

Data Fusion in Ubiquitous Networked Robot Systems for Urban Services

Luis Merino · Andrew Gilbert · Jesús Capitán · Richard Bowden · John Illingworth · Aníbal Ollero

the date of receipt and acceptance should be inserted later

Abstract There is a clear trend in the use of robots to accomplish services that can help humans. In this paper, robots acting in urban environments are considered for the task of person guiding. Nowadays, it is common to have ubiquitous sensors integrated within the buildings, such as camera networks; and wireless communications, like 3G or WiFi. Such infrastructure can be directly used by robotic platforms. The paper shows how combining the information from the robots and the sensors allows tracking failures to be overcome, by being more robust under occlusion, clutter and lighting changes. The paper describes the algorithms for tracking with a set of fixed surveillance cameras and the algorithms for position tracking using the signal strength received by a Wireless Sensor Network (WSN). Moreover, an algorithm to obtain estimations on the positions of people from cameras on board robots is described. The estimate from all these sources are then combined using a decentralised data fusion algorithm to provide an increase in performance. This scheme is scalable and can handle communication latencies and failures. We present results of the system operating in real time on a large outdoor environment, including 22 non-overlapping cameras, WSN and several robots.

1 Introduction

There is an increasing interest in service robotics, that is, robot systems that provide services to human users. The EU Project called URUS (Ubiquitous Net-

L. Merino
School of Engineering, Pablo de Olavide University, 41013, Seville, Spain

A. Gilbert · R. Bowden · J. Illingworth
Centre for Vision Speech and Signal Processing, University of Surrey, Guildford, GU2 7XH, UK

J. Capitán
Institute for Systems and Robotics, Instituto Superior Tecnico, Lisbon, Portugal

A. Ollero
School of Engineering, University of Seville, Spain
Centre for Advanced Aerospace Technology, Parque Tecnológico y Aeronáutico de Andalucía, C. Wilbur y Orville Wright 17-19-21, 41309, La Rinconada, Spain

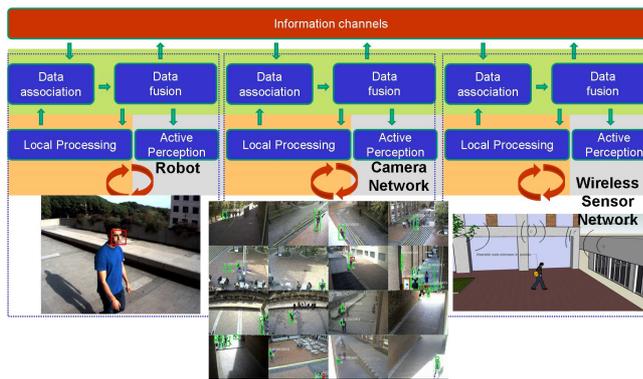


Fig. 1: Overview of the URUS perception system. Information from robots, fixed cameras and a WSN is considered.

working **R**obotics in **U**rban **S**ettings) [43] considers a team of mobile robots, a set of static cameras and a Wireless Sensor Network (WSN) in a urban environment to offer urban services. All these elements can communicate through wireless links (using 3G and WiFi), and constitute an example of a system with Ubiquitous Networked Robots (UNR). The system integrates robots, sensors, communications and mobile devices in a cooperative way, which means not only a physical interconnection between these elements, but also, for example, the development of novel intelligent methods of cooperation for task-oriented purposes.

The system is used, among other tasks, for person guiding by robots. This scenario consists of having a robot guiding a person towards a given destination. This person is to be detected and then must be tracked continuously to allow the robot to accomplish its task. Real scenarios involve dynamic environments and varying conditions. The robustness and reliability of autonomous perception in these scenarios are critical. In most cases, a single autonomous entity (i.e. a robot or a static surveillance camera) is not able to acquire all the information required for the application because of the characteristic of the particular task, i.e. loss of visibility. Thus, the cooperation among robots and between robots and heterogeneous sensors embedded in the environment through information fusion is relevant. The set of fixed cameras can obtain global views of the scene; however, as they are static, they cannot react to non-covered zones and illumination changes such as shadows can affect the system. Robots carry local cameras and can move to suitable positions, reacting to the changing conditions. However, their field of view is limited and they can lose the person they are tracking. Wireless devices can also help to localize the people, estimating their positions by measuring the signal strength from different static receivers, but the resolution obtained is usually low, and depends on the density of anchored receivers.

This paper presents the perception system developed in URUS for the UNR (see Fig. 1). The objective was to create a cooperative perception system, in which the different elements of the UNR system collaborate to obtain more precise information for the task assigned. In order to cope with scalability, a decentralised data fusion algorithm is employed, in which only local estimations and local communications are used. The main novelty of this work is the fusion of the various

elements into a single system. Furthermore, some aspects of the work have been only demonstrated in simplified settings. Here we extend the techniques and apply the approaches to a real outdoors application, using the fusion to overcome the limitations of any single approach. This paper builds on the preliminary work of the authors [15]; however, it has been extensively enhanced in terms of description and evaluation. The different techniques employed are described more thoroughly, including new parts of the system and the approach is thoroughly evaluated with extensive complex outdoor tests.

The next section presents related work. After an overview of the full system in Section 3, the paper will present the individual input sensor algorithms. Thus, Section 4 firstly describes the process to extract information from a set of fixed cameras. In Section 5 the main ideas of person tracking from on-board robot cameras is presented, and Section 6 describes the use of the signal strength from wireless sensors for tracking. Finally, the results of the tracking from all sensors are used to infer the position of the person in a global coordinate system through a data fusion process. This system is described in Section 7. The paper ends showing results obtained during the experiments of the URUS project, in an urban scenario involving 22 fixed cameras, a WSN of 30 wireless Mica2 nodes, and several robots.

2 Related Work

2.1 Tracking with camera networks

There has been many attempts to track people and other moving objects inter camera. The early tracking algorithms [4, 27, 7, 33] require both camera calibration and overlapping fields of view to compute the handover of objects of interest between cameras. Others [28, 22] can work with non-overlapping cameras but still require calibration. The use of non-overlapping cameras is more realistic and reduces hardware limitations, though, as the handover of objects between cameras cannot be explicitly observed, reasoning must be used. Probabilistic approaches have been therefore proposed; [22] presents an approach to track cars on a highway modelling appearance and transition times as Gaussian distributions, though this was within a relatively controlled environment. Through the use of supervised off-line learning period it has been possible to model the camera topology and path probabilities of objects [48, 8]. However, often probabilistic solutions for non-overlapping cameras are used under restrictive assumptions or within limited applications. Illumination changes between cameras can be a challenge therefore a number of approaches [39, 16, 23] adjust the colour appearance of cameras to improve performance.

Approaches have been proposed that do not require *a priori* correspondences to be explicitly stated [25, 26, 9]; instead they use the observed motion over time to establish reappearance periods. Nevertheless, batch processing was initially performed on the data. Therefore, if the environment changes significantly, the system must be “rebooted” and correspondences re-learned. The use of local invariant features has also seen an increase in popularity over the last few years [11, 17].

Within this paper, the fixed camera tracking builds on the limited work of [14]. The camera transition relationships were incrementally learnt, to model both the

colour variations and posterior probability distributions of the spatio-temporal transition links between cameras.

2.2 Radio-signal based tracking

There is an increasing interest in systems that use the signal strength received by wireless devices for localization purposes. Many systems are devoted to the localization of static devices by using the signal received from a small set of very well-localized static devices (called beacons) [40], or the signal received by a well-localized mobile node, usually on board a robot [3].

There is also work devoted to the tracking of mobile nodes by using radio signals, which is the problem of estimating the position of a mobile node from the signal received by a set of static devices whose position are known. A tutorial on the main issues and approaches for the problem is presented in [20]. Many algorithms use, beside signal strength, additional information to obtain range estimates or even direction of arrival estimates. For instance, [37] considers the use of particle filters for tracking a mobile node using Time of Arrival, Difference of Time of Arrival and power measurements, presenting results in simulation. The works [30,29] use the Doppler shift of interference signals to estimate the velocity and position of mobile nodes. These approaches require the precise synchronization of the emission of signals. In our approach, only signal strength is used, through a calibrated model for radio propagation. Particle filters are used in [44,34] for localization in indoors scenarios. Here, a similar approach is used, but outdoors urban scenarios are considered. Moreover, the previous works require a full calibration of a signal map model. Here, a simple model of radio propagation is calibrated, and map information is used just in the prediction phase of the filter. Also, the tracking is benefited from the data fusion with other sensor networks. It is worth to mention that there are approaches in which the signal strength model is learnt [12,21].

