

Multi-Tracker Partition Fusion

ObaidUllah Khalid, Juan C. SanMiguel, and Andrea Cavallaro

Abstract—We propose a decision-level approach to fuse the output of multiple trackers based on their estimated individual performance. The proposed approach is divided into three steps. First, we group trackers into clusters based on the spatio-temporal pair-wise correlation of their short-term trajectories. Then, we evaluate performance based on reverse-time analysis with an adaptive reference frame and define the cluster with trackers that appear to be successfully following the target as the *on-target* cluster. Finally, the state estimations produced by trackers in the *on-target* cluster are fused to obtain the target state. The proposed fusion approach uses standard tracker outputs and can therefore combine various types of trackers. We tested the proposed approach with several combinations of state-of-the-art trackers, and also compared it with individual trackers and other fusion approaches. The results show that the proposed approach improves the state estimation accuracy under multiple tracking challenges.

Index Terms—Visual tracking, Fusion, Online performance evaluation, Tracker correlation

I. INTRODUCTION

VISUAL tracking is widely used in applications such as video surveillance, human-computer interaction, activity recognition and video indexing. A tracker faces several challenges such as occlusions, clutter, changes in target scale or appearance and variations in scene illumination. Because no individual tracker can still provide accurate results for all challenges [1], fusing complementary trackers whose expected failures are uncorrelated can increase robustness.

Fusion can be performed at feature or decision-level [2]. Feature-based approaches fuse multiple features in a single tracking framework to adapt to appearance changes [2], [3], [4], [5], [6]. When the features have variable dimensionality and range, adaptation methods are needed to integrate new features [7]. Decision-level fusion combines the output of multiple trackers [8], [9], [10], [7], [11], [12], [13], [14]. Fusion can happen sequentially [15], [11], [16], using outputs from specific trackers [8] or employing likelihood-based fusion [13]; in parallel [17] or with a hybrid approach [18].

The online evaluation of the quality of current tracking results uses current and past information only, and may

help the fusion process [19]. In order to weight trackers or features prior to fusion, online performance evaluation identifies the trackers that follow the target and estimate the accuracy of their outputs at run time [20]. Fusion approaches estimate performance using target velocity [21], democratic integration [6], filter uncertainty [22], likelihood [13] and tracker correlation [17]. Existing performance evaluators provide different score ranges for each tracker [2], require specific trackers [20] or are computationally expensive [23].

In this paper, we propose a decision-level fusion framework that combines the outputs of selected trackers over time based on the spatio-temporal relationships of their results. The main novelties of the proposed approach are a method to identify the trackers that are expected to be on the correct target and the definition of an adaptive reference frame for online performance analysis. We group trackers hierarchically based on their agreement in estimating the target state in terms of spatial location and direction of movement. Using this spatio-temporal agreement, we determine which groups (clusters) of trackers are in the same region and identify the one that is *on-target*. This identification is achieved using an adaptive time-reversed performance evaluation. This evaluation compares the results of trackers running in the reverse temporal direction with the results of the fused output at a specific frame. This specific frame is adaptively determined via online performance evaluation and motion analysis. The final output is then generated by fusing the outputs of trackers within the *on-target* cluster and the selected *on-target* cluster is propagated over time until a split or merge is detected.

The paper is organised as follows: Section II discusses the related work. The overview of the proposed framework is given in Section III, while the tracker clustering and the reverse-analysis evaluation are described in Section IV and Section V, respectively. Experimental results are presented in Section VI. Finally, Section VII concludes the paper.

II. RELATED WORK

In this section we discuss methods for online evaluation of trackers and features, and their combination.

Tracker performance can be evaluated using trajectories, observation likelihood or the spatial uncertainty of the target hypotheses. Comparing target state properties such as target velocity [21] to empirical thresholds limits the approach to specific data. By the reversibility property of Markov chains [23], tracker performance is evaluated using another tracker running in reverse temporal direction until a reference frame in which the tracker under evaluation (forward) is assumed to be correct. Then, reverse-forward

O. Khalid and A. Cavallaro are with the Centre for Intelligent Sensing, Queen Mary University of London, UK (e-mail: {o.khalid, a.cavallaro}@qmul.ac.uk). Juan C. SanMiguel is with University Autonoma of Madrid, Spain (email: juancarlos.sanmiguel@uam.es). Juan C. SanMiguel performed part of this work when he was Post-Doctoral Research Assistant at Queen Mary University of London.

O. Khalid was supported by the Erasmus Mundus Joint Doctorate in Interactive and Cognitive Environments, which is funded by the Education, Audiovisual & Culture Executive Agency (EACEA). Juan C. SanMiguel acknowledges the support of the Spanish Government (HA-Video TEC2014-5317-R). A. Cavallaro acknowledges the support of the Artemis JU and by the UK Technology Strategy Board through the COPCAMS Project, under Grant 332913.

results are compared using the Mahalanobis distance at the reference frame [23]. However, this reverse-analysis considerably increases the computational cost. The tracker likelihood can also be used as performance indicator [13]. However, distractors (i.e. objects with similar features to those of the target) may produce a high likelihood thus generating misleading measurements.

Performance measures can take advantage of multi-hypothesis trackers such as the spatial uncertainty of the hypotheses (particles) of particle filter (PF) to weight feature contributions [2], to compute a feature rejection probability [5] or to be employed jointly with reverse analysis [20]. When multiple trackers or features are used, performance weights can be estimated as the distance to the fused output [24], [6], the average change of multiple features [25] or correlation among point trajectories [26]. Such weighting often requires multiple features whose average result is assumed to be accurate.

Fusion at feature level considers multiple visual features to be combined in a single-tracker framework. An example is the sum of feature likelihoods (colour histograms and intensity gradients) [27]. Results can be improved using a performance-based feature weighting, such as the fused-feature distance of democratic integration [24], [6]. Appearance models based on sparse coding [28] have also been used for feature-level fusion, where the weighting is determined by the contribution of the feature template to track the target [29], while other approaches discard features far away from the target model [30].

Decision-level fusion combines the output of multiple trackers in cascade or in parallel. A *cascade* for fusion defines an execution order where each tracker output is used by the next tracker. Examples include the combination of two trackers (region and shape) and two detectors (head and motion for people tracking) [15]; the sequential execution of the template-based Mean Shift (MS) and appearance-based trackers [11]; and the integration of three PFs and one Kalman filter (KF) [16]. Moreover, trackers can be *integrated within the framework of another tracker* [8], [9]. For instance, a head tracker uses MS to improve the PF tracker predictions [8].

In *parallel* tracker fusion, two trackers may be combined using target motion [21] or probability density functions [7]. Moreover, tracker performance within a parallel framework can be measured as disagreement with other trackers [17], [31] or can be used to select the best tracker [32]. Other approaches may use tracker correlation to improve the overall tracking performance by correcting PFs [18] and KFs [26]. These approaches determine the accuracy as the spatial uncertainty of hypotheses whose value may vary across trackers, thus making tracker fusion difficult.

Learning-based methods have also been proposed [33] where labels (foreground/background) are assigned to image patches. A Bayesian approach is employed for fusion where tracker accuracy is the distance between the fused output and the output of each tracker. Using likelihood as the per-

TABLE I
DECISION-LEVEL FUSION APPROACHES. KEY - C: CASCADE; P: PARALLEL; S: SPATIAL; T: TEMPORAL; PDF: PROBABILITY DENSITY FUNCTION; GLAD: GENERATIVE MODEL OF LABELS, ABILITIES AND DIFFICULTIES [34].

Ref.	Type	Reliability evaluator		Tracker correlation		Fused trackers	Fusion method
		S	T	S	T		
[11]	C	-	-	-	-	All	Correction
[8], [9]		-	-	-	-		Kernel-Bayesian
[7]		-	-	-	-		Product of PDFs
[14]		✓	-	-	-		Weighted sum
[33]		✓	-	-	-		GLAD
[21]	P	✓	✓	-	-		Correction
[17]		✓	✓	✓	-		Mixture of Gaussians
[18]		✓	-	-	-		Weighted sum
[26]		✓	-	✓	-		Correction
[31]		✓	✓	✓	-		Correction + Average
[35]		✓	✓	✓	-	Selected	Average
[12], [13]		✓	-	-	-		Interaction
Proposed		✓	✓	✓	✓		Average

formance estimator within *tracker interaction* and sampling based approaches, the tracker with the highest likelihood is chosen [12]. Similarly, multiple motion and appearance models can be used to form a single compound tracker [13].

Decision-level fusion approaches are summarized in Table I and compared with the proposed approach.

III. OVERVIEW

We propose a framework to cluster trackers over time and to select the best-performing ones for fusion to improve the overall accuracy of target state estimation. The proposed approach is inspired by the *test and select framework* [36] for ensemble combination where accurate classifiers are fused assuming that their errors are diverse. Considering trackers as classifiers, we extend this framework to video tracking by introducing spatio-temporal correlation and adaptive online performance evaluation (Fig. 1).

Let $\mathbf{I} = \{I_t\}_{t=1}^T$ be a video sequence of T frames and $\mathbf{F} = \{F^k\}_{k=1}^K$ be a set of K trackers. Let the target state x_t^k be a bounding box, defined by a four dimensional vector $[u_t^k, v_t^k, w_t^k, h_t^k]$, where u_t^k, v_t^k are the target position; w_t^k and h_t^k are its width and height, respectively. Each tracker F^k uses the observation z_t^k and the target model at frame I_{t-1}, ϕ_{t-1}^k , to estimate the target state at time t :

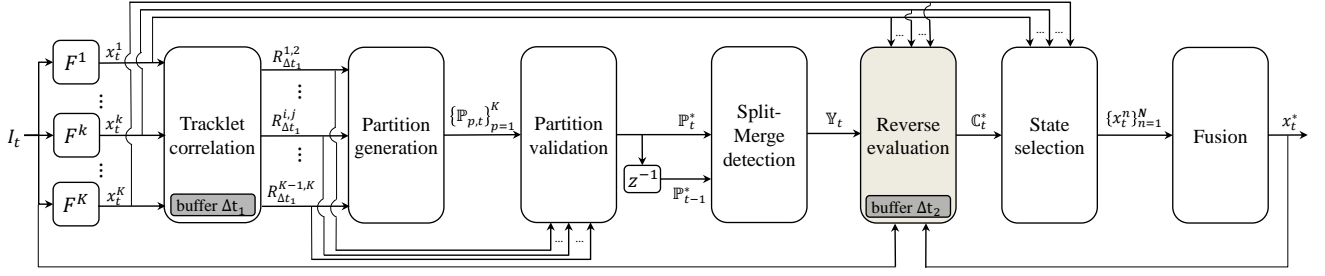
$$x_t^k = F^k(x_{t-1}^k, z_t^k, \phi_{t-1}^k), \quad (1)$$

where x_{t-1}^k is the state estimate (i.e. the tracker output) at the previous time step.

