# ONLINE MULTI-PERSON TRACKING VIA ROBUST COLLABORATIVE MODEL

*Mohamed A. Naiel, M. Omair Ahmad,* FIEEE,
*and M.N.S. Swamy,* FIEEE
Dept. of Electrical and Computer Engineering
Concordia University
Montreal, Quebec, Canada H3G 1M8
{m_naiel, omair, swamy}@ece.concordia.ca

*Yi Wu, and Ming-Hsuan Yang,* SMIEEE
Dept. of Electrical Engineering and Computer Science
University of California at Merced
Merced, CA 95344, USA
{ywu29, mhyang}@ucmerced.edu

## ABSTRACT

The past decade has witnessed significant progress in object detection and tracking in videos. In this paper, we present a model for collaboration between a pre-trained object detector and multiple single object trackers in the particle filter tracking framework. For each frame, we construct an association between the trackers and the detections, and when a tracker is successfully associated to a detection, we treat this detection as the key-sample for this tracker. We present a dual motion model that incorporates the associated detections with the object dynamics. Then, a likelihood function provides different weights for the propagated and the newly created particles, reducing the effect of false positives and missed detections in the tracking process. In addition, we use generative and discriminative appearance models to maximize the appearance variation among the targets. The performance of the proposed algorithm compares favorably with that of the state-of-the-art approaches on three public sequences.

## 1. INTRODUCTION

Multi-object tracking in videos is one of the challenging problems in computer vision. It has applications in automatic analysis of surveillance videos, behavior analysis, road traffic management etc. The advancement in object detection promotes collaboration between the processes of detection and tracking in the tracking-by-detection scheme [1].

In the multi-object tracking field, offline approaches based on the global optimization of all object trajectories usually achieve better results than online counterparts [2, 3, 4] do. However, such offline methods are time-delayed. On the other hand, the online methods [1, 5, 6] have been developed within the tracking-by-detection framework, and have applied data association between detections and trackers in an online manner. These approaches perform well for real-time applications.

Online multi-object tracking can be achieved by using joint state-space model for multi-targets [7, 8, 9, 10]. For instance, a mixture particle filter has been proposed in [8] to obtain the posterior probability by using 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* (PHD) filter [11] has been incorporated for visual multi-target tracking in [9, 12]. The particle PHD filter has a linear complexity with respect to the number of targets. However, it does not maintain the target identity, and as a result, requires an online clustering method to detect the peaks of the particle weights and applies data association for each cluster.

The degeneracy problem of the particle filter [13] has been addressed by several researchers [6, 14, 15, 16]. To overcome this problem, it is required to enhance the proposal distribution and/or the re-sampling step, in order to increase the number of effective particles and the diversity of the particles. For instance, Rui and Chen [15] have used the unscented Kalman Filter (UKF) for generating the proposal distribution of the particle filter, and the resulting scheme is known as the unscented particle filter (UPF). An immune genetic algorithm for visual tracking was subsequently introduced [17]. Recently, in [6] the Metropolis Hastings algorithm has been used to sample particles from associated detections in the tracking-by-detection framework. However, these methods do not exploit the collaboration between detectors and trackers [15, 17], or consider the effect of false positive detections on the trackers [6].

In this paper, we propose an online multi-object tracking algorithm by using a robust collaborative model for interaction between multiple single-object trackers and a pre-trained 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. We propose a dual motion model that incorporates the associated detections with the object dynamics. Furthermore, the likelihood function provides different weights for the propagated and the newly created particles
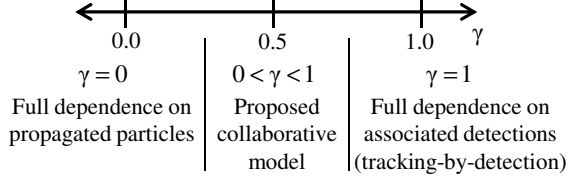
Fig. 1. Effect of changing the collaborative factor $\gamma$.

sampled from the associated detections, resulting a reduction in the effect of false positives and missed detections in the tracking process. The appearance model of the trackers is based on the sparsity-based representation [18, 19], and the two dimensional principal component analysis (2DPCA) [20] to maximize the appearance variation among targets.

