

Opportunistic Image Acquisition of Individual and Group Activities in a Distributed Camera Network

Chong Ding, *Student Member, IEEE*, Jawadul H. Bappy, *Student Member, IEEE*, Jay A. Farrell, *Fellow, IEEE*, and Amit K. Roy-Chowdhury, *Senior Member, IEEE*

Abstract—The decreasing cost and size of video sensors has led to camera networks becoming pervasive in our lives. However, the ability to analyze these images effectively is very much a function of the quality of the acquired images. In this paper we consider the problem of automatically controlling the fields of view of individual pan, tilt, zoom (PTZ) cameras in a camera network leading to improved situation awareness (e.g. where and what are the critical targets and events) in a region of interest. The network of cameras attempts to observe the entire region of interest at some minimum resolution while opportunistically acquiring high resolution images of critical events in real time. Since many activities involve groups of people interacting, an important decision that the network needs to make is whether to focus on individuals or groups of them. This is achieved by understanding the performance of video analysis tasks and designing camera control strategies to improve a metric that quantifies the quality of the source imagery. Optimization strategies, along with a distributed implementation, are proposed, and their theoretical properties analyzed. The proposed methods bring together computer vision and network control ideas. Performance of the proposed methodologies discussed herein, has been evaluated on a real life wireless network of pan, tilt and zoom capable cameras.

Index Terms—Camera networks, distributed optimization, opportunistic sensing and tracking.

1 INTRODUCTION

The focus of most work in video analysis has been on improving the performance of detection, tracking and recognition algorithms in a variety of settings. Large-scale experimental analysis has consistently demonstrated the limitations in the improvements that can be obtained by focusing solely on the processing side. Pose and image resolution remain two of the hardest challenges to be overcome for robustness in scene understanding. However, research focused on video acquisition strategies driven by the need to maximize performance goals have been quite limited.

Our goal is to develop ‘optimal’ sensing strategies in a network of visual sensors, with a focus on tightly integrating the sensing and processing tasks in order to image, track and identify targets effectively. We consider wide-area scene understanding problems where the sensors are pan-tilt-zoom (PTZ) cameras. In keeping with normal behavior observed widely, we assume that the scene will be populated with people acting individually or in groups. We envision that these visual sensors will have the capability to

analyze their own data and perform collaborative decision making in coordination with other sensors. This would enable the sensors to maneuver themselves optimally, i.e., change the pan, tilt and zoom (PTZ) parameters of the cameras, so as to obtain image sequences that can be analyzed with a high degree of reliability. A specific challenge that this network will have to address in the process is to decide whether to image individuals separately or in groups.

1.1 Related Work

Historically, camera networks [1] have consisted of mostly static cameras, where the layout is designed to satisfy some predefined objective such as covering an area. Optimal camera placement strategies proposed in [2] and [3] were solved by using a camera placement metric that captures occlusion in 3-D environments. The solution to the problem of optimal camera placement given coverage constraints was presented in [4], and can also be used to come up with an initial camera configuration.

Currently most of data collected by these networks is analyzed manually, a task that is extremely tedious and limits the potential of the installed network. This has sparked a huge interest in automated analysis in camera networks. The authors in [5] considered distributed estimation in camera networks, while [6] handled distributed tracking in camera networks. A review of automated analysis in camera networks can

- C. Ding is with the Department of Computer Science and Engineering, University of California, Riverside, Riverside, CA 92521.
E-mail: cding@cs.ucr.edu
- J. H. Bappy, A. K. Roy-Chowdhury and J. A. Farrell are with the Department of Electrical and Computer Engineering, University of California, Riverside, Riverside, CA 92521.
E-mail: {mbappy, amitrc, farrell}@ee.ucr.edu

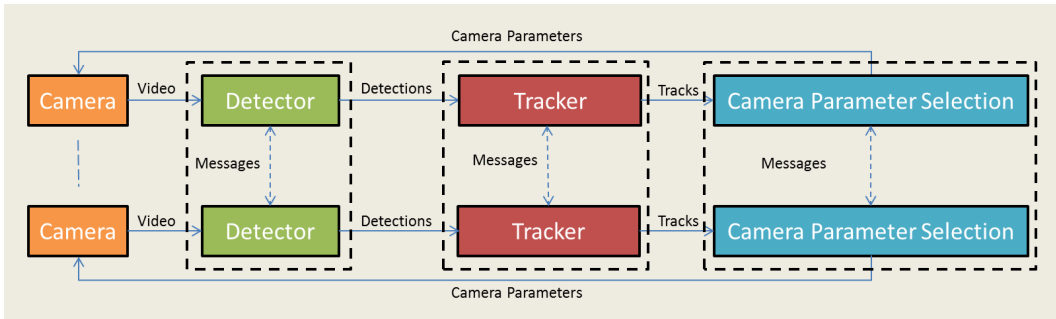


Fig. 1. Diagram depicting the framework for integrating scene analysis and PTZ control.

be found in [1]. A solution to improve performance within the realm of static cameras is to increase the total number of cameras observing a region and selectively use a subset depending on the scene [7]. Another approach is to use a system with camera parameters that can be dynamically altered according to the scene. In [8], a master-slave configuration is used to acquire high resolution images and scheduling is considered as a dynamic discrete optimization problem. Early work done in active vision [9] looked at the problem of moving a camera to improve imagery. Performing active vision in a distributed camera network, where the cameras coordinate among themselves, remains relatively unexplored.

Much of the research in controllable camera networks in the last decade was focused on centralized solutions to master-slave systems where static cameras directed the PTZ cameras. The path planning inspired approach proposed by [10] used static cameras to track all targets in a virtual environment while the PTZ cameras were each assigned to obtain high resolution video from a particular target. This approach showed that given the predicted tracks of all the targets, a set of one-to-one mappings between cameras and targets can be formed to acquire high resolution videos. A method for determining good sensor configurations that would maximize performance measures was introduced in [11]. The configuration framework allowed for the presence of random occluding objects. In [12], an information-theoretic approach was incorporated based on tracker uncertainty for active scene analysis where no movement cost is considered for selecting camera parameters.

A more recent approach in [13] and [14] uses the Expectation-Maximization (EM) algorithm to find the optimal configuration of PTZ cameras given a map of activities and the value of each discretized ground coordinate is determined using this map. The game theoretic approach in [15] and [16] showed how local value functions could be designed to constantly improve the tracking accuracy of all targets observed by a self-configuring camera network. Prior work in target tracking with the ability to obtain high resolution shots was shown in [17] using a fixed cost function. While the approaches in these papers were decentralized and could achieve high accuracy, they

were focused largely on designing specific imaging functions on a per-target basis, and not on the overall system optimization required for more complex scenes involving interactions between multiple targets. Also, the scalability of the system was limited, a fact that we address in Sec. (3). More discussion on how to design games in general for distributed optimization can be found in [18].

