

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/44676013>

Tracking and Activity Recognition Through Consensus in Distributed Camera Networks

Article in IEEE Transactions on Image Processing · October 2010

DOI: 10.1109/TIP.2010.2052823 · Source: PubMed

CITATIONS

92

READS

104

6 authors, including:



B. Song

University of California, Riverside

39 PUBLICATIONS **1,185** CITATIONS

[SEE PROFILE](#)



Tuba Kamal

University of Karachi

98 PUBLICATIONS **964** CITATIONS

[SEE PROFILE](#)



Cristian Soto

Universidad de la Amazonia

6 PUBLICATIONS **267** CITATIONS

[SEE PROFILE](#)



Chong Ding

HRL Laboratories, LLC

17 PUBLICATIONS **367** CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Self-organizing Approximation Based Adaptive and Optimal Control [View project](#)



Digital Media Forensic [View project](#)

Tracking and Activity Recognition Through Consensus in Distributed Camera Networks

Bi Song, *Member, IEEE*, Ahmed T. Kamal, *Student Member, IEEE*, Cristian Soto, Chong Ding, *Student Member, IEEE*, Amit K. Roy-Chowdhury, *Senior Member, IEEE*, and Jay A. Farrell, *Fellow, IEEE*

Abstract—Camera networks are being deployed for various applications like security and surveillance, disaster response and environmental modeling. However, there is little automated processing of the data. Moreover, most methods for multi-camera analysis are centralized schemes that require the data to be present at a central server. In many applications, this is prohibitively expensive, both technically and economically. In this paper, we investigate distributed scene analysis algorithms by leveraging upon concepts of consensus that have been studied in the context of multi-agent systems, but have had little applications in video analysis. Each camera estimates certain parameters based on its own sensed data which is then shared locally with the neighboring cameras in an iterative fashion, and a final estimate is arrived at in the network using consensus algorithms. We specifically focus on two basic problems - tracking and activity recognition. For multi-target tracking in a distributed camera network, we show how the Kalman-Consensus algorithm can be adapted to take into account the directional nature of video sensors and the network topology. For the activity recognition problem, we derive a probabilistic consensus scheme that combines the similarity scores of neighboring cameras to come up with a probability for each action at the network level. Thorough experimental results are shown on real data along with a quantitative analysis.

I. INTRODUCTION

Networks of video cameras are being installed in many applications, *e.g.*, surveillance and security, disaster response, environmental monitoring, etc. Currently, most of the data collected by such networks is analyzed manually, a task that is extremely tedious and reduces the potential of the installed networks. Therefore, it is essential to develop tools for analyzing the data collected from these cameras and summarizing the results in a manner that is meaningful to the end user. Tracking and activity recognition are two fundamental tasks in this regard. In this paper, we develop methods for tracking and activity recognition in a distributed network of cameras.

For many applications, for a number of reasons it is desirable that the video analysis tasks be decentralized. For example, there may be constraints of bandwidth, secure transmission, and difficulty in analyzing a huge amount of data centrally. In such situations, the cameras would have to act

as autonomous agents making decisions in a decentralized manner. At the same time, however, the decisions of the cameras need to be coordinated so that there is a consensus on the state (*e.g.*, position, activity) of the target even if each camera is an autonomous agent. Thus, the cameras, acting as autonomous agents, analyze the raw data locally, exchange only distilled information that is relevant to the collaboration, and reach a shared, global analysis of the scene.

Although there are a number of methods in video analysis that deal with multiple cameras, and even camera networks, *distributed* processing in camera networks has received very little attention. In Sec. II, we will review the current state of the art in camera networks and will see that very few methods are capable of distributed analysis of video. On the other hand, distributed processing has been extensively studied in the multi-agent systems and cooperative control literature [29]. Methods have been developed for reaching consensus on a state observed independently by multiple sensors. However, there is very little study on the applicability of these methods in camera networks.

In this paper, we show how to develop methods for tracking and activity recognition in a camera network where processing is distributed across the cameras. For this purpose, we show how consensus algorithms can be developed that are capable of converging to a solution, *i.e.*, target state, based on local decision making and exchange of these decisions (not sensed data) among the cameras. We focus on two problems. For distributed tracking, we show how the Kalman consensus algorithm [28] can be adapted to camera networks taking into account issues like network topology, handoff and fault tolerance. For activity recognition, we derive a new consensus algorithm based on the recognized activity at each camera and the transition probabilities between various activities. Experimental results and quantitative evaluation for both these methods are presented. Note that here we assume ideal communication between cameras which are connected, *i.e.*, communication is not a bottleneck. This proposed work is a proof-of-concept study in using distributed processing algorithms for video analysis. In the future, the practical constraints of using consensus algorithms in camera networks should be considered.

We start with a review of consensus algorithms for distributed estimation. Thereafter, in Sec. IV, we present a variant of the Kalman-Consensus approach for distributed tracking in the camera network and show experimental results, that are analyzed quantitatively. In Sec. V, we study the problem of activity recognition in a consensus framework. For this

Copyright (c) 2010 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

The authors gratefully acknowledge support from NSF grants ECS-0622176 and CNS-0551741, ONR grant N00014-09-1-0666, and ARO grant W911NF-07-1-0485.

B. Song, A. T. Kamal, C. Ding, A. K. Roy-Chowdhury and J. A. Farrell are with University of California, Riverside, CA, 92521 USA (e-mail: {bsong, akamal, amitrc, farrell}@ee.ucr.edu, cding@cs.ucr.edu).

C. Soto is with Digital Western, the work was done when the author was at University of California, Riverside.

purpose, we derive a completely new algorithm that shows how local decisions at each camera node can be combined to come up with a consensus on the state representing the activity. Again, experimental results are shown and analyzed.

II. PAST WORK ON SCENE ANALYSIS IN CAMERA NETWORKS

Our review of scene analysis algorithms will be limited to those directly related to the application domain of camera networks.

There have been a few papers in the recent past that deal with networks of video sensors. Particular interest has been focused on learning a network topology [21], [40], i.e., configuring connections between cameras and entry/exit points in their view. Some of the existing methods on tracking over the network, include [34], [37]. Other interesting problems in camera networks, like object/behavior detection and matching across cameras, camera handoff and camera placement have been addressed in [1], [10], [16], [39], [46]. There has also been recent work on tracking people in a multi-camera setup [8], [17]. However, these methods do not address the issue of distributed processing.

In [22], a distributed target tracking approach using a cluster-based Kalman filter was proposed. Here, a camera is selected as a cluster head which aggregates all the measurements of a target to estimate its position using a Kalman filter and sends that estimate to a central base station. Our proposed tracking system differs from this method in that each camera has a consensus-based estimate of the target's state and thus there is no need for additional computation and communication to select a cluster head. As will be described in Section IV, we apply in a special way the distributed Kalman-Consensus filter [28] which has been shown to be more effective than other distributed Kalman filter schemes. Consensus schemes have been gaining popularity in computer vision applications involving multiple cameras [41]. A related work that deals with tracking targets in a camera network with PTZ cameras is [33]. Here, the authors proposed a mixture between a distributed and a centralized scheme using both static and PTZ cameras in a virtual camera network environment. Our approach to tracking in the camera network, however, is completely distributed using consensus algorithms. Another problem that has received some attention in this context is the development of distributed embedded smart cameras [3]. The focus of this paper, however, is on the algorithm side, rather than building a specific smart camera architecture.

The problem of multi-view activity recognition have been addressed in many papers, e.g., [44], [45], but the information of multiple views is fused centrally. Our proposed framework is decentralized: each camera determines a probabilistic measure of similarity of its own observed activities to a pre-defined dictionary and information is dispersed to compute a consensus-based estimate. A preliminary framework for distributed tracking and control in camera network was presented in [38]. However, instead of only considering target-based network topology [38], in this paper we also define a network topology based on communication constraints which is more

important in practice. Besides tracking through consensus, we also address another fundamental task of distributed activity recognition, derive a probabilistic consensus scheme, and show experimental results on real data with a quantitative analysis.

III. CONSENSUS ALGORITHMS FOR DISTRIBUTED ESTIMATION

In the multi-agent systems literature, *consensus* means that the agents reach an agreement regarding a certain quantity of interest that depends on the measurements of all sensors in a network. The network may not be fully connected, so there is no central unit that has access to all the data from the sensors. Consequently, a *consensus algorithm* is an interaction rule that specifies information exchange between a sensor and its neighbors that guarantees that all the nodes reach a consensus. The interaction topology of a network of sensors is represented using a graph $G = (V, E)$ with the set of nodes $V = \{1, 2, \dots, n\}$ and edges $E \subseteq V \times V$. Each sensor node $i = 1, \dots, n$ maintains an estimate $\mathbf{x}_i \in \mathbb{R}^m$ of a quantity $\mathbf{x} \in \mathbb{R}^m$. Consensus is achieved when $\mathbf{x}_1 = \mathbf{x}_2 = \dots = \mathbf{x}_n$, which is an n -dimensional subspace of \mathbb{R}^{mn} . A thorough review of consensus in networked multi-agent systems can be found in [29]. Here we briefly review some of the basic approaches needed for this paper.

A. Brief Review

In a network of agents, consensus can be defined as reaching an agreement through cooperation regarding a certain quantity of interest that depends on the information available to measurements from all agents. An interaction rule that specifies the information exchange between an agent and all of its neighbors in the network and the method by which the information is used, is called a consensus algorithm (or protocol). Cooperation means giving consent to providing one's state and following a common protocol that serves group objective.

For example, in a network of temperature sensor, the sensors' estimates of temperature could be different due to sense noise and local variation. The sensors then interchange information with their neighboring sensors, and use the information to refine their local estimates. Consensus is reached when all sensors agree on a single value.

Distributed computing [20] has been a challenging field in computer science for the last few decades. A lot of work has been done on consensus algorithms which formed the baseline for distributed computing. Formally the study of consensus originated in management science and statistics in 1960s (see [6]). The work in [42] on asynchronous asymptotic agreement problems in distributed decision making systems and parallel computing [2] were the initial works in systems and control theory on a distributed network. A theoretical framework for defining and solving consensus problems for networked dynamic systems was introduced in [30] building on the earlier work of [11]. Consensus algorithms for reaching an agreement without computing any objective function appeared in the work of [15]. Further theoretical extensions of this work were presented in [35] with a focus towards treatment of directed

information flow in networks. In [15], a formal analysis was provided for emergence of alignment. The setup in [30] was originally created with the vision of designing agent-based amorphous computers for collaborative information processing in networks. Later, [30] was used in development of flocking algorithms with guaranteed convergence and the capability to deal with obstacles and adversarial agents [27]. Recent works related to multi agent networked systems include consensus [19], collective behavior of flocks and swarms [27], sensor fusion [28], random networks [13], synchronization of coupled oscillators [32], algebraic connectivity of complex networks [26], asynchronous distributed algorithms [23], formation control for multi robot systems [9], dynamic graphs [24], and complexity of coordinated tasks [14].