Let *on-target* and *off-target* be the labels that indicate whether a tracker follows the target successfully or not, respectively. The goal is to identify the successful trackers given the outputs x_t^k by labeling them as:

$$x_t^k \rightarrow l_t^k \in \{\text{on-target}, \text{off-target}\}. \quad (2)$$

We determine l_t^k by recognizing groups of trackers (clusters) following the same region in the frame and identifying the cluster with the *on-target* trackers $\mathbb{C}_t^* = \{F^n\}_{n=1}^N \subseteq \mathbf{F}$

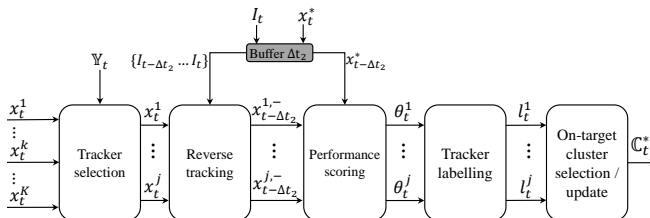
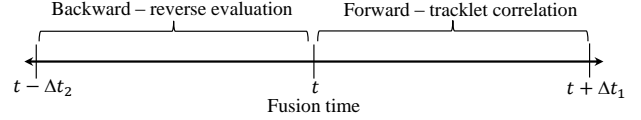
Fig. 1. Block diagram of the proposed approach to fuse the output of K trackers.

($N \leq K$). Assuming that all the trackers are initialized with a ground truth, the approach starts with a single *on-target* cluster. When trackers fail, they split into different clusters, of which only one (i.e. \mathbb{C}_t^*) or none correctly tracks the target. For each frame I_t we compute $\binom{K}{2}$ scores $R_{\Delta t_1}^{i,j}$ to determine spatio-temporal relationships between pairs of trackers, measured as similarity of spatial location and direction of movement of short-term trajectories (tracklet correlation) over a temporal window Δt_1 . These spatio-temporal scores are then employed to generate partition hypotheses $\{\mathbb{P}_{p,t}\}_{p=1}^K$ to divide the K trackers into clusters. After validating the best partition \mathbb{P}_t^* by exploring the correlations among tracklet data, the *on-target* cluster \mathbb{C}_t^* is determined by online performance evaluation of the trackers that are expected to be following the target. Such evaluation uses reverse-tracking [23] over a sliding temporal window Δt_2 (Fig. 2), which requires standard tracker outputs (e.g. bounding boxes), thus providing a generic evaluator across trackers.

The proposed approach employs two temporal windows Δt_1 and Δt_2 (Fig. 3), during which data is buffered from future and past time instants, respectively. The temporal window used for tracklet correlation makes the proposed approach suitable for applications that can tolerate a short latency Δt_1 .

\mathbb{C}_t^* is propagated until the detection of a split or a merge, which happens when trackers leave or join the cluster \mathbb{C}_t^* , respectively. A split or merge indicates that some or all of the *on-target* trackers may have failed. When such changes occur, all trackers are re-evaluated to determine the new *on-target* cluster \mathbb{C}_t^* in the partition \mathbb{P}_t^* .

Only the trackers belonging to the *on-target* cluster \mathbb{C}_t^*

Fig. 2. Block diagram of the reverse evaluation that identifies the *on-target* cluster \mathbb{C}_t^* .Fig. 3. Temporal windows Δt_1 and Δt_2 employed by the proposed approach to account for forward and backward data, respectively. Forward data are used to determine the relationships among trackers via their trajectories (Section IV). Backward data are used to check tracker performance via a time-reversed evaluator (Section V).

are used to compute the final target state x_t^* :

$$x_t^* = \frac{1}{N} \sum_{n=1}^N x_t^n. \quad (3)$$

One advantage of the proposed clustering is that it helps to reduce the computational load by avoiding to apply reverse-evaluation over all trackers when they maintain their spatio-temporal relationships over time.

IV. TRACKER CLUSTERING

A. Tracklet correlation

We combine spatial and temporal features of the short-term trackers' trajectories (tracklets) to obtain a set of pairwise correlation scores $R_{\Delta t_1}^{i,j}$, for $1 \leq i, j \leq K$ with $i \neq j$; for pairs of trackers F^i and F^j over a temporal window Δt_1 . The spatial agreement for F^i and F^j is based on their outputs x_t^i and x_t^j at frame I_t :

$$O_t^{i,j} = \frac{2|A_t^i \cap A_t^j|}{|A_t^i| + |A_t^j|}, \quad (4)$$

where A_t^i and A_t^j are the set containing the pixels of the bounding boxes generated by tracker F^i and F^j , respectively, and $|\cdot|$ is the cardinality of a set. $O_t^{i,j} \in [0, 1]$ and a value of 1 (0) represents a full agreement (disagreement). The spatial agreement over time is computed by averaging $O_t^{i,j}$ over Δt_1 :

$$O_{\Delta t_1}^{i,j} = \frac{1}{\Delta t_1} \sum_t^{t+\Delta t_1} O_t^{i,j}. \quad (5)$$

In order to estimate the agreement for motion direction, we compute a score $r_{\Delta t_1}^{i,j}$ using the directional feature \vec{d}^k of each F^k [37] over Δt_1 that encodes the trajectory direction:

$$\vec{d}^k = (u_{t+\Delta t_1}^k - u_t^k, v_{t+\Delta t_1}^k - v_t^k). \quad (6)$$

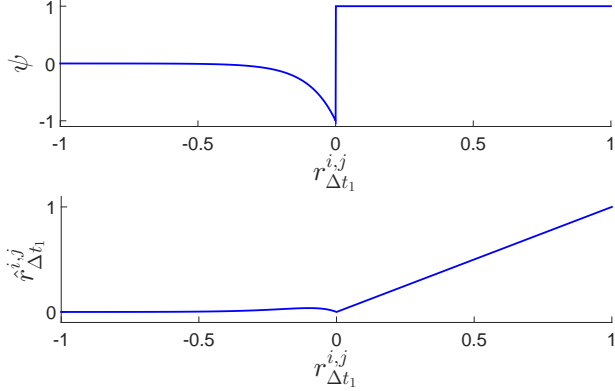


Fig. 4. Weighting function ψ (top) and weighted directional score $\hat{r}_{\Delta t_1}^{i,j}$ (bottom) for directional feature normalization using $\lambda = 10$.

The directional similarity score $r_{\Delta t_1}^{i,j}$ is computed between F^i and F^j using the cosine similarity:

$$r_{\Delta t_1}^{i,j} = \cos \left((\vec{d}_t^i \cdot \vec{d}_t^j) / (|\vec{d}_t^i| \cdot |\vec{d}_t^j|) \right), \quad (7)$$

where $r_{\Delta t_1}^{i,j} \in [-1, 1]$ and negative values represent (estimated) targets moving in opposite directions.

The desired correlation score $R_{\Delta t_1}^{i,j}$ is obtained by combining $O_{\Delta t_1}^{i,j}$ and $r_{\Delta t_1}^{i,j}$ after normalization of $r_{\Delta t_1}^{i,j}$ to $[0, 1]$. Since we are interested in an agreement on the direction of motion, we seek values $r_{\Delta t_1}^{i,j} \in [0, 1]$ and define a weighted directional similarity score $\hat{r}_{\Delta t_1}^{i,j}$ as:

$$\hat{r}_{\Delta t_1}^{i,j} = \psi(\lambda, r_{\Delta t_1}^{i,j}) \cdot r_{\Delta t_1}^{i,j}, \quad (8)$$

where $\hat{r}_{\Delta t_1}^{i,j} \in [0, 1]$ and $\psi \in [-1, 1]$ is a weighting function that assigns a constant weight to $r_{\Delta t_1}^{i,j} \in [0, 1]$, while $r_{\Delta t_1}^{i,j} \in [-1, 0]$ are given low weights. Such a function is defined as:

$$\psi(\lambda, r_{\Delta t_1}^{i,j}) = \begin{cases} 1 & \text{if } 0 \leq r_{\Delta t_1}^{i,j} \leq 1, \\ -e^{\lambda \cdot r_{\Delta t_1}^{i,j}} & \text{if } -1 \leq r_{\Delta t_1}^{i,j} < 0, \end{cases} \quad (9)$$

where $\lambda \in (0, \infty)$ is the decay rate of ψ . Values of λ close to zero give smooth transitions for $\psi \in [-1, 0]$ turning into high $\hat{r}_{\Delta t_1}^{i,j}$ values when $r_{\Delta t_1}^{i,j} \in [-1, 0]$. High values of λ give abrupt transitions for ψ turning into $\hat{r}_{\Delta t_1}^{i,j}$ values close to zero. Fig. 4 shows the relations between $r_{\Delta t_1}^{i,j}$ and $\hat{r}_{\Delta t_1}^{i,j}$ for $\lambda = 10$ (bottom graph) and between $r_{\Delta t_1}^{i,j}$ and ψ (top graph).

$R_{\Delta t_1}^{i,j}$ is finally computed as follows:

$$R_{\Delta t_1}^{i,j} = \omega \cdot O_{\Delta t_1}^{i,j} + (1 - \omega) \cdot \hat{r}_{\Delta t_1}^{i,j}, \quad (10)$$

where $\omega \in [0, 1]$. High (low) values of ω prioritize the spatial overlap (trajectory direction), which can be useful for short (long) Δt_1 .

B. Partition generation

A single partition $\mathbb{P}_{p,t}$ of \mathbf{F} is a collection of non-empty clusters $\mathbb{C}_{p,t}^a$ ($a = 1, \dots, |\mathbb{P}_{p,t}|$) such that each tracker in \mathbf{F} is in exactly one $\mathbb{C}_{p,t}^a$, i.e. all $\mathbb{C}_{p,t}^a$ are mutually disjoint.

At each time step, K trackers can be grouped into clusters $\mathbb{C}_{p,t}^a$, forming a single partition $\mathbb{P}_{p,t}$, where $|\mathbb{P}_{p,t}| \in [1, K]$. Let $[\cdot]$ represent a partition. For example, $[\{F^1, \dots, F^k\}]$ means that all trackers are clustered together ($|\mathbb{P}_{p,t}| = 1$, initial condition) and $[\{F^1\}, \dots, \{F^k\}]$ means that each tracker is a single cluster ($|\mathbb{P}_{p,t}| = K$).

Our aim is to hypothesize a set of partitions $\{\mathbb{P}_{p,t}\}_{p=1}^B$ to cluster the trackers. All possible partitions $\mathbb{P}_{p,t}$ can be systematically enumerated with an exhaustive search [38]. The set size is given by the Bell number \mathcal{B} [39], which increases exponentially with K . For example, with $K = 8$ trackers $\mathcal{B} = 4140$ partitions are generated.

To reduce the computational complexity, we use a greedy search that determines the most plausible partitions for a given number of clusters. Since the optimum number and composition of clusters is unknown, we take advantage of the hierarchical structure of the tracker relationships to generate a set of partitions whose cardinality ranges from 1 (i.e. a single cluster) to K (i.e. each tracker is a cluster), with $K \ll \mathcal{B}$.