## 2. PROPOSED ALGORITHM

The proposed multi-object tracking system consists of three main components: pre-trained object detector, multiple single-object trackers and data association module. Specifically, the *object detector* is applied on every frame and supports the data association module with a set of detections at time $t$, $\mathcal{D}^t$. We use FPDW [21] as the baseline pedestrian detector. The *object tracker* adopts a hybrid motion model, and a particle filter with a robust collaborative model is used to find the best estimate of the target location. Further, it constructs the target appearance model, which consists of sparsity-based generative model (SGM) [19] with local features, 2DPCA-based generative model (PGM) with holistic features and sparsity-based discriminative classifier (SDC) with holistic features. Finally, the *data association* module is used to construct the similarity matrix $S$ to find the association between existing trackers, $b_t \in \mathcal{B}_e^t$, and detections, $d_t \in \mathcal{D}^t$, at time $t$. Furthermore, it controls the initialization and termination status of the trackers, and supports the tracker with key-samples from the target trajectory.

### 2.1. Tracking Approach

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 $\mathbf{x}_t^i$. The state $\mathbf{x}$ consists of translation $(x, y)$, average velocity $(v_x, v_y)$, scale $\hat{s}$, rotation angle $\theta$, aspect ratio $\eta$, and skew direction $\phi$. Since in the current frame, 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, it can be shown that given the observation $\mathbf{z}_t$ at time $t$, the weight is proportional to the likelihood, namely, $\mathbf{w}_t^i \propto p(\mathbf{z}_t|\mathbf{x}_t^i)$. Let $\mathcal{I}_{b_t}$ be a gate function representing the state of the tracker $b_t$ associated with the detection $d_t$ at time $t$.

**Motion Model:** In our approach, we adopt a hybrid motion model that depends on the first-order Markov motion model
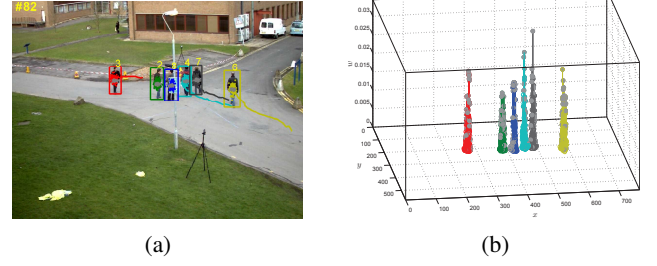


(a)                            (b)

Fig. 2. 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).

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 particles of the tracker. 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 $X^{b_t, d_t}$ as

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

where $Q$ and $P$ are the Gaussian noise vectors, $F = I + B$, $I$ being an identity matrix of size $8 \times 8$, $B = \begin{bmatrix} C & 0_{2,4} \\ 0_{6,4} & 0_{6,4} \end{bmatrix}$, $C = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$, and $0_{n_1, n_2}$ representing a zero matrix of size $(n_1 \times n_2)$.

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

$$p(\mathbf{z}_t|\mathbf{x}_t^i) = \pi^i H(\mathbf{x}_t^i) \quad (2)$$

where,

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

and $\gamma \in [0, 1]$ is the collaborative factor. In (2), the particles from the associated detections and previous propagated particles are weighted differently. Figure 1 shows the effect of changing the collaborative factor value. Figures 2 (a) and (b) show 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 weighting term, $\pi^i$, 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 positive detections, so the tracker should not depend completely on the detector. From our experiments, we find that value of $\gamma$ between 0.5 to 0.6 gives the best results. If the detector suffers from missing detections, i.e., $\mathcal{I}_{b_t} = 0$, the likelihood function in (2) will only depend on the previously propagated particles $\mathbf{x}_t^i \in X^{b_t}$, which represents the bootstrap particle filter [13]. Our collaborative model is based on the hybrid motion model that incorporates associated detections with the dynamic motion model. On the contrary, the motion model adopted in [1] depends only on propagated particles, and the likelihood function depends on tracker appearance model and the detector confidence density. In [8] the collaborative model only exists in the proposal distribution and the likelihood is without weighting collaborative factor.