The goals of most consensus algorithms usually include [12]:

- 1. Validity:** The final answer that achieves consensus is a valid answer.
- 2. Agreement:** All processes agree as to what the agreed upon answer was by the end of the process.
- 3. Termination:** The consensus process eventually ends with each process contributing.
- 4. Integrity:** Processes vote only once.

Many consensus algorithms contain a series of events (and related messages) during a decision-making round. Typical events include Proposal and Decision. Here, proposal typically means the communication of the state of each agent and decision is the process of an agent deciding on proposals received from its neighbors after which it is not going to receive any proposal from the neighbors to come a different conclusion. In our application domain of camera networks, the agents are the cameras and the state vector we are trying to estimate are the position and velocity of a set of targets and the ID of an activity based on a learned dictionary of activities.

B. Consensus in Distributed Camera Networks

In distributed camera networks, the cameras act as autonomous agents. Each camera determines its own estimate of the object's state (e.g., position, activity label). The cameras then share local estimates with their neighboring cameras in an iterative fashion, and a final estimate is arrived at in the network using consensus algorithms [29].

1) *Distributed Tracking:* There have been recent attempts to achieve dynamic state estimation in a consensus-like manner. In contrast to a central Kalman filter where state information coming from several sensors is fused in a central station, Distributed Kalman Filters (DKF) compute a consensus-based estimate on the state of interest with only point-to-point communication between the sensors [28]. A distributed Kalman filtering (DKF) strategy that obtains consensus on state estimates was presented in [28]. The overall performance of this so-called Kalman-Consensus filter has been shown to be superior to other distributed approaches. It is on this DKF strategy that we base our distributed tracking algorithm. The mathematical details are presented in Section IV-A.

2) *Distributed Activity Recognition:* There have been methods on multi-view activity recognition [44], [45], but the

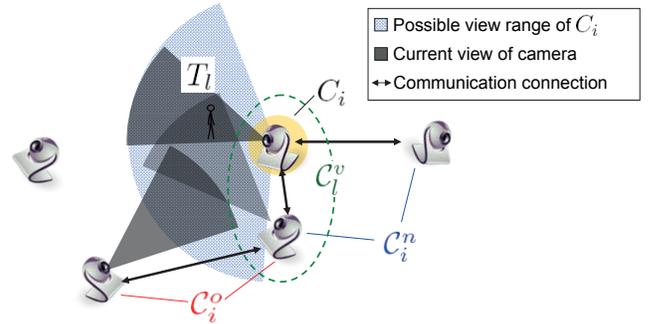


Fig. 1. Conceptual illustration of camera network topologies. $C_i^v \subset \mathcal{C}$ is the subset of all cameras viewing target T_l and the rest of the cameras are $C_i^v \subset \mathcal{C}$. $C_i^n \subset \mathcal{C}$ is the set of *neighboring cameras* of C_i and defined as all the cameras with which C_i is able to communicate. $C_i^o \subset \mathcal{C}$ is the set of *overlapping cameras* of C_i , and is defined as all the cameras with which C_i can *potentially* have an overlapping field of view.

information of multiple views is fused centrally. In this paper, we propose a framework for distributed activity recognition. Each camera determines a probabilistic measure of similarity of its own observed activities to a pre-defined dictionary, and then disperses this information to compute a consensus-based estimate with only point-to-point communication between the cameras. We show mathematically how to compute this consensus based on the similarity score computed at each camera and the transition probabilities between activities (can be uniform if no prior information is available).

IV. DISTRIBUTED TARGET TRACKING USING KALMAN CONSENSUS FILTERING

In this section, we present the first major result of this paper - how to track multiple targets in a camera network using a consensus algorithm that relies on the tracks obtained at individual cameras. For this purpose, we leverage upon the Kalman-Consensus algorithm in the distributed processing and multi-agent systems literature [28], [29]. However, there are some major differences due to the nature of cameras, and we show how to handle them.

Cameras are directional sensors and thus geographically neighboring cameras may be viewing very different portions of the scene. On the other hand, cameras that are geographically far away may be observing the same target. Therefore, we can define a target-based network topology, where the neighborhood structure is defined with respect to each target. Since targets are dynamic, this target-based topology changes over time. However, the communication constraints due to bandwidth limitation or physical network connection, which is most important in practice, naturally determine the communication-based topology of network. The communication-based topology is somewhat static, since the bandwidth limitation or physical connection won't change in a short time period. The distributed tracking is achieved by considering both the communication and target-based network topologies. In the next section, we will describe this process in more detail. Also, we will show how to take into account the handoff of targets as they move between cameras.

A. Problem Formulation

Let \mathcal{C} be the set of all cameras in the network. We can then define the subset of all cameras viewing target T_l as $\mathcal{C}_l^v \subset \mathcal{C}$ and the rest of the cameras as $\mathcal{C}_l^{v-} \subset \mathcal{C}$. Each camera C_i will also have its set of *neighboring cameras* $\mathcal{C}_i^n \subset \mathcal{C}$. Based on the communication constraints due to bandwidth limitation and network connection, we define the set \mathcal{C}_i^n as all the cameras with which C_i is able to communicate directly. In other words, C_i can assume that no other cameras other than its neighbors \mathcal{C}_i^n exist as no information flows directly from non-neighboring cameras to C_i . Note that the set of neighbors need not be geographical neighbors. We also define the set of *overlapping cameras* of C_i as $\mathcal{C}_i^o \subset \mathcal{C}$; since all the cameras can change their PTZ parameters and have therefore several possible fields of view, we define the set \mathcal{C}_i^o as all the cameras with which C_i can *potentially* have an overlapping field of view. By definition, it becomes clear then that for each $C_i \in \mathcal{C}_l^v$, it is true that $\mathcal{C}_l^v \subset \{\mathcal{C}_i^o \cup C_i\}$. We define $\mathcal{C}_i^c \subset \mathcal{C}$ as the connected component that C_i is in. We assume $\mathcal{C}_i^o \subset \mathcal{C}_i^c$, that is to say, C_i is able to exchange information with its overlapping cameras directly or via other cameras (**Assumption ***). An example of the camera network is shown in Figure 1.

As mentioned earlier, we propose a special application of the Kalman-Consensus Filter presented in [28] to solve the problem of finding a consensus on the state vectors of multiple targets in a camera network. We consider the situation where targets are moving on a ground plane and a homography between each camera's image plane and the ground plane is known. We will show how the state vector estimation for each target by each camera (i.e., each camera's estimates based on its individual measurements) can be combined together through the consensus scheme. This method is independent of the tracking scheme employed in each camera, which may be ~~Kalman-Consensus~~ **not based on the Kalman filter**.

To model the motion of a target T_l on the ground plane as observed by camera C_i , we consider a linear dynamical system with time propagation and observation models:

$$\mathbf{x}^l(k+1) = \mathbf{A}^l(k)\mathbf{x}^l(k) + \mathbf{B}^l(k)\mathbf{w}^l(k); \quad \mathbf{x}^l(0) \quad (1)$$

$$\mathbf{z}_i^l(k) = \mathbf{F}_i^l(k)\mathbf{x}^l(k) + \mathbf{v}_i^l(k) \quad (2)$$

where $\mathbf{w}^l(k)$ and $\mathbf{v}_i^l(k)$ are zero mean white Gaussian noise ($\mathbf{w}^l(k) \sim \mathcal{N}(0, \mathbf{Q}^l)$, $\mathbf{v}_i^l(k) \sim \mathcal{N}(0, \mathbf{R}_i^l)$) and $\mathbf{x}^l(0) \sim \mathcal{N}(\mathbf{x}_0^l, \mathbf{P}_0)$ is the initial state of the target. We define the state of the target at time step k as $\mathbf{x}^l(k) = (x^l(k), y^l(k), \dot{x}^l(k), \dot{y}^l(k))^T$ where $(x^l(k), y^l(k))$ and $(\dot{x}^l(k), \dot{y}^l(k))$ are the position and velocity of target T_l in the x and y directions respectively. The vector \mathbf{x}_i^l is the state of target T_l by C_i based on the measurements in C_i only. The vector $\mathbf{z}_i^l(k)$ is the noisy measurement at camera C_i . $\mathbf{z}_i^l(k)$ can be measured on either ground plane or image plane. We consider both cases and can show that it doesn't affect the performance of distributed Kalman-Consensus tracking algorithm, i.e., these two different cases of $\mathbf{z}_i^l(k)$ give equivalent tracking results.

Case 1: $\mathbf{z}_i^l(k)$ is the sensed target position $(x_i^l(k), y_i^l(k))$ on the ground plane based on the pre-computed homography between image plane and ground plane. We use a subscript

g to represent this case, i.e.,

$$\begin{aligned} (\mathbf{z}_i^l)_g(k) &= (\mathbf{F}_i)_g \mathbf{x}^l(k) + (\mathbf{v}_i^l)_g(k) \\ \text{and } (\mathbf{F}_i)_g &= \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}. \end{aligned} \quad (3)$$

Case 2: $\mathbf{z}_i^l(k)$ is the observed target position on the image plane of C_i . To differentiate with Case 1, we use a subscript c , i.e.,

$$\begin{aligned} (\mathbf{z}_i^l)_c(k) &= (\mathbf{F}_i)_c \mathbf{x}^l(k) + (\mathbf{v}_i^l)_c(k) \\ \text{and } (\mathbf{F}_i)_c &= \begin{bmatrix} (f_{11})_i & (f_{12})_i & 0 & 0 \\ (f_{21})_i & (f_{22})_i & 0 & 0 \end{bmatrix} = [\tilde{F}_i \quad \mathbf{0}], \end{aligned} \quad (4)$$

where $\tilde{F}_i = \begin{bmatrix} (f_{11})_i & (f_{12})_i \\ (f_{21})_i & (f_{22})_i \end{bmatrix}$, $\tilde{F}_i \in \{\mathbb{F} : \mathbb{R}^2 \rightarrow \mathbb{R}^2\}$ denotes the mapping from ground plane to the image plane of C_i ¹. \square

Our special implementation of the Kalman-Consensus distributed tracking algorithm is presented in Algorithm 1. We describe it for the general system model of Equations (1) and (2) and is applicable for the two special cases described above. This algorithm is performed in a distributed fashion by each camera node C_i . At each time step k and for each target T_l , we assume we are given the prior estimated target state $\bar{\mathbf{x}}_i^l$ and the error covariance matrix \mathbf{P}_i^l at k using measurements up to and including time $(k-1)$. At time step $k=0$, the Kalman-Consensus filter is initialized with $\mathbf{P}_i^l = \mathbf{P}_0$ and $\bar{\mathbf{x}}_i^l = \bar{\mathbf{x}}_0^l =$ average of $(\mathbf{z}^l)_g(0)$'s of cameras viewing T_l .