We use hierarchical clustering (HC) [40] to determine the relationships between trackers based on the pair-wise tracker correlation scores ($R_{\Delta t_1}^{i,j}$). Based on the distance between the trackers' outputs, we obtain a dendrogram, which is inspected by a divisive (top-down) approach to determine each partition $\mathbb{P}_{p,t}$. The search starts with the partition that groups all trackers in one cluster $\mathbb{P}_{1,t}$. Recursively moving down the tree, a different $\mathbb{P}_{p,t}$ is generated at each level, with the final partition having each tracker in a separated cluster $\mathbb{P}_{K,t}$. A partition $\mathbb{P}_{p,t}$ is obtained as:

$$\mathbb{P}_{p,t} = f(\beta^*(p)), \quad (11)$$

where $p = 1, \dots, K$ and $f(\beta^*)$ is a HC-based function that provides a cluster partition given an optimum distance threshold β^* , which is computed as:

$$\beta^*(p) = \underset{\beta}{\operatorname{argmin}} \{ |f(\beta)| - p \}, \quad (12)$$

where $\beta = 0, \dots, \max\{R_{\Delta t_1}^{i,j}\}$.

The proposed greedy search has a linear relationship between the size of $\{\mathbb{P}_{p,t}\}_{p=1}^K$ and K , which significantly speeds up the search. Fig. 5(a) and Fig. 5(c) show an example for four trackers and the scores for their spatio-temporal relations, which are used to compute the dendrogram illustrated in Fig. 5(b).

C. Partition validation

After generating the set of partitions $\{\mathbb{P}_{p,t}\}_{p=1}^K$, the objective is to select the partition \mathbb{P}_t^* that best represents the spatio-temporal relations among trackers. We therefore define the score $S(\mathbb{P}_{p,t})$ as:

$$S(\mathbb{P}_{p,t}) = \frac{1}{|\mathbb{P}_{p,t}|} \sum_{a=1}^{|\mathbb{P}_{p,t}|} Q(\mathbb{C}_{p,t}^a), \quad (13)$$

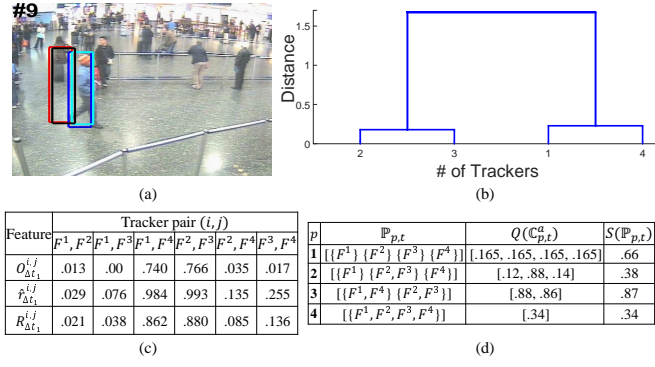


Fig. 5. (a) Tracking results for frame 9 of the MCTTR0205a sequence (TRECVID) —: F^1 ; —: F^2 ; —: F^3 ; —: F^4 . (b) The dendrogram obtained by hierarchical clustering. (c) Pair-wise tracker correlations scores. (d) Hypothesized partitions and cluster scores, where $\mathbb{P}_{3,t}$ has the highest score.

where $Q(\mathbb{C}_{p,t}^a)$ is the score for a single cluster $\mathbb{C}_{p,t}^a \in \mathbb{P}_{p,t}$. $S(\mathbb{P}_{p,t})$ determines the partition \mathbb{P}_t^* as:

$$\mathbb{P}_t^* = \underset{p}{\operatorname{argmax}} \{S(\mathbb{P}_{p,t})\}, \quad (14)$$

where $p = 1, \dots, K$ and $Q(\mathbb{C}_{p,t}^a)$ is dependent upon the pair-wise relationship score between trackers F^i and F^j in $\mathbb{C}_{p,t}^a$, which is obtained in the *tracklet correlation* block (Section IV-A).

The score $Q(\mathbb{C}_{p,t}^a)$ is computed as:

$$Q(\mathbb{C}_{p,t}^a) = \begin{cases} \frac{1}{\nu} \sum_{i=1}^{|\mathbb{C}_{p,t}^a|} \sum_{j=1}^{|\mathbb{C}_{p,t}^a|} R_{\Delta t}^{i,j} & \text{if } |\mathbb{C}_{p,t}^a| > 1, \\ 1 - \max_{b \in 1, \dots, |\mathbb{P}_{p,t}|} (Q(\mathbb{C}_{p,t}^a \cup \mathbb{C}_{p,t}^b)) & \text{if } |\mathbb{C}_{p,t}^a| = 1, \end{cases} \quad (15)$$

where $\nu = \binom{|\mathbb{C}_{p,t}^a|}{2}$ is the total number of tracker-pair combinations within the cluster. Since a pair-wise score for a single tracker in a cluster, $|\mathbb{C}_{p,t}^a| = 1$, cannot be obtained, we compute its pair-wise scores with trackers in other clusters. Therefore $\mathbb{C}_{p,t}^a \cup \mathbb{C}_{p,t}^b$ indicates the hypothetical case where the tracker in $\mathbb{C}_{p,t}^a$ becomes part of cluster $\mathbb{C}_{p,t}^b$, and b is any of the remaining clusters within $\mathbb{P}_{p,t}$ ($b \neq a$).

When each tracker is a single cluster, i.e. $|\mathbb{P}_{p,t}| = K$, each cluster score $Q(\mathbb{C}_{p,t}^a)$ is computed as:

$$Q(\mathbb{C}_{p,t}^a) = 1 - Q(\mathbb{C}_{p,t}^b), \quad (16)$$

where $\mathbb{C}_{p,t}^b$ is the cluster containing all trackers. Fig. 5(d) shows the computed cluster and partition scores, where $\mathbb{P}_{3,t}$ achieves the highest score.

D. Split-Merge detection

After determining \mathbb{P}_t^* , the *split-merge detection* step identifies changes between \mathbb{P}_t^* and the previous partition \mathbb{P}_{t-1}^* . Such changes may occur due to trackers leaving or joining the *on-target* cluster in the previous time step $\mathbb{C}_{t-1}^* \in \mathbb{P}_{t-1}^*$, hence modifying the structure of \mathbb{P}_{t-1}^* . Thus, \mathbb{C}_{t-1}^* cannot be propagated to the current time and reverse-evaluation

is required to identify the current \mathbb{C}_t^* . We apply reverse-evaluation over a set of trackers \mathbb{Y}_t , selected as:

$$\mathbb{Y}_t = \begin{cases} \mathbb{C}_{t-1}^* & \text{if } \mathbb{P}_{t-1}^* \equiv \mathbb{P}_t^*, \\ \mathbb{P}_t^* & \text{otherwise,} \end{cases} \quad (17)$$

where the condition $\mathbb{P}_{t-1}^* \equiv \mathbb{P}_t^*$ checks the similarity between the number of clusters and their members (i.e. trackers). When this condition is satisfied, an existing cluster $\mathbb{C}_t^* \in \mathbb{P}_t^*$ equivalent to \mathbb{C}_{t-1}^* is used at the current time. However, when a split or merge occurs, all the trackers in \mathbb{P}_t^* are evaluated.

V. ‘On-target’ CLUSTER IDENTIFICATION

We evaluate the performance of each tracker in the set \mathbb{Y}_t (*tracker selection* block in Fig. 2). This performance evaluation either determines the *on-target* cluster \mathbb{C}_t^* of the valid partition \mathbb{P}_t^* or validates the *on-target* cluster from the previous time step \mathbb{C}_{t-1}^* . We cast this problem as an online tracker evaluation and use the time-reversibility of target motion to assess the performance of the trackers. We first review reverse-based evaluation methods and then we present our proposed improvements.

A. Reverse-based online evaluation

Reverse-based evaluation [23] measures the performance of a tracker during runtime using the generated results. For each frame where we evaluate the tracker, a reversed-tracker (i.e. the same tracker operating in reverse-time) is applied. Using the tracker output x_t^k , as the reverse-tracker initialization $x_{t-1}^{k,-}$, the reverse-tracker obtains its output as:

$$x_{t-1}^{k,-} = F^k(x_t^{k,-}, z_t^k, \phi_t^k), \quad (18)$$

where $x_{t-1}^{k,-}$ is the reverse-tracker output at time $t-1$. Then the result of the reverse-tracker and that of the tracker are compared to obtain a similarity score θ_t^k by means of the Mahalanobis distance between the likelihood distributions of the forward and reverse target estimations. This comparison is performed at a certain time instant t_{ref} , known as reference frame I_{ref} . I_{ref} is a frame where the tracker is known to be *on-target* and it is usually $I_{ref} = I_1$ [23], i.e. the frame where the target is initialized.

This approach has two major limitations. First, the forward-reverse similarity uses the Mahalanobis distance that returns unbounded scores $\theta_t^k \in [0, +\infty)$, which can have different range of values depending on the trackers employed in the fusion framework. Hence θ_t^k may be inappropriate to compare the trackers to be combined. Second, running the reverse tracker until the first frame implies an exponential growth in computational time. A faster approximation is proposed where I_{ref} is moved ahead in time. However, tracker errors are accumulated over time by the reverse-tracker, thus leading to drift [41]. For example, if the tracker loses the target and gets locked on the background the forward-reverse similarity may give high scores θ_t^k , since the

reverse-tracker is incorrectly initialized by the wrong tracker estimations.

We address these shortcomings for reverse-tracking evaluation as described next.

B. Performance score and reference frame update

For the limitation of unbounded θ_t^k scores, we compare the reverse-tracker and fused outputs to obtain a $\theta_t^k \in [0, 1]$:

$$\theta_t^k = G(x_{ref}^{k,-}, x_{ref}^{*,k}), \quad (19)$$

where $x_{ref}^{k,-}$ and $x_{ref}^{*,k}$ are the reverse-tracker and fused outputs at I_{ref} , respectively; G defines the output similarity and is computed using (4) where A_t^i and A_t^j are replaced by $A_{ref}^{*,k}$ and $A_{ref}^{k,-}$, respectively. $A_{ref}^{*,k}$ and $A_{ref}^{k,-}$ are the sets containing the pixels of the bounding boxes of $x_{ref}^{*,k}$ and $x_{ref}^{k,-}$, respectively.

For the limitation of the exponential growth of computational time when $I_{ref}=1$, we update I_{ref} over time so that the computational cost is bounded and reverse-evaluation can be applied to long sequences. We implement such update assuming that the fused output is *on-target* and that the target has changed position from I_{ref} to the current frame, thus making the motion information useful for reverse-analysis.

