-	-	
	Proposed	74.54%	69.15%	13.02%	<u>12</u>	<u>158</u>	
	Yoon et al. (2015)*	<u>73</u>	70.16%	17.23%	9.49%	126	
	Bae and Yoon (2014)*	70.65%	<u>69</u>	<u>16</u>	13.07%	320	
	Poiesi et al. (2013)	54.60%	63.70%	23.80%	21.70%	285	
LISA10 Urban	Benfold and Reid (2011)	45.40%	50.80%	29.00%	26.20%	-	
	Proposed	100.00%	82.68%	0.00%	0.00%	0	
	Yoon et al. (2015)*	<u>99</u>	81.98%	<u>0</u>	<u>0</u>	0	
	Bae and Yoon (2014)*	98.33%	<u>82</u>	1.67%	0.00%	0	
LISA10 Sunny	Proposed	97.22%	78.28%	0.78%	1.98%	0	
	Yoon et al. (2015)*	92.89%	77.20%	6.89%	0.24%	0	
	Bae and Yoon (2014)*	<u>97</u>	<u>77</u>	<u>2</u>	<u>0</u>	0	

Note: * denotes the results obtained by utilizing the code provided by the authors of the paper, where the detection results and GT annotations that have been used with the proposed scheme are used. The best and the second best results on each dataset are represented in boldface and underscored, respectively.

in the data association step as in Breitenstein et al. (2009, 2011), where the gate function provides higher weight for detections located in the direction of motion of the target.

5. Conclusion

In this paper, we have presented a robust collaborative model that enhances the interaction between a pre-trained object detector and a number of single-object online trackers in the particle filter framework. The proposed scheme is based on incorporating the associated detections with the motion model, in addition to the likelihood function providing different weights for the propagated and the newly created particles sampled from the associated detections, providing a reduction on the effect of the detector errors on the tracking process. We have exploited sparse representation and 2DPCA to construct diverse features that maximize the appearance variation among the trackers. Furthermore, we have presented a conservative sample selection scheme to update the appearance model of every tracker. Experimental results on benchmark datasets have shown that the proposed scheme outperforms state-of-the-art multi-object tracking methods in most of the cases.

Acknowledgments

The authors would like to thank Dr. Y. Wu for his helpful discussions and suggestions. They would also like to thank all the authors that made their codes available for comparison of the proposed algorithm with theirs and the anonymous reviewers for their constructive comments and suggestions. M.A. Naiel would

like to acknowledge the support from Concordia University to conduct this research. This work is supported by research grants from the Natural Sciences and Engineering Research Council (NSERC) of Canada and the Regroupement Stratégique en Microsystèmes du Québec (ReSMiQ) awarded to M.O. Ahmad and M.N.S. Swamy. J. Lim is supported by the National Research Foundation (NRF) of Korea grant #2014R1A1A2058501. M.-H. Yang is supported in part by the National Science Foundation (NSF) CAREER grant #1149783 and a gift from Panasonic.

References

- Andriyenko, A., Schindler, K., 2011. Multi-target tracking by continuous energy minimization. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1265–1272.
- Andriyenko, A., Schindler, K., Roth, S., 2012. Discrete-continuous optimization for multi-target tracking. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1926–1933.
- Bae, S.H., Yoon, K.J., 2014. Robust online multi-object tracking based on tracklet confidence and online discriminative appearance learning. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1218–1225.
- Benfold, B., Reid, I., 2011. Stable multi-target tracking in real-time surveillance video. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 3457–3464.
- Bernardin, K., Stiefelwagen, R., 2008. Evaluating multiple object tracking performance: the CLEAR MOT metrics. J. Image Video Process. 2008, 1–10.
- Breitenstein, M.D., Reichlin, F., Leibe, B., Koller-Meier, E., Gool, L.V., 2009. Robust tracking-by-detection using a detector confidence particle filter. In: Proc. IEEE International Conference on Computer Vision, pp. 1515–1522.
- Breitenstein, M.D., Reichlin, F., Leibe, B., Koller-Meier, E., Gool, L.V., 2011. Online multi-person tracking-by-detection from a single, uncalibrated camera. IEEE Trans. Pattern Anal. Mach. Intell. 33 (9), 1820–1833.