## 2.2. Appearance Model

In our approach SGM and SDC are used in a way different from that in [19]. First, we do not use the collaboration between SGM and SDC [19], whereas we use the SGM with the PGM to compute the similarity matrix of the data association module to take advantage of the occlusion handling scheme (6), and the modified SDC model is used to obtain the likelihood of the particle filter (2). Therefore, the resulting tracker is more efficient. Second, our SDC confidence measure depends on the sparsity concentration index (SCI), given by (4). Finally, we update the SDC tracker with high confidence keysamples.

**Sparsity-based Discriminative Classifier:** The SDC is used to construct the target discriminative appearance model, in order to evaluate the confidence score in (2). Every SDC tracker is initialized by using $N_p$ positive training samples taken 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, the training samples are transformed to a fixed size $m \times n$, and then vectorized, normalized and stacked together to form a matrix $A \in \mathbb{R}^{r \times N^t}$, where $r = m \times n$ and $N^t = N_p + N_n$. Let the measurement corresponding to the candidate location $\mathbf{x}_t^i$ be denoted by $\mathbf{z}_t^i \in \mathbb{R}^r$. Then, we obtain the sparse coefficients $\alpha^i$ for the $i^{th}$ candidate by solving the optimization problem, $\min_{\alpha^i} \left\| \mathbf{z}_t^i - A\alpha^i \right\|_2^2 + \lambda \left\| \alpha^i \right\|_1$. Then, we can obtain the classifier confidence score as

$$H(\mathbf{x}_t^i) = \exp\left(-\frac{(\varepsilon_+^i - \varepsilon_-^i)}{\sigma}\right)\Omega(\alpha^i) \quad (3)$$

where $\varepsilon_+^i = \left\| \mathbf{z}_t^i - A_+\alpha_+^i \right\|_2^2$, $\varepsilon_-^i = \left\| \mathbf{z}_t^i - A_-\alpha_-^i \right\|_2^2$, $\sigma$ is used to adjust the confidence measure, and $\Omega(\alpha^i)$ represents the sparsity concentration index (SCI) [18] defined as

$$\Omega(\alpha^i) = \frac{J \cdot \max_j \|\delta_j'(\alpha^i)\|_1 / \|\alpha^i\|_1 - 1}{J - 1} \in [0, 1] \quad (4)$$

$\delta_j'$ being a function that selects the coefficients corresponding to the $j^{th}$ class and suppresses the rest, and $J$ being the number of classes ($J = 2$ in our case). The SCI checks the validity of the candidate so that it can be represented by a linear combination of the training samples in one class.

**Sparsity-based Generative Model:** We use the SGM for contributing to the similarity function of the data association module in (6). The SGM is concerned with representing the appearance of the positive class of the tracker by using $M$ local patches of the initial object or candidate location $c$, where each candidate is represented by a sparse histogram feature vector $\rho$. In order to handle occlusion, the patch reconstruction error is used to suppress the coefficients of the occluded patches. Let $\psi_j$ be the non occlusion indicator; the final histogram can be represented by $\varphi = \rho \odot \psi$, where $\odot$ denotes the element-wise multiplication. The resulting histogram, $\varphi$, taking the spatial representation into consideration, can handle occlusion effectively. The generative model similarity, $G_{SGM}$, between the candidate $\varphi_c$ and the model $\varphi$ is measured by using the histogram intersection kernel. We refer the reader to [19] for more details about SGM.

**2DPCA-based Generative Model:** We utilize the 2DPCA [20] as a generative model to compute the data association similarity in (6). Unlike SGM, PGM is based on holistic representation for the object. For each tracker $b_t$, we use $N$ positive samples, $\{Y_j\}_{j=1}^N$ each of size $m \times n$. Then, we evaluate the optimal orthonormal matrix $V$ that maximizes the total scatter in the learned subspace. For each candidate location, $Y^c$, the similarity between the testing and the training features can be computed as

$$G_{PGM} = \exp(-\varepsilon_{PGM}/\hat{\sigma}^2) \quad (5)$$

where $\varepsilon_{PGM}$ represents the reconstruction error between the candidate sample and the training example with the minimum $l_2$-norm to the test example in the 2DPCA feature space.