In recent years, group information such as track splitting and merging, has been studied in several tracking methods ([19], [20], [21], [22], [23], [24], [25]). In [19], [20] and [22], group information is utilized as a constraint in order to improve individual tracking performance. In [23], particle filter is used to model both individual and group information jointly. In [24] and [25], the authors address the group tracking problem with a descriptor of appearance features. In [26], the authors propose a person re-identification method that finds the one to one correspondences between targets in different cameras. However, none of these works consider image acquisition with individual or group information.

Main Contribution. The main contribution of this paper is to present a general approach where interesting events that the targets are involved in are detected, and the system automatically decides *when, where and how* to image these targets. An aspect that requires particular attention in this regard is that it is often difficult to maintain long tracks on objects, especially when they are interacting with each other. For example, if an individual merges with a group, we may lose track of that individual. However, if the group later splits, we may be able to get the track back. This is predicated upon the availability of high quality images at opportune moments of time so that the tracks can be reconstructed. It is a major motivation for this work, and the experimental results presented later will highlight this aspect further.

1.2 Solution Overview

Configurations of cameras with large fields-of-view (FOV) can monitor a large area of the environment but may not be able to supply images for reliable recognition tasks. Configurations where the cameras are zoomed in on specific areas of interest can gather

TABLE 1
Notation Summary

Parameter	Variable
i -th camera	c_i
Set of all PTZ settings available to c_i	A_i
The PTZ setting for all cameras	\mathbf{a}
The PTZ setting for c_i , all cameras except c_i	a_i, \mathbf{a}_{-i}
The PTZ setting for c_i and its neighbors	\mathbf{a}^{c_i}
Measurement vector, measurement covariance	\mathbf{u}, \mathbf{C}
Rotation matrix from frame a to frame b	\mathbf{R}_a^b
State vector for track j	\mathbf{x}_j
Position vector for track j in world frame	\mathbf{p}_j^w
Area coverage utility	$U_{area}(\mathbf{a})$
Imaging utility	$U_{image}(\mathbf{a})$
Local area coverage utility	$u_{area}(\mathbf{a}^{c_i})$
Local imaging utility	$u_{image}(\mathbf{a}^{c_i})$

images useful for recognition, but result in a very limited view of the scene.

In this paper we will show how to automatically acquire high resolution images, based on the state of the scene and video surveillance system, while covering a region of interest. Ideally the system would be able to acquire detailed images of targets before and after any events that could potentially affect the ability to create long term tracks and recognize the targets that are being tracked.

This framework is shown in Fig. (1) and can find application in any visual surveillance system containing PTZ cameras. The raw video from each of the cameras is first processed through a detector. Any resulting detections are then associated and used to update the target tracks by the tracker. These tracks are then used to decide on the next set of PTZ settings for each camera. In a traditional completely centralized system, the raw video would be sent from each camera to a central server where the detector, tracker and parameter selection modules would all be located. In the fully distributed case, each camera would communicate only necessary information, e.g., the state estimates of the targets. The local objective function of each camera must be aligned with the global objective function such that any local decisions, using only locally available information, improve the value of the global objective function.

2 MODELING OF SCENE DYNAMICS

The purpose of this section is to define the time propagation and measurement models of our system, which will form the basis for evaluating the performance of the video analysis tasks. Additional notation is summarized in Table 1. The camera network assumes that all the PTZ cameras are calibrated.

Time Propagation Model: With time step $T = 1$, as track j moves throughout the area, its trajectory in discrete time is modeled in state space as,

$$\mathbf{x}_j(k+1) = \Phi \mathbf{x}_j(k) + \omega, \quad (1)$$

where, $\mathbf{x}_j = [\mathbf{p}_j^w, \mathbf{v}_j^w]^\top$ with $\mathbf{p}_j^w = [x_j^w, y_j^w, z_j^w]^\top$ is the position in the world frame, $\mathbf{v}_j^w = [v_{x,j}^w, v_{y,j}^w, v_{z,j}^w]^\top$ is the velocity in the world frame, and $j = 1, \dots, N_t$ is the track number. Also, $\omega \sim \mathcal{N}(\mathbf{0}_{6 \times 1}, \mathbf{Q})$ is the process noise and

$$\Phi = \begin{bmatrix} \mathbf{I} & \mathbf{I} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}, \quad (2)$$

is the state transition matrix where $\mathbf{0}$ and $\mathbf{I} \in \mathbb{R}^{3 \times 3}$. The state estimate and its error covariance matrix are propagated between sampling instants using recursive Kalman Filter equations [27].

Measurement Model: This section derives the measurement model for track j by camera i . The position of the j -th track in the i -th camera frame is $\mathbf{p}_j^{c_i} = [x_j^{c_i}, y_j^{c_i}, z_j^{c_i}]^\top$, the rotation from world to the i -th camera frame is denoted by $\mathbf{R}_w^{c_i}$ and the position of the camera in the world frame is $\mathbf{p}_{c_i}^w$.

The position of track j in the i -th camera frame is related to it's position in the world frame by

$$\begin{bmatrix} \mathbf{p}_j^{c_i} \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R}_w^{c_i} & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\mathbf{p}_{c_i}^w \\ \mathbf{0}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{p}_j^w \\ 1 \end{bmatrix} \quad (3)$$

where $\mathbf{0} \in \mathbb{R}^{3 \times 1}$ and $\mathbf{I} \in \mathbb{R}^{3 \times 3}$ is the identity matrix.

Using the homogeneous co-ordinate representation, the i -th camera's image plane co-ordinates for the j -th track can be represented as \mathbf{p}_j^i . The result of performing feature detection on the image from Camera i is a two-dimensional projection of the centroid position of the j -th track on the image plane in pixel coordinates \mathbf{u}_j^i and it's associated error covariance \mathbf{C}_j^i . Accounting for noise, the measurement from the i -th camera can be modeled as

$$\mathbf{u}_j^i = h(\mathbf{p}_j^i) + \boldsymbol{\eta}_j^i = \begin{bmatrix} \frac{x_j^i}{z_j^i} \\ \frac{y_j^i}{z_j^i} \end{bmatrix} + \boldsymbol{\eta}_j^i, \quad (4)$$

where we assume that $\boldsymbol{\eta}_j^i \sim \mathcal{N}(\mathbf{0}_{2 \times 1}, \mathbf{C}_j^i)$. The linearized measurement model for every measurement at time-step k is thus given by,

$$\mathbf{u}_j^i(k) = \mathbf{H}_j^i(k) \mathbf{x}_j^i(k) + \boldsymbol{\eta}_j^i(k). \quad (5)$$

where \mathbf{H}_j^i is the observation matrix [28].