- Brendel, W., Amer, M.R., Todorovic, S., 2011. Multiobject tracking as maximum weight independent set. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1273–1280.
- Butt, A.A., Collins, R.T., 2012. Multiple target tracking using frame triplets. In: Proc. Asian Conference on Computer Vision, pp. 163–176.
- Dalal, N., Triggs, B., 2005. Histograms of oriented gradients for human detection. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 886–893.
- Dollár, P., Piotr's Image and Video Matlab Toolbox (PMT). <https://pdollar.github.io/toolbox/>. Last retrieved, January 18, 2015.
- Dollár, P., Appel, R., Belongie, S., Perona, P., 2014. Fast feature pyramids for object detection. *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (8), 1532–1545.
- Dollár, P., Belongie, S., Perona, P., 2010. The fastest pedestrian detector in the west. In: Proc. British Machine Vision Conference.
- Duffner, S., Odobez, J., 2013. Track creation and deletion framework for long-term online multi-face tracking. *IEEE Trans. Image Process.* 22 (1), 272–285.
- Eiselein, V., Arp, D., Patzold, M., Sikora, T., 2012. Real-time multi-human tracking using a probability hypothesis density filter and multiple detectors. In: Proc. IEEE International Conference on Advanced Video and Signal-Based Surveillance, pp. 325–330.
- Everingham, M., Gool, L.V., Williams, C.K., Winn, J., Zisserman, A., 2010. The Pascal visual object classes (VOC) challenge. *Int. J. Comput. Vis.* 88 (2), 303–338.
- Ferryman, J., 2009. In: Proc. IEEE Workshop Performance Evaluation of Tracking and Surveillance.
- Gomez, D.G., Lerasle, F., Peña, A.M.L., 2012. State-driven particle filter for multi-person tracking. In: Proc. Advanced Concepts for Intelligent Vision Systems. In: Lecture Notes in Comput. Sci., 7517, pp. 467–478.
- Gordon, N.J., Salmond, D.J., Smith, A.F.M., 1993. Novel approach to nonlinear and non-Gaussian Bayesian state estimation. *IEE Proc. F (Radar Signal Process.)* 140 (2), 107–113.
- Han, H., Ding, Y.-S., Hao, K.-R., Liang, X., 2011. An evolutionary particle filter with the immune genetic algorithm for intelligent video target tracking. *Comput. Math. Applicat.* 62 (7), 2685–2695.
- Huang, Y., Djuric, P., 2004. A hybrid importance function for particle filtering. *IEEE Signal Process. Lett.* 11 (3), 404–406.
- Izadinia, H., Saleemi, I., Li, W., Shah, M., 2012. (MP)²T: multiple people multiple parts tracker. In: Proc. European Conference on Computer Vision, pp. 100–114.
- Jin, Y., Mokhtarian, F., 2007. Variational particle filter for multi-object tracking. In: Proc. IEEE International Conference on Computer Vision, pp. 1–8.
- Jinxia, Y., Yongli, T., Jingmin, X., Qian, Z., 2012. Research on particle filter based on an improved hybrid proposal distribution with adaptive parameter optimization. In: Proc. International Conference on Intelligent Computation Technology and Automation, pp. 406–409.
- Kuo, C.-H., Huang, C., Nevatia, R., 2010. Multi-target tracking by on-line learned discriminative appearance models. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 685–692.
- Leal-Taixé, L., Milan, A., Reid, I., Roth, S., Schindler, K., 2015. MOTChallenge 2015: towards a benchmark for multi-target tracking, 2015. [arXiv:1504.01942](https://arxiv.org/abs/1504.01942).
- Leal-Taixé, L., Pons-Moll, G., Rosenhahn, B., 2011. Everybody needs somebody: modeling social and grouping behavior on a linear programming multiple people tracker. In: Proc. IEEE International Conference on Computer Vision Workshops, pp. 120–127.
- Maggio, E., Piccardi, E., Regazzoni, C., Cavallaro, A., 2007. Particle PHD filtering for multi-target visual tracking. In: Proc. International Conference on Acoustics, Speech, and Signal Process., 1, pp. 1–1101–1–1104.
- Maggio, E., Taj, M., Cavallaro, A., 2008. Efficient multitarget visual tracking using random finite sets. *IEEE Trans. Circuits Syst. Video Technol.* 18 (8), 1016–1027.
- Mahler, R., 2003. Multitarget Bayes filtering via first-order multitarget moments. *IEEE Trans. Aerosp. Electron. Syst.* 39 (4), 1152–1178.
- Naiel, M.A., Ahmad, M.O., Swamy, M.N.S., 2014. Vehicle detection using TD2DHOG features. In: Proc. IEEE New Circuits and System Conference, pp. 389–392.
- Okuma, K., Taleghani, A., Freitas, N.D., Little, J.J., Lowe, D.G., 2004. A boosted particle filter: multitarget detection and tracking. In: Proc. European Conference on Computer Vision, pp. 28–39.