The motion of bounding boxes is minimal when the tracker is *on-target* and the target is static; or when the tracker drifts from the target and gets locked onto a static background region. Because it is difficult to differentiate between these two situations, we analyze significant motion changes of the trackers compared with their average motion. The maximum motion M^k is computed over temporal window Δt_2 using F^k trajectory (Fig. 3), and the top-left ($u_{1,t}^k, v_{1,t}^k$) and bottom-right ($u_{2,t}^k, v_{2,t}^k$) coordinates of the bounding box. Motion for u_1^k over Δt_2 is computed as:

$$M_{u_1^k} = \frac{1}{\Delta t_2} \sum_{t'=t-\Delta t_2}^t (M_{u_{1,t'}^k} - M_{u_{1,t'-1}^k}). \quad (20)$$

Motion for $v_{1,t}^k, u_{2,t}^k$ and $v_{2,t}^k$ is computed using (20), where $M_{u_1^k}$ is replaced by $M_{v_{1,t}^k}, M_{u_{2,t}^k}$ and $M_{v_{2,t}^k}$, respectively. $M^k = \max(M_{u_1^k}, M_{v_{1,t}^k}, M_{u_{2,t}^k}, M_{v_{2,t}^k})$ returns the maximum motion for F^k . $M^* = \frac{1}{\Delta t_2} \sum_{t'=1}^t (M_{x_{t'}}^* - M_{x_{t'-1}}^*)$ determines the motion of the fused output, which is used as a common threshold to compare the motion of all trackers in the framework.

The performance of each tracker is computed using (19). To determine a single I_{ref} for all trackers, we use $\max(M^k)$ and $\max(\theta_t^k)$ to select the best performing tracker for that temporal window.

We adaptively estimate and update I_{ref} by combining motion analysis and performance of the tracker as:

$$I_{ref} = \begin{cases} I_{t-\Delta t_2} & \text{if } \max(M^k) \geq M^* \text{ and } \max(\theta_t^k) \geq \tau_1, \\ I_{ref} & \text{otherwise,} \end{cases} \quad (21)$$

where $\tau_1 = 0.5$ is the minimum tracker accuracy [1].

TABLE II
SEQUENCES USED IN THE EXPERIMENTS. KEY - BC: BACKGROUND CLUTTER; P: POSE CHANGES; O: OCCLUSIONS; I: ILLUMINATION CHANGES; S: SCALE CHANGES; M: MOTION CHANGES; BS: BACKGROUND SIMILARITY.

Dataset	Sequence name	Target Class	Size		Total Frames	Challenges
			Target	Frame		
Students	P1	University Students	Person	21 x 75	720 x 576	250 BC, P, O, I
	P2		Person	22 x 69		250 BC, P, I
	P3		Person	25 x 61		165 P, I
CAVIAR	P4	Browse_WhileWaiting1	Person	50 x 24	384 x 288	200 P, I
	P5	OneLeaveShopReenter1Front	Person	16 x 56		195 P, I, S, BS
	P6	OneLeaveShopReenter2Front	Person	14 x 50		300 P, I, S, BS
	P7	ThreePastShop2cor	Person	56 x 142		170 P, I, S
PETS	P8	S2.L2 walking	Person	14 x 50	768 x 576	140 P, I, S, BS
	P9	PETS2001 Dataset 1	Vehicle	56 x 142		150 P, I, S
	P10	S2.L1 walking	Person	22 x 68		150 BC, I, S
	P11	PETS2001 Dataset 1	Vehicle	72 x 56		180 S, M
	P12	S2.L1 walking	Person	72 x 56		90 O, M
	P13		Person	72 x 56		150 M, O, P
	P14		Person	72 x 56		110 S, M
	P14		Person	72 x 56		110 S, M
LTD	P15	NissanSkylineChase	Person	37 x 21	640 x 275	300 I, M, S
David	P16	David Indoor	Head	91 x 116	320 x 240	130 I, S, M, BS
AVSS2007	P17	Abandoned baggage	Person	60 x 240	720 x 576	200 P, I
TRECVID	P18	MCTTR0205a	Person	72 x 226	720 x 576	50 P, I, BS
	P19		Person	64 x 204		40 P, I, BS
	P20		Person	64 x 204		40 I, M, BS
MIT Traffic	P21	MV2_001	Vehicle	34 x 26	720 x 480	160 I, S, M
	P22	MV2_006	Vehicle	70 x 36	720 x 576	160 O, I, M

C. On-target cluster selection/update

Reverse-evaluation identifies the *on-target* trackers by using the individual performance scores θ_t^k of trackers in \mathbb{Y}_t . Trackers with $\theta_t^k \geq \tau_1$ are labeled as *on-target*, enabling the method to select \mathbb{C}_t^* as the cluster \mathbb{C}_t^a with all *on-target* trackers:

$$\mathbb{C}_t^* = \{\mathbb{C}_t^a \in \mathbb{P}_t^* : l_t^k = \text{on-target}, \forall F_t^k \in \mathbb{C}_t^a\}. \quad (22)$$

VI. EXPERIMENTAL RESULTS

A. Setup

1) *Dataset*: For evaluating the proposed approach, multi-Tracker Partition Fusion (TPF), we consider the following datasets: Students¹, CAVIAR², PETS (2009³ and 2001⁴), LTD⁵, TRECVID2009⁶, MIT Traffic⁷, David⁸ and AVSS2007⁹. We have selected 22 sequences (3580 frames) to cover indoor and outdoor scenarios containing tracking challenges such as occlusions, background clutter, pose, motion and illumination changes. Table II describes the selected sequences and Fig. 6 shows the target initializations.

¹<http://graphics.cs.ucy.ac.cy/research/downloads/crowd-data>

²<http://groups.inf.ed.ac.uk/vision/CAVIAR/CAVIARDATA1/>

³<http://www.cvg.reading.ac.uk/PETS2009>

⁴<http://www.cvg.reading.ac.uk/slides/pets.html>

⁵<http://www.micc.unifi.it/LTDT2014>

⁶<http://trecvid.nist.gov/trecvid.data.html#tv09>

⁷<http://www.ee.cuhk.edu.hk/~xgwang/HBM.html>

⁸<http://www.cs.toronto.edu/~dross/ivt/>

⁹<http://www.avss2007.org/>

2) *Trackers*: We apply TPF to combinations of up to eight trackers using publicly available authors' implementations.

The first tracker is the *Sparse features based Tracker* (ST) [42], which is PF-based and uses sparse (intensity) features to generate the target appearance model. The Maximum a Posteriori criterion is employed to estimate the target state. The second tracker is the *Adaptive Fragments-based Tracker* (AFT) [43] that models the target appearance with various fragments. Fragment reliability is based on colour similarity between the current and previous fragment, to integrate highly-reliable fragments within a PF framework. The third tracker is the *Locally Orderless Tracker* (LOT) [44] that divides the target into *superpixels* using the HSV colour space and employs a PF to track the target. The fourth tracker is the *Incremental Visual Tracker* (IVT) [45] that performs on-line updating to account for appearance changes, and a PF to track the target over time. The fifth tracker is the *scale and orientation Adaptive MS Tracker* (AMS) [46], that estimates the changes in scale and orientation of the target using the MS framework by employing Gaussian kernels and image moments. The sixth tracker is the *Fast Compressive Tracker* (FCT) [47] that projects the original image to a low-dimensional space. The projected features are then used to formulate tracking as a binary classification task via a naive Bayesian classifier. The seventh tracker is the *L1 Tracker* (L1T) [48], which is based on PF and models the target by sparse linear combinations of target and trivial templates (set of unit vectors). Assuming an affine motion model, tracking is performed by solving the L_1 minimization problem. The eighth tracker is the *Least Soft-Threshold Squares Tracker* (LSST) [49], which is based on PF and performs linear regression via least-soft threshold squares distance between the observation and the target model.

We have implemented six TPF configurations: TPF_3 (ST, AFT, LOT), TPF_4 (ST, AFT, LOT, IVT), TPF_5 (ST, AFT, LOT, IVT, FCT), TPF_6 (ST, AFT, LOT, IVT, FCT, AMS), TPF_7 (ST, AFT, LOT, IVT, FCT, AMS, L1T) and TPF_8 (ST, AFT, LOT, IVT, FCT, AMS, L1T, LSST). We use TPF_3 for Section VI-C and VI-D, while Sections VI-E, VI-F and VI-G use all six configurations.

We compare TPF with the eight selected trackers, two recent trackers: STRUCK (STR) [50] and Kernelized Correlation Filters (KCF) [3]; and three state-of-the-art decision-level fusion approaches: Average fusion (AvgF), Symbiotic Tracker (SymT) [17] and Visual Tracker Sampler (VTS) [13]. STRUCK [50] is a tracking-by-detection approach using SVMs with Gaussian kernels. Three features have been tested (Haar, raw pixels and histograms) and we report the results for histograms as they outperformed Haar features and raw pixels. KCF [3] employs correlation among filters based on histograms of oriented gradients features. AvgF combines the eight trackers by assigning equal weights to each tracker. SymT estimates trackers' relationships based on their spatial agreement only, and the tracker performance is based on displacements between consecutive frames. SymT has been re-implemented as described in [17]. VTS

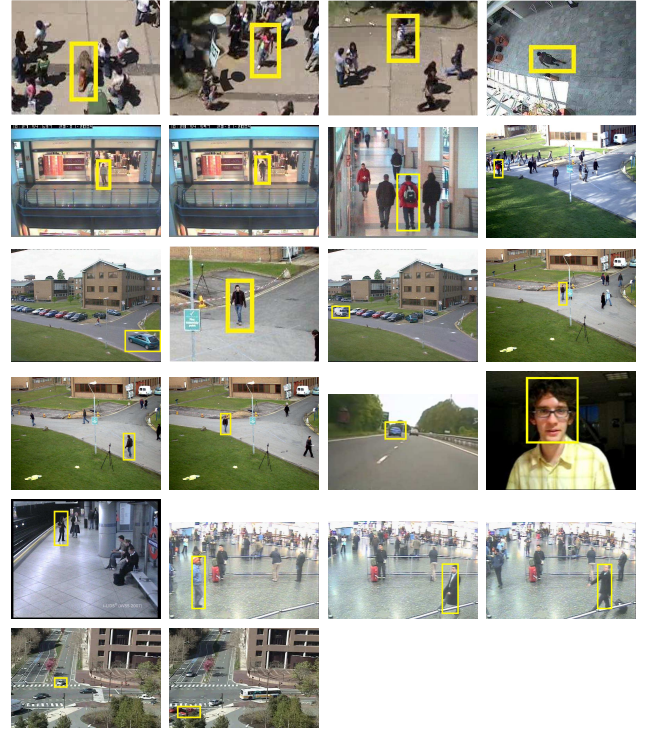


Fig. 6. Target initializations. The order from top-left to bottom-right along the rows corresponds to the row order in Table II.

combines two motion and four appearance models to get eight trackers, using a likelihood-based tracker performance. For STRUCK, KCF and VTS we use the authors' implementation.