Detection, Association & Tracking: In our framework we assume that the detector for each camera provides a centroid position \mathbf{u}_j^i and error covariance \mathbf{C}_j^i on the image plane for each detection. The resulting detections can then be associated to existing tracks using any existing data association technique. A review of object tracking methods can be found in [27]. It will typically be the case that each camera only detects some of the targets.

There are several different events that affect our tracker. When the probability that a detection is associated to an existing track is below a threshold, a new track is initialized by transforming the centroid position and error covariance of the detection to the world frame as the targets are on the ground plane. Since

the relation between the image and ground planes is non-linear, we use Eqn. (5) to perform measurement updates for the state estimates and error covariances for tracks using an Extended Kalman filter.

3 DISTRIBUTED SOLUTION

One of the challenges in vision networks is scalability with the size of the network. This motivates us to provide a distributed solution to the problem, where each camera is able to communicate with its neighbors, as defined by a communication graph described below. In the proposed distributed network, the cameras are only able to communicate directly with their immediate neighbors on the communication graph and only have access to information provided by these neighbors. A major benefit of such a framework is that new cameras with varying capabilities can be easily added or removed from the network at any time. We do not address the problem of distributed tracking and data association, as these have been presented in recent work [29], [30] and [31] can be used directly.

In a distributed camera network, the camera communications are modeled by a communication graph. To ease computation we quantize the entire area B into a set of blocks, $B = \{b_1, b_2, \dots, b_{N_b}\}$. $B^i \subseteq B$ is the set of blocks for which there exist a valid parameter setting so that the block is in the Field of View (FOV) of camera c_i .

Definition 1. The communication graph has cameras as nodes and there exist an edge from camera i' to camera i'' when those cameras can directly receive information from each other.

Definition 2. In vision graph, cameras are represented as nodes. The graph has an edge from camera i' to camera i'' when the set of blocks $B^{i'} \cap B^{i''} \neq \emptyset$.

These graphs are important to our problem as the communication graph determines *who knows what and when*, while the vision graphs determines which cameras can share operation regions (i.e., have overlapping FOV's when their PTZ parameters are appropriately chosen). The example in Fig. (2(a)) shows 4 cameras each covering a portion of the area. For the regions B^i depicted in Fig. (2(a)), the corresponding vision graph is represented by Fig. (2(b)).

3.1 Design of Local Utility Function

In a distributed camera network, one possible solution is to design the objective functions so that they implement a potential game [18], which guarantees that the cameras will converge to a Nash equilibrium (N.E.) (see Sec. 3.4 for the proof). While this is a distributed solution there are a few undesirable characteristics such as allowing only one camera to change its proposed parameter settings at each iteration step [16]). The existing literature only accounts for the communication graph connectivity and the amount of communication increases with the number of cameras.

In the following, we incorporate information from the vision graph by considering the cameras that can affect the value of the local camera's PTZ settings.

We present a design of the value functions for the different objectives of our system. In the following sections let $C_i \subseteq C$ be the set of cameras containing c_i and its neighbors in communication graph, $A^{c_i} = \times_{i \in C_i} A_i$ be the set of possible PTZ settings of the cameras in C_i .

Local Area Coverage: Instead of considering all cameras in the network, we write the local area coverage utility to only consider the pan tilt and zoom settings $\mathbf{a}^{c_i} \in A^{c_i}$ as,

$$u_{area}(\mathbf{a}^{c_i}) = \sum_{p \in B_j} \left(1 - \prod_{q \in C_i} (1 - \beta_p^q) \right). \quad (6)$$

where B_j is the set of blocks observed by the cameras in C_i and β_p^q is defined as

$$\beta_p^q = \begin{cases} 1 - e^{-\lambda_1 \frac{r_p^q}{r_p^{max}}} & \text{if } r_p^{max} > r_p^q > r_p^{min} \\ 0 & \text{otherwise} \end{cases}, \quad (7)$$

r_p^{min} is the minimum acceptable resolution in terms of pixel height at which p -th block should be viewed, r_p^{max} is the maximum height of the target in pixels of c_i 's image plane, and r_p^q is the resolution at which block p would be imaged for the proposed action a^{c_i} . The conditions under which this value function increases are when area covered with an acceptable resolution increases or the area can be viewed at a higher resolution.

Eqn. (6) implies that only cameras who are neighbors of c_i in communication graph can affect this utility value. By defining the local utility in such a way, the change in value of the local utility, denoted by Δu_{area} , and global utility ΔU_{area} , by a change in the PTZ parameters of camera in c_i from a_i to b_i (where $a_i, b_i \in A_i$) are related by:

$$\begin{aligned} \Delta u_{area} &= u_{area}(b_i, \mathbf{a}_{-i}^{c_i}) - u_{area}(a_i, \mathbf{a}_{-i}^{c_i}) \\ &= \sum_{l \in B^i} \left(1 - \prod_{j \in C_i} (1 - \beta_l^j) \right) \Big|_{b_i} - \\ &\quad \sum_{l \in B^i} \left(1 - \prod_{j \in C_i} (1 - \beta_l^j) \right) \Big|_{a_i} \\ &= U_{area}(b_i, \mathbf{a}_{-i}) - U_{area}(a_i, \mathbf{a}_{-i}) = \Delta U_{area}. \end{aligned} \quad (8)$$

This is true because B^i contains all blocks that can be affected by a change in a_i in U_{area} and all the cameras that can affect the value of blocks in B^i are in the set $C_i \subseteq C$.