2.3 Robot person tracking

Tracking from mobile platforms like robots in outdoor scenarios is a hard problem affected by clutter, illumination changes in the case of vision approaches, occlusions, etc. Most of the approaches combine people detection and people tracking modules for this task. The people detection module tries to obtain person hypotheses analyzing the sensor data, and is usually computationally demanding. Many classification techniques are used for this task, like boosting [32], SVM [35], etc. The tracking module is usually a feature tracking algorithm applied to the initial hypothesis given by the detection module, which can be run at a higher rate than the detection algorithm, like CamShift [2]. In most cases, both modules support each other, so when the tracker is lost new hypotheses from the detector can be used. More complex combinations, including what is called cognitive feedback are also considered [10,13].

In the work presented here, a combination of state of the art algorithms for detection [49] and tracking [2] from a single robot camera is used. This algorithm works relatively well, although outdoor scenarios pose difficulties to it. The key

issue in the paper is to show how the combination of the local information obtained by the robot with the information received from other elements of the UNR system can improve the results.

2.4 Bayesian decentralised data fusion

Fusion of data gathered from a network of heterogeneous sensors is a highly relevant problem in robotics that has been widely addressed in the literature. Most of those works are based on Bayesian approaches, where the sensors are modeled like uncertain sources. However, there are alternative methods for dealing with uncertainties.

Possibility theory, which is built on the arithmetic of fuzzy sets, has been used for uncertain reasoning. For instance, in [5] possibility theory is used for cooperative localisation and ball position estimation within the framework of Robocup. The decentralization of possibility-based systems system is not clear though. Moreover, fuzzy techniques are mainly suitable for control systems, in which there are small sets of rules and no chains of inferences. In consensus theory [41, 38], the idea is to reach agreements among different sources, so it can also be used for sensor fusion. These techniques are typically used to achieve collective coordinated dynamics in multi-agent systems. They are, however, less adequate for applications in which the state evolves with time (e.g. tracking). For typical consensus algorithms, convergence to a same value among the sources is proved, but it cannot be assured that this value is the correct one [41].

In multi-sensor data fusion, Bayesian approaches provide a sound mathematical framework and allow for a better modelling when the uncertain sources are complex. Moreover, in some cases (like the one proposed in this paper), they can be decentralised in an efficient manner. Even though fusing all the information in a central node is simpler, decentralised systems are scalable and more robust under communication failures, since only local information and communication are used.

The main issues and problems with decentralised information fusion in Bayesian settings can be traced back to the work of [18], where the Information Filter (IF, dual of the Kalman Filter) is used as the main tool for data fusion for process plant monitoring. The works [18, 46, 36] demonstrate that, for the case of static states (for instance, in mapping applications, when estimating the location of a set of static objects), the decentralised implementation of the IF allows a local estimation that is the same as the one obtained by a centralised IF with access to all the information (provided that sufficient information is exchanged). In the case of dynamic states, for instance in tracking applications (like the one considered here), it was noticed in [42, 1] that if only an estimation about the last state is exchanged between the decentralised nodes, information will be missed with respect to a centralised node. The problem is due to the fact that there are some information not taken into account when performing the prediction steps in each fusing node.

The idea of the Channel Filters in order to fuse the information in a consistent manner (non-overconfident) is considered in the works [31, 19]. These works require a fixed topology between nodes with no loops. Other options are conservative fusion rules that achieve a consistent estimation without the need for Channel Filters when no assumptions can be made about the network topology, like the *Covariance*



Fig. 2: 16 of the 22 cameras within the experiment system. The cameras can provide overall information about the complete scenario.

Intersection algorithm [24]. Moreover, [47] presents the *Covariance Union* method, which tries to deal with disagreement in a Gaussian decentralised fusion setup.

In this paper, an Decentralized Information Filter over the state trajectory is proposed as the main algorithm for scalable data fusion. It will be seen that the exact centralised estimation can be recovered due to the use of delayed states. In previous works like [31], this fusion was not possible for dynamic states without missing some information. Besides, the filter proposed deals with these issues and communication delays in an efficient manner.

3 UNR System Overview

The Ubiquitous Networked Robots system developed in the URUS Project consists of a team of mobile robots, equipped with cameras and other sensors for localisation, navigation and perception; a fixed camera network for environment perception (see Fig. 2); and a Wireless Sensor Network that uses the signal strength of the received messages from a mobile device to determine the position of a person carrying it. An architecture for urban robots networked with the environment has been developed [43]. This architecture provides a decisional layer and communication capabilities. The robots can switch between WiFi and 3G to communicate with the control station, other robots and the camera and sensors networks. This paper is focused on the perception part of the system.

From the perception point of view, the information obtained by the fixed camera network or the Wireless Sensor Network can be shared with the local information from each robot to improve the perception. That way, each robot obtains a better picture of the world than it would do alone. In this case, the tracks on the image plane obtained by the camera network will be fused with the information from the other systems (robots and WSN) to improve the tracking of the person being guided. Thus, it is possible to cope with occlusions and obtain better tracking capabilities, since information of different modalities is employed; and cope with non-covered zones, since the robots can move to cover these zones.

The system consists of a set of fusion nodes which implements a decentralised data fusion algorithm. Each fusion node only employs local information (data from

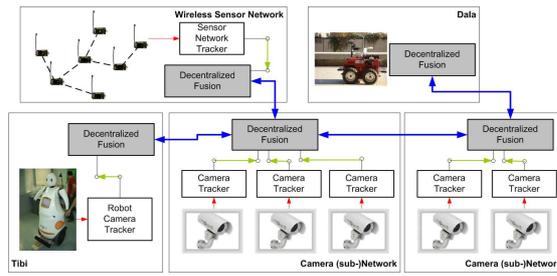


Fig. 3: A block description of the URUS perception system. The different subsystems are integrated in a decentralised manner through a set of decentralised data fusion nodes. Locally, each system can process and integrate its data in a central way (like the WSN) or in a distributed way (like the camera network). Some systems can obtain information from the rest of the network even in the case they do not have local sensors.

local sensors; for instance, a camera subnet, or the sensors on board the robot) to obtain a local estimation of the variables of interest (in this case, the position of the person being tracked). Then, these nodes share their local estimations among themselves if they are within communication range. As the nodes only use local data and communications, the system is scalable. Also, as each node accumulates information from its local sensors, temporal communication failures can be tackled without losing information.

Figure 3 shows a block description of the system, for the camera network, each fusion node considers information from a small subset of cameras, which are processed in a distributed way, with a separate tracker obtaining estimations from each camera. The WSN process messages from all the network in a gateway to localize the mobile node using the signal strength. Similarly the on-board robot camera tracks nearby people. Then, the local estimations of the different nodes are fused in a decentralised way.

4 Fixed Camera Tracking

The fixed cameras cover a wide area of the experiment site and therefore, in many cases, they are the foundation for the fusion of the other sensors within the UNR; they are able to track objects of interest both on and across different cameras without explicit calibration periods.

Figure 4 gives a general overview of the processes for a single camera. Each camera is a self-contained node connected to others via a network, meaning that it can easily be distributed over multiple processors or machines. There are two inputs to the node, the camera image pixel values at (a) and also the estimated location of objects from previous frames at (f). The estimated location of previous frames allows the Kalman filter to use data from other sources to overcome occlusion that will occur with the foreground objects. There is a single output (e), this contains the location in camera coordinates of all the detected and tracked objects in the frame.

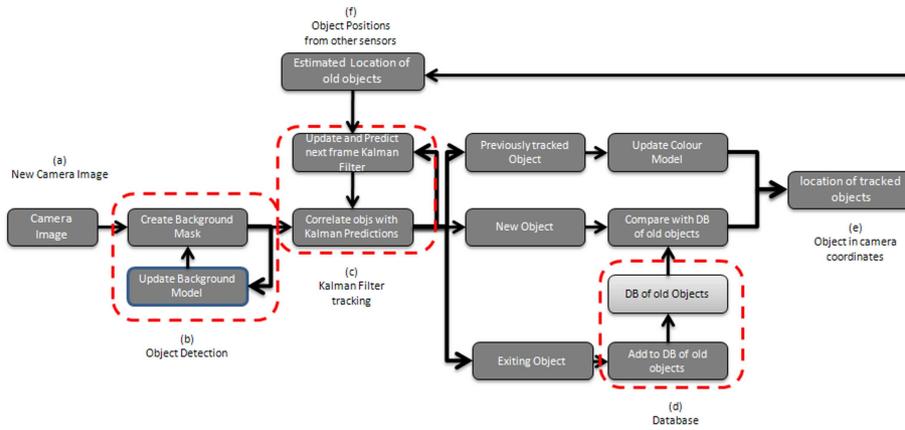


Fig. 4: System Overview of tracking objects using a single camera

To detect moving objects within an image, the static background is modelled in a similar fashion to that originally presented by [25]. The foreground objects are identified from the background mask through connected component analysis. This provides a bounding box centred over each object. A path track of each object over time is created using correlation between frames. A Kalman Filter is used to provide temporal correspondence between the detected foreground objects inter frame. A histogram is used as an objects descriptor as it is spatially invariant and, through quantisation, a degree of invariance to illumination can be achieved. Each object is then given a unique label for identification. Figure 5 shows an example of tracking multiple moving objects on a single camera at the experimental site.