Comparing with the Kalman filter with centralized fusion (i.e., all the cameras send their measurements to a central processor, and tracking is preformed centrally, see Appendix A), we can see the fundamentals of Kalman-Consensus tracking algorithm described in Algorithm 1. If C_i is viewing a target T_l , it obtains T_l 's measurement \mathbf{z}_i^l and computes the corresponding information vector \mathbf{u}_i^l and matrix \mathbf{U}_i^l . Similar to [31], we define the information matrix and vector of $C_i \in \mathcal{C}_l^{v-}$ as $\mathbf{U}_i^l = 0$ and $\mathbf{u}_i^l = 0$ by assuming that their output matrices are zero, i.e., $\mathbf{F}_i^l = 0$ for all $C_i \in \mathcal{C}_l^{v-}$ to avoid any ambiguity arising from the lack of measurements in these cameras. If $C_i \in \mathcal{C}_l^v$ and the communication graph for \mathcal{C}_l^v is fully connected, such that C_i can receive information from all the other cameras viewing the same target, by fusing information vectors and matrixes, the local state estimation at C_i is the same as central estimation. However, in the more typical situation, the neighbors of each cameras are different; therefore, at each time instant the information each camera receives to fuse may also be different. There is no guarantee that the state estimates at different cameras remain cohesive. Thus a consensus step is implemented right as part of the estimation step. By comparing the fusion step (5) and Kalman-consensus state estimation step (6) in Algorithm 1 with the centralized state estimation (26) in Appendix A, it can be seen that our Kalman-consensus filter is essentially a distributed

¹As homography is applied on homogeneous coordinates, the mapping from ground plane to the image plane is nonlinear, and F is a linear approximation. Since $F \in \{\mathbb{F} : \mathbb{R}^2 \rightarrow \mathbb{R}^2\}$, \tilde{F} is invertible.

Algorithm 1 Distributed Kalman-Consensus tracking algorithm performed by every C_i at discrete time step k . The state estimate of T_l by C_i is represented by \mathbf{x}_i^l with error covariance matrix \mathbf{P}_i^l (see Sec. IV-A).

Input: $\bar{\mathbf{x}}_i^l$ and \mathbf{P}_i^l valid at k using measurements from time step $k-1$
for each T_l that is being viewed by $\{C_i^c \cup C_i\}$ **do**
 Obtain measurement \mathbf{z}_i^l with covariance \mathbf{R}_i^l
 Compute information vector and matrix

$$\begin{aligned}\mathbf{u}_i^l &= \mathbf{F}_i^{lT} (\mathbf{R}_i^l)^{-1} \mathbf{z}_i^l \\ \mathbf{U}_i^l &= \mathbf{F}_i^{lT} (\mathbf{R}_i^l)^{-1} \mathbf{F}_i^l\end{aligned}$$

Send messages $\mathbf{m}_i^l = (\mathbf{u}_i^l, \mathbf{U}_i^l, \bar{\mathbf{x}}_i^l)$ to neighboring cameras C_i^n
 Receive messages $\mathbf{m}_j^l = (\mathbf{u}_j^l, \mathbf{U}_j^l, \bar{\mathbf{x}}_j^l)$ from all cameras $C_j \in C_i^n$
 Fuse information matrices and vectors

$$\mathbf{y}_i^l = \sum_{j \in (C_i \cup C_i^n)} \mathbf{u}_j^l, \quad \mathbf{S}_i^l = \sum_{j \in (C_i \cup C_i^n)} \mathbf{U}_j^l \quad (5)$$

Compute the Kalman-Consensus state estimate

$$\begin{aligned}\mathbf{M}_i^l &= ((\mathbf{P}_i^l)^{-1} + \mathbf{S}_i^l)^{-1} \\ \hat{\mathbf{x}}_i^l &= \bar{\mathbf{x}}_i^l + \mathbf{M}_i^l (\mathbf{y}_i^l - \mathbf{S}_i^l \bar{\mathbf{x}}_i^l) + \gamma \mathbf{M}_i^l \sum_{j \in C_i^n} (\bar{\mathbf{x}}_j^l - \bar{\mathbf{x}}_i^l)\end{aligned} \quad (6)$$

$$\gamma = 1 / (|\mathbf{M}_i^l| + 1), \quad |\mathbf{X}| = (\text{tr}(\mathbf{X}^T \mathbf{X}))^{\frac{1}{2}}$$

Propagate the state and error covariance matrix from time step k to $k+1$

$$\begin{aligned}\mathbf{P}_i^l &\leftarrow \mathbf{A}^l \mathbf{M}_i^l \mathbf{A}^{lT} + \mathbf{B}^l \mathbf{Q}^l \mathbf{B}^{lT} \\ \bar{\mathbf{x}}_i^l &\leftarrow \mathbf{A}^l \hat{\mathbf{x}}_i^l\end{aligned} \quad (7)$$

end for

implementation of the centralized case with the consideration of communication constraint by adding a consensus term in (6). It is proved in [28] that all estimators asymptotically reach an unbiased consensus, i.e., $\hat{\mathbf{x}}_1 = \dots = \hat{\mathbf{x}}_n = \mathbf{x}$.

As shown in Algorithm 1, the information vector \mathbf{u}_i and \mathbf{U}_i exchanged between camera nodes are computed with measurement \mathbf{z}_i^l , covariance matrix \mathbf{R}_i^l and output matrix \mathbf{F}_i^l . Consider the two cases of measurement $(\mathbf{z}_i^l)_g$ and $(\mathbf{z}_i^l)_c$ as in (3) and (4). We denote their corresponding information vector and matrix as $(\mathbf{u}_i)_g, (\mathbf{U}_i)_g$ and $(\mathbf{u}_i)_c, (\mathbf{U}_i)_c$ respectively. The following shows that $(\mathbf{u}_i)_g = (\mathbf{u}_i)_c$ and $(\mathbf{U}_i)_g = (\mathbf{U}_i)_c$.

Recall that $(\mathbf{z}_i^l)_g$ and $(\mathbf{z}_i^l)_c$ are the measurements on ground plane and on the image plane of C_i respectively and \tilde{F}_i^l is the mapping from ground plane to the image plane. It is obvious that

$$\begin{aligned}(\mathbf{z}_i^l)_c &= \tilde{F}_i^l (\mathbf{z}_i^l)_g \\ \Rightarrow (\mathbf{F}_i^l)_c \mathbf{x}^l + (\mathbf{v}_i^l)_c &= \tilde{F}_i^l (\mathbf{F}_i^l)_g \mathbf{x}^l + \tilde{F}_i^l (\mathbf{v}_i^l)_g.\end{aligned} \quad (8)$$

– from (3) and (4)

Then

$$\begin{aligned}(\mathbf{F}_i^l)_c &= \tilde{F}_i^l (\mathbf{F}_i^l)_g, \\ (\mathbf{v}_i^l)_c &= \tilde{F}_i^l (\mathbf{v}_i^l)_g \Rightarrow (\mathbf{R}_i^l)_c = \tilde{F}_i^l (\mathbf{R}_i^l)_g \tilde{F}_i^{lT}.\end{aligned} \quad (9)$$

– from (8) and definition of covariance matrix

So the information vector and matrix are

$$\begin{aligned}(\mathbf{u}_i)_c &= (\mathbf{F}_i^l)_c^T (\mathbf{R}_i^l)_c^{-1} (\mathbf{z}_i^l)_c \\ &= (\tilde{F}_i^l (\mathbf{F}_i^l)_g)^T (\tilde{F}_i^l (\mathbf{R}_i^l)_g \tilde{F}_i^{lT})^{-1} \tilde{F}_i^l (\mathbf{z}_i^l)_g \\ &\quad \text{– substituting (8) and (9)} \\ &= (\mathbf{F}_i^l)_g^T \tilde{F}_i^{lT} (\tilde{F}_i^l)^{-1} (\mathbf{R}_i^l)_g^{-1} \tilde{F}_i^{-1} \tilde{F}_i^l (\mathbf{z}_i^l)_g \\ &= (\mathbf{F}_i^l)_g^T (\mathbf{R}_i^l)_g^{-1} (\mathbf{z}_i^l)_g \\ &= (\mathbf{u}_i)_g\end{aligned} \quad (10)$$

and

$$\begin{aligned}(\mathbf{U}_i)_c &= (\mathbf{F}_i^l)_c^T (\mathbf{R}_i^l)_c^{-1} (\mathbf{F}_i^l)_c \\ &= (\tilde{F}_i^l (\mathbf{F}_i^l)_g)^T (\tilde{F}_i^l (\mathbf{R}_i^l)_g \tilde{F}_i^{lT})^{-1} \tilde{F}_i^l (\mathbf{F}_i^l)_g \\ &\quad \text{– substituting (8) and (9)} \\ &= (\mathbf{F}_i^l)_g^T \tilde{F}_i^{lT} (\tilde{F}_i^l)^{-1} (\mathbf{R}_i^l)_g^{-1} \tilde{F}_i^{-1} \tilde{F}_i^l (\mathbf{F}_i^l)_g \\ &= (\mathbf{F}_i^l)_g^T (\mathbf{R}_i^l)_g^{-1} (\mathbf{F}_i^l)_g \\ &= (\mathbf{U}_i)_g.\end{aligned} \quad (11)$$

Since the information message exchanged between cameras are the same for both cases of \mathbf{z}_i^l , whether the measurement \mathbf{z}_i^l is measured on ground plane or image plane does not affect the tracking algorithm; these two cases give the same result.

C. Handoff and Fault Tolerance

Through this algorithm, each C_i has a consensus-based ground plane state estimate of each target that is being viewed by the cameras with which C_i can exchange information directly or indirectly, even if C_i has never seen some of the targets. Since we are assuming that the network of cameras as a whole is always covering the entire area under surveillance, each target will always be seen by at least one camera. Also, by our definition of overlapping cameras, a target T_l will always move from one camera C_i 's FOV to the FOV of an overlapping camera $C_j \in C_i^o$. Moreover, by **Assumption ***, C_i can exchange information with its overlapping cameras, C_i^o , directly or via other cameras. Therefore, C_j can take over the tracking of T_l and find the target correspondence in a seamless way since it had knowledge of T_l 's ground plane position through the consensus-tracking before it even entered its FOV. Additional target features could be used to find the target correspondences in a cluttered scene.

Another advantage of the fact that cameras have knowledge of all the targets in their neighborhood is that in the event of a sudden failure of camera node C_i , the targets that were viewed by C_i are not suddenly lost by the camera network.

We have also considered the fact that a camera may take a short amount of time to change its parameters to a new position in a non-static camera network. If no camera is viewing the target for the short amount of time it takes for the cameras to come to a new set of parameters to cover the entire area, the target state estimate and covariance continue to propagate by (7). This does not translate to a significant decrease in tracking performance as seen in our experiments.