3) *Parameters of the proposed approach*: For TPF, the temporal window for reverse-analysis is initially set to $\Delta t_2 = 10$ since it provides a good speed-accuracy trade-off as shown in [23]. This value is updated if the motion or the performance of trackers is below the thresholds (see Section V-B). The temporal window for tracklet correlation is set to $\Delta t_1 = 10$ to keep an initial forward-backward symmetry for analysis, since no prior information is available to define the importance of one analysis over the other. $\omega = 0.5$ ensures equal weighting for the features. For (8), we heuristically found that $\lambda \in [5, 15]$ gives the desired ψ behaviour so we used $\lambda = 10$. Finally, $\tau_1 = \tau_2 = 0.5$ as previously set [1].

B. Evaluation measures

We measure the deviation from ground-truth data (GT) as the overlap score $O_t^{k,GT}$ between the output of tracker F^k and GT annotation using (4). A_t^i and A_t^j are replaced by A_t^{GT} and A_t^k , the sets of pixels contained in the GT and F^k target estimations, respectively. Values close to 1 (0) indicate high (low) tracking performance. The mean of $O_t^{k,GT}$ is computed for each sequence.

We also measure the performance of TPF when assigning the *on-target* and *off-target* labels to trackers and clusters from the valid partition \mathbb{P}_t^* . GT information is used to compute the overlap score for each cluster $O_t^{c,GT}$ by taking the average of $O_t^{k,GT}$ for the trackers within the cluster.

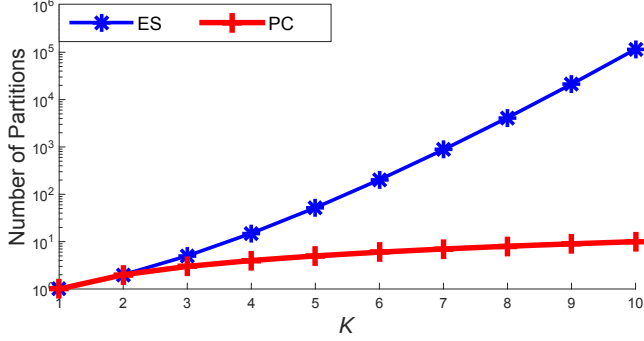


Fig. 7. Comparison of methods to generate partitions $\{\mathbb{P}_{p,t}\}_{p=1}^B$ of K trackers. With exhaustive search (ES) B grows exponentially, while the proposed clustering (PC) is bounded by $B = K$.

The *on-target* trackers are defined for $O_t^{k,GT} \geq \tau_2$ and $O_t^{C^*,GT} \geq \tau_2$ corresponds to the *on-target* cluster \mathbb{C}_t^* . To simplify notation in the remaining sections, we denote the mean ground-truth overlap $O_t^{k,GT}$ for each sequence as O_G .

Using n_{TP} , n_{FP} , n_{TN} and n_{FN} as number of true positives (TP), false positives (FP), true negatives (TN) and false negatives (FN), we compute the precision, $Pr = \frac{n_{TP}}{n_{TP}+n_{FP}}$, the recall, $Re = \frac{n_{TP}}{n_{TP}+n_{FN}}$ and $F\text{-score} = 2 \cdot \frac{Pr \cdot Re}{Pr + Re}$, [1]. Values for the F -score close to 1 (0) indicate high (low) accuracy. n_{TP} (n_{FP}) and n_{TN} (n_{FN}) are the number of clusters or trackers correctly (incorrectly) labeled as *on-target* and *off-target*, respectively.

C. Tracker clustering

We evaluate the performance of the partition generation approach and the features employed to cluster the trackers.

1) *Comparison of the proposed clustering (PC) with exhaustive search (ES)*: Fig. 7 compares the generated set of partitions $\{\mathbb{P}_{p,t}\}_{p=1}^B$ with an increasing number of trackers K for both approaches. The accuracy of their results is equal as PC and ES select the same valid partition. However the size grows exponentially for ES with an increasing K whereas PC keeps the size of $\{\mathbb{P}_{p,t}\}_{p=1}^B$ bounded with respect to $B = K$.

2) *Performance analysis of features*: TPF combines the features $O_{\Delta t_1}^{i,j}$ and $\hat{r}_{\Delta t_1}^{i,j}$ to get $R_{\Delta t_1}^{i,j}$. An accuracy comparison at cluster and tracker-level is presented in Table III. The results indicate that combining both features outperforms using single features. At tracker-level, using both features improves the F -score by 5% (7%) compared to using only the overlap (direction) feature. Similarly at cluster-level an improvement of 6% (7%) is observed in comparison to the overlap (direction) feature.

Fig. 8 shows the tracking accuracy O_G for the three features, where $\omega = 0.5$ improves results globally in 60% of the sequences. Individual features do not always increase performance since no feature is optimum for all situations.

D. On-target cluster identification

1) *Performance analysis for motion*: Table IV compares the proposed approach with and without motion analysis

TABLE III
COMPARISON OF FEATURE COMBINATIONS FOR THE PROPOSED APPROACH. RESULTS SHOW THE F -SCORE AT TRACKER-LEVEL AND CLUSTER-LEVEL, WITH DIFFERENT FEATURE WEIGHTS ω IN (10).

	Tracker-level			Cluster-level		
	$\omega=0$	$\omega=0.5$	$\omega=1$	$\omega=0$	$\omega=0.5$	$\omega=1$
P1	.92	.90	.87	.94	.92	.89
P2	.98	.98	.97	.99	.99	.98
P3	.91	.92	.97	.87	.90	.98
P4	.87	.90	.94	.80	.87	.92
P5	.87	.90	.94	.93	.93	.99
P6	.44	.48	.45	.38	.42	.40
P7	.86	.78	.84	.99	.87	.87
P8	.88	.94	.93	.90	.97	.96
P9	.87	.81	.82	.99	.95	.84
P10	.46	.82	.37	.57	.90	.34
P11	.92	1	.94	1	1	.98
P12	.63	.77	.83	.60	.75	.83
P13	.88	.95	.87	.87	.98	.90
P14	.95	.96	.98	1	1	1
P15	.96	.97	.98	1	1	.99
P16	.81	.81	.93	.81	.88	1
P17	.52	.55	.56	.53	.67	.47
P18	.51	.93	.88	.60	.99	.96
P19	.98	.98	.98	1	1	1
P20	.95	1	.98	.97	1	1
P21	.95	.96	.96	.95	.98	.99
P22	.91	.95	.52	.99	1	.60
Mean	.82	.88	.84	.85	.91	.86

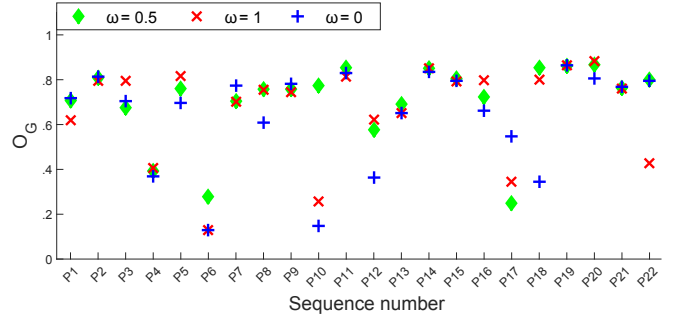


Fig. 8. Tracker accuracy O_G using individual features, ($\omega = 1$ for overlap and $\omega = 0$ for direction) and their equal combination ($\omega = 0.5$).

(TPF and TPF', respectively) to update I_{ref} , in terms of the F -score for selecting the *on-target* trackers. Results for P4-P7 indicate the case when trackers might lose the target due to background clutter and get fixed on the background. Since TPF' is unable to detect this situation in P4, it determines the trackers to be always *on-target*. TPF improves TPF' by 10%, 22% and 25% for ST, AFT and LOT, respectively (19% mean improvement). For P5, P6 and P7 TPF improves performance by 8%, 20% and 20%, respectively. P12 remains occluded between frames 29-39 where ST loses the target and becomes locked on to foreground objects being labeled as *on-target* by TPF', whereas AFT and LOT are labeled as *off-target*. Using motion, TPF improves by 44%, 8% and 85% (40% mean improvement). For P18, ST loses the target at frame 7 due to similar background. TPF' assumes ST to be *on-target*, while AFT and LOT are labeled as *off-target*. TPF uses motion to correctly label AFT and LOT as *on-target* achieving an overall improvement of 420% in

TABLE IV

F-SCORE WITH (TPF) AND WITHOUT (TPF') MOTION ANALYSIS FOR THE THREE FUSED TRACKERS WITH THE REFERENCE FRAME I_{ref} UPDATED USING MOTION ANALYSIS (SECTION V). KEY - ST: SPARSE TRACKER, AFT: ADAPTIVE FRAGMENTS BASED TRACKER; LOT: LOCALLY ORDERLESS TRACKER.

	ST [42]		AFT [43]		LOT [44]		Sequence mean for all trackers	
	TPF'	TPF	TPF'	TPF	TPF'	TPF	TPF'	TPF
P1	.92	.88	.96	.92	.94	.90	.94	.90
P2	1	1	1	.97	1	.97	1	.98
P3	.95	.99	.95	.91	1	.87	.97	.92
P4	.85	.94	.78	.95	.59	.74	.74	.88
P5	.81	.87	.94	.92	.61	.76	.79	.85
P6	.26	.41	.41	.60	.52	.43	.40	.48
P7	.53	.67	.74	.82	.63	.78	.63	.76
P8	.88	.87	.99	.94	1	.98	.95	.93
P9	.90	.72	.90	.78	.89	.91	.89	.80
P10	.88	.75	.91	.83	.88	.82	.89	.80
P11	1	1	1	1	1	1	1	1
P12	.57	.82	.61	.66	.46	.85	.55	.77
P13	.79	.98	.83	.96	.94	.92	.86	.95
P14	1	.98	.96	.95	.91	.94	.96	.96
P15	.99	.98	.96	.96	.95	.99	.97	.98
P16	.91	.85	.99	.86	.88	.73	.93	.81
P17	.51	.54	.80	.40	.88	.48	.73	.47
P18	.29	.33	.05	.87	.08	1	.14	.73
P19	1	1	1	1	1	.95	1	.98
P20	1	1	.97	.99	1	1	.99	1
P21	.94	.88	1	1	1	.99	.98	.96
P22	.92	.96	.81	.97	.92	.86	.88	.93
Tracker mean	.81	.83	.84	.87	.82	.86	.83	.85

comparison to TPF'. ST remains *on-target* for the first 6 frames of the sequence; where TPF incorrectly labels it *off-target* in 5 out of the 6 frames, hence resulting in lower values for TPF. For P17, the target does not move for most of the sequence. ST and AFT lose the target at frame 45 due to similar background, and form a cluster. Due to the stationary target, TPF assumes the ST-AFT cluster to be *on-target* resulting in incorrect labels for all trackers, hence decreasing performance by 35%. Globally, TPF improves TPF' by 2%, 4% and 5% for ST, FT and LOT, respectively.