Local High Resolution Imaging: The local value function for high resolution imaging of the tracked objects can then be written as,

$$u_{image}(\mathbf{a}^{c_i}) = \sum_{j \in T^i} \left(1 - \prod_{k \in C_i} (1 - \zeta_j^k) \right) s_j, \quad (9)$$

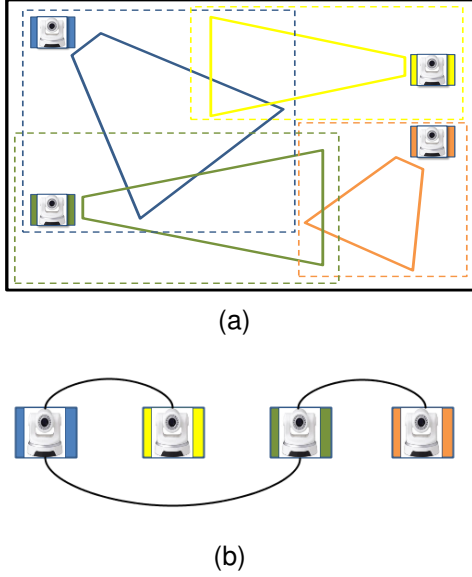


Fig. 2. (a) Shows Field of View(FOV) for 4 cameras as closed trapezoids. The dotted rectangles show the possible area that can be covered given all the settings available to each camera. (b) Shows the corresponding vision graph, connecting cameras that have overlap in their fields of view.

where T^i are the targets being viewed by the cameras in the set C_i . ζ_j^k is defined as

$$\zeta_j^k = \begin{cases} 1 - e^{-\lambda_2 \frac{r_j^k}{r_j^{max}}} & \text{if } r_j^{max} > r_j^k > r_j^{min} \\ 0 & \text{otherwise} \end{cases}, \quad (10)$$

Here r_j^k is the resolution weighted by the probability $p_{j,l}$ of a target j is at block l which is computed using $\hat{\mathbf{P}}_j^w(k)^+$ and $\hat{\mathbf{P}}_j^w(k)^+$ and can be written as,

$$r_j^k = \begin{cases} \sum_{b_l \in B_j} (p_{j,l} r_l^k), & \text{if } \forall b_l \in B_j, r_l^k > 0 \\ 0 & \text{otherwise} \end{cases}, \quad (11)$$

The scaling parameter s_j in eqn. (9) is an importance factor to weight the value of getting high resolution images of target j by changing it dynamically over time. We populate this variable using a function monotonically increasing in time, initialized to a low value when a new track j is added to the system. Once a high resolution image of the track has been acquired $s_j = 0$ for the remaining lifetime of the track.

3.2 Communication & Vision Graphs

The change in value of the local and global utility for high resolution imaging defined in eqn. (9), denoted by Δu_{image} and ΔU_{image} respectively. The relation-

ship between Δu_{image} and ΔU_{image} are given by,

$$\begin{aligned} \Delta u_{image} &= u_{image}(b_i, \mathbf{a}_{-i}^{c_i}) - u_{image}(a_i, \mathbf{a}_{-i}^{c_i}) \\ &= \sum_{j \in T^{c_i}} \left(1 - \prod_{i \in C_i} (1 - \beta_l^i) \right) \Big|_{b_i} - \\ &\quad \sum_{j \in T^{c_i}} \left(1 - \prod_{i \in C_i} (1 - \beta_l^i) \right) \Big|_{a_i} \\ &= U_{image}(b_i, \mathbf{a}_{-i}) - U_{image}(a_i, \mathbf{a}_{-i}) = \Delta U_{image} \end{aligned} \quad (12)$$

This is true because only targets who have a probability of being in B^i can be affected by a change in a_i , and the only cameras that can affect blocks in B^i are the neighboring cameras C_i .

Local Utility Function: Given the new local value functions for area coverage and high resolution imaging defined in eqns. (6) and (9). The utility function representing the local value of the system can be represented as

$$u_{local}(\mathbf{a}^{c_i}) = u_{area}(\mathbf{a}^{c_i}) + u_{image}(\mathbf{a}^{c_i}), \quad (13)$$

and the change in $\Delta u_{local}^{c_i}$, global utility ΔU_{global} and the PTZ parameters of camera in c_i from a_i to b_i (where $a_i, b_i \in A_i$) are related by:

$$\begin{aligned} \Delta u_{local}^{c_i} &= u_{local}(b_i, \mathbf{a}_{-i}^{c_i}) - u_{local}(a_i, \mathbf{a}_{-i}^{c_i}) \\ &= \Delta U_{global}. \end{aligned} \quad (14)$$

Since a change in the local utility $\Delta u_{local}^{c_i}$ is equivalent to the change in the global utility ΔU_{global} for any a_i . We can increase the global utility by increasing any camera's local utility.

3.3 Distributed Optimization Algorithm

The alignment between the local utility and the global utility is necessary to ensure that local decisions are coordinated with the global objective of the system. Given the local utility defined in eqn. (13) and its alignment to the global utility of the camera network we modify the algorithm run at each camera i so that more than one camera can change PTZ parameters at each time step.

Step 1: c_i computes the best PTZ setting,

$$\max_{a_i \in A_i} (\Delta u_{local}^{c_i}), \quad (15)$$

for camera i and send the proposed setting a_i and utility $\max_{a_i \in A_i} (\Delta u_{local}^{c_i})$ to all neighboring cameras.

Step 2: Receive the proposed setting a'_i and utility $\max_{a'_i \in A'_i} (\Delta u_{local}^{c'_i})$ from each other camera in the network.

Step 3: If c_i can provide the greatest improvement,

$$\max_{c'_i \in C_i} (\max_{a'_i \in A'_i} (\Delta u_{local}^{c'_i})) = \max_{a_i \in A_i} (\Delta u_{local}^{c_i}), \quad (16)$$

then execute the proposed settings. In the case that multiple cameras have the same improvement any tie

breaker such that $g(c_i, c'_i) = g(c'_i, c_i)$ can be used.

Step 4: Send current settings to all neighboring cameras.

Step 5: Recieve current settings from all neighboring cameras and update \mathbf{a}^{c_i} .

3.4 Deterministic Equilibrium

We now analyze the convergence characteristics of the distributed optimization algorithm addressed in Sec. (3.3) and show that the cameras can arrive at a Nash equilibrium (N.E.) with one or more cameras deciding to change parameters simultaneously.

Definition 3. A path in the parameter space S of the camera network is a sequence of camera network parameters $(\mathbf{a}^1, \mathbf{a}^2, \dots, \mathbf{a}^{N_s})$ such that each two consecutive camera network states differs in at least one camera PTZ setting.

Definition 4. An improvement path is a path in which $U_{global}(\mathbf{a}^k) < U_{global}(\mathbf{a}^{k+1})$, where \mathbf{a}^k and \mathbf{a}^{k+1} differ in at least one camera state, for all $k = 1, 2, \dots$

Definition 5. The camera network is at a Nash equilibrium if $\forall \mathbf{a}_{-i} \in A_{-i}$ and $\forall a_i, b_i \in A_i$, $\Delta U_{global} = U_{global}(b_i, \mathbf{a}_{-i}) - U_{global}(a_i, \mathbf{a}_{-i}) \leq 0$ and $\Delta u_{local}^{c_i} = u_{local}(b_i, \mathbf{a}_{-i}^{c_i}) - u_{local}(a_i, \mathbf{a}_{-i}^{c_i}) \leq 0$.