D. Experimental Results

We tested our approach for tracking in a real camera network composed of 10 PTZ cameras looking over an outdoor area of approximately 10000 sq. feet. In the area under surveillance, there were 8 targets in total that were to be tracked using our distributed Kalman-Consensus filtering approach. In our experiment, the measurements (i.e., the observed positions of targets) are obtained using histogram of gradient (HOG) human detector [5]. The association of measurements to targets is achieved based on appearance (color) and motion information. Figure 2 shows the tracking results as viewed by each camera at 4 time instants.

The results are shown on a non-static camera network. The cameras are controlled to always cover the entire area under surveillance through a game theoretic control framework we proposed in [38]. As explained above, the change of camera settings does not affect the procedure of the Kalman-consensus filter. Figure 2 (a) shows the initial settings of the camera network that covers the entire area. As the targets are observed in this area, the single-view tracking module in each camera determines the ground plane position of each target in its FOV and sends that information to the Kalman-Consensus filter which processes it together with the information received from the Kalman-Consensus filters of neighboring cameras as described in Section IV.

Figure 2(b) shows the instant when a camera C_6 is focused on a target T_1^h . Figures 2(b) and (c) show the dynamics of the targets in the camera network. All targets are tracked using the Kalman-consensus scheme, although we show the marked track for only one target. The handoff of T_1^h is clearly shown in Figure 2(d) from C_6 to C_3 . It is to be noted that every time a target goes from one camera's FOV into another one, or when a camera changes its parameters, the network topologies for the targets, i.e., C_l^v and C_l^{v-} , also change.

Figure 3(a) shows the distributed Kalman-Consensus tracks for the 8 targets. The measurements of the different cameras are shown in a light gray color. As can be seen, the Kalman-Consensus filter in each camera comes to a smooth estimate of the actual state for each target.

Figure 3(b) shows the distributed tracking results on the ground plane for one of the targets, T_5 . The dots correspond to the ground plane measurements from different cameras viewing the target while the solid line is the consensus-based estimate. As can be expected, the individual positions are different for each camera due to calibration and single-view tracking inaccuracies. As can be seen clearly, even though C_5^v is time varying, the Kalman-Consensus filter estimates the target's position seamlessly at all times.

In Figure 3(a) and (b), the cameras that are viewing the same target can communicate with each other directly, i.e., $\forall l, C_l^v$ is a fully connected graph. As shown in Sec. IV-B, the results are exactly the same as a centralized case similar to each cluster of [22]. We denote the results of this fully connected case as KCF1. In order to show the effect of the network communication topology on the Kalman-consensus tracking, we consider an example of a partially connected network, which is shown on the right-top of Figure 3(c). Compared to

the fully connected one, direct communication does not exist between camera 1 and camera 3, neither between camera 4 and camera 8. Figure 3(c) shows the KCF tracking results at Camera 1 for this case, which is denoted as KCF2. It is slightly different with KCF1, due to the difference of fused information. The consensus method is guaranteed to have the same result as centralized case if there are no limitations on the communication capabilities. In the case of partial connection between cameras, KCF will converge to the same estimate centralized result as the number of consensus iterations goes to infinity [28]. However, the limited communication will result in differences from the centralized result for finite steps (as shown in Figure 3(c)). However, even in this case, the consensus result is better than that obtained at each individual camera, as shown in Figure 3(d) and explained below.

In order to measure tracking performance, we compare the tracking results with the groundtruth trajectory, which is shown in Figure 3(c). In the table at the bottom, we show the minimum, maximum and average distances to the groundtruth of KCF1, KCF2 and individual camera tracks. It can be seen that KCF1 performs best and KCF2 is better than individual camera tracks. We also look at the output error covariance matrix \mathbf{P} of the Kalman filter. The higher the trace of \mathbf{P} is, the lower the tracking accuracy is. Figure 3(d) shows the traces of the covariance matrix of the tracking error for the same target as in Figure 3(b) and (c). The colored lines with symbols correspond to tracking results from different cameras using their own measurements only (as each camera runs an independent Kalman filter), while the solid black line is the result of consensus-based estimate for the fully connected case (which will be the same for the centralized case) and dashed purple line is for the partially connected one. As can be seen clearly, the Kalman-Consensus filter with full connection performs the best, and partially connected one does better than individual Kalman filters without consensus.

V. DISTRIBUTED ACTIVITY RECOGNITION THROUGH CONSENSUS

In this section, we consider the problem of activity recognition in a camera network where processing power is distributed across the network and there is no central processor accumulating and analyzing all the data. Each camera computes a similarity measure of the observed activities in its views against a dictionary of pre-defined activities. Also, the transition probability between activities is known. This is a common assumption used in many activity recognition approaches and can be learned *a priori* from training data [4], [7], [18], [25], [36]. If no such information is available, the transition matrix can be assumed to be uniform. Based on the computed similarities at each camera node and the learned transition matrices, we show how to compute the consensus estimate in a probabilistic framework. Essentially, the consensus is a probability of similarity of the observed activity against the dictionary taking into account the decisions of the individual cameras.



Fig. 2. Each sub-figure shows 10 cameras at one of four time instants denoted by k . The track of one target, marked with a box, is shown. All targets are tracked using the Kalman-Consensus filtering approach, but are not marked for clarity.

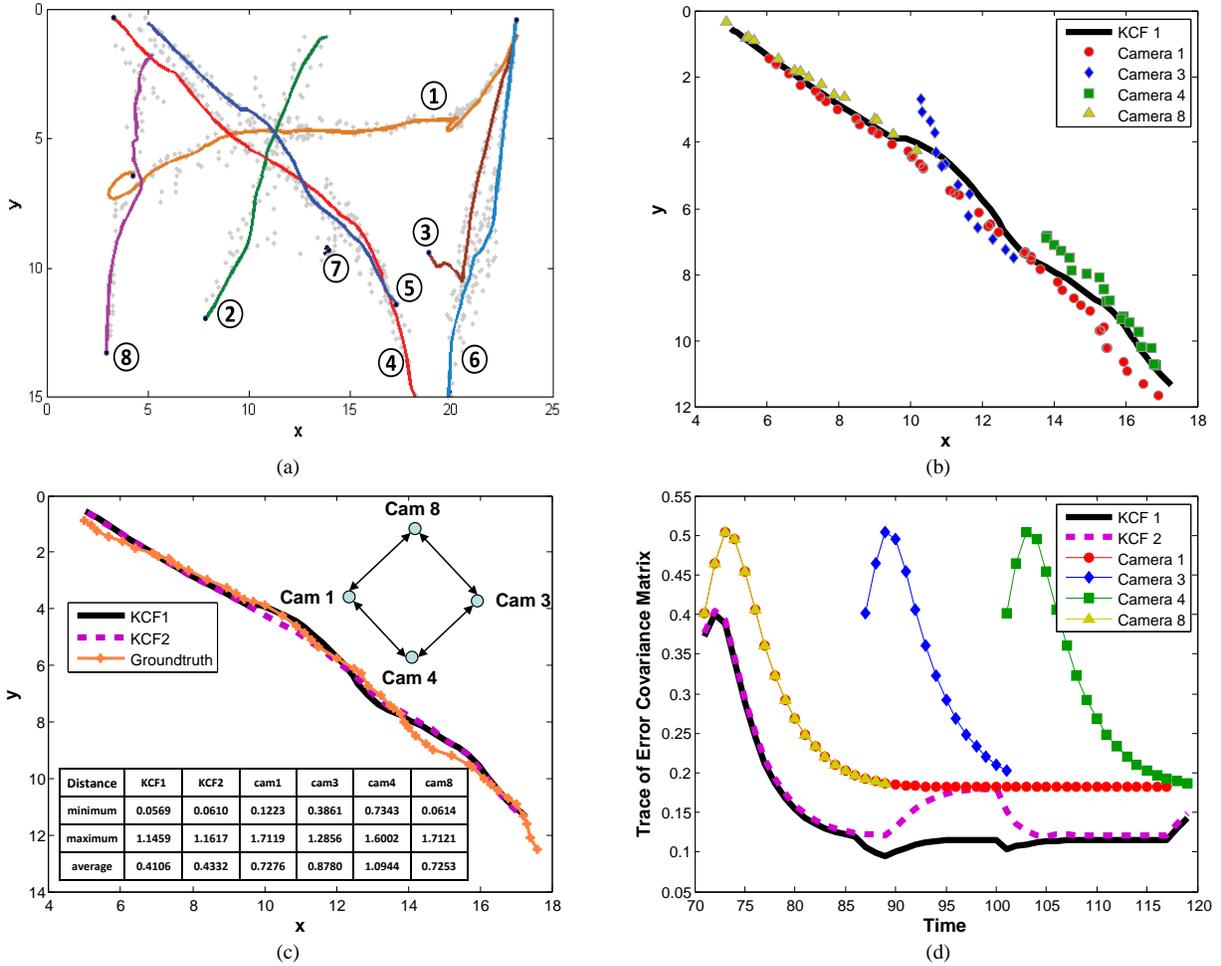


Fig. 3. Tracking results. (a): Distributed Kalman-Consensus tracking trajectories for 8 targets. Measurements from all cameras are shown in a light gray color. (b): Tracking results on the ground plane for one of the targets T_5 . In (a) and (b), the cameras that are viewing the same target can communicate with each other directly, i.e., $\forall l, C_l^v$ is a fully connected graph. The results are exactly same as centralized case. We denote the results of this full connection as KCF1. (c): KCF tracking results at Camera 1 given an example of a partially connected camera network, which is shown on the top-right. This case is denoted as KCF2. We can see that Cam (1,3) and Cam (4,8) cannot communicate. The groundtruth trajectory is also marked. The comparison of tracking performances (minimum, maximum and average distances to the groundtruth) of KCF1, KCF2 and individual camera tracks are shown in the table at the bottom. (d): Trace of the error covariance of the tracking results for the same target shown in (b) and (c).

A. Problem Formulation and Main Result

Let us assume that there are N_c cameras viewing a person performing some actions. The observation of camera C_i in the k^{th} time interval is denoted as $O_i(k), i = 1, \dots, N_c$. Let $\mathbf{O}(k)$ be the collection of observations from all the cameras, i.e., $\mathbf{O}(k) = \{O_1(k), \dots, O_{N_c}(k)\}$. Its history is $\mathcal{O}^k = \{\mathbf{O}(1), \dots, \mathbf{O}(k)\}$. The problem of activity recognition can be formulated so as to estimate the conditional probability, $P(y(k)|\mathcal{O}^k)$, where $y(k) \in \{1, \dots, Y\}$ is the label of the class of activity in a dictionary of Y activities with history $\mathcal{Y}^k = \{y(1), \dots, y(k)\}$.