2) *Evaluation of fast approximation:* Fig. 9 compares the proposed update for the reference frame I_{ref} with the original approach [23]. The average result for three trackers (ST, AFT, LOT) is presented in terms of the overlap score O_G between the GT and the existing forward estimation in I_{ref} obtained by [23] and TPF. TPF improves [23] in 16 out of 22 sequences. I_{ref} is updated only when tracker(s) are found to be *on-target*. For instance, all three trackers fail between frames 95-110 for P4. TPF detects and does not update I_{ref} after frame 110, whereas [23] keeps moving I_{ref} forward, thus accumulating tracker errors. For P14-P16, TPF achieves similar tracking accuracy to [23], while for P9 and P21 the three trackers are able to track the target throughout the sequences.

E. Combining trackers

Performance comparisons are based on the overlap score O_G to measure the area overlap between the final target estimate and the ground truth data.

Table V compares the six TPF configurations (TPF_3, \dots, TPF_8) showing that TPF_3 is the best and

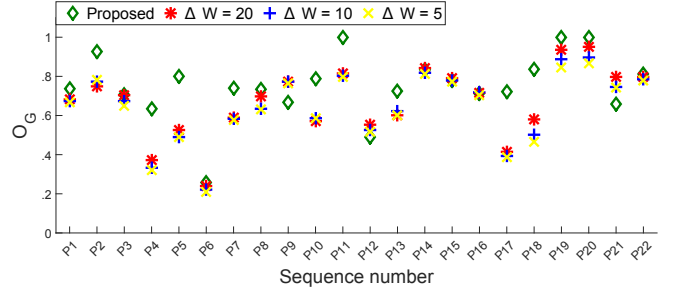


Fig. 9. Comparison of I_{ref} selected by the proposed approach and the original approach based on fixed temporal windows $\Delta W = 5, 10, 20$.

average tracking accuracy decreases with increasing number of trackers. There are two main reasons for this accuracy drop. First, low performing trackers in the *on-target* cluster decrease the overall tracking accuracy when fused using the average. This is indicated by results for P7, P9 and P15, where all trackers are *on-target* for most of the sequence. This can be further validated by a comparison of results with AvgF (Table V) for these sequences. Second, the number of splitting and merging of clusters naturally increases as we include more trackers, thus increasing the chances of wrong reverse-analysis evaluations. For P1, the target undergoes occlusions between frames 75-95. FCT, L1T and LSST lose the target due to occlusion, however FCT (a deterministic tracker) achieves the best performance score during this interval, reducing the overall accuracy. Results for TPF_5 and TPF_8 indicate this scenario. Similarly for P10, a drop in accuracy of TPF_5 , TPF_6 and TPF_7 occurs when the target undergoes occlusions between frames 20-35. All trackers lose the target at frame 20. However, LOT and AMS regain the target. The target remains stationary from frame 45 till the end of the sequence. This scenario allows failed trackers to achieve a higher performance score during reverse-analysis, hence reducing the tracking accuracy. TPF removes low-performing trackers to improve the overall tracking accuracy. Results for P3, P10, P13, P18 and P19 indicate that TPF outperforms all trackers. Furthermore TPF has similar performance to the best performing tracker(s) for the other sequences except for P6, P12, P17 and P22. TPF_3 achieves an overall improvement of 23%, 15%, 8%, 23%, 21%, 27%, 13% and 17% in O_G in comparison to the individual trackers ST, AFT, LOT, IVT, FCT, AMS, L1T and LSST respectively. Moreover, all other TPF configurations (TPF_4, \dots, TPF_8) also achieve better results as compared to all 8 trackers.

Fig. 10 compares tracker accuracy using O_G values for selected sequences. The target in P2 changes its pose, causing AMS and L1T to lose the target between frames 80-90. Both failing trackers at this point are discarded by TPF. Performance of FCT and LSST drops gradually after frame 130 due to background clutter. The performance of TPF_7 drops at frame 140 where the output is corrupted by low performing trackers (FCT and LSST), which are incorrectly determined as *on-target*, while other TPF configurations make use of the best performing trackers. For

TABLE V

MEAN OVERLAP SCORE COMPARISON IN TERMS OF O_G (TEN INDEPENDENT RUNS). KEY - ST: SPARSE TRACKER [42], AFT: ADAPTIVE FRAGMENTS BASED TRACKER [43], LOT: LOCALLY ORDERLESS TRACKER [44], IVT: INCREMENTAL VISUAL TRACKER [45], FCT: FAST COMPRESSIVE TRACKER [47], AMS: MEAN SHIFT TRACKER [46], L1T: L1 TRACKER [48], LSST: LEAST SOFT-THRESHOLD SQUARES TRACKER [49], AVGF: AVERAGE FUSION; SYMT: SYMBIOTIC TRACKER [17]; VTS: VISUAL TRACKER SAMPLER [13]; STR: STRUCK [50]; KCF: KERNELIZED CORRELATION FILTER TRACKER [3].

	Fused trackers								Proposed approaches								Selected state-of-the-art				
	ST	AFT	LOT	IVT	FCT	AMS	LIT	LSST	TPF_3	TPF_4	TPF_5	TPF_6	TPF_7	TPF_8	AvgF	SymT	VTS	STR	KCF		
P1	.65	.70	.64	.64	.14	.72	.20	.22	.72	.70	.53	.67	.67	.50	.25	.35	.24	.75	.21		
P2	.80	.76	.79	.82	.58	.23	.82	.50	.81	.80	.78	.74	.74	.73	.40	.65	.77	.68	.87		
P3	.40	.70	.73	.57	.40	.47	.79	.75	.70	.68	.74	.70	.82	.82	.78	.81	.77	.78	.85		
P4	.35	.43	.25	.26	.41	.30	.23	.31	.41	.41	.41	.42	.42	.42	.30	.33	.37	.39	.64		
P5	.81	.82	.31	.64	.39	.39	.38	.75	.76	.71	.51	.70	.47	.57	.64	.72	.59	.68	.85		
P6	.12	.14	.10	.13	.10	.29	.31	.10	.13	.12	.12	.12	.13	.12	.12	.12	.14	.13	.10		
P7	.91	.79	.78	.91	.81	.81	.92	.93	.85	.86	.81	.85	.85	.89	.89	.89	.85	.80	.83		
P8	.63	.72	.78	.43	.04	.60	.55	.41	.77	.77	.71	.65	.63	.58	.47	.62	.78	.74	.67		
P9	.75	.74	.89	.81	.74	.75	.89	.83	.77	.83	.81	.80	.81	.82	.84	.84	.92	.82	.74		
P10	.11	.15	.74	.11	.44	.70	.12	.12	.77	.75	.55	.48	.37	.76	.12	.12	.14	.91	.11		
P11	.75	.74	.91	.81	.75	.12	.90	.92	.86	.85	.83	.77	.80	.82	.69	.90	.84	.54	.76		
P12	.36	.52	.25	.35	.63	.33	.25	.35	.55	.40	.39	.40	.37	.35	.30	.33	.36	.28	.27		
P13	.60	.63	.42	.55	.67	.48	.51	.52	.71	.63	.70	.68	.61	.67	.53	.55	.56	.69	.44		
P14	.87	.68	.79	.85	.78	.73	.78	.87	.86	.83	.82	.76	.82	.76	.82	.83	.85	.82	.82		
P15	.74	.80	.78	.74	.77	.61	.90	.89	.81	.76	.79	.78	.81	.82	.82	.83	.78	.69	.79		
P16	.78	.79	.74	.56	.78	.80	.71	.92	.74	.85	.83	.82	.82	.89	.70	.77	.86	.78	.79		
P17	.26	.22	.65	.37	.20	.31	.53	.28	.26	.37	.31	.30	.40	.36	.43	.41	.24	.16	.84		
P18	.42	.33	.87	.10	.85	.79	.69	.86	.86	.89	.87	.81	.80	.85	.28	.38	.87	.84	.12		
P19	.86	.81	.77	.87	.87	.56	.78	.86	.89	.87	.87	.90	.89	.89	.87	.88	.80	.76	.89		
P20	.90	.79	.85	.89	.89	.58	.84	.79	.88	.85	.86	.85	.86	.85	.83	.85	.87	.78	.90		
P21	.75	.73	.84	.67	.73	.78	.83	.83	.77	.79	.77	.73	.84	.74	.78	.79	.85	.75	.74		
P22	.52	.42	.36	.52	.68	.80	.81	.13	.43	.54	.50	.51	.52	.50	.43	.54	.37	.13	.13		
Mean	.57	.61	.65	.57	.58	.55	.62	.60	.70	.69	.66	.66	.66	.67	.56	.61	.62	.63	.61		

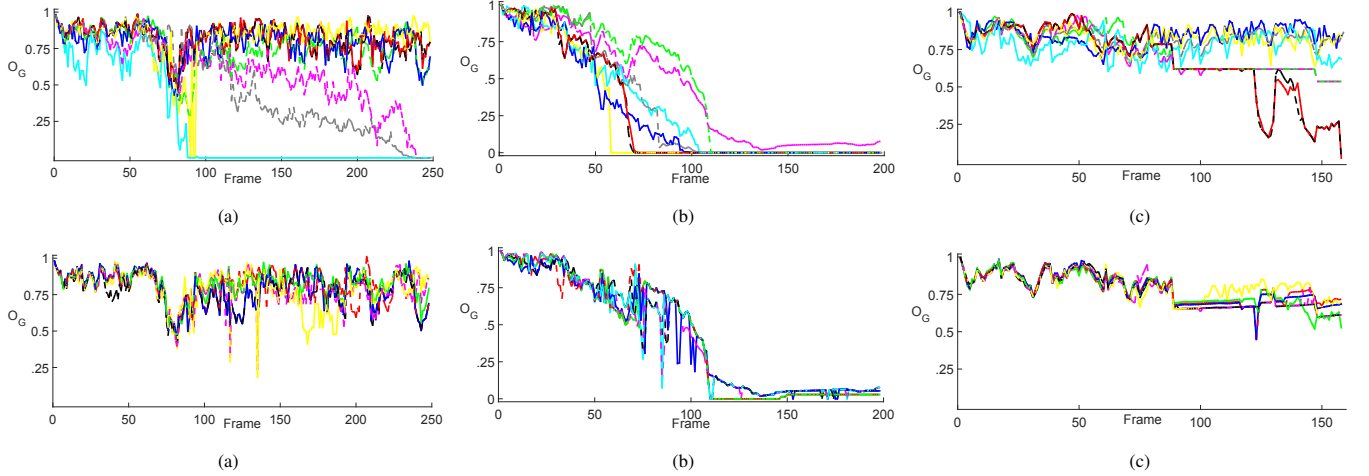


Fig. 10. O_G scores for trackers and TPF configurations under analysis for selected sequences. (a) Students-P2, (b) CAVIAR-P4, (c) MITTraffic-P21. Top row: Trackers; —: ST; - - -: AFT; —: LOT; - - -: IVT; —: AMS; - - -: FCT; —: L1T; —: LSST. Bottom row: TPF configurations; —: TPF_3 ; - - -: TPF_4 ; —: TPF_5 ; - - -: TPF_6 ; —: TPF_7 ; - - -: TPF_8 .