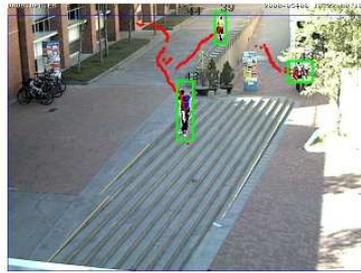


Fig. 5: The track paths of three objects over time

When the object of interest enters a new camera, the transfer of the object's label to the new camera is a challenge as the cameras have no overlapping fields of view, making many traditional image plane calibration techniques impossible. In addition, the large number of cameras means traditional time consuming calibration is infeasible. Therefore the approach needs to learn the relationships between the cameras automatically. This is achieved by way of two individually weak cues,

modelling the colour, and movement of objects inter camera. These two weak cues are then fused to allow the technique to determine if objects have been previously tracked on another camera or are new object instances. By incrementally learning the cues over time, the accuracy is able to increase without any supervised input.

4.1 Forming temporal movement links inter camera

We make use of the key assumption that, given time, objects (such as people) will follow similar routes inter camera and that the repetition of the routes will form marked and consistent trends in the overall data. These temporal transition links are used to link regions of the cameras together, producing a probabilistic distribution of object movements between cameras.

Initially the system is divided so that each camera is defined as a single region. This coarse detail allows immediate operation and tracking of people, unlike the approaches of [25,9], that require batch processing. While incoming data is stored to allow refinement and subdivision of the entry and exit regions using people tracked inter camera. All newly detected objects are compared to previously tracked objects within a set time window. The colour similarity is calculated and used to increment a probability distribution with respect to the reappearance period, as shown in Fig. 6. Regularly, the noise floor level is measured for each link, if the maximum peak of the distribution is found to exceed the noise floor level, this indicates a possible correlation between the two regions. Figure 6 shows the probability distribution for two regions with a distinct link at around 13 seconds.

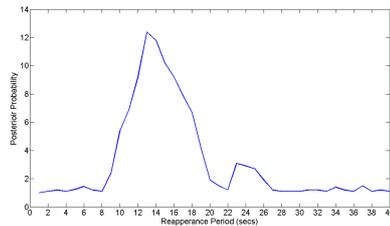


Fig. 6: An example of a probability distribution, $f^{x|y}$ showing a distinct link between two regions x and y

When a link is found between two regions, they can be subdivided to create four new equal sized sub-regions, this aims to increase the detail level of the entry and exit areas. The previous data is then reused and incorporated with future evidence to form links in the newly subdivided regions.

It is likely that many of the subdivided regions will not form coherent links; therefore, if a link between two regions has no data in it, it is removed to minimise the number of links maintained. In addition, if a region is found to have no links to any other region, the region is also removed. This policy of removing unused and invalid regions improves system scalability. As the process proceeds, the regions start to visually represent entry and exit points of the cameras.

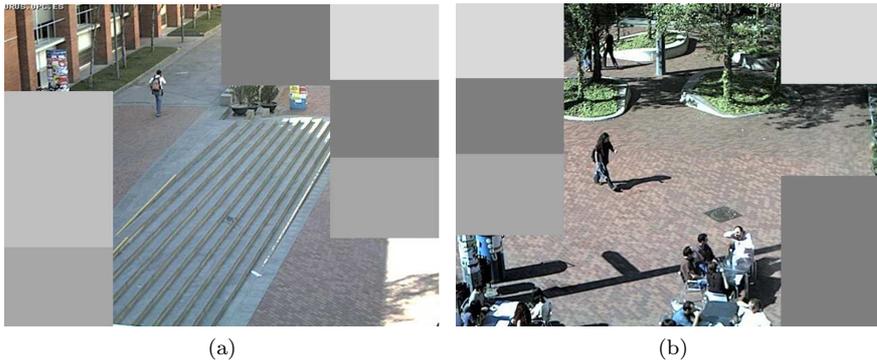


Fig. 7: Entry and exit regions of two of the cameras, the darker the region the greater the quantity of track entry and exits

Figure 7 shows the entry and exit regions of two cameras in an experimental site. Even though regions cover initially the whole image, over time they start to represent the main entry and exit area of the cameras. The regions are continuously subdivided if a link to another region is found. However when the regions are subdivided a number of times it is likely that many neighbouring regions contain similar links to other neighbouring regions. Therefore the correlation of neighbouring region’s link distributions is examined using Bhattacharyya coefficient. If the distributions from two neighbouring regions are found to be highly correlated, the regions spatial areas are combined to form a single region, to further increase the overall number of samples.

4.2 Modelling colour variations

The colour quantisation descriptor used to form temporal reappearance links in the previous section assumes a similar colour response between cameras. However this is seldom the case, especially on outdoor environments. Therefore, a colour calibration of the cameras is proposed that can be learnt incrementally simultaneously with the temporal relationships discussed in the section above. The idea of using a colour transformation matrix to calibrate the cameras has been proposed before, however the experiments are often inside and limited [23, 39, 14]. we propose to use the colour transformation matrix to calibrate the multiple outdoor cameras. The people tracked inter camera are used as the calibration objects, and a transformation matrix is formed incrementally to model the colour changes between specific cameras.

Between each set of cameras a colour transformation matrix is formed. Initially, this is an identity matrix. This assumes the ideal case of a uniform prior of colour variation between cameras. When a person is tracked inter camera and is identified as the same object, the difference between the two colour descriptors is modelled by a transform matrix. The matrix is calculated by computing the transformation that maps the person’s descriptor from the previous camera to the person’s current descriptor. This transformation is computed via SVD. The

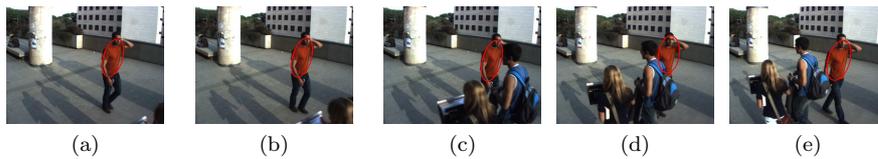


Fig. 8: By using the local people tracking and detection module, the robot can obtain estimations of the person on the image plane. The employed tracking algorithm is able to handle temporal occlusions.

matrix is then averaged with the appropriate camera transformation matrix, and repeated with other tracked people to gradually build a colour transformation between cameras. This method will introduce small errors, however it is in keeping with the incremental theme of the work. This is especially important when the environment is outdoors, as the lighting between cameras is constantly changing slightly, and this approach allows the method to continually update and adapt to the colour changes between cameras over time.

With the weak cues learnt, the reappearance probability of an object is then used to weight the observation likelihood obtained through colour similarity to obtain a posterior probability of a match. Tracking objects is then achieved by maximising the posterior probability within a set time window.

5 Robot Camera Tracking

The robots carry on-board cameras that are used for person guiding. These cameras can be used to obtain local estimations on the position of the person to be guided. A combination of state of the art algorithms for person detection and tracking is used.

The person detection algorithm applied to the image is the one of [49]. This detection module is launched when the robot is requested to guide a person and it is close to the location where the person is waiting. Once the person is detected, it is tracked using an algorithm based on the CamShift technique [2]. While the algorithm is able to handle temporal occlusions (see Fig. 8), due to changes in illumination, the changing field of view of the camera when the robot moves, or even the person going out of the field of view, the tracking system is not sufficient to maintain the track on the person continuously. Therefore, the results from the tracking and the detection applications are combined, so that the robot employs the person detector whenever the tracker is lost to recover the track. The algorithm determines that the person is lost employing some heuristics, like the track going out to the limits of the image or size restrictions on the blob. As a result, the robots can obtain estimations of the pose of the person on the image plane.

More complex algorithms like [13] could be used and included into the system. However, it is out of the scope of the paper to develop a robust person tracking system based on a single camera on-board the robot. Nevertheless, the claim is that, in general, a system based only on local information will not be robust enough to be able to guide one person through the whole scenario. Moreover, information

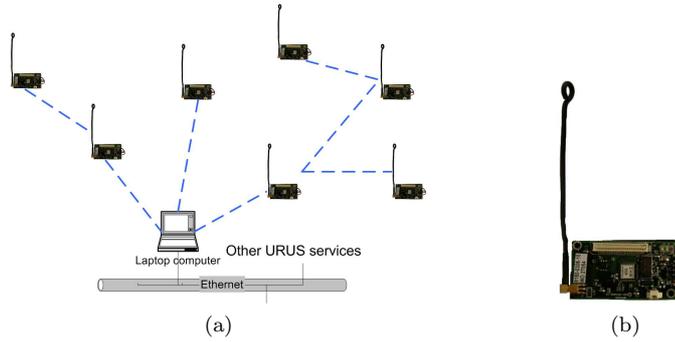


Fig. 9: Left: WSN architecture: the nodes can establish an ad-hoc network to relay information to a gateway. This information (messages indicating the power received from a mobile node) is processed in the gateway to obtain an estimation of the position of the mobile node. Right: WSN Mica2 node.

from one camera alone is not sufficient to estimate the full 3D position of the person. The following sections will show how the combination of the local camera information and the information from the other subsystems (camera network and WSN) can overcome these problems.