## 2.3. Data Association

The goal of the data association is to find association between existing trackers and detections every frame. Furthermore, new trackers may be created by using un-associated detections. The similarity matrix, $S$, is used to measure the relation between the trackers, $b_t \in \mathcal{B}_e^t$, and the detections, $d_t \in \mathcal{D}^t$. The similarity between $b_t$ and $d_t$ is defined as

$$S(b_t, d_t) = G(b_t, d_t)O(b_t, d_t) \quad (6)$$

where $G(b_t, d_t) = G_{SGM}(b_t, d_t) + G_{PGM}(b_t, d_t)$ measures 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 confused detections. The overlap ratio is based on the PASCAL VOC criterion [22], which is defined as the area of intersection divided by the area of union. The association is performed using the Hungarian algorithm to match a detection to a tracker, similar to the approach adopted in [1, 5], and is carried out online.
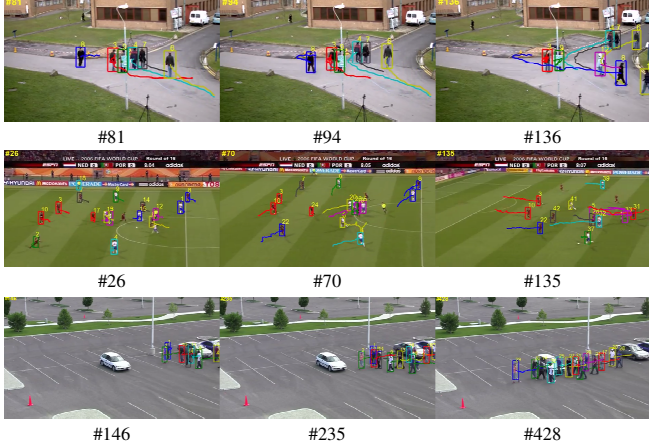
**Fig. 3**. Sample tracking results for three sequences. 1st row: PETS09 S2L1 view1; 2nd row: Soccer; 3rd row: UCF-PL.

**Table 1**. Performance measures of CLEAR MOT metrics.

| Sequence | Method | MOTA | MOTP | FNR | FPR | IDSW |
|---|---|---|---|---|---|---|
| PETS09 S2L1 | proposed + FPDW | 82.17% | 72.11% | **8.43%** | 9.25% | **4** |
| | Zhang *et al.* [27] | **93.27%** | 68.17% | - | - | 19 |
| | Breitenstein *et al.* [28] | 79.70% | 56.30% | - | - | - |
| | Gerónimo *et al.* [29] | 51.1% | **75.0%** | 45.2% | - | 0 |
| Soccer | proposed + FPDW | 71.46% | 70.80% | 16.81% | 12.16% | 7 |
| UCF-PL | proposed + FPDW | 79.26% | 73.91% | 13.11% | 7.60% | 13 |
| | proposed + [5][1] | **84.54%** | 73.24% | **8.58%** | **6.86%** | **4** |
| | Shu *et al.* [5] | 79.30% | **74.10%** | 18.30% | 8.70% | - |

## 3. EXPERIMENTAL RESULTS

We evaluate the proposed multi-person tracking algorithm on three challenging datasets: the PETS09 S2L1 view1 [23], UCF Parking Lot (UCF-PL) dataset [5], and Soccer dataset [24]. In all datasets we use comparable parameter settings without extensive parameter tuning, where we use $\gamma = 0.54$ as the default value. In order to measure the performance of our technique and compare it to that of several state-of-the-art methods, we use the CLEAR MOT metrics [25]: multiple object tracking accuracy (MOTA), multiple object tracking precision (MOTP), false negative rate (FNR), false positive rate (FPR), and identity switches (IDSW). We use an overlap threshold of 0.5 to compute the evaluation metrics in all experiments. For PETS09, we use the ground truth data available in [26], while for Soccer and UCF-PL we use the ground truth data provided by the authors. The Matlab implementation of the proposed algorithm on a PC with 2.9 GHz CPU takes on average 1.997 seconds per frame[2] to track 16 players in the Soccer sequence of frame size $960 \times 544$.

Figure 3 illustrates the tracking results on the testing sequences[3] and Table 1 shows the performance of our algorithm compared to that of some recent online methods. On PETS09

---

[1]Using the detection results for the part based detector proposed in [5].