P4, all trackers lose the target between frames 60-110. All TPF configurations identify and achieve accuracy close to the best performing tracker (AFT). However TPF fails when all trackers are off-target. P21 undergoes scale changes as it moves away from the camera. O_G for all TPF configurations drops after frame 80, since all trackers remain *on-target* and form a single cluster. After frame 130, O_G for ST and IVT drops significantly since they cannot handle scale changes. However these trackers are discarded by TPF, while the performance for TPF_7 further improves as a better performing tracker (L1T) is added in the framework.

F. Comparison with the state-of-the-art (SOA) approaches

Table V compares the TPF configurations and the related SOA. AvgF and SymT have been tested using the eight trackers. STRUCK is the best for P1 and P10, achieving the best average results among the selected SOA approaches. KCF achieves the best results for P2, P3, P4, P5, P17 and P20. However, it is unable to handle occlusions as shown for P10, P12 and P22. SymT fails to determine a low performing tracker, hence reducing the overall tracking accuracy. It achieves good performance when most of the trackers are accurate as indicated by results for P7 and P15.

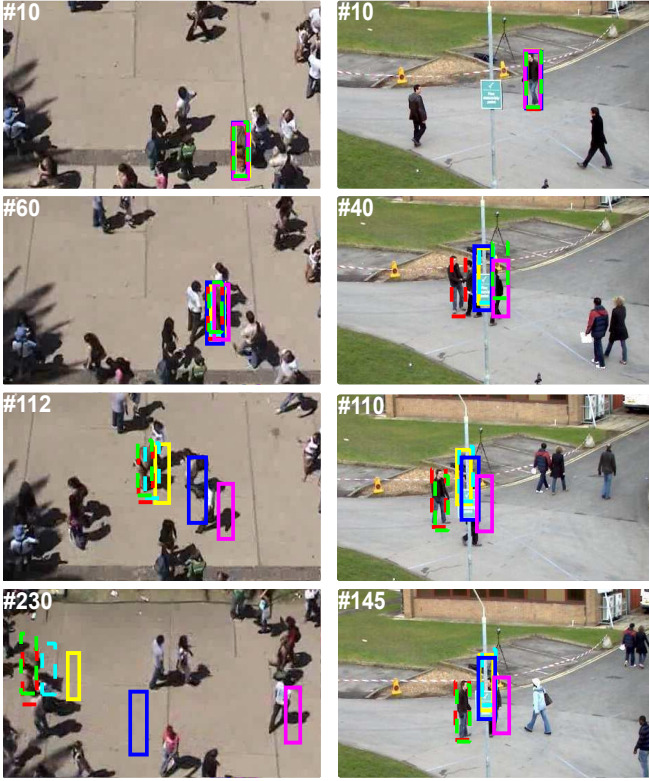


Fig. 11. Sample tracking results for Students-P1 (left column) and PETS-P10 (right column). - - -: TPF; - - -: STRUCK; —: VTS; - - -: SymT; —: AvgF; —: KCF.



Fig. 12. Sample tracking results for CAVIAR-P4 (left column) and PETS-P12 (right column). - - -: TPF; - - -: STRUCK; —: VTS; - - -: SymT; —: AvgF; —: KCF.

TPF on the other hand is able to use the best performing trackers and the overall accuracy is not dependent on the percentage of good trackers. VTS performs relatively well and shows the best results for P8, P9 and P21. However, it fails for P1 and P10 due to occlusions, and for P17 and P22 due to similarly coloured background. Although the state-of-the-art approaches outperform some employed trackers (ST, IVT, FCT and AMS, see Table V), TPF_3 shows an overall improvement of 23%, 15%, 13%, 11% and 15% in O_G in comparison to AvgF, SymT, VTS, STRUCK and KCF, respectively.

Sample tracking results for some sequences are shown in Fig. 11 and Fig. 12 where it can be seen that TPF correctly discards wrong trackers as they start to fail due to tracking challenges. For clarity we only present comparisons between TPF_3 and the SOA. Examples in Fig. 11 show that all trackers correctly follow the target at the beginning of the sequence. As target occlusions are more frequent, only STRUCK is able to perform similarly to TPF_3 as seen in frame 230 for P1 and frame 145 for P10. The right column of Fig. 12 depicts the situation where only TPF_3 is able to adapt to changes in target scale and occlusions whereas all the compared trackers fail, as seen in frame 85. The left column of Fig. 12 shows an example where none of trackers obtain accurate position estimations after an illumination change (frames 117 and 192) and the best trackers (KCF, STRUCK and TPF_3) achieve low accuracy.



Fig. 13. Average computational time for the stages of the proposed approach. For each configuration, the average is computed over the complete dataset and the total number of trackers.

G. Computational Cost

Fig. 13 presents the cost for the *trackers*, *tracker clustering* (Section IV) and *on-target cluster identification* (Section V) in terms of average computational time. The cost of the fusion stage is negligible and therefore ignored. The cost of the *trackers* considers running in parallel the trackers to fuse and depends on the employed approaches, being heavily influenced by the slowest tracker (LOT). The computational time for *tracker clustering* slightly increases with the number of trackers. Since *on-target cluster identification* uses reverse-analysis, the computational time becomes dependent on the trackers in the *on-target* cluster \mathbb{C}_t^* and the tracking challenges present in the sequence. This trend is also highlighted by the overall cost for the

TABLE VI

COMPUTATIONAL COST OF THE TRACKERS AND THE PROPOSED APPROACH (TPF) MEASURED AS FRAMES PER SECOND (FPS). KEY - ST: SPARSE TRACKER, AFT: ADAPTIVE FRAGMENTS BASED TRACKER; LOT: LOCALLY ORDERLESS TRACKER; IVT: INCREMENTAL VISUAL TRACKER; FCT: FAST COMPRESSIVE TRACKER; AMS: MEAN SHIFT TRACKER; LIT: L1 TRACKER; LSST: LEAST SOFT-THRESHOLD SQUARES TRACKER.

	Fused Trackers								Proposed approaches					
	ST [42]	AFT [43]	LOT [44]	IVT [45]	FCT [47]	AMS [46]	LIT [48]	LSST [49]	TPF_3	TPF_4	TPF_5	TPF_6	TPF_7	TPF_8
FPS	50.1±0.1	5.9±4.0	0.3±0.1	48.1±1.4	8.0±6.2	108.0±31.2	7.7±1.7	3.4±0.6	2.5±4.8	0.3±0.4	0.2±0.3	0.2±0.1	0.1±0.1	0.1±0.1

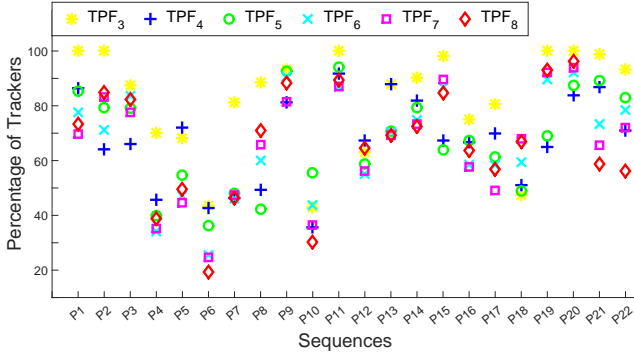


Fig. 14. Percentage of trackers used by the proposed approach for different tracker combinations.

TPF configurations presented in Table VI, where TPF_3 achieves the best computational cost. Fig. 14 shows the average number of trackers used by TPF, highlighting its advantage to cluster trackers, and using only the ones *on-target* for the various TPF combinations.

VII. CONCLUSIONS

We presented an approach to dynamically select and combine the results of successful (i.e. *on-target*) trackers in a decision-level fusion framework. The proposed approach determines relationships between trackers by analyzing the position and direction of movement of their estimated states. These spatio-temporal features are combined to estimate pair-wise tracker correlation scores that determine clusters of similarly performing trackers over time. An adaptive online evaluator identifies the trackers that are *on-target* and propagates them over time until a split or merge of this group (cluster) of trackers is detected. The final target state is estimated by fusing the outputs from the trackers that are in the *on-target* cluster. Experimental results show that the proposed approach outperforms state-of-the-art methods and the combined trackers. Moreover, the proposed time-reversed evaluation improves the original approach by using motion analysis and tracker performance to temporally update the reference frame.

As future work, we will include the performance weight of each tracker in the fusion stage.