6 Wireless Sensor Network Tracking

Another element considered in the UNR system is a network of wireless Mica2 sensor nodes. These Mica2 nodes are able to sense different quantities, like pressure, temperature, humidity, etc. Moreover, they have wireless communication devices, and are able to form networks and relay the information they gather to a gateway (see Fig. 9).

In addition, the signal strength received by the set of static nodes (Received Signal Strength Indicator, RSSI) can be used to infer the position of a mobile object or a person carrying one of the nodes (the emitter). In the application considered here, the user that wants to be guided carries one of the nodes.

The algorithm to estimate and track the node position is based on particle filtering. In the particle filter, the current belief about the position of the mobile node is represented by a set of particles $\{\mathbf{x}_t^{(i)}\}$, which represent hypotheses about the current position of the person that carries the node (see Fig. 10).

In each iteration of the filter, kinematic models of the motion of the person and map information are used to predict the future position of the particles. The likelihood of these particles is updated when new messages are received from the static network. The technique is summarized in Algorithm 1, where \mathbf{z}_t^j is the measurement provided by each static node j , consisting of its position \mathbf{x}^j and the strength $RSSI_t^j$ of the received signal from the mobile node. Following subsections describe the main steps in this algorithm.

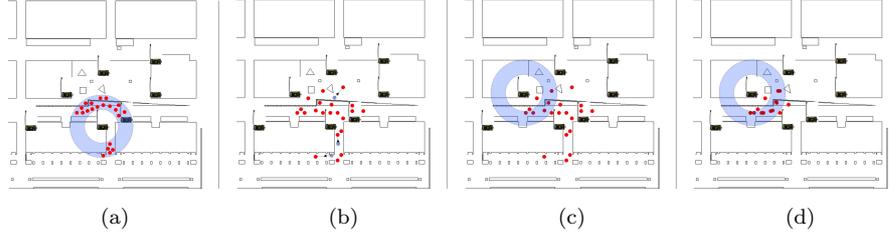


Fig. 10: Particles (red) are used to represent person hypotheses. (a) The filter is initiated when the first message is received by sampling uniformly from a spherical annulus around the receiver. (b) Particles are predicted at each iteration; map information is also taken into account. (c) New measurements are used to weight the particles. (d) Re-sampling maintains the number of particles bounded.

Algorithm 1 $\{\mathbf{x}_t^{(i)}, \omega_t^{(i)}; i = 1, \dots, L\} \leftarrow \text{Particle_filter}(\{\mathbf{x}_{t-1}^{(i)}, \omega_{t-1}^{(i)}; i = 1, \dots, L\}, \mathbf{z}_t^j = \{\mathbf{x}^j, RSSI_t^j\})$

```

1: for  $i = 1$  to  $L$  do
2:    $\mathbf{x}_t^{(i)} \leftarrow \text{sample\_kinematic\_model}(\mathbf{x}_{t-1}^{(i)})$ 
3: end for
4: if Message from network  $\mathbf{z}_t^j$  then
5:   for  $i = 1$  to  $L$  do
6:     Compute  $d_t^{(i)} = \|\mathbf{x}_t^{(i)} - \mathbf{x}^j\|$ 
7:     Determine  $\mu(d_t^{(i)})$  and  $\sigma(d_t^{(i)})$ 
8:     Update weight  $\omega_t^{(i)} = p(RSSI_t^j | \mathbf{x}_t^{(i)}) \omega_{t-1}^{(i)}$  with  $p(RSSI | \mathbf{x}_t^{(i)}) = \mathcal{N}(\mu(d_t^{(i)}), \sigma(d_t^{(i)}))$ 
9:   end for
10: end if
11: Normalize weights  $\{\omega_t^{(i)}\}, i = 1, \dots, L$ 
12: Compute  $N_{eff} = \frac{1}{\sum_{i=1}^L (\omega_t^{(i)})^2}$ 
13: if  $N_{eff} < N_{th}$  then
14:   Re-sample with replacement  $L$  particles from  $\{\mathbf{x}_t^{(i)}, \omega_t^{(i)}; i = 1, \dots, L\}$ , according to the weights  $\omega_t^{(i)}$ 
15: end if

```

6.1 Prior prediction and importance functions

As a prior, the filter is initialised with the first message received from the mobile node, considering an uniform distribution on a spherical annulus around the receiver. The map of the scenario is taken into account when sampling from this prior (see Fig. 10a), considering that the person is not inside any building.

Each time step, the position of the particles are predicted from their previous position (line 2 of Algorithm 1). No further information is assumed, and similarly to [44] the prediction function employed is a Brownian motion model, in which new particles are sampled from a Gaussian distribution centered at the previous particle position [45] (Fig. 10b). However, this model also considers map information to discard infeasible motions (like going through walls).

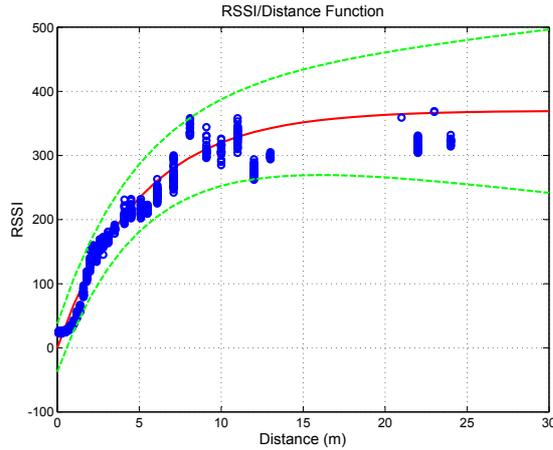


Fig. 11: RSSI-Distance functions, $\mu(d_k)$ and $\sigma(d_k)$. These functions relate the distance between two nodes and the RSSI received in mean and std. deviation. It has been experimentally computed using a large set of RSSI/distance couples. The RSSI representation is the one used in the Mica2 nodes, 0 is the maximum signal strength and 375 the minimum. Dots: A sub-set of the experimental set of data. Solid line: Estimated mean $\mu(d_k)$. Dashed lines: standard deviation confidence interval based on $\sigma(d_k)$.

6.2 The likelihood function

The likelihood function $p(RSSI_t|\mathbf{x}_t)$ plays a very important role in the estimation process, since each time a message is received, this likelihood is used to update the particles weights (lines 5 to 9, Fig. 10c). The likelihood models the correlation that exists between the distance that separate two nodes and the *RSSI* value. Figure 11 shows experimental data on the RSSI values for given distances. It can be seen that the correlation between RSSI and distance decreases with the distance between the two nodes, transmitter and receiver [3]. This is mainly caused by radio-frequency effects such as radio reflection, multi-path or antenna polarisation.

The model used here considers that the conditional density $p(RSSI_t^j|\mathbf{x}_t)$ can be approximated as a Gaussian distribution for a given distance $d_t^j = \|\mathbf{x}_t - \mathbf{x}^j\|$ between the mobile node and static node j , as follows:

$$RSSI_t^j = \mu(d_t^j) + \mathcal{N}(0, \sigma(d_t^j)) \quad (1)$$

From the experimental data, it can be seen that this model can represent adequately the relations for distances below 15 meters. At the same time, this function allows an efficient evaluation within the particle filter. Please notice that the functions $\mu(d_t^j)$ and $\sigma(d_t^j)$ are themselves non-linear functions of the distance (which itself is a non-linear function of the state), so the full model is non-linear. These functions are estimated during a calibration procedure (Figure 11 shows also the estimated functions for the data set).

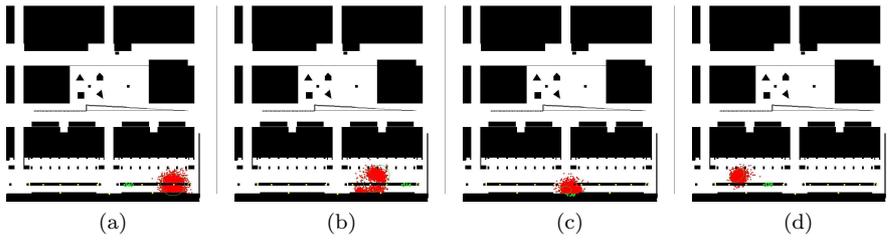


Fig. 12: A sequence of the 500 particles employed in the filter for an experiment. Red points represent the particles. Yellow points represent the static nodes, being the green one the emitter at each frame.