[2]It does not include the time for object detection.

[3]Supplementary material shows several qualitative tracking results https://www.youtube.com/watch?v=1rtz3MVrX2I
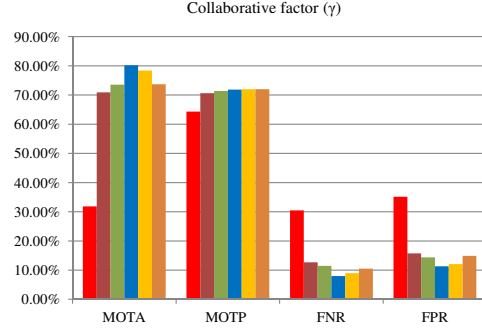


**Fig. 4**. Performance of the proposed approach on PETS09 S2L1 as the value of the collaborative factor is varied.

S2L1 view1, we achieve the second highest MOTA performance. For the Soccer sequence, we achieve a good performance considering that the sequence has high similarity among targets of the same team, and high missing detections. For the UCF-PL sequence we achieve a better MOTA performance than that offered by the online method [5], when using the same detection results.

**Collaborative factor effect:** To measure the effect of the proposed collaborative model, we vary the value of the collaborative factor $\gamma$ in the interval $[0, 1]$ in steps of $0.2$. Figure 4 illustrates the performance of the proposed approach for different values of $\gamma$. When $\gamma = 0$, the likelihood function of the particle filter is based completely on the propagated particles, and hence, the tracker suffers from the degeneracy problem, whereas when $\gamma = 1$, the likelihood function is based on the associated detections, and the tracker suffers from false positives, and missed detections. The highest performance is achieved when $\gamma = 0.6$; this is due to the fact that the tracker achieves a balance between the two particle sets in this case.

## 4. CONCLUSIONS

We have presented a robust collaborative model that can enhance the interaction between a pre-trained object detector and multiple single-object trackers in a particle filter framework. The proposed algorithm is based on incorporating the associated detections with the motion model of the particle filter, 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 used sparsity-based representation and the 2DPCA to construct discriminative features that maximize the appearance variation among the trackers. The proposed algorithm has been evaluated on three public sequences and the performance compares favorably with that of the state-of-the-art approaches.

# References

[1] M. D. Breitenstein, F. Reichlin, B. Leibe, E. Koller-Meier, and L. V. Gool, "Robust tracking-by-detection using a detector confidence particle filter," in *Proc. ICCV*, 2009, pp. 1515–1522.

[2] A. R. Zamir, A. Dehghan, and M. Shah, "GMCP-tracker: Global multi-object tracking using generalized minimum clique graphs," in *Proc. ECCV*, 2012, pp. 343–356.

[3] W. Brendel, M. Amer, and S. Todorovic., "Multiobject tracking as maximum weight independent set," in *Proc. CVPR*, 2011, pp. 1273–1280.

[4] C.-H. Kuo, C. Huang, and R. Nevatia, "Multi-target tracking by on-line learned discriminative appearance models," in *Proc. CVPR*, 2010, pp. 685–692.

[5] G. Shu, A. Dehghan, O. Oreifej, E. Hand, and M. Shah, "Part-based multiple-person tracking with partial occlusion handling," in *Proc. CVPR*, 2012, pp. 1815–1821.

[6] S. Santhoshkumar, S. Karthikeyan, and B. S. Manjunath, "Robust multiple object tracking by detection with interacting markov chain monte carlo," in *Proc. ICIP*, 2013.

[7] J. Vermaak, A. Doucet, and P. Perez, "Maintaining multi-modality through mixture tracking," in *Proc. ICCV*, 2003, pp. 1110–1116.

[8] K. Okuma, A. Taleghani, N. D. Freitas, J. J. Little, and D. G. Lowe, "A boosted particle filter: Multitarget detection and tracking," in *Proc. ECCV*, 2004, pp. 28–39.

[9] E. Maggio, M. Taj, and A. Cavallaro, "Efficient multitarget visual tracking using random finite sets," *IEEE Trans. on CSVT*, vol. 18, no. 8, pp. 1016–1027, 2008.