Claim 1: Given a fixed environment and an arbitrary set of initial camera settings, the network of cameras will improve upon the global utility U_{global} until it reaches a Nash equilibrium.

Proof: As there are a finite number of camera network configurations, it is easy to see that every improvement path is finite and arrives at a Nash equilibrium. If the camera network is not at a Nash equilibrium then there is a set of cameras such that $\max_{a_i \in A_i} (\Delta u_{local}^{c_i}) > 0$. After applying the condition in Eqn. (16) the set of cameras that decide to move is $C_{move} \subseteq C$ resulting in the camera network state \mathbf{a}_{move} . Also, $\forall c'_i, c''_i \in C_{move}$, $c'_i \notin C''_i$ and $c''_i \notin C'_i$. This means that the contribution of c'_i to the global utility is independent of the contribution of c''_i . The resulting change in the global value can then be written as

$$\Delta U_{global} |_{\mathbf{a}_{move}} = \sum_{i \in C_{move}} (\Delta u_{local}^{c_i} |_{\mathbf{a}_{move}}). \quad (17)$$

Since $\Delta u_{local}^{c_i} |_{\mathbf{a}_{move}} > 0, \forall i \in C_{move}$, then $\Delta U_{global} |_{\mathbf{a}_{move}} > 0$. Thus by Def. (4) and (5) the camera network is either at a N.E. or moves along an improvement path. \square

4 EXPERIMENTAL ANALYSIS & RESULTS

In this section we will explain the details of our experimental setup and our results. The experiments were performed on a wireless PTZ camera network consisting of Axis 215 PTZ-E cameras over a $20m \times 30m$ region. As we do not currently have the capability to access the video from the cameras directly and commands must be sent across the wireless network through the provided VAPIX API resulting in poor

response time due to latency. In our physical system we noticed a delay ranging from $80ms$ to $200ms$ for each message to arrive over TCP. This limits the rate at which we can change the parameters of the camera network as there is significant delay in the video response after sending a command.

The detection method we used in our experiments was very simple in design. First the area was uniformly divided into a number of overlapping blocks, each $1m \times 1m \times 1.8m$. The frames from each camera are processed using a motion subtraction [32] algorithm to generate a motion image. The blocks in the cameras operational region are then mapped to the image plane using the homography.

In this implementation we make no assumptions as to the number of people present in each detection or track. Two people walking close together may result in a single detection with a covariance encompassing both. This means that if two people stay close together for the duration they are in the region they will be tracked as a group rather than as two individual tracks. Since we do not explicitly count the number of people in each track, a few special events such as track split and track merge can occur.

Split track: When two or more people who are walking together start drifting apart, the detector will return multiple detections that associate to their existing track. The existing track is removed and one new track for each detection is created.

Merge track: When two or more people start walking ever closer to each other, the detector will eventually return a single detection for all of the people. This event occurs when one detection associates to more than one track. In such cases the existing tracks are removed and a new track is created for the group.

One thing to note here is that as people form groups or split apart there is a short time period where multiple track split and track merge events may occur rapidly until they are close enough to form a group or far enough apart to be tracked separately. The scaling parameter s_j is tuned such that high resolution images are not prioritized until after a track has existed for a period of time so as to avoid capturing multiple high resolution images during merge and split events.

Parameter Selection: There are many ways to speed up the search for the PTZ setting that maximizes the local utility, the easiest approach is to increase the quantization of the parameter space. For our experiments the available pan settings were quantized into 5 degree increments and bounded based on the position of the camera and the area under surveillance. The tilt settings were quantized into one degree increments and was restricted to a maximum tilt of 30 degrees. Each camera also was restricted to three different zoom settings. Some additional techniques such as a lookup table were also used to improve the computation speed such that each camera can detect, track and decide on parameters within a $33ms$ time window

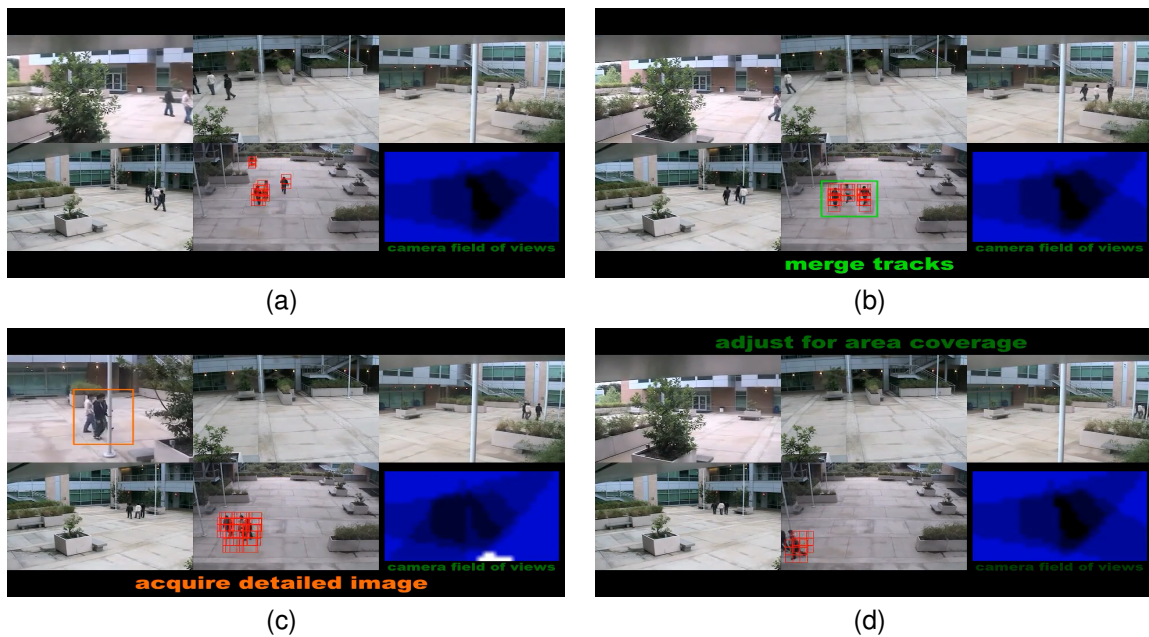


Fig. 3. Scene involving a pair of people and two individual people. The images show the result when two of the tracks merge into a new group. (a) Two individuals and a group of two people are distinguishable by the detections, shown in red boxes, in the bottom middle camera. (b) Shows that the detections of the group of two and an individual in the scene can no longer be separated and a new group is formed. (c) A high resolution image is acquired by the upper left camera. (d) The cameras reconfigure to cover the area.

on a 2.66GHz Intel Core i560m. The communication overhead of our physical setup reduces the rate at which our cameras change their parameters.