It is a somewhat general assumption that the state transitions of activity class y are governed by the transition matrix for a 1st order Markov chain [36]:

$$\begin{aligned} P(y(k) = a|y(k-1) = a', \mathcal{Y}^{k-2}) \\ = P(y(k) = a|y(k-1) = a') \\ = m(a', a). \end{aligned} \quad (12)$$

$m(a', a)$ can be learned *a priori* from training data; if no such

information is available, the transition matrix can be assumed to be uniform. Given $y(k)$, observation $\mathbf{O}(k)$ is assumed to be independent of other observations and states, i.e.,

$$P(\mathbf{O}(k)|\mathcal{Y}^k, \mathcal{O}^{k-1}) = P(\mathbf{O}(k)|y(k)). \quad (13)$$

Based on Bayes' rule and above Markov chain assumption, we can show that the following relationship holds (see Appendix B for proof):

Result 1.

$$\begin{aligned} P(y(k)|\mathcal{O}^k) = \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \\ \cdot \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \\ \cdot \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right). \end{aligned} \quad (14)$$

where $\sum_{y(k)}$ mean summing over all values of $y(k) = 1, \dots, Y$. \square

Analysis of Result 1: By observing the righthand side of equation (14), we notice that $P(O_j(k)|y(k)), j = 1, \dots, N_c$ is the likelihood of camera C_j 's observation. The first term of the righthand side is a constant with respect to $y(k)$, so that it can be treated as a normalization factor and denoted by $\gamma(k)$, i.e.,

$$\gamma(k) \triangleq \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})}.$$

So we rewrite (14) as

$$P(y(k)|\mathcal{O}^k) = \gamma(k) \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \cdot \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right). \quad (15)$$

We define the state of the activity at camera C_i as $\mathbf{w}_i = [w_i^1, w_i^2, \dots, w_i^Y]^T$, where

$$w_i^a \triangleq P(y(k) = a|\mathcal{O}^k), a = 1, \dots, Y.$$

The likelihood of camera C_i 's observation is denoted by $\mathbf{v}_i = [v_i^1, v_i^2, \dots, v_i^Y]^T$, where

$$v_i^a \triangleq P(O_i(k)|y(k) = a), a = 1, \dots, Y.$$

Thus,

$$w_i^a(k) = \gamma(k) \prod_{j=1}^{N_c} P(O_j(k)|y(k) = a) \cdot \left(\sum_{y(k-1)} P(y(k) = a|y(k-1) = a')P(y(k-1) = a'|\mathcal{O}^{k-1}) \right) = \gamma(k) \prod_{j=1}^{N_c} v_j^a(k) \left(\sum_{a'=1}^Y m(a', a)w_i^{a'}(k-1) \right) \quad (16)$$

Based on the above argument, we have the activity recognition algorithm described in Algorithm 2 for each camera in the network.

Regarding the normalization factor $\gamma(k)$, we have the following result (see Appendix C for details).

Result 2.

$$\gamma(k) = \left[\sum_{y(k)} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \cdot \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \right]^{-1}.$$

□

Result 3. The local activity recognition procedure for node i based on fusion of the recognition results in all the cameras

Algorithm 2 Distributed Consensus based activity recognition algorithm performed by every C_i at step k .

Input: $\bar{\mathbf{w}}_i(k-1)$

for each person that is being viewed by $\{C_i^c \cup C_i\}$ **do**

Obtain observations $O_i(k)$

Compute local likelihood

$$\mathbf{v}_i(k) = \begin{bmatrix} v_i^1(k) \\ \vdots \\ v_i^Y(k) \end{bmatrix} = \begin{bmatrix} P(O_i(k)|y(k) = 1) \\ \vdots \\ P(O_i(k)|y(k) = Y) \end{bmatrix}$$

Send $\mathbf{v}_i(k)$ to neighboring cameras C_i^n

Receive $\mathbf{v}_j(k)$ from all cameras $C_j \in C_i^n$

Fuse information to estimate activity state

$$\begin{aligned} \mathbf{w}_i(k) &= \begin{bmatrix} w_i^1(k) \\ \vdots \\ w_i^Y(k) \end{bmatrix} \\ &= \begin{bmatrix} \gamma(k) \prod_{j \in (C_i \cup C_i^n)} v_j^1(k) \left(\sum_{a'=1}^Y m(a', 1) \bar{w}_i^{a'}(k-1) \right) \\ \vdots \\ \gamma(k) \prod_{j \in (C_i \cup C_i^n)} v_j^Y(k) \left(\sum_{a'=1}^Y m(a', Y) \bar{w}_i^{a'}(k-1) \right) \end{bmatrix} \\ &= \gamma(k) \prod_{j \in (C_i \cup C_i^n)} \Lambda(\mathbf{v}_j(k)) \mathbf{M}^T \bar{\mathbf{w}}_i(k-1), \end{aligned}$$

$$\begin{aligned} \gamma(k) &= \left(\sum_{a=1}^Y \prod_{j \in (C_i \cup C_i^n)} v_j^a(k) \left(\sum_{a'=1}^Y m(a', a) \bar{w}_i^{a'}(k-1) \right) \right)^{-1} \\ &= \left(\mathbf{1}_Y^T \cdot \prod_{j \in (C_i \cup C_i^n)} \Lambda(\mathbf{v}_j(k)) \mathbf{M}^T \bar{\mathbf{w}}_i(k-1) \right)^{-1}, \end{aligned}$$

where \mathbf{M} is a $Y \times Y$ matrix with $(i, j)^{th}$ element to be $m(i, j)$,

$$\Lambda(\mathbf{v}_j(k)) = \begin{bmatrix} v_j^1(k) & & \\ & \ddots & \\ & & v_j^Y(k) \end{bmatrix},$$

and $\mathbf{1}_Y = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$ with Y elements.

repeat

Send $\mathbf{w}_i(k)$ to neighboring cameras C_i^n

Receive $\mathbf{w}_j(k)$ from all cameras $C_j \in C_i^n$

Compute the Consensus state estimate

$$\bar{\mathbf{w}}_i(k) = \mathbf{w}_i(k) + \epsilon \sum_{j \in C_i^n} (\mathbf{w}_j(k) - \mathbf{w}_i(k))$$

until either a predefined iteration number is reached or $\sum_{j \in C_i^n} (\mathbf{w}_j(k) - \mathbf{w}_i(k))$ is smaller than a predefined small value

end for

is

$$\begin{aligned} v_i^a(k) &= P(O_i(k)|y(k) = a), a = 1, \dots, Y, \\ w_i^a(k) &= \gamma(k) \prod_{j=1}^{N_c} v_j^a(k) \left(\sum_{a'=1}^Y m(a', a) w_i^{a'}(k-1) \right), \quad (17) \\ a &= 1, \dots, Y, \\ \gamma(k) &= \left(\sum_{a=1}^Y \prod_{j=1}^{N_c} v_j^a(k) \left(\sum_{a'=1}^Y m(a', a) w_i^{a'}(k-1) \right) \right)^{-1}. \end{aligned}$$

□

The proof of Result 3 follows directly from Results 1 and 2.

Based on the network topology defined in Section IV-A, each camera can only communicate with its neighbors. According to this local activity recognition algorithm, there is no guarantee that the estimates remain cohesive among nodes. We use an ad hoc approach by implementing a consensus step right after the estimation step to reduce the disagreement regarding the estimates obtained in Result 3, from which Algorithm 2 can be inferred. This consensus approach is similar to the one proposed in [28] for the Kalman-Consensus filtering. However, a number of iterations are done in each time segment so as to converge to a consensus estimate.

The cameras that exchange information in the consensus stage are defined based on the communication constraints; therefore, it is possible that a camera involved in the consensus does not view the activity. In this case, such a camera transmits a value of $\mathbf{v}_i = \frac{1}{Y} \mathbf{1}_Y$, i.e., by assuming equal likelihood for all possible action classes.

B. Experimental Evaluation

To validate our proposed consensus approach for activity recognition, we carried out an experimental evaluation. We did activity recognition using multiple cameras and came to a consensus about the actions taking place using the theory of Section V-A.

For this, we used the IXMAS dataset [45]. In the dataset, there are sequences of images of different people doing several actions. The extracted silhouettes of the people in those actions are also given in the dataset. Five cameras were used to capture the whole activity which were placed at pan and tilt angles of $(120^\circ, 10^\circ)$, $(90^\circ, 10^\circ)$, $(30^\circ, 30^\circ)$, $(0^\circ, 10^\circ)$ and $(30^\circ, 90^\circ)$, where 0° pan angle means looking at a person from the front and 90° means to look at him from the left. A 3-dimensional motion-model of each person doing the actions is also given which has approximately 3500 voxels on a person.

We used the 3-dimensional motion-model as our training data and the silhouettes extracted from each camera as our test data. To build our training database, we took the orthographic projection of the 3-dimensional voxels on an image plane by rotating our virtual camera around the model with pan angles from 0° to 330° in increments of 30° and for each pan angle we used tilt angles of 10° and 30° . The actions we used in our experiments from the dataset are: looking at watch, scratching head, sit, wave hand, punch, kick and pointing a gun. These are later referred to as Actions 1 through 7. For each action and each camera viewpoint, we extracted the shape silhouette using 40 landmarks, i.e. 40 uniformly distributed contour points per shape in each frame. In a similar fashion we extracted the shape sequences of the test data, i.e. the silhouettes from different camera views.

For matching two shape sequences, we used a shape-based activity recognition algorithm based on work in [43]. The distance between two shape sequences is measured by comparing the mean shapes of the two sequence. Then we took the reciprocal of the distance measure to get a similarity measure between two shape sequences and normalized the similarity measures to convert them to probabilities. A block diagram of the overall activity recognition process is given in Figure 4.

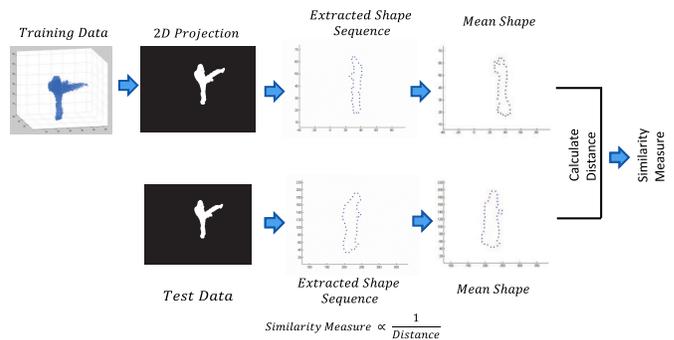


Fig. 4. Block diagram of the activity recognition process. For training using the 3d action models, orthographic projections were taken for different viewing angles. From the projections, shape sequences were extracted from which the mean shape was calculated for each viewing angle. For testing, in similar way the mean shapes were extracted. The Euclidean distance between the mean shapes were computed by comparing the test mean shape to all the training mean shapes. The ones with the lowest distance was selected for each action in the dictionary. Taking the reciprocal of the distance measure and normalizing it so that the sum of all the similarities is 1 gave the similarity measure of the actions.