REFERENCES

- [1] A. Smeulders, D. Chu, R. Cucchiara, S. Calderara, A. Dehghan, and M. Shah, "Visual tracking: An experimental survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 7, pp. 1442–1468, July 2014.
- [2] J. Wu, S. Hu, and Y. Wang, "Adaptive multifeature visual tracking in a probability-hypothesis-density filtering framework," *Signal Process.*, vol. 93, no. 11, pp. 2915 – 2926, April 2013.
- [3] J. Henriques, R. Caseiro, P. Martins, and J. Batista, "High-speed tracking with kernelized correlation filters," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 3, pp. 583–596, March 2015.
- [4] E. Erdem, S. Dubuisson, and I. Bloch, "Visual tracking by fusing multiple cues with context-sensitive reliabilities," *Pattern Recogn.*, vol. 45, no. 5, pp. 1948–1959, May 2012.
- [5] V. Badrinarayanan, P. Perez, F. Le Clerc, and L. Oisel, "Probabilistic color and adaptive multi-feature tracking with dynamically switched priority between cues," in *IEEE Int. Conf. Comput. Vis.*, Oct 2007, pp. 1–8.
- [6] M. Spengler and B. Schiele, "Towards robust multi-cue integration for visual tracking," *Machine Vis. and Appl.*, vol. 14, no. 1, pp. 50–58, April 2003.
- [7] I. Leichter, M. Lindenbaum, and E. Rivlin, "A general framework for combining visual trackers: the black boxes approach," *Int. J. Comput. Vision*, vol. 67, pp. 343–363, May 2006.
- [8] X. Zhang, W. Hu, H. Bao, and S. Maybank, "Robust head tracking based on multiple cues fusion in the kernel-bayesian framework," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 23, no. 7, pp. 1197–1208, July 2013.
- [9] E. Maggio and A. Cavallaro, "Hybrid particle filter and mean shift tracker with adaptive transition model," in *IEEE Int. Conf. on Acoustics, Speech, and Signal Process.*, vol. 2, March 2005, pp. 221–224.
- [10] C. Chang, R. Ansari, and A. Khokhar, "Multiple object tracking with kernel particle filter," in *IEEE Conf. on Comput. Vis. and Pattern Recogn.*, vol. 1, June 2005, pp. 566–573.
- [11] J. Santner, C. Leistner, A. Saffari, T. Pock, and H. Bischof, "Prost: Parallel robust online simple tracking," in *IEEE Conf. on Comput. Vis. and Pattern Recogn.*, June 2010, pp. 723–730.
- [12] J. Yoon, D. Kim, and K.-J. Yoon, "Visual tracking via adaptive tracker selection with multiple features," in *Eur. Conf. on Comput. Vis.*, Oct 2012, pp. 28–41.
- [13] J. Kwon and K. Lee, "Tracking by sampling and integrating multiple trackers," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 7, pp. 1428–1441, July 2014.
- [14] M. Heber, M. Godec, M. Ruther, P. Roth, and H. Bischof, "Segmentation-based tracking by support fusion," *Comput. Vis. and Image Understanding*, vol. 117, no. 6, pp. 573–586, June 2013.
- [15] N. Siebel and S. Maybank, "Fusion of multiple tracking algorithms for robust people tracking," in *Eur. Conf. on Comput. Vis.*, May 2002, pp. 1–15.
- [16] F. Moreno-Noguer, A. Sanfeliu, and D. Samaras, "Dependent multiple cue integration for robust tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 4, pp. 670–685, April 2008.
- [17] Y. Gao, R. Ji, L. Zhang, and A. Hauptmann, "Symbiotic tracker ensemble towards a unified tracking framework," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24, no. 7, pp. 1122 – 1131, July 2014.
- [18] T. Biresaw, A. Cavallaro, and C. Regazzoni, "Tracker-level fusion for robust bayesian visual tracking," in *IEEE Trans. Circuits Syst. Video Technol.*, vol. 25, no. 5, May 2015, pp. 776–789.
- [19] Y. Wu, J. Lim, and M.-H. Yang, "Online object tracking: A benchmark," in *IEEE Conf. on Comput. Vis. and Pattern Recogn.*, June 2013, pp. 2411–2418.
- [20] J. SanMiguel, A. Cavallaro, and J. Martinez, "Adaptive online performance evaluation of video trackers," *IEEE Trans. Image Process.*, vol. 21, no. 5, pp. 2812–2823, May 2012.
- [21] K. Shearer, K. Wong, and S. Venkatesh, "Combining multiple tracking algorithms for improved general performance," *Pattern Recogn.*, vol. 34, no. 6, pp. 1257–1269, June 2001.

- [22] E. Maggio, F. Smerladi, and A. Cavallaro, "Adaptive multifeature tracking in a particle filtering framework," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 17, no. 10, pp. 1348–1359, Oct 2007.
- [23] H. Wu, A. Sankaranarayanan, and R. Chellappa, "Online empirical evaluation of tracking algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 8, pp. 1443–1458, Aug 2010.
- [24] J. Triesch and C. Malsburg, "Democratic integration: Self-organized integration of adaptive cues," *Neural Comput.*, vol. 13, no. 9, pp. 2049–2074, Sept 2001.
- [25] Z. Kalal, K. Mikolajczyk, and J. Matas, "Forward-backward error: Automatic detection of tracking failures," in *Int. Conf. on Pattern Recogn.*, Aug 2010, pp. 2756–2759.
- [26] T. Biresaw, A. Cavallaro, and C. Regazzoni, "Correlation-based self-correcting tracking," *Neurocomputing*, vol. 152, no. 1, pp. 345–358, March 2015.
- [27] S. Birchfield, "Elliptical head tracking using intensity gradients and color histograms," in *IEEE Comput. Society Conf. on Comput. Vis. and Pattern Recogn.*, June 1998, pp. 232–237.
- [28] S. Zhang, H. Yao, X. Sun, and X. Lu, "Sparse coding based visual tracking: Review and experimental comparison," *Pattern Recogn.*, July 2013.
- [29] T. Zhang, B. Ghanem, S. Liu, and N. Ahuja, "Robust visual tracking via multi-task sparse learning," in *IEEE Conf. on Comput. Vis. and Pattern Recogn.*, June 2012, pp. 2042–2049.
- [30] Z. Hong, X. Mei, D. Prokhorov, and D. Tao, "Tracking via robust multi-task multi-view joint sparse representation," in *IEEE Int. Conf. Comput. Vis.*, Dec 2013, pp. 649–656.
- [31] Q. Li, X. Wang, W. Wang, Y. Jiang, Z. Zhou, and Z. Tu, "Disagreement-based multi-system tracking," in *Asian Conf. on Comput. Vis.*, vol. 7729, Nov 2012, pp. 320–334.
- [32] B. Stenger, T. Woodley, and R. Cipolla, "Learning to track with multiple observers," in *IEEE Conf. on Comput. Vis. and Pattern Recogn.*, June 2009, pp. 2647–2654.
- [33] B. Zhong, H. Yao, S. Chen, R. Ji, T. Chin, and H. Wang, "Visual tracking via weakly supervised learning from multiple imperfect oracles," *Pattern Recogn.*, vol. 47, no. 3, pp. 1395 – 1410, March 2014.
- [34] J. Whitehill, P. Ruvolo, T. Wu, J. Bergsma, and J. Movellan, "Whose vote should count more: Optimal integration of labels from labelers of unknown expertise," in *Proc. of Neural Information Process. Systems*, Dec 2009, pp. 2035–2043.
- [35] C. Bailer, A. Pagani, and D. Stricker, "A superior tracking approach: Building a strong tracker through fusion," in *Eur. Conf. on Comput. Vis.*, vol. 8695, Sept 2014, pp. 170–185.
- [36] A. Sharkey, N. Sharkey, U. Gerecke, and G. Chandroth, "The test and select approach to ensemble combination," in *Mult. Classifier Sys.*, June 2000, vol. 1857, pp. 30–44.
- [37] N. Anjum and A. Cavallaro, "Multifeature object trajectory clustering for video analysis," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 18, no. 11, pp. 1555–1564, Nov 2008.
- [38] J. Lucas, *Introduction to Abstract Mathematics*. Rowman and Littlefield, 1990.
- [39] J. Conway and R. Guy, *The Book of Numbers*. Springer Science and Business Media, 1996.
- [40] F. Jiang, Y. Wu, and A. Katsaggelos, "A dynamic hierarchical clustering method for trajectory-based unusual video event detection," *IEEE Trans. Image Process.*, vol. 18, no. 4, pp. 907–913, April 2009.
- [41] M. Yang, Y. Wu, and G. Hua, "Context-aware visual tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 7, pp. 1195–1209, July 2009.
- [42] D. Wang, H. Lu, and M.-H. Yang, "Online object tracking with sparse prototypes," *IEEE Trans. Image Process.*, vol. 22, no. 1, pp. 314–325, Jan 2013.
- [43] E. Erdem, S. Dubuisson, and I. Bloch, "Fragments based tracking with adaptive cue integration," *Comput. Vis. and Image Understanding*, vol. 116, no. 7, pp. 827 – 841, July 2012.
- [44] S. Oron, A. Bar-Hillel, D. Levi, and S. Avidan, "Locally orderless tracking," *Int. J. Comput. Vision*, vol. 111, no. 2, pp. 213–228, Jan 2014.
- [45] D. Ross, J. Lim, R. Lin, and M. Yang, "Incremental learning for robust visual tracking," *Int. J. Comput. Vision*, vol. 77, no. 1-3, pp. 125–141, May 2008.
- [46] J. Ning, L. Zhang, D. Zhang, and C. Wu, "Scale and orientation adaptive mean shift tracking," *IET Comput. Vis.*, vol. 6, no. 1, pp. 52–61, Jan 2012.
- [47] K. Zhang, L. Zhang, and M. Yang, "Fast compressive tracking," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 10, pp. 2002–2015, Oct 2014.
- [48] X. Mei and H. Ling, "Robust visual tracking and vehicle classification via sparse representation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 11, pp. 2259–2272, Nov 2011.
- [49] D. Wang, H. Lu, and M. Yang, "Least soft-threshold squares tracking," in *IEEE Conf. on Comput. Vis. and Pattern Recogn. (CVPR)*, June 2013, pp. 2371–2378.
- [50] S. Hare, A. Saffari, and P. Torr, "Struck: Structured output tracking with kernels," in *IEEE Int. Conf. Comput. Vis.*, Nov 2011, pp. 263–270.



ObaidUllah Khalid received his MS degree in Network Engineering from Illinois Institute of Technology (IIT), Chicago in 2008 and BS degree in Computer Engineering from SirSyed University of Engineering and Technology (SSUET) in 2004. He served as a lecturer in National University of Sciences and Technology (NUST), Pakistan. Since December 2012, he has been with Queen Mary University of London, UK, and Alpen-Adria Universitat Klagenfurt, Austria, as a PhD researcher under the supervision of Prof. A. Cavallaro and Prof. B. Rinner. He was awarded National University of Sciences and Technology (NUST) scholarship for his MS studies and Erasmus Mundus fellowship for his Double Doctorate in Interactive and Cognitive Environments. His research interests include tracker-level fusion, online performance evaluation and video tracking.



Juan Carlos SanMiguel received the PhD degree in Computer Science and Telecommunication in 2011 at Universidad Autnoma de Madrid (Spain). From June 2013 to June 2014, he was a postdoctoral researcher at Queen Mary University of London (UK) under a Marie Curie IAPP fellowship. He is currently assistant professor at Universidad Autnoma de Madrid (<http://www.uam.es>) and researcher at the Video Processing and Understanding Lab (<http://www-vpu.eps.uam.es>). His current research interests are in the field of computer vision, focused on online performance evaluation and multi-camera activity understanding for video segmentation and tracking.



Andrea Cavallaro is Professor of Multimedia Signal Processing and Director of the Centre for Intelligent Sensing at Queen Mary University of London, UK. He received his Ph.D. in Electrical Engineering from the Swiss Federal Institute of Technology (EPFL), Lausanne, in 2002. He was a Research Fellow with British Telecommunications (BT) in 2004/2005 and was awarded the Royal Academy of Engineering teaching Prize in 2007; three student paper awards on target tracking and perceptually sensitive coding at IEEE ICASSP in 2005, 2007 and 2009; and the best paper award at IEEE AVSS 2009. Prof. Cavallaro is Associate Editor for the IEEE Transactions on Circuits and Systems for Video Technology and member of the editorial board of IEEE Multimedia. He is a past Area Editor for IEEE Signal Processing Magazine and a past Associate Editor for the IEEE Transactions on Image Processing, IEEE Transactions on Multimedia, IEEE Transactions on Signal Processing, and IEEE Signal Processing Magazine. He has published over 160 journal and conference papers, one monograph on Video tracking (2011, Wiley) and three edited books: Multi-camera networks (2009, Elsevier); Analysis, retrieval and delivery of multimedia content (2012, Springer); and Intelligent multimedia surveillance (2013, Springer).