4.1 Experimental Results

We ran a number of different test scenarios through our system. A video of all the scenes is available at <http://www.ee.ucr.edu/~amitrc/CameraNetworks.php>. Here, we will explain two of them (scene 4 and 5 of the video) that will demonstrate all possible cases such as track splitting, merging and individual tracking. The objective is to cover the entire area while opportunistically acquiring high resolution shots of people in the area.

First Scenario. In this scenario, two pairs of people enter the scene from opposite sides of the courtyard. One pair splits up shortly upon entering the area, while the other pair stays together. Some time after the first group split, one of its members walks towards and joins together with other pair forming a group of three as they proceed to exit the area.

The system is expected to form a track and acquire a high resolution image for each of the initial pairs of people. After the first group splits, two new tracks are expected to be created and high resolution images of each should be taken. When the track of the second pair and one of the members of the first pair merge, the system is expected to create a new track and capture a corresponding high resolution image.

Fig. (3) shows the behavior of the camera network

in response to the merging of the tracks. In Fig. (3(a)) we can see three tracks, two people walking alone and a pair of people walking together. As two of the tracks start getting closer, it becomes more difficult for our basic detector to separate the pair of people from the person walking alone. This can be seen in Fig. (3(b)) as the tracks merge to form a group of three people. A high resolution image is taken in Fig. (3(c)) a short while after the formation of the new track.

Second Scenario. The final scenario is the most complicated and shows how the system behaves when it does not have the time to complete all the desired tasks. In this scenario four people enter the courtyard from each of the four corners within seconds of each other. They each walk towards the opposite end of the courtyard resulting in the four of them meeting in the center of the courtyard before reaching their final destinations.

As the paths of the people in the scene start intersecting (shown in Fig. 4(a)) it becomes more difficult to segment the detections of the individuals and results in the merging of all four people into a new single track as in Fig. (4(b)). Since this is a large group covering a significant portion of the area, very little area is needed to be sacrificed to take a high resolution image of the entire group resulting in the image shown in Fig. (4(c)). During the time period where the four tracks are becoming merged into one, many short lived tracks are created as the detector begins to have difficulty separating

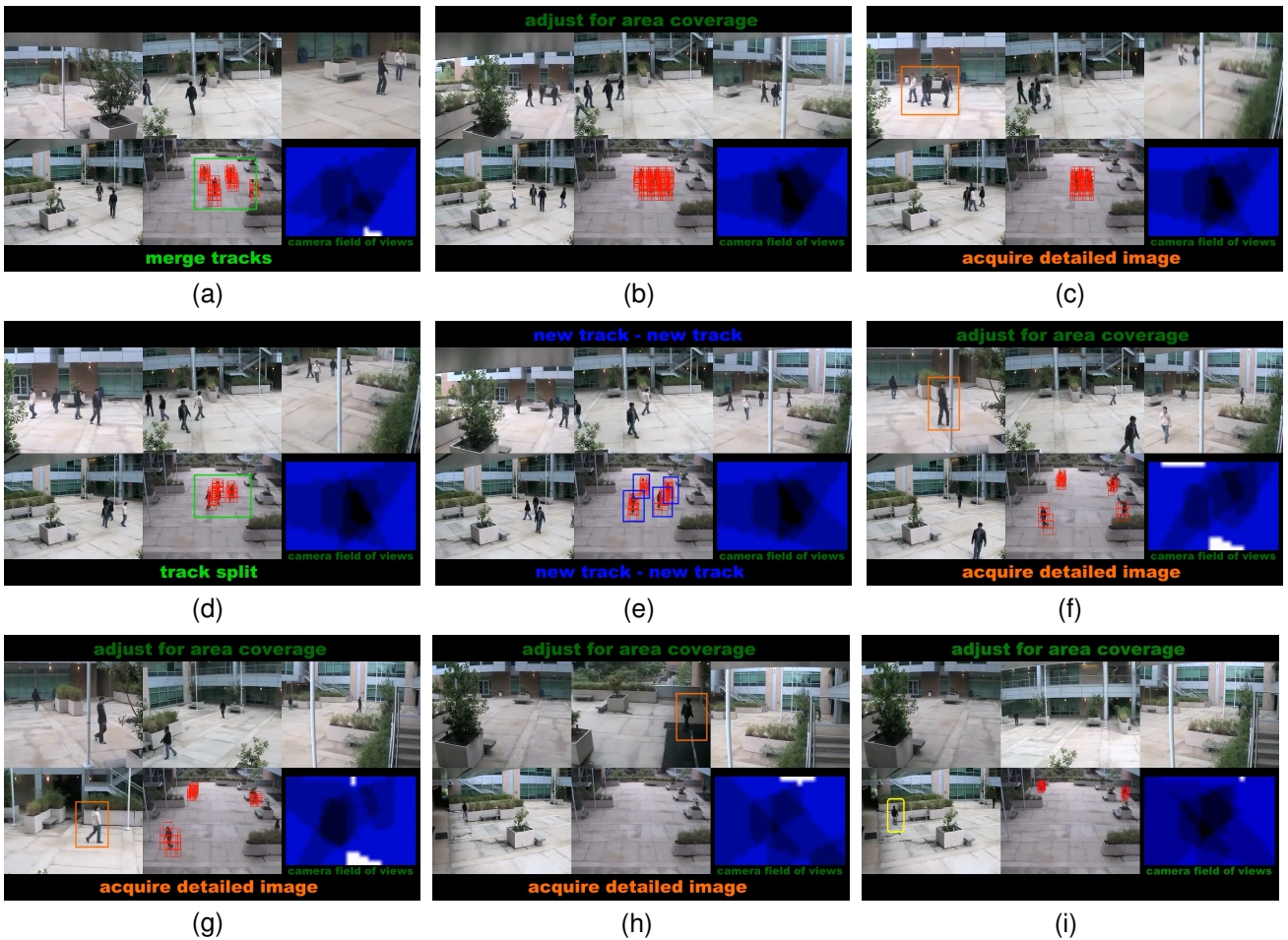


Fig. 4. Shows the sequence of images taken in response to merging and splitting tracks. (a) As the targets move closer together their tracks begin to merge into a single track. (b) The 4 people have grouped together to create a new track. (c) A high resolution image of the entire group is taken by the upper left camera. (d) The targets are beginning to split apart. (e) A new track is created for each individual as there is enough separation to generate stable tracks. (f), (g) and (h) show high resolution images are acquired for 3 of the individual tracks. (i) The 4th target is leaving the surveillance area before the network could take a high resolution shot (yellow box).

individual tracks. If the targets stay in a position where the detector cannot consistently separate or group them together, a long term track will not exist and no high resolution images will be acquired. As the scene progresses, the people start to separate as can be seen in Fig. (4(d)). A short while after the split, distinct tracks for each person can again be formed in Fig. (4(e)). A sequence of high resolution images for each of the targets are then taken as shown in fig. (4c-e). It is important to note that the third high resolution image is only taken after the other two zoomed in cameras have returned to covering a significant portion of the area. As more cameras zoom in for high resolution images, the amount of area covered by overlapping FOVs decreases. Due to the trade-off between the value gained from area coverage and high resolution imaging, the less cameras responsible for covering the area, the less likely a camera will zoom in for a detailed image.