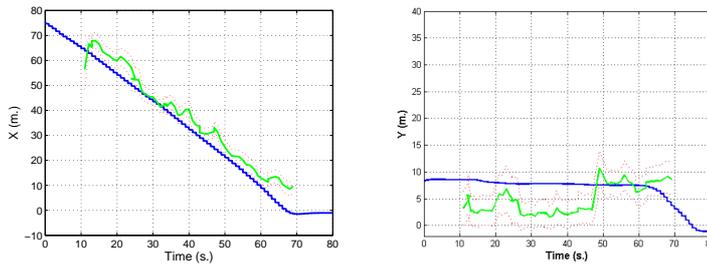


Fig. 13: Estimated position of the person by the WSN (green) and position of the guiding robot (blue) estimated by using its navigation modules. Dashed lines represent sigma intervals.

6.3 Filter evolution

Although Section 8 will show additional results, Figure 12 presents the evolution of the particles for a particular tracking experiment performed at the experimental site. 500 particles are employed, and the algorithm runs at more than 1 Hz. Figure 13 shows the estimated position of the mobile node carried by a person estimated by the WSN. It is compared to that of the guiding robot, which is some meters ahead. In this experiment, this means that the robot has a lower X and higher Y coordinates.

When the filter converges to a Gaussian distribution, the estimated mean and covariance can be fed to the decentralised fusion system that will be explained in the next section. The convergence is determined by analyzing the Kullback-Leibler divergence between the particle distribution and a Gaussian distribution with the same mean and covariance as the particles.

7 Decentralised Data Fusion for Person Tracking

Using the trackers described above, the camera network, the robots and the WSN are able to obtain local estimations of the position of the people on the image

plane or in a 3D coordinate system. This information, characterised as Gaussian distributions (mean and covariance matrix), can be fused in order to obtain a more accurate estimation of the 3D position of the person.

As commented in Section 3, the idea is to implement a decentralised fusion approach, in which each fusion node only employs local information (data from local sensors; for instance, a camera subnet, or the sensors on board the robot), and then *shares* its estimation with neighbouring nodes (see Fig. 3, right). Thus, scalability and robustness are improved and bandwidth requirements alleviated.

The question is how to integrate measurements from all the sources and deal with communication delays without losing any information with respect to a centralised solution. A novel Bayesian filter that keeps track of delayed states is proposed to recover exactly the centralised estimation. Moreover, it is shown how in the case of Gaussian distributions, these state trajectories can be maintained in an efficient manner.

7.1 Delayed-State Information Filter

The Information Filter (IF), which corresponds to the dual implementation of the Kalman Filter (KF), is a suitable approach for decentralised multi-robot estimation. Whereas the KF represents the distribution using its first $\boldsymbol{\mu}$ and second $\boldsymbol{\Sigma}$ order moments, the IF employs the so-called *canonical representation*. The fundamental elements are the *information vector* $\boldsymbol{\xi} = \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}$ and the *information matrix* $\boldsymbol{\Omega} = \boldsymbol{\Sigma}^{-1}$. Prediction and updating equations for the (standard) IF can also be derived from the standard KF [6]. In the case of non-linear prediction or measurement models, first order linearisation leads to the Extended Information Filter (EIF). Even though the prediction stage becomes more complex for the IF, the update stage is simpler. Hence, the use of the IF for multi-robot applications is justified by the additive nature of its updating steps. This simplifies a lot the filter when there is a single prediction step but multiple updates for the different data sources.

Formally, in a Delayed-State Information Filter, the belief over the full trajectory of the state up to the current time step t , is denoted by $\boldsymbol{\Omega}^t$ and $\boldsymbol{\xi}^t$. Thus, delayed states are also considered instead of just estimating the current state \mathbf{x}_t . Let us consider the system:

$$\mathbf{x}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \boldsymbol{\nu}_t \quad (2)$$

$$\mathbf{z}_t = \mathbf{g}_t(\mathbf{x}_t) + \boldsymbol{\varepsilon}_t \quad (3)$$

where \mathbf{x}_t is the person's position and velocity at time t , \mathbf{z}_t represents the estimations obtained by the camera network, robots or WSN at time t , \mathbf{A}_t is the prediction model, \mathbf{g}_t the corresponding measurement model, and $\boldsymbol{\nu}_t$ and $\boldsymbol{\varepsilon}_t$ are Gaussian noises. Knowing the information matrix and vector for the person trajectory up to time $t-1$, $\boldsymbol{\Omega}^{t-1}$ and $\boldsymbol{\xi}^{t-1}$, the estimation of this trajectory can be updated up to time t incorporating the local measurements. This is done according to Algorithm 2 [6], where $\mathbf{M}_t = \nabla \mathbf{g}_t(\bar{\boldsymbol{\mu}}_t)$, \mathbf{R}_t is the covariance of the additive noise for the prediction model (2) and \mathbf{S}_t is the covariance matrix of the noise in the measurement (3). **Add_M** adds a block row and a block column of zeros to

Algorithm 2 $(\xi^t, \Omega^t) \leftarrow$ Delayed State Information Filter $(\xi^{t-1}, \Omega^{t-1}, \mathbf{z}_t)$

- 1: $\bar{\Omega}^t = \mathbf{Add_M}(\Omega^{t-1}) + \begin{pmatrix} \mathbf{I} & \mathbf{0}^T \\ -\mathbf{A}_t^T & \mathbf{R}_t^{-1}(\mathbf{I} - \mathbf{A}_t) \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$
 - 2: $\bar{\xi}^t = \mathbf{Add_Row}(\xi^{t-1})$
 - 3: $\Omega^t = \bar{\Omega}^t + \begin{pmatrix} \mathbf{M}_t^T \mathbf{S}_t^{-1} \mathbf{M}_t & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$
 - 4: $\xi^t = \bar{\xi}_t + \begin{pmatrix} \mathbf{M}_t^T \mathbf{S}_t^{-1}(\mathbf{z}_t - \mathbf{g}_t(\bar{\mu}_t) + \mathbf{M}_t \bar{\mu}_t) \\ \mathbf{0} \end{pmatrix}$
-

the previous information matrix and **Add_Row** adds a block row of zeros to the previous information vector. For the movement of the person, a linear model $(\mathbf{A}_t, \mathbf{R}_t)$ is used like in [6]. Here, an initial estimation of the person position is also assumed in order to initialise the filter.

In general, for a state trajectory of n steps, the storage required would be $O(n^2)$. However, in the IF filter proposed, when the trajectory grows the matrix structure is block tridiagonal and symmetric at any time, what leads to a storage $O(n)$. Moreover, the computational complexity of the algorithm itself is $O(1)$, as the prediction and updating computations at each time instant only involve the previous block. These advantages allow the proposed approach to deal with delayed states more efficiently than a normal KF would do.

If several measurements from different sensors \mathbf{z}_t^i arrive at time t , and these measurements are conditionally independent given the state \mathbf{x}_t , with measurement functions \mathbf{g}_t^i and noise \mathbf{S}_t^i , then step 3 of Algorithm 2 becomes (equivalently for step 4):

$$\Omega^t = \bar{\Omega}^t + \sum_i \begin{pmatrix} \mathbf{M}_t^{iT} (\mathbf{S}_t^i)^{-1} \mathbf{M}_t^i & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \quad (4)$$

$$\xi^t = \bar{\xi}^t + \sum_i \begin{pmatrix} \mathbf{M}_t^{iT} (\mathbf{S}_t^i)^{-1} (\mathbf{z}_t^i - \mathbf{g}_t^i(\bar{\mu}_t) + \mathbf{M}_t^i \bar{\mu}_t) \\ \mathbf{0} \end{pmatrix} \quad (5)$$

Moreover, maintaining delayed states allows the filter to incorporation of delayed and asequent data, by adding their contribution to the corresponding elements of the information vector and matrix of the state trajectory.

7.1.1 Measurement functions

Regarding the measurements functions $\mathbf{g}_t(\mathbf{x}_t)$ considered in the system, the following considerations must be noted:

- the camera network, as described in Section 4, obtains measurements on the image plane. The position of the tracked objects can be transformed into the world coordinate system through a set of homographies that are obtained beforehand through a calibration process (although these homographies are not used by the camera trackers themselves).
- the robot obtains observations on the image plane, and then $\mathbf{g}_t(\mathbf{x}_t)$ is in this case the camera pin-hole model.
- the WSN provides 3D estimations on the position of the person in the world coordinate system (see Section 6).

7.2 Decentralised Information Filter

The main interest of the IF is that it can be easily decentralised. In a decentralised approach, each node i of the network employs only its local data \mathbf{z}_t^i to obtain a local estimation of the person trajectory (given by $\boldsymbol{\xi}^{i,t}$ and $\boldsymbol{\Omega}^{i,t}$) and then shares its belief with its neighbours. The received information $\boldsymbol{\xi}^{j,t}$ and $\boldsymbol{\Omega}^{j,t}$ from other nodes is locally fused in order to improve the local perception of the world. Ideally, the decentralised fusion rule should produce the same result locally as that obtained by a central node employing a centralised filter.