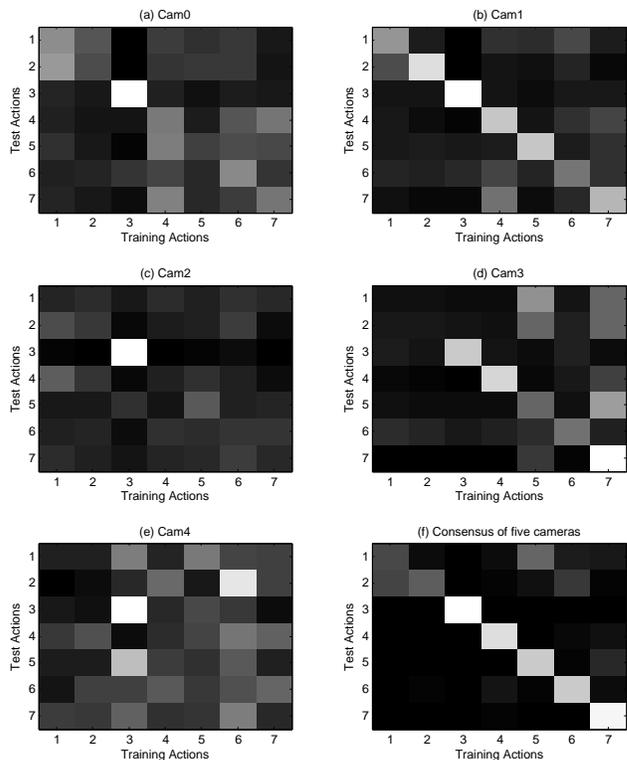
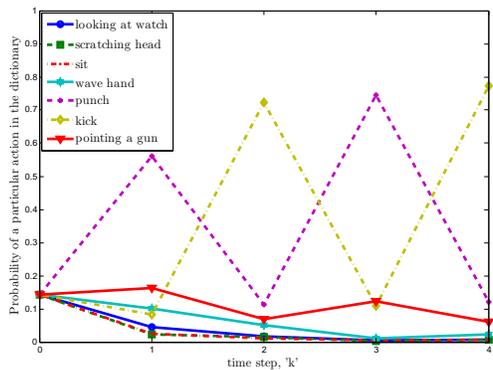


Fig. 5. (a-e) Similarity matrices of the activities for the cameras cam0,cam1,cam2,cam3 and cam4; (f) Similarity matrix of the activities for the consensus of all these cameras. Actions 1 through 7 are looking at watch, scratching head, sit, wave hand, punch, kick and pointing a gun, respectively, all from the IXMAS dataset.

The activity recognition is performed individually in each of the cameras depending on its own current observation. In our experiment, we have five cameras, i.e. cam0, cam1, cam2, cam3 and cam4. We consider a network topology where the network is not a full mesh, rather each camera is connected to two other cameras only. So, after the activity recognition stage, each camera shares the detection result with its immediate neighbor. Each camera fuses the detection results of itself and its neighbors, final detection result from the previous time step $k-1$, and the transition probabilities between different actions,



(a)

action	looking at watch	scratching head	sit	wave hand	punch	kick	point
looking at watch	0.1804	0.1854	0.1790	0.1825	0.0871	0.1033	0.0823
scratching head	0.1742	0.1758	0.1896	0.1736	0.0991	0.0908	0.0969
sit	0.1914	0.1889	0.1118	0.1949	0.0992	0.1066	0.1071
wave hand	0.2172	0.1055	0.1151	0.2050	0.1093	0.1212	0.1266
punch	0.0807	0.0943	0.0801	0.0984	0.2500	0.2396	0.1569
kick	0.0810	0.0842	0.0818	0.0865	0.2485	0.2481	0.1699
point	0.1480	0.1483	0.1387	0.1364	0.1500	0.1317	0.1469

(b)

Fig. 6. (a): Graphical representation of the final detection result, i.e. the result of the consensus stage in each time step, for the sequence punch-kick-punch-kick. The vertices in each line at each time step represent the probability of a particular action. It was assumed that in time step $k = 0$, all the activities were equally likely. We use a non-uniform transition matrix, as shown in (b), where there are high transition probability between the punch and kick, and there is also some moderately high transition probability between looking at watch, scratch, sit, wave hand and point.

and gets a new probability distribution of the actions. After this stage, the cameras initiate the consensus algorithm and try to converge to the same detection results.

In Figure 5, we show the similarity matrices, i.e. the probability of match for each test activity (the row of a matrix). The more white the cell block is, the test data it refers to is detected with more probability as that action. Five of the images represent the similarity matrix for the test data captured by each camera and the sixth image shows the similarity matrix of the consensus for all of these cameras. The similarity scores of correct matching are the diagonal values of the similarity matrix. Comparing with other values in the matrix, the higher the diagonal values (brighter in the image) are, the less confusing the recognition result is. By comparing the similarity matrix of consensus with the test data captured by each camera (compare (f) with (a)-(e)), it is clear that the recognition result after consensus has less confusion than others.

Next, in Figure 6(a), we show a graphical representation of the final detection result for a sequence of punch-kick-punch-kick, by plotting the result of the consensus stage in each time step. The vertices in each line at each time step, represent the probability of a particular action in the dictionary. It was assumed that in time step $k = 0$, all the activities were equally likely. As an example, we use a non-uniform transition matrix where there is high transition probability between punch and kick, and there is also some moderately high transition probability between looking at watch, scratch, sit, wave hand and point. The transition matrix between different actions is shown in Figure 6(b). In practice, if some prior knowledge is available, the transition matrices can be learned/manually set. As the transition probability between punch and kick is high, it can be seen that the recognition result (after consensus) keeps on improving.

Finally, we generate a statistics to observe the performance of the probability of correct match for individual cameras versus their consensus. We use every possible subset of the five

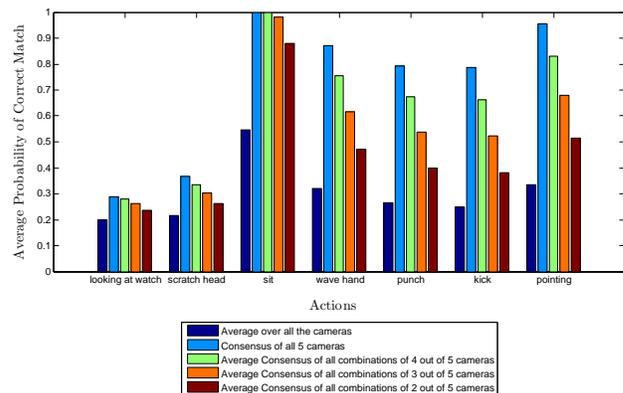


Fig. 7. Comparison of average probability of correct match for individual camera and their consensus for all the activities. There are seven sets of bars for seven different actions and in each set, there are five bars where the leftmost one (blue) is the average probability of correct match for individual cameras and the next four bars the average probability of correct match of the consensus over all the combinations of cameras taking respectively five, four, three and two out of five cameras.

cameras by considering five, four, three and two cameras to determine their consensus and show that the consensus result is better than an individual camera, on average. This result shows the fault tolerance aspect of the consensus process. The result is shown in Figure 7.

C. Discussion

- *Experimental Setup*: We did the experiments by running the algorithms on independent threads, one for each camera, and communication between the threads using existing protocols. We assume here that communication is not a bottleneck. This proposed work is a proof-of-concept study in using distributed processing algorithms for video analysis. Future work should consider the practical constraints of using consensus algorithms in camera networks.

- *Temporal Segmentation for Activity Recognition*: In our distributed activity recognition procedure, the video sequence is divided into segments, where each segment is treated as

a observed variable (image features) and associated with a hidden variable (activity state). In our experiments, in order to provide a clear comparison of our results with ground truth, the observed sequence from each camera is temporally segmented based on the ground truth. In practice, such a precise segmentation is not required; the observed sequence can be uniformly divided into short segments (e.g., 4 seconds each). Since our activity recognition results are represented in the form of probabilities, the non-dominant recognition results on the segments where activity transitions happen won't affect the recognition on their subsequent segments.

- *Synchronization*: The cameras in the network have been pre-synchronized, however, the frame synchronization may not be perfect due to slight frame rate difference between cameras. So the transmitted information between cameras includes a time stamp. In the distributed tracking framework, when a camera fuses the information (e.g., state estimations) from its neighboring cameras, it will do interpolation of the information vector \mathbf{u} (in Algorithm 1) as necessary. This will ensure that the information being fused is synchronized. While the activity recognition is done on each segment, unlike the frame based Kalman-consensus tracking, a precise synchronization of the cameras is not needed; precision of pre-synchronization is enough.

- *Selection of Parameters*: We can see that the consensus step in Algorithm 2 is a gradient descent algorithm that minimizes the cost function $g(\mathbf{w}_i) = \frac{1}{2} \sum_{j \in \mathcal{C}_i^n} (\mathbf{w}_i - \mathbf{w}_j)^2$. The step-size $\epsilon > 0$ should be a small number. The choice of ϵ is based on reasoning similar to what is used for gradient descent. The simplest way is to set ϵ a fixed small number, while some suggest using an adaptive step-size. In our experiments, the step-size ϵ is fixed at 0.01.

- *Integration of tracking and activity recognition*: Since the distributed tracking and activity recognition can be achieved through analogous frameworks (though the detailed fundamentals are different) by estimating locally and fusing through consensus, it is possible to integrate these two by designing integrated local estimation and fusion schemes. We address the integration as a future work.

VI. CONCLUSION AND FUTURE WORK

We investigated in this paper distributed scene analysis algorithms by leveraging upon concepts of consensus. We addressed two fundamental tasks - tracking and activity recognition in a distributed camera network. We proposed a robust approach to distributed multi-target tracking in a network of cameras. A distributed Kalman-Consensus filtering approach was used together with a dynamic network topology for persistently tracking multiple targets across several camera views. A probabilistic consensus scheme for activity recognition was provided, which combines the similarity scores of neighboring cameras to come up with a probability for each action at the network level. In the future, we will look at the integration of tracking and activity recognition into a single framework and more complex activities that span a larger area.