4.2 Analysis of Active Selection

We analyze the active vision scenario (high resolution shots) in terms of number of Pixels. We determine average and maximum number of pixels with and without active selection of cameras. Number of pixels is calculated from every frame within a specific time span for both with and without zoom and averaged over total frames to find average number of pixels. In Table 2, average and maximum number of pixels of whole target and face region of the targets are provided to demonstrate the effect of active control in terms of resolution. Since the operations in aforementioned scenarios are similar, we choose 4 targets from the final scenario to compute the number of pixels. In final scenario, all targets have been zoomed in twice- before track merging and after track splitting. Only target 1 has been zoomed in once- before track

TABLE 2
Quantitative Analysis of Active Control

Targets	Avg. Nb of Pixels (with No control)	Avg. Nb of Pixels (with control)	Max. Nb of Pixels (with No control)	Max. Nb of Pixels (with control)
1. Face	171.70	1058.27	231.00	4228.20
Whole Body	515.11	3386.48	693.00	12896.00
2. Face	176.81	735.15	480.00	4318.94
Whole Body	530.44	2278.96	1440.00	13000.00
3. Face	317.67	852.57	480.00	4201.32
Whole Body	953.01	2728.21	1440.00	12688.00
4. Face	321.00	785.21	783.33	4819.74
Whole Body	963.00	2434.14	2350.00	14652.00

merging; after track splitting it is lost. High resolution shot acquired through active selection strategy gives us detailed information about targets.

4.2.1 Network Delay

In the second scenario as shown in Fig. 4(i), we can see that a high resolution image of one of the targets was not acquired before he exited the area after the group split. This is because the time span from when the individual track was created until the time he left the area was too short for any camera to gain utility from capturing a high resolution image. If the delay of the network was very small or almost instantaneous, it would have been possible to capture the high resolution shot of the target.

4.2.2 Communication Cost

In the centralized scheme data needs to be present at a central server. As the number of cameras increases the communication cost becomes high. There are many distributed frameworks, some are consensus based and some are not. The non-consensus based distributed network presented in [16] allows only one camera to change their PTZ setting at each iteration time. When we scale the size of the camera network this can quickly become a problem. To overcome this, we design a distributed system in a way that only camera nodes that are part of the vision graph will exchange information with each other since they sense the same target. The number of such cameras looking at a particular target will usually be small in most applications. Thus the amount of data that is exchanged between these cameras is much lower than a centralized case since central server needs to access all the information from all the cameras in the network. For example, if there are a total of N_c cameras in the network and N_v in a particular clique of a vision graph (i.e., cameras that sense the same target and thus need to exchange information about that target), the relative amount of data exchanged in the distributed case is N_v/N_c of the centralized case (approximately). For typical values of N_v and N_c , we can see that the distributed case provides significant

TABLE 3
Data Transfer Rate for Distributed (with Different Clique Sizes) and Centralized Scheme.

Distributed Scheme		Centralized Scheme
Clique-2	Clique-3	
5.750 MB/s	8.775 MB/s	15.125 MB/s

reduction in communication cost. In Table 3, data transfer rate for distributed and central network has been provided. Here, $N_c = 5$ and $N_v = \{2, 3\}$ are used in vision graph.

5 CONCLUSION

We described an active camera parameter selection method for capturing high resolution images of targets, when certain events occur, while maintaining coverage of a large area using a camera network. A key contribution of this article is showing how the parameter selection can be distributed in a network so that multiple cameras can change parameters simultaneously while prior approaches allowed only one camera to update its parameters at a time. Future work should consider co-design methodologies whereby the effect of communication limitations or processing capabilities are modeled in the camera control objective functions.

ACKNOWLEDGMENTS

This work was partially supported by NSF grant 1544969.

REFERENCES

- [1] A. K. Roy-Chowdhury and B. Song, "Camera networks: the acquisition and analysis of videos over wide areas," *Synthesis Lectures on Computer Vision*, vol. 3, no. 1, pp. 1–133, 2012.
- [2] 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.
- [3] Z. Tu and P. Bhattacharya, "Game-theoretic surveillance over arbitrary floor plan using a video camera network," *Signal, Image and Video Processing*, pp. 1–17, 2013. [Online]. Available: <http://dx.doi.org/10.1007/s11760-013-0484-8>
- [4] U. M. Erdem and S. Sclaroff, "Automated camera layout to satisfy task-specific and floor plan-specific coverage requirements," *Comput. Vis. Image Underst.*, vol. 103, no. 3, pp. 156–169, 2006.
- [5] R. Tron and R. Vidal, "Distributed Algorithms for Camera Sensor Networks," *IEEE Signal Processing Magazine*, vol. 3, pp. 32–45, May 2011.
- [6] M. Taj and A. Cavallaro, "Distributed and Decentralized Multicamera Tracking," *IEEE Signal Processing Magazine*, vol. 3, pp. 46–58, May 2011.
- [7] Dieber, B. and Micheloni, C. and Rinner, B., "Resource-Aware Coverage and Task Assignment in Visual Sensor Networks," *Circuits and Systems for Video Technology, IEEE Transactions on*, vol. 21, no. 10, pp. 1424–1437, 2011.
- [8] A. D. Bagdanov, A. del Bimbo, and F. Pernici, "Acquisition of high-resolution images through on-line saccade sequence planning," in *Proceedings of the third ACM international workshop on Video surveillance & sensor networks*, 2005, pp. 121–130.
- [9] A. Blake and A. Yuille, Eds., *Active Vision*. MIT Press, 1992.