If a node i runs Algorithm 2 using only its local information, the one step equations are:

$$\boldsymbol{\Omega}^{i,t} = \underbrace{\begin{pmatrix} \mathbf{0} & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{0} & \boldsymbol{\Omega}^{i,t-1} & \\ \mathbf{0} & & \end{pmatrix}}_{\text{prior}} + \underbrace{\begin{pmatrix} \mathbf{R}_t^{-1} & -\mathbf{R}_t^{-1} \mathbf{A}_t & \mathbf{0}^T \\ -\mathbf{A}_t^T \mathbf{R}_t^{-1} & \mathbf{A}_t^T \mathbf{R}_t^{-1} \mathbf{A}_t & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}}_{\text{prediction}} + \underbrace{\begin{pmatrix} \mathbf{M}_t^{iT} (\mathbf{S}_t^i)^{-1} \mathbf{M}_t^i & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}}_{\text{update}} \quad (6)$$

$$\boldsymbol{\xi}^{i,t} = \underbrace{\begin{pmatrix} \mathbf{0} \\ \boldsymbol{\xi}^{i,t-1} \end{pmatrix}}_{\text{prior}} + \underbrace{\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}}_{\text{prediction}} + \underbrace{\begin{pmatrix} \mathbf{M}_t^{iT} (\mathbf{S}_t^i)^{-1} (\mathbf{z}_t^i - \mathbf{g}_t^i(\bar{\boldsymbol{\mu}}_t) + \mathbf{M}_t^i \bar{\boldsymbol{\mu}}_t) \\ \mathbf{0} \end{pmatrix}}_{\text{update}} \quad (7)$$

When this node i is within communication range of another node j , they can share their beliefs, represented by their information vectors $\boldsymbol{\xi}^{i,t}$ and $\boldsymbol{\xi}^{j,t}$, and matrices $\boldsymbol{\Omega}^{i,t}$ and $\boldsymbol{\Omega}^{j,t}$. Assume, without loss of generality, that nodes i and j have common priors $\boldsymbol{\Omega}^{j,t-1} = \boldsymbol{\Omega}^{i,t-1}$ and $\boldsymbol{\xi}^{i,t-1} = \boldsymbol{\xi}^{j,t-1}$ (for instance, due to previous communications). This node j will apply the same equations (6) and (7), but with its own measurement \mathbf{z}_t^j and covariance \mathbf{S}_t^j . Then, it can be seen that the the next fusion rule, proposed by the authors [6]:

$$\boldsymbol{\Omega}^{i,t} \leftarrow \boldsymbol{\Omega}^{i,t} + \boldsymbol{\Omega}^{j,t} - \boldsymbol{\Omega}^{ij,t} \quad (8)$$

$$\boldsymbol{\xi}^{i,t} \leftarrow \boldsymbol{\xi}^{i,t} + \boldsymbol{\xi}^{j,t} - \boldsymbol{\xi}^{ij,t} \quad (9)$$

where

$$\boldsymbol{\Omega}^{ij,t} = \underbrace{\begin{pmatrix} \mathbf{0} & \mathbf{0}^T & \mathbf{0}^T \\ \mathbf{0} & \boldsymbol{\Omega}^{i,t-1} & \\ \mathbf{0} & & \end{pmatrix}}_{\text{prior}} + \underbrace{\begin{pmatrix} \mathbf{R}_t^{-1} & -\mathbf{R}_t^{-1} \mathbf{A}_t & \mathbf{0}^T \\ -\mathbf{A}_t^T \mathbf{R}_t^{-1} & \mathbf{A}_t^T \mathbf{R}_t^{-1} \mathbf{A}_t & \mathbf{0}^T \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}}_{\text{prediction}} \quad (10)$$

$$\boldsymbol{\xi}^{ij,t} = \underbrace{\begin{pmatrix} \mathbf{0} \\ \boldsymbol{\xi}^{i,t-1} \end{pmatrix}}_{\text{prior}} + \underbrace{\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}}_{\text{prediction}} \quad (11)$$

allows node i to recover locally the same estimation as the one a central entity receiving the information from i and j would obtain. The equations mean that each node must sum up the information received from other nodes. The additional terms $\boldsymbol{\Omega}^{ij,t}$ and $\boldsymbol{\xi}^{ij,t}$ represent the common information between the nodes. This

common information is due to previous communications between nodes (previous priors) and common prediction information, and should be removed to avoid double counting of information (i.e. rumour propagation). After eliminating this common information, the remaining information from j is due to its local data (or data gathered from other nodes connected to j).

These fusion equations can be applied also to the general case in which the nodes exchange information at arbitrary instants. The bottom line is that each node computes part of the update that a central node would compute (which is a sum for all the information received, see (4) and (5)). The fusion rule allows a fusing node to recover the centralised solution if the common information is properly removed. As long as a tree-shaped logical topology in the perception system (no cycles or duplicated paths of information) is assumed, this common information can be maintained by a separated EIF, a so-called channel filter [31].

It is important to remark that, using these fusion equations and considering trajectories (delayed states), the local filter can obtain an estimation that is equal to that obtained by a centralised system [6], which is not the case when using just the last state [36, 42] (unless all measurements arrive in order and no measurement is lost). Another advantage of using delayed states is that the belief states can be fused asynchronously without missing information. Each node in the UNR system can accumulate evidence, and send it whenever it is possible. Also, as commented before, asequent and delayed measurements can be incorporated in the filter. However, as the state grows over time, the size of the message needed to communicate its belief also does. For the normal operation of the system, only the state trajectory over a time interval is needed, so these belief trajectories can be bounded by marginalizing out old states (which is a cheap operation due to the block diagonal nature of the information matrix). Note that the trajectories should be longer than the maximum expected delay in the network in order not to miss any measurement information.

Finally, when no assumptions about the network topology can be made (e.g. due to the existence of mobile robots, possible losses of communication links, etc.), another option to remove the common information is to employ a conservative fusion rule, which ensures that the system does not become overconfident even in presence of duplicated information. For the case of the IF, there is an analytic solution for this, given by the Covariance Intersection algorithm of [24]. Therefore, the conservative rule to combine the local belief of a robot i with that received from another robot j is given by:

$$\mathbf{\Omega}^{i,t} \leftarrow \omega \mathbf{\Omega}^{i,t} + (1 - \omega) \mathbf{\Omega}^{j,t} \quad (12)$$

$$\boldsymbol{\xi}^{i,t} \leftarrow \omega \boldsymbol{\xi}^{i,t} + (1 - \omega) \boldsymbol{\xi}^{j,t} \quad (13)$$

for $\omega \in [0, 1]$. It can be demonstrated that the estimation is consistent (in the sense that no overconfident estimations are done) for any ω . The value of ω can be selected following some criteria, such as maximizing the obtained determinant of $\mathbf{\Omega}^{i,t}$ (minimizing the entropy of the final distribution). Here, the option chosen is to use ω as a fixed weight that determines the system confidence in its own estimation and the neighbour's ones.

Although employing the CI formula avoids the need to maintain an estimation of the common information transmitted to the neighbour systems, as these

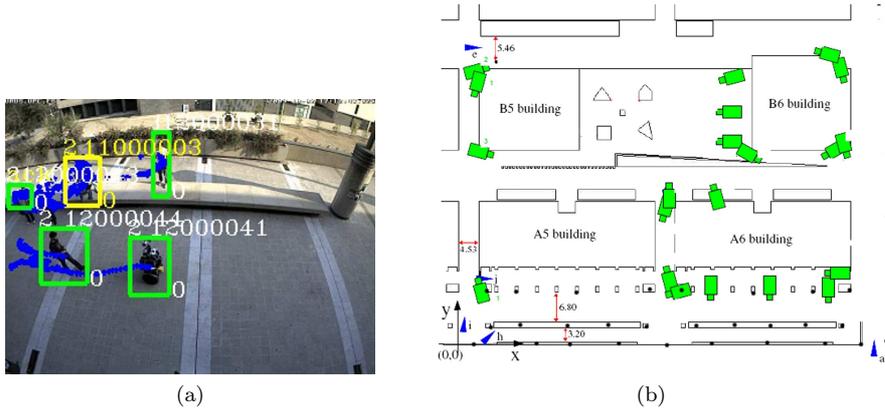


Fig. 14: (a) One of the robots (on the image marked as 12000041) guiding a person (12000044). (b) Scenario. The dimension is 100 by 100 meters, approximately. Cameras in green and Mica2 nodes as black dots.

fusion rules are conservative, some information is lost with respect to the purely centralised case.

7.3 Data association

Each fusion node of the system should be able to associate its local observations to the current tracks. In the case of the camera network, this is done by combining the inter camera information and geometric information. As commented in Section 4, the system is able to handle inter camera tracking without calibration, using as weak cues reappearance probabilities and colour information. Therefore, the system uses this information for data association. As this scheme may fail, the non-associated observations are also passed through a data association procedure based on the Mahalanobis distance, using the estimated 3D position obtained from the homographies.

The data association in the case of the WSN node is straightforward, as the messages from the WSN are tagged with an ID. The image tracker in the case of the robot maintains the identity of the tracked persons while they are on the image plane. The Mahalanobis distance is also used to associate new measurements with previous tracks.