- Poiesi, F., Mazzon, R., Cavallaro, A., 2013. Multi-target tracking on confidence maps: an application to people tracking. *Comput. Vis. Image Under.* 117 (10), 1257–1272.
- Rui, Y., Chen, Y., 2001. Better proposal distributions: Object tracking using unscented particle filter. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 786–793.
- Salti, S., Cavallaro, A., Stefano, L.D., 2012. Adaptive appearance modeling for video tracking: Survey and evaluation. *IEEE Trans. Image Process.* 21 (10), 4334–4348.
- Santhoshkumar, S., Karthikeyan, S., Manjunath, B.S., 2013. Robust multiple object tracking by detection with interacting Markov chain Monte Carlo. In: Proc. IEEE International Conference on Image Process.
- Schumann, A., Bäuml, M., Stiefelwagen, R., 2013. Person tracking-by-detection with efficient selection of part-detectors. In: Proc. IEEE International Conference on Advanced Video and Signal-Based Surveillance, pp. 43–50.
- Shitrit, H.B., Berclaz, J., Fleuret, F., Fua, P., 2011. Tracking multiple people under global appearance constraints. In: Proc. IEEE International Conference on Computer Vision, pp. 137–144.
- Shu, G., Dehghan, A., Oreifej, O., Hand, E., Shah, M., 2012. Part-based multiple-person tracking with partial occlusion handling. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1815–1821.
- Sivaraman, S., Trivedi, M., 2010. A general active-learning framework for on-road vehicle recognition and tracking. *IEEE Trans. Intell. Transp. Syst.* 11 (2), 267–276.
- Vermaak, J., Doucet, A., Perez, P., 2003. Maintaining multimodality through mixture tracking. In: Proc. IEEE International Conference on Computer Vision, pp. 1110–1116.
- Wang, T., Gu, I., Shi, P., 2007. Object tracking using incremental 2D-PCA learning and ML estimation. In: Proc. International Conference on Acoustics, Speech, and Signal Process., 1, pp. 933–936.
- Wright, J., Yang, A.Y., Ganesh, A., Sastry, S.S., Ma, Y., 2009. Robust face recognition via sparse representation. *IEEE Trans. Pattern Anal. Mach. Intell.* 31 (2), 210–227.
- Wu, Y., Lim, J., Yang, M.-H., 2013. Online object tracking: a benchmark. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 2411–2418.
- Wu, Y., Tong, X., Zhang, Y., Lu, H., 2008. Boosted interactively distributed particle filter for automatic multi-object tracking. In: Proc. IEEE International Conference on Image Process., pp. 1844–1847.
- Wu, Z., Thangali, A., Sclaroff, S., Betke, M., 2012. Coupling detection and data association for multiple object tracking. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1948–1955.
- Yang, B., Nevatia, R., 2012. Multi-target tracking by online learning of non-linear motion patterns and robust appearance models. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1918–1925.
- Yang, J., Zhang, D., Frangi, A.F., Yang, J.-Y., 2004. Two dimensional PCA: a new approach to appearance-based face representation and recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 26 (1), 131–137.
- Yang, M., Lv, F., Xu, W., Gong, Y., 2009. Detection driven adaptive multi-cue integration for multiple human tracking. In: Proc. IEEE International Conference on Computer Vision, pp. 1554–1561.
- Yoon, J.H., Yang, M.H., Lim, J., Yoon, K.J., 2015. Bayesian multi-object tracking using motion context from multiple objects. In: Proc. IEEE Winter Conference on Applications of Computer Vision, pp. 33–40.
- Zamir, A.R., Dehghan, A., Shah, M., 2012. GMCP-tracker: global multi-object tracking using generalized minimum clique graphs. In: Proc. European Conference on Computer Vision, pp. 343–356.
- Zhang, D., Zhou, Z.-H., 2005. (2D)²PCA: Two-directional two-dimensional PCA for efficient face representation and recognition. *Neurocomputing* 69 (1–3), 224–231.
- Zhang, J., Presti, L.L., Sclaroff, S., 2012. Online multi-person tracking by tracker hierarchy. In: Proc. IEEE International Conference on Advanced Video and Signal-Based Surveillance, pp. 379–385.
- Zhong, W., Lu, H., Yang, M.-H., 2012. Robust object tracking via sparsity-based collaborative model. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition, pp. 1838–1845.
- Zhou, X., Li, Y., He, B., Bai, T., 2014. GM-PHD-based multi-target visual tracking using entropy distribution and game theory. *IEEE Trans. Ind. Informat.* 10 (2), 1064–1076.