- [10] F. Qureshi and D. Terzopoulos, "Planning Ahead for PTZ Camera Assignment and Handoff," in *IEEE/ACM Intl. Conf. on Distributed Smart Cameras*, Como, Italy, Aug-Sep 2009, pp. 1-8.
- [11] A. Mittal and L. Davis, "A general method for sensor planning in multi-sensor systems: Extension to random occlusion," *International Journal of Computer Vision*, vol. 76, pp. 31-52, 2008.
- [12] E. Sommerlade and I. Reid, "Information theoretic Active Scene Exploration," in *IEEE Conf. on Computer Vision and Pattern Recognition*, 2008.
- [13] C. Piciarelli, C. Micheloni, and G. Foresti, "PTZ camera network reconfiguration," in *IEEE/ACM Intl. Conf. on Distributed Smart Cameras*, Como, Italy, Aug. 2009, pp. 1-8.
- [14] C. Micheloni, B. Rinner, and G. L. Foresti, "Video Analysis in Pan-Tilt-Zoom Camera Networks," *IEEE Signal Processing Magazine*, vol. 27, no. 5, pp. 78-90, September 2010.
- [15] B. Song, C. Ding, A. T. Kamal, J. A. Farrell, and A. K. Roy-Chowdhury, "Distributed Wide Area Scene Analysis in Reconfigurable Camera Networks," *IEEE Signal Processing Magazine*, vol. 3, pp. 20-31, May 2011.
- [16] C. Ding, B. Song, A. A. Morye, J. A. Farrell, and A. K. Roy-Chowdhury, "Collaborative Sensing in a Distributed PTZ Camera Network," *Image Processing, IEEE Transactions on*, vol. 21, no. 7, pp. 3282-3295, July 2012.
- [17] A. A. Morye and C. Ding and B. Song and A. K. Roy-Chowdhury and J. A. Farrell, "Optimized Imaging and Target Tracking within a Distributed Camera Network," in *American Control Conference*, 2011.
- [18] N. Li and J. Marden, "Designing games for distributed optimization," in *IEEE Conf. on Decision and Control*, Florida, USA, Dec. 2011.
- [19] S. Pellegrini, A. Ess, and L. Van Gool, "Improving data association by joint modeling of pedestrian trajectories and groupings," in *Computer Vision-ECCV 2010*. Springer, 2010, pp. 452-465.
- [20] Z. Qin and C. R. Shelton, "Improving multi-target tracking via social grouping," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 1972-1978.
- [21] A. Chakraborty, A. Das, and A. K. Roy-Chowdhury, "Network consistent data association," *IEEE transactions on pattern analysis and machine intelligence*, 2016.
- [22] X. Chen, Z. Qin, L. An, and B. Bhanu, "An online learned elementary grouping model for multi-target tracking," in *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*. IEEE, 2014, pp. 1242-1249.
- [23] L. Bazzani, M. Cristani, and V. Murino, "Decentralized particle filter for joint individual-group tracking," in *CVPR*. IEEE, 2012, pp. 1886-1893.
- [24] W.-S. Zheng, S. Gong, and T. Xiang, "Associating groups of people," in *BMVC*, vol. 2, 2009, p. 6.
- [25] Y. Cai, V. Takala, and M. Pietikäinen, "Matching groups of people by covariance descriptor," in *ICPR*. IEEE, 2010, pp. 2744-2747.
- [26] R. Vezzani, D. Baltieri, and R. Cucchiara, "People reidentification in surveillance and forensics: A survey," *ACM Computing Surveys (CSUR)*, vol. 46, no. 2, p. 29, 2013.
- [27] Yilmaz, Alper and Javed, Omar and Shah, Mubarak, "Object tracking: A survey," *ACM Comput. Surv.*, vol. 38, no. 4, Dec. 2006. [Online]. Available: <http://doi.acm.org/10.1145/1177352.1177355>
- [28] T. Ueshiba and F. Tomita, "Plane-based Calibration Algorithm for Multi-camera Systems via Factorization of Homography Matrices," in *IEEE Intl. Conf. on Computer Vision*, 2003.
- [29] A. T. Kamal, J. Bappy, J. Farrell, and A. Roy-Chowdhury, "Distributed multi-target tracking and data association in vision networks," *PAMI*, 2015.
- [30] A. T. Kamal, J. Farrell, A. K. Roy-Chowdhury et al., "Information weighted consensus filters and their application in distributed camera networks," *Automatic Control, IEEE Transactions on*, vol. 58, no. 12, pp. 3112-3125, 2013.
- [31] A. T. Kamal, J. Farrell, A. K. Roy-Chowdhury et al., "Information consensus for distributed multi-target tracking," *CVPR*, pp. 2403-2410, 2013.
- [32] Z. Zivkovic, "Improved adaptive gaussian mixture model for background subtraction," in *ICPR*, vol. 2, 2004, pp. 28-31.



time systems.



Chong Ding received B.S. (2008) and Ph.D. (2013) degrees in computer science from the University of California, Riverside. He is a Research Staff member in the Information and System Sciences Laboratory at HRL Laboratories, LLC located in Malibu, California. He is the author of several technical publications on intelligent camera networks and his primary research interests include autonomous vehicles, mobile ad-hoc networks, network security, and distributed and real-

Jawadul H. Bappy received the B.S. degree in Electrical and Electronic Engineering from the Bangladesh University of Engineering and Technology, Dhaka in 2012. He is currently pursuing his Ph.D. degree in Electrical and Computer Engineering at University of California, Riverside. His main research interests include wide area scene analysis, scene understanding, object recognition and machine learning.



and as President in 2014. He is a Fellow of the IEEE, a Fellow of AAAS, a Distinguished Member of IEEE CSS, and author of over 200 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). He is interested in applied research related to estimation, planning, and control of intelligent autonomous agents.



Amit K. Roy-Chowdhury received the Bachelor's degree in electrical engineering from Jadavpur University, Calcutta, India, the Master's degree in systems science and automation from the Indian Institute of Science, Bangalore, India, and the Ph.D. degree in electrical engineering from the University of Maryland, College Park.

He is a Professor of electrical engineering and a Cooperating Faculty in the Department of Computer Science, University of California, Riverside. His broad research interests include the areas of image processing and analysis, computer vision, and video communications and statistical methods for signal analysis. His current research projects include intelligent camera networks, wide-area scene analysis, motion analysis in video, activity recognition and search, video-based biometrics (face and gait), biological video analysis, and distributed video compression. He is a coauthor of two books *Camera Networks: The Acquisition and Analysis of Videos over Wide Areas and Recognition of Humans and Their Activities Using Video*. He is the editor of the book *Distributed Video Sensor Networks*. He has been on the organizing and program committees of multiple computer vision and image processing conferences and is serving on the editorial boards of multiple journals.