Moreover, the decentralised nodes should be able to associate the received tracks to the local tracks. For this track-to-track fusion, the Mahalanobis distance is used again.

8 Experimental Results

The techniques described above have been tested during the experimental sessions of the URUS European Project. The experiments were carried out outdoors at the

Seq	Time of Day	Length	Weather	Num of People
1	11:00	80mins	Sunny	1200
2	12:00	180mins	Cloudy	3750
3	16:00	200mins	Sunny	1750

Table 1: Details of the fixed camera test sequences

UPC Campus North site, in Barcelona, Spain. The experimental setup consisted of twenty two fixed colour cameras with mostly non-overlapping fields of view. Moreover, a network of 30 Mica2 nodes was deployed in the campus. A set of different robots were involved in the experiments. Figure 14 shows a moment of these experiments and the final deployment of sensors. Usually, the camera network is used as the initial point for the person guiding task. First, some partial results will be described concerning the camera network. Then, the results obtained with the full system will be shown.

8.1 Fixed camera sensor results

A series of partial experiments concerning just tracking of people using the fixed camera sensors were performed. The experiments were conducted in October, with cloudy and sunny weather. The cameras are wall mounted around 5 to 10 meters high. There are large gaps between many of the cameras of up to 30 seconds. The twenty two time synchronised IP video feeds are fed into three quad-core PC with real-time person tracking. There is no calibration of the camera environment with no *a priori* information provided. Over time, additional information is incorporated into the system to learn links between regions on the cameras and improve tracking accuracy. The experimental data was accumulated from 9am for 5 days tracking a total of around 140,000 people in total or around 1,200 people per camera per day. Often the same person was tracked on multiple cameras as they moved around the site, therefore the number of unique people was around 500 per day per camera or 55,000 unique people over the complete system. Figure 15 shows resultant temporal likelihoods for a number of inter camera links at a single subdivision level. The evaluation of the tracking was performed using three unique sequences taken at different times of day on different days.

The black vertical line indicates a reappearance of zero seconds. It can be seen that there are strong links between cameras 3 and 4 and between 3 and 5, whereas there are no visible links between 3 and 6 and between 3 and 14. This is due to the increased distance and people will rarely reappear on cameras 6 and 14 after they were tracked on camera 3. Table 1 shows the details of the three test sequences used to evaluate the approach. Table 2 shows the results of tracking people inter camera on the three sequences. A subdivisions of 0 means using no region link cues, only basic colour correlation to match and track people inter camera. A subdivision of 1 is a single region per camera, i.e. the initial camera to camera linking, whereas a subdivision of 2 is where any suitable single camera regions are subdivided into 4 equally sized news regions as described in Section 4.1. For the test sequence, all people that occurred on multiple cameras were ground-truthed and a true positive occurred when a person was assigned the same ID that they were assigned in a previous region. A false positive indicates when a person was

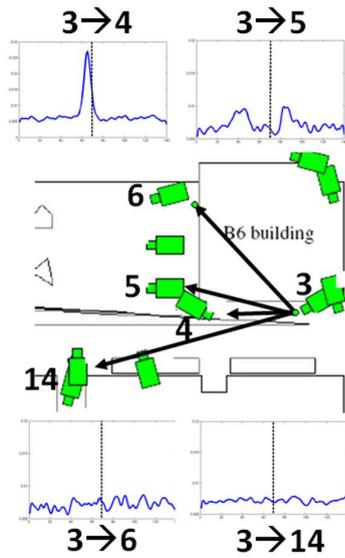


Fig. 15: Inter camera temporal likelihoods

Seq	Subdivisions	True Positive	False Positive	False Negative
1	0	5%	35%	60%
	1	50%	26%	24%
	2	52%	24%	23%
2	0	5%	46%	49%
	1	47%	22%	31%
	2	52%	24%	23%
3	0	3%	17%	80%
	1	58%	21%	21%
	2	62%	21%	17%

Table 2: Fixed camera tracking of people

assigned an incorrect ID, and a false negative is when a person who has moved inter camera is given a new ID instead of the ID from their previous region.

The column for 0 subdivisions indicates performance without learning the temporal and colour relationships between regions. It is generally poor because of the large colour variations inter camera as well as shadow and lighting changes. The subdivision level 1 performs far better, with the additional detail of 2 subdivisions providing a further improvement. The reason for the greater performance on sequence 3 is due to the time of day of the experiment. Being later in the day there was less simultaneous traffic on the system, this meant there were less possible correlation options for people tracked cross camera. Figure 16 gives example frames of tracking inter camera for two separate people. Figure 17 shows the estimated position of the person using only information from the camera network (the cameras are homography calibrated, although this is not used by the intra and inter camera tracking algorithms).



Fig. 16: Cross camera tracking (a) Person 11000001 on camera 11, (b) Person 11000001 correctly identified on camera 12 (c) Person 13000027 on camera 13 (d) Person 13000027 correctly identified on camera 12.

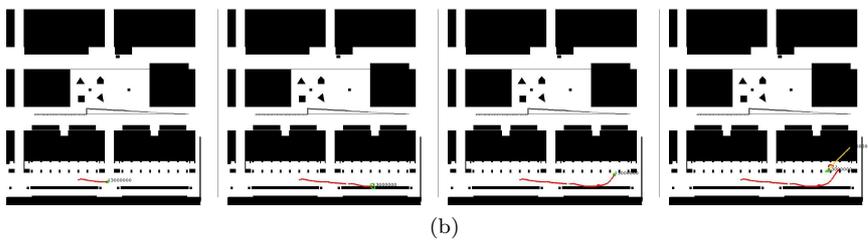
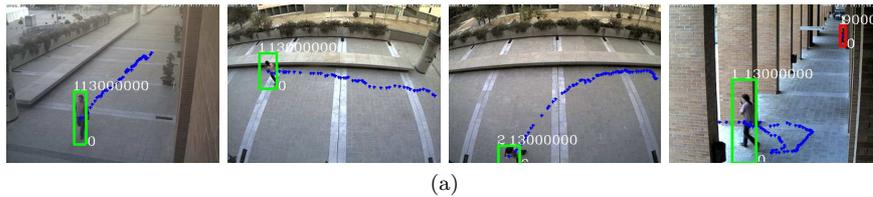


Fig. 17: (a) Tracks on the image plane of 4 different cameras. The identity is correctly handed over the cameras using the weak cues described in Section 4. (b) Estimated position of the person on the UPC campus.

8.2 Robot and WSN

In order to illustrate the benefits from the data fusion process, a first setup is presented here. This setup considers information from one camera on board the robot Romeo (4-wheel vehicle) and the WSN. The objective was to track one person cooperatively. In this case, just two nodes of the decentralised fusion are used: one on board the robot and one for the WSN. These nodes locally integrate information from a monocular camera (see Fig. 18) and from the signal strength-based estimations (Section 6, see Fig. 18a), respectively.

Figure 19 shows the X and Y estimations obtained by the robot alone and when the robot combines its information with the one provided by the WSN. For these outdoor and urban experiments, obtaining the real position of the person was not easy, since the tests were run on-line and there were some areas without GPS coverage. Therefore, the person moved together with the robot, and the robot position was used as ground truth to check the estimation. The trajectory of the robot was measured accurately by its navigation software (laser, map knowledge,

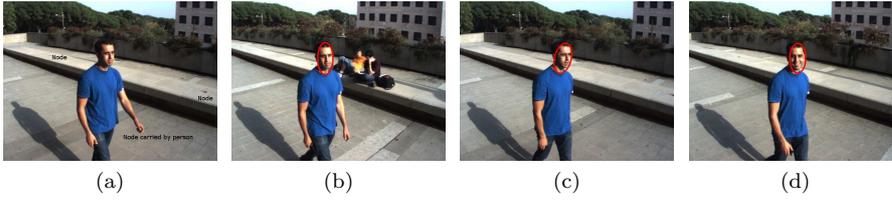


Fig. 18: (a) The person is carrying a Mica2 node during the experiment. (b,c,d) The robot is able to obtain local observations on the image plane of the face of the person.

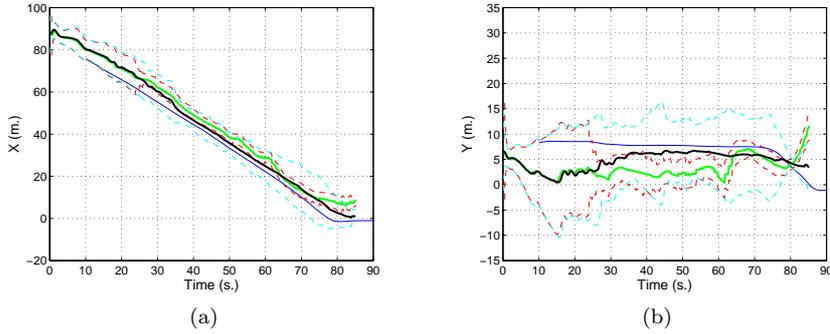


Fig. 19: Tracking using one on-board camera and the WSN. Black: robot alone. Green: robot and WSN. Dashed lines are the sigma intervals and the blue solid line represents the robot trajectory.