APPENDIX A KALMAN FILTER WITH CENTRALIZED INFORMATION FUSION

Consider a Kalman filter with centralized information fusion, i.e., each camera sends its observation to a central processor, and tracking (i.e. state estimation) is performed centrally. As in (2), the sensing model at camera C_i of target T_l is $\mathbf{z}_i^l = \mathbf{F}_i^l \mathbf{x}^l + \mathbf{v}_i^l$. Thus the central measurement, observation noise and observation matrix are defined as

$$\mathbf{z}^l = \begin{bmatrix} \mathbf{z}_1^l \\ \mathbf{z}_2^l \\ \vdots \\ \mathbf{z}_{N_c}^l \end{bmatrix}, \mathbf{v}^l = \begin{bmatrix} \mathbf{v}_1^l \\ \mathbf{v}_2^l \\ \vdots \\ \mathbf{v}_{N_c}^l \end{bmatrix}, \mathbf{F}^l = \begin{bmatrix} \mathbf{F}_1^l \\ \mathbf{F}_2^l \\ \vdots \\ \mathbf{F}_{N_c}^l \end{bmatrix}, \quad (18)$$

where N_c is the total number of cameras. Then we get

$$\mathbf{z}^l = \mathbf{F}^l \mathbf{x}^l + \mathbf{v}^l, \quad (19)$$

where $\mathbf{x}^l = (x^l, y^l, \dot{x}^l, \dot{y}^l)^T$ is the same as in Sec. IV-B.

By assuming \mathbf{v}_i^l 's are uncorrelated, the covariance matrix of \mathbf{v}^l is

$$\mathbf{R}^l = \text{diag}(\mathbf{R}_1^l, \mathbf{R}_2^l, \dots, \mathbf{R}_{N_c}^l). \quad (20)$$

Thus, the Kalman filter iterations in the information form are

$$\begin{aligned} \mathbf{M}^l(k) &= \left[(\mathbf{P}^l(k))^{-1} + \mathbf{F}^l(k)^T \mathbf{R}^l(k)^{-1} \mathbf{F}^l(k) \right]^{-1} \\ &= \left[(\mathbf{P}^l(k))^{-1} + \sum_{i=1}^{N_c} \mathbf{F}_i^l(k)^T \mathbf{R}_i^l(k)^{-1} \mathbf{F}_i^l(k) \right]^{-1}, \end{aligned} \quad (21)$$

$$\mathbf{K}^l(k) = \mathbf{M}^l(k) \mathbf{F}^l(k)^T \mathbf{R}^l(k)^{-1}, \quad (22)$$

$$\begin{aligned} \hat{\mathbf{x}}^l(k) &= \bar{\mathbf{x}}^l(k) + \mathbf{K}^l(k) (\mathbf{z}^l(k) - \mathbf{F}^l(k) \bar{\mathbf{x}}^l(k)) \\ &= \bar{\mathbf{x}}^l(k) + \mathbf{M}^l(k) \left[\mathbf{F}^l(k)^T \mathbf{R}^l(k)^{-1} \mathbf{z}^l(k) \right. \\ &\quad \left. - \mathbf{F}^l(k)^T \mathbf{R}^l(k)^{-1} \mathbf{F}^l(k) \bar{\mathbf{x}}^l(k) \right] \\ &= \bar{\mathbf{x}}^l(k) + \mathbf{M}^l(k) \left[\sum_{i=1}^{N_c} \mathbf{F}_i^l(k)^T \mathbf{R}_i^l(k)^{-1} \mathbf{z}_i^l(k) \right. \\ &\quad \left. - \left(\sum_{i=1}^{N_c} \mathbf{F}_i^l(k)^T \mathbf{R}_i^l(k)^{-1} \mathbf{F}_i^l(k) \right) \bar{\mathbf{x}}^l(k) \right], \end{aligned} \quad (23)$$

$$\mathbf{P}^l(k+1) = \mathbf{A}^l \mathbf{M}^l(k) \mathbf{A}^{lT} + \mathbf{B}^l \mathbf{Q}^l \mathbf{B}^{lT}, \quad (24)$$

$$\bar{\mathbf{x}}^l(k+1) = \mathbf{A}^l \hat{\mathbf{x}}^l(k). \quad (25)$$

Denoting the information vector and matrix at camera C_i as $\mathbf{u}_i^l(k) = \mathbf{F}_i^l(k)^T \mathbf{R}_i^l(k)^{-1} \mathbf{z}_i^l(k)$ and $\mathbf{U}_i^l(k) = \mathbf{F}_i^l(k)^T \mathbf{R}_i^l(k)^{-1} \mathbf{F}_i^l(k)$, (23) can be rewritten as

$$\hat{\mathbf{x}}^l(k) = \bar{\mathbf{x}}^l(k) + \mathbf{M}^l(k) \left[\sum_{i=1}^{N_c} \mathbf{u}_i^l(k) - \left(\sum_{i=1}^{N_c} \mathbf{U}_i^l(k) \right) \bar{\mathbf{x}}^l(k) \right] \quad (26)$$

By comparing (26) with Algorithm 1 where each camera fuses its information vector and matrix and those from its neighbors, it is clearly shown that Kalman-consensus filter is

a distributed implementation. If $\forall l, C_l^v$ is a fully connected graph, i.e., all cameras that are viewing the same target can communicate with each other directly, the Kalman-consensus filter will provide exactly the same result as the Kalman filter with centralized fusion.

APPENDIX B PROOF OF RESULT 1

Assuming there are N_c cameras viewing a person performing some actions, the observations of camera C_i in k^{th} time interval are denoted as $O_i(k), i = 1, \dots, N_c$. Let $\mathbf{O}(k)$ be the collection of observations from all the cameras, i.e., $\mathbf{O}(k) = \{O_1(k), \dots, O_{N_c}(k)\}$ and its history is $\mathcal{O}^k = \{\mathbf{O}(1), \dots, \mathbf{O}(k)\}$. Then the statement

$$P(y(k)|\mathcal{O}^k) = \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \cdot \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right)$$

holds $\forall N_c \geq 1$.

Proof:

$$P(y(k)|\mathcal{O}^k) = \frac{P(y(k), \mathcal{O}^k)}{P(\mathcal{O}^k)} \quad \text{– from Bayes' rule} \quad (27)$$

We notice that $P(y(k), \mathcal{O}^k)$ is the forward variable in hidden Markov model. According to the recursion of the Forward Algorithm, i.e.,

$$P(y(k), \mathcal{O}^k) = \left[\sum_{y(k-1)} P(y(k-1), \mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \cdot P(O(k)|y(k)), \quad (28)$$

(27) becomes

$$\begin{aligned} & P(y(k), \mathcal{O}^k) \\ &= \frac{P(\mathbf{O}(k)|y(k))}{P(\mathcal{O}^k)} \\ & \cdot \left[\sum_{y(k-1)} P(y(k-1), \mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \\ & \quad \text{– substituting (28)} \\ &= \frac{P(\mathbf{O}(k)|y(k))}{P(\mathcal{O}^k)} \\ & \cdot \left[\sum_{y(k-1)} P(y(k-1)|\mathcal{O}^{k-1})P(\mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \\ & \quad \text{– from Bayes' rule} \\ &= \frac{P(\mathcal{O}^{k-1})P(\mathbf{O}(k)|y(k))}{P(\mathcal{O}^k)} \\ & \cdot \left[\sum_{y(k-1)} P(y(k-1)|\mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \end{aligned} \quad (29)$$

$$\begin{aligned} &= \frac{P(O(k)|y(k))}{P(\mathcal{O}^k)/P(\mathcal{O}^{k-1})} \\ & \cdot \left[\sum_{y(k-1)} P(y(k-1)|\mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \\ &= \frac{P(O(k)|y(k))}{P(\mathbf{O}(k), \mathcal{O}^{k-1})/P(\mathcal{O}^{k-1})} \\ & \cdot \left[\sum_{y(k-1)} P(y(k-1)|\mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \\ & \quad \text{– expanding } P(\mathcal{O}^k) \\ &= \frac{P(O(k)|y(k))}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})P(\mathcal{O}^{k-1})/P(\mathcal{O}^{k-1})} \\ & \cdot \left[\sum_{y(k-1)} P(y(k-1)|\mathcal{O}^{k-1})P(y(k)|y(k-1)) \right] \\ & \quad \text{– from Bayes' rule} \\ &= \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} P(\mathbf{O}(k)|y(k)) \\ & \cdot \left[\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right] \end{aligned} \quad (30)$$

The observation of Camera C_j , $O_j(k)$, is determined by the activity being performed and the view point of C_j . If the activity is known, i.e., $y(k)$ is given, $O_j(k)$ only depends on the view point of C_j and is independent of observations of other cameras, i.e.,

$$\begin{aligned} P(\mathbf{O}(k)|y(k)) &= P(O_1(k), \dots, O_{N_c}(k)|y(k)) \\ &= \prod_{j=1}^{N_c} P(O_j(k)|y(k)). \end{aligned} \quad (31)$$

So we get

$$P(y(k)|\mathcal{O}^k) = \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \cdot \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right),$$

which is the statement of Result 1. \blacksquare

APPENDIX C PROOF OF RESULT 2

Given that

$$P(y(k)|\mathcal{O}^k) = \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \cdot \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \quad (*)$$

where $y(k) \in \{1, \dots, Y\}$, then

$$\begin{aligned} \gamma(k) &\triangleq \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \\ &= \left[\sum_{y(k)} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \right. \\ &\quad \cdot \left. \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \right]^{-1} \end{aligned}$$

Proof: Since $y(k) \in \{1, \dots, Y\}$, it can be inferred that

$$\sum_{a=1}^Y P(y(k)|\mathcal{O}^k) = 1.$$

By substituting (*), we have

$$\begin{aligned} 1 &= \sum_{y(k)} \left[\frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \right. \\ &\quad \cdot \left. \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \right] \\ \Rightarrow 1 &= \frac{1}{P(\mathbf{O}(k)|\mathcal{O}^{k-1})} \\ &\quad \cdot \left[\sum_{y(k)} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \right. \\ &\quad \cdot \left. \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \right] \\ \Rightarrow 1 &= \gamma(k) \left[\sum_{y(k)} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \right. \\ &\quad \cdot \left. \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \right] \\ \Rightarrow \gamma(k) &= \left[\sum_{y(k)} \prod_{j=1}^{N_c} P(O_j(k)|y(k)) \right. \\ &\quad \cdot \left. \left(\sum_{y(k-1)} P(y(k)|y(k-1))P(y(k-1)|\mathcal{O}^{k-1}) \right) \right]^{-1} \end{aligned}$$