[10] V. Eiselein, D. Arp, M. Patzold, and T. Sikora, "Real-time multi-human tracking using a probability hypothesis density filter and multiple detectors," in *Proc. IEEE Int. Conf. on AVSS*, 2012, pp. 325–330.

[11] R. Mahler, "Multitarget bayes filtering via first-order multi-target moments," *IEEE Trans. on Aerosp. and Electron. Syst.*, vol. 39, no. 4, pp. 1152–1178, 2003.

[12] E. Maggio, E. Piccardo, C. Regazzoni, and A. Cavallaro, "Particle phd filtering for multi-target visual tracking," in *Proc. ICASSP*, vol. 1, 2007, pp. I–1101–I–1104.

[13] N. J. Gordon, D. J. Salmond, and A. F. M. Smith, "Novel approach to nonlinear and non-gaussian bayesian state estimation," *IEE Proc. F (Radar and Signal Process.)*, vol. 140, no. 2, pp. 107–113, 1993.

[14] Y. Jinxia, T. Yongli, X. Jingmin, and Z. Qian, "Research on particle filter based on an improved hybrid proposal distribution with adaptive parameter optimization," in *Proc. Int. Conf. on Intell. Computation Technol. and Automation*, 2012, pp. 406–409.

[15] Y. Rui and Y. Chen, "Better proposal distributions: Object tracking using unscented particle filter," in *Proc. CVPR*, 2001, pp. 786–793.

[16] Y. Huang and P. Djuric, "A hybrid importance function for particle filtering," *IEEE Signal Process. Letters*, vol. 11, no. 3, pp. 404–406, 2004.

[17] H. Han, Y.-S. Ding, K.-R. Hao, and X. Liang, "An evolutionary particle filter with the immune genetic algorithm for intelligent video target tracking," *Comput. and Mathematics with Applicat.*, vol. 62, no. 7, pp. 2685–2695, 2011.

[18] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *IEEE Trans. on PAMI*, vol. 31, no. 2, pp. 210–227, 2009.

[19] W. Zhong, H. Lu, and M.-H. Yang, "Robust object tracking via sparsity-based collaborative model," in *Proc. CVPR*, 2012, pp. 1838–1845.

[20] J. Yang, D. Zhang, A. F. Frangi, and J.-Y. Yang, "Two dimensional pca: A new approach to appearance-based face representation and recognition," *IEEE Trans. on PAMI*, vol. 26, no. 1, pp. 131–137, 2004.

[21] P. Dollár, S. Belongie, and P. Perona, "The fastest pedestrian detector in the west," in *Proc. British Machine Vis. Conf.*, 2010.

[22] M. Everingham, L. V. Gool, C. K. Williams, J.Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *Int. J. Comput. Vis.*, vol. 88, no. 2, pp. 303–338, 2010.

[23] J. Ferryman, in *Proc. IEEE Workshop Performance Evaluation of Tracking and Surveillance*, 2009.

[24] Y. Wu, X. Tong, Y. Zhang, and H. Lu, "Boosted interactively distributed particle filter for automatic multi-object tracking," in *Proc. ICIP*, 2008, pp. 1844–1847.

[25] K. Bernardin and R. Stiefelhagen, "Evaluating multiple object tracking performance: the clear mot metrics," *J. on Image and Video Process.*, pp. 1–10, 2008.

[26] http://www.gris.informatik.tu-darmstadt.de/~aandriye/data.html, last retrieved Feb. 13th, 2014.

[27] J. Zhang, L. L. Presti, and S. Sclaroff, "Online multi-person tracking by tracker hierarchy," in *Proc. IEEE Int. Conf. on AVSS*, 2012, pp. 379–385.

[28] M. D. Breitenstein, F. Reichlin, B. Leibe, E. Koller-Meier, and L. V. Gool, "Online multiperson tracking-by-detection from a single, uncalibrated camera," *IEEE Trans. on PAMI*, vol. 33, no. 9, pp. 1820–1833, 2011.

[29] D. Gernimo, F. Lerasle, and A. M. López, "State-driven particle filter for multi-person tracking," in *Proc. Advanced Concepts for Intell. Vis. Syst.*, ser. Lecture Notes in Comput. Sci., vol. 7517, 2012, pp. 467–478.