REFERENCES

- [1] A. Alahi, D. Marimon, M. Bierlaire, and M. Kunt. A master-slave approach for object detection and matching with fixed and mobile cameras. In *Intl. Conf. on Image Processing*, 2008.
- [2] D. P. Bertsekas and J. Tsitsiklis. *Parallel and Distributed Computation*. Upper Saddle River, NJ: Prentice-Hall, 1989.
- [3] M. Bramberger, A. Doblender, A. Maier, B. Rinner, and H. Schwabach. Distributed Embedded Smart Cameras for Surveillance Applications. *IEEE Computer*, 2006.
- [4] T.-J. Cham and J. M. Rehg. A multiple hypothesis approach to figure tracking. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 1999.
- [5] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2005.
- [6] M. H. DeGroot. Reaching a consensus. *Journal of the American Statistical Association*, 69(345):118–121, March 1974.
- [7] A. Doucet, N. d. Freitas, and N. Gordon. *Sequential Monte Carlo Methods in Practice*. Springer, 2001.
- [8] W. Du and J. Piater. Multi-camera People Tracking by Collaborative Particle Filters and Principal Axis-Based Integration. In *Asian Conf. on Computer Vision*, 2007.
- [9] M. Egerstedt and X. Hu. Formation control with virtual leaders and reduced communications. *IEEE Trans. Robot. Autom.*, 17(6):947–951, 2001.
- [10] E. Ermis, V. Saligrama, P. Jodoin, and J. Konrad. Abnormal behavior detection and behavior matching for networked cameras. In *IEEE/ACM Intl. Conf. on Distributed Smart Cameras*, 2008.
- [11] J. A. Fax. Optimal and cooperative control of vehicle formations. *Ph.D. dissertation, Control Dynamical Syst., California Inst. Technol., Pasadena, CA*, 2001.
- [12] R. Guerraoui, M. Hurfin, A. Mostefaoui, R. Oliveira, M. Raynal, and A. Schiper. Consensus in asynchronous distributed systems: A concise guided tour. In *LNCS 1752*, pages 33–47. Springer-Verlag, 1999.
- [13] Y. Hatano and M. Mesbahi. Agreement over random networks. *IEEE Trans. on Automatic Control*, Nov 2005.
- [14] D. Hristu and K. Morgansen. Limited communication control. *Syst. Control Lett.*, 37:193–205, Jul 1999.
- [15] A. Jadbabaie, J. Lin, and A. S. Morse. Coordination of groups of mobile autonomous agents using nearest neighbor rules. *IEEE Trans. on Automatic Control*, 48(6):988–1001, Jun 2003.
- [16] S. Khan, O. Javed, Z. Rasheed, and M. Shah. Camera Handoff: Tracking in Multiple Uncalibrated Stationary Cameras. In *IEEE Workshop on Human Motion*, 2000.
- [17] S. Khan and M. Shah. A Multiview Approach to Tracking People in Crowded Scenes Using a Planar Homography Constraint. In *Euro. Conf. on Computer Vision*, 2006.
- [18] L. Liao, D. Fox, and H. Kautz. Location-based activity recognition using relational markov networks. In *Proc. of the International Joint Conference on Artificial Intelligence*, 2005.
- [19] Z. Lin, M. Brouke, and B. Francis. Local control strategies for groups of mobile autonomous agents. *IEEE Trans. on Automatic Control*, 49(4):622–629, Apr 2004.
- [20] N. Lynch. *Distributed Algorithms*. Morgan Kaufmann Publishers, San Mateo, CA, 1996.
- [21] D. Markis, T. Ellis, and J. Black. Bridging the Gap Between Cameras. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2004.
- [22] H. Medeiros, J. Park, and A. Kak. Distributed object tracking using a cluster-based kalman filter in wireless camera networks. *IEEE Journal of Selected Topics in Signal Processing*, 2(4):448–463, Aug. 2008.
- [23] M. Mehyar, D. Spanos, J. Pongsjapan, S. Low, and R. M. Murray. Distributed averaging on asynchronous communication networks. In *IEEE Conf. on Decision and Control and European Control Conference*, pages 7446–7451, Dec 2005.
- [24] M. Mesbahi. On state-dependent dynamic graphs and their controllability properties. *IEEE Trans. Autom. Control*, 50(3):387–392, Mar 2005.
- [25] B. North, A. Blake, M. Isard, and J. Rittscher. Learning and classification of complex dynamics. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 22(9):1016–1034, 2000.
- [26] R. Olfati-Saber. Ultrafast consensus in small-world networks. in *Proc. Am. Control Conf.*, 2005.
- [27] R. Olfati-Saber. Flocking for multi-agent dynamic systems: Algorithms and theory. *IEEE Trans. Automatic Control*, 51(3):401–420, Mar 2006.
- [28] R. Olfati-Saber. Distributed kalman filtering for sensor networks. *IEEE Conf. on Decision and Control*, 2007.
- [29] R. Olfati-Saber, J. Fax, and R. Murray. Consensus and Cooperation in Networked Multi-Agent Systems. *Proceedings of the IEEE*, 95(1):215–233, Jan. 2007.
- [30] R. Olfati-Saber and R. M. Murray. Consensus problems in networks of agents with switching topology and time-delays. *IEEE Trans. on Automatic Control*, 49(9):1520–1533, Sep 2004.
- [31] R. Olfati-Saber and N. F. Sandell. Distributed tracking in sensor networks with limited sensing range. *Proceedings of the American Control Conference*, June 2008.
- [32] V. M. Preciado and G. C. Verghese. Synchronization in generalized Erdős-Rénye networks of nonlinear oscillators. In *IEEE Conf. on Decision and Control and European Control Conference*, 2005.
- [33] F. Qureshi and D. Terzopoulos. Surveillance in Virtual Reality: System Design and Multi-Camera Control. *IEEE Conf. on Computer Vision and Pattern Recognition*, 2007.
- [34] A. Rahimi and T. Darrell. Simultaneous calibration and tracking with a network of non-overlapping sensors. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2004.
- [35] W. Ren and R. W. Beard. Consensus seeking in multi-agent systems under dynamically changing interaction topologies. *IEEE Trans. on Automatic Control*, 50(5):655–661, May 2005.
- [36] J. Rittscher and A. Black. Classification of human body motion. In *Intl.*

Conf. on Computer Vision, 1999.

- [37] B. Song and A. Roy-Chowdhury. Stochastic Adaptive Tracking in a Camera Network. In *Intl. Conf. on Computer Vision*, 2007.
- [38] C. Soto, B. Song, and A. Roy-Chowdhury. Distributed multi-target tracking in a self-configuring camera network. In *IEEE Conf. on Computer Vision and Pattern Recognition*, 2009.
- [39] B. Stancil, C. Zhang, and T. Chen. Active Multicamera Networks: From Rendering to Surveillance. *IEEE Journal on Selected Topics in Signal Processing Special Issue on Distributed Processing in Vision Networks*, August 2008.
- [40] K. Tieu, G. Dalley, and W. Grimson. Inference of Non-Overlapping Camera Network Topology by Measuring Statistical Dependence. In *Intl. Conf. on Computer Vision*, 2005.
- [41] R. Tron, R. Vidal, and A. Terzis. Distributed pose averaging in camera networks via consensus on SE(3). *IEEE/ACM Intl. Conf. on Distributed Smart Cameras*, Sept. 2008.
- [42] J. N. Tsitsiklis, D. P. Bertsekas, and M. Athans. Distributed asynchronous deterministic and stochastic gradient optimization algorithms. *IEEE Trans. on Automatic Control*, 31(9):803–812, Sep 1986.
- [43] A. Veeraraghavan, A. Roy-Chowdhury, and R. Chellappa. Matching shape sequences in video with applications in human motion analysis. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2005.
- [44] Y. Wang, K. Huang, and T. Tan. Multi-view gymnastic activity recognition recognition with fused hmm. In *Asian Conf. on Computer Vision*, 2007.
- [45] D. Weinland, E. Boyer, and R. Ronfard. Action recognition from arbitrary views using 3d exemplars. In *Intl. Conf. on Computer Vision*, 2007.
- [46] J. Zhao, S. C. Cheung, and T. Nguyen. Optimal Camera Network Configurations for Visual Tagging. *IEEE Journal on Selected Topics in Signal Processing Special Issue on Distributed Processing in Vision Networks*, August 2008.



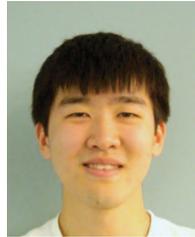
Bi Song received the B.S. and M.S. degrees in electronic engineering and information science from the University of Science and Technology of China, Hefei, in 2001 and 2004, respectively. She received the Ph.D. degree in 2009 from the Department of Electrical Engineering at the University of California, Riverside, where she is currently working as a postdoctoral scholar. Her main research interests include image processing and analysis, computer vision.



Ahmed T. Kamal received the B.S. degree in Electrical and Electronic Engineering from the Bangladesh University of Engineering and Technology, Dhaka in 2008. He received the M.S. degree in Electrical Engineering from the University of California, Riverside in 2010. He is currently a Ph.D. candidate in the Department of Electrical Engineering in the same University. His main research interests include intelligent camera networks, wide-area scene analysis, activity recognition and search and distributed information fusion.



Cristian Soto received the M.S. degrees from the Department of Electrical Engineering at the University of California, Riverside in 2008. He is currently with Western Digital, CA, USA.



Chong Ding received the B.S. degree in Computer Science from the University of California, Riverside in 2008. He is currently a Ph.D. candidate in the Department of Computer Science in the same University. His main research interests include intelligent camera networks, wide-area scene analysis and distributed and real-time systems.



Amit K. Roy-Chowdhury is an Associate Professor of Electrical Engineering and a Cooperating Faculty in the Dept. of Computer Science at the University of California, Riverside. He received his Bachelor's degree in Electrical Engineering from Jadavpur University, Calcutta, India, his Masters in Systems Science and Automation from the Indian Institute of Science, Bangalore, India, and his PhD in Electrical Engineering from the University of Maryland, College Park, USA. His broad research interests are in the areas of image processing and analysis, computer vision, video communications and statistical methods for signal analysis. His current research projects include intelligent camera networks, wide-area scene analysis, physics-based mathematical modeling of images, activity recognition and search, video-based biometrics (face and gait), biological video analysis and distributed video compression. The work is supported by the National Science Foundation, Office of Naval Research, Army Research Office, DARPA and private industries. Dr. Roy-Chowdhury has over eighty papers in peer-reviewed journals, conferences and edited books. He is an author of the book titled "Recognition of Humans and Their Activities Using Video".



Jay A. Farrell received B.S. degrees (1986) in physics and electrical engineering from Iowa State University, and M.S. (1988) and Ph.D. (1989) degrees in electrical engineering from the University of Notre Dame. At Charles Stark Draper Lab (1989-1994), he was principal investigator on projects involving intelligent and learning control systems for autonomous vehicles. Dr. Farrell received the Engineering Vice President's Best Technical Publication Award in 1990, and Recognition Awards for Outstanding Performance and Achievement in 1991 and 1993. He is a Professor and former Chair of the Department of Electrical Engineering at the University of California, Riverside. He has served as Vice President Finance and Vice President of Technical Activities for the IEEE Control Systems Society (CSS). He is a Fellow of the IEEE (2008), a Distinguished Member of IEEE CSS, and author of over 160 technical publications. He is author of the book "Aided Navigation: GPS with High Rate Sensors" (McGraw-Hill 2008). He is also co-author of the books "The Global Positioning System and Inertial Navigation" (McGraw-Hill, 1998) and "Adaptive Approximation Based Control: Unifying Neural, Fuzzy and Traditional Adaptive Approximation Approaches" (John Wiley 2006).