GPS, etc.). The person is following behind the robot (see Fig. 18) (which in this trajectory means that the X coordinates of the person are larger than that of the robot) and some meters beside the robot (a lower Y coordinate).

Regarding the accuracy of the method, it can be seen that the error in the estimation (which is determined by the standard deviation) is enough for tracking people in an area of around 2,500 square meters. Moreover, the key point is that the fusion system improves the accuracy of independent sensors. Thus, it can be seen how the introduction of the WSN reduces the uncertainty; as we have a monocular camera, the uncertainty on the person position is quite big in both axes when the robot is alone. In this case, the initial position of the person is computed assuming a known height of the face.

8.3 UNR experiments with decentralised fusion

In this setup, one robot, the WSN of 30 nodes and 7 IP cameras are used. In the experiment, one person was following the robot, which was manually controlled. The setup of the perception system is a decentralised node on the robot, one for

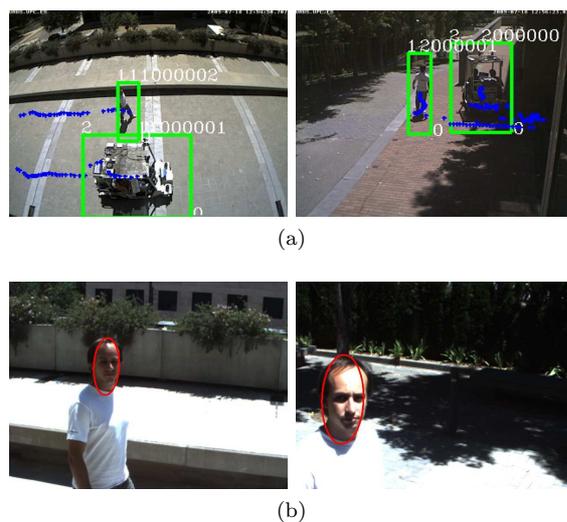


Fig. 20: (a) Tracks obtained by the camera network. (b) Tracks obtained by the camera on-board Romeo.

the WSN and 2 for the IP cameras, one integrating measurements from 3 cameras and the other from 4 cameras.

Figure 20 shows some examples of the tracks obtained by the on-board camera and the camera network. Along the trajectory, of more than 350 meters, there were gaps in the camera coverage. Moreover, the robot lost track several times due to the changes in illumination. Finally, the WSN coverage was limited to a certain part of the campus.

Figure 21a shows the estimated position of the person with the full system running. The total length of the experiment is around 350 meters and 5 minutes. The person is usually besides the robot (which means that the X or Y coordinates are the same) and the robot position is used as ground truth again. The system is able to maintain the estimation on the person position for the full trajectory. There are WSN coverage between 0 and 150 seconds, approximately. Figure 21b shows an interval of the trajectory. In this part, only WSN and robot information is available. Sensor fusion is beneficial because, although the WSN measurements have lower precision, they bound the error from the monocular camera. At 75 seconds, the person enters under coverage of the camera network. This leads to a big reduction in uncertainty.

8.3.1 Other issues

The same kind of experiment was repeated several times. The communication between the fusion nodes on board the robots and the fusion nodes related to the camera network and the WSN was done using WiFi and 3G. Software running on the robot was able to measure the quality of the WiFi link, and to switch to 3G whenever this quality dropped below a certain threshold.

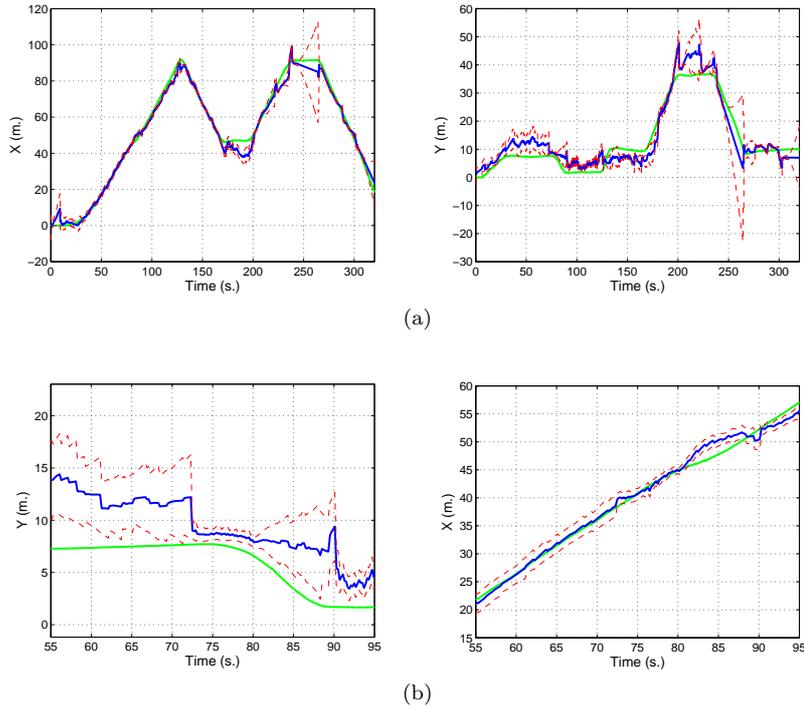


Fig. 21: Estimated position of the person (blue) compared to the position of the robot (green). Dashed lines represent the standard deviation of the estimation. (a) Complete trajectory. (b) A section of the trajectory. The person is following the robot with the same X coordinate up to time 80 seconds. Then the robot changes orientation. The person is separated from the robot around 3-4 meters.

The switching between communication networks created from time to time gaps of several seconds. Moreover, although 3G had more stable coverage in the scenario, it had lower bandwidth and higher latencies. The use of decentralised nodes allowed the system to cope with communication gaps, as in the meantime, the local nodes were accumulating information. When the communication links were recovered, the nodes exchanged their estimations. Moreover, as delayed states were considered, this delayed information (and also information delayed due to the latencies) could be fused in a correct way, and no information was lost.

Figure 22 compares the decentralised estimation with a centralised off-line implementation. In this experiment, 4 different cameras and the WSN were running. For the centralised implementation, all the information is received and fused in a single node (and no information is missed; it can be considered as a gold standard estimation); whereas for the decentralised case three fusion nodes were used (two of them processing locally information from 2 cameras; the remaining one processing locally information from the WSN). The estimation obtained by one of the fusion nodes is shown in Fig. 22. It can be seen how, with some latencies depending on the conditions, the decentralised node obtains an estimation quite

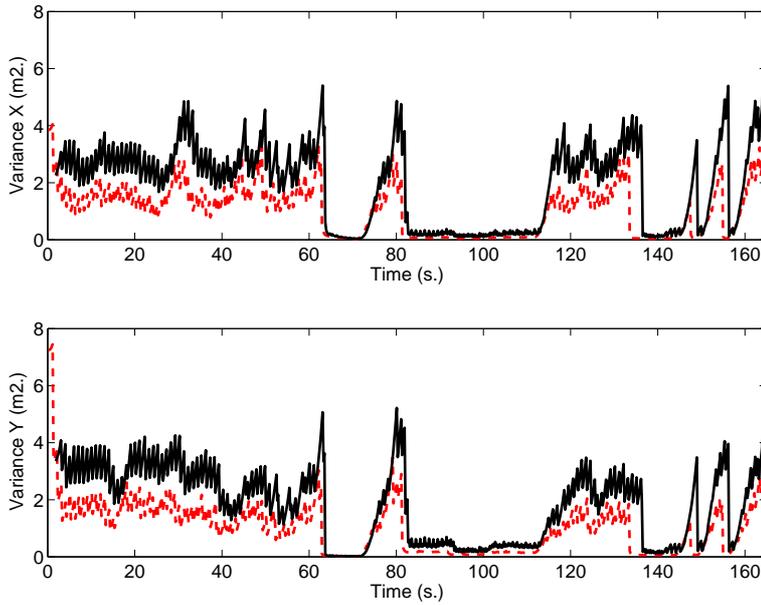


Fig. 22: Estimated variance by a central node receiving all the information (dashed and red) compared to the estimation in a decentralised node (solid and black).

close to the centralised one, even though some information is lost due to the conservative fusion rule employed in this case (Covariance Intersection, see Section 7. Moreover, the decentralised estimation is consistent, in the sense that it never accumulates more information than the one obtained in the ideal centralised filter.

9 Conclusions

The combination of robots and ambient intelligence (like embedded sensors and camera networks) seems a clear trend in the near future. This paper has presented a decentralised system that aims to use multiple sensors to accurately track people within a surveillance context. The system makes extensive use of data fusion procedures to incorporate all the information available.

The algorithms are real-time and operate on realistic outdoor environments. The fixed camera tracking provides high accuracy by incrementally learning the colour and temporal relationships between regions on non-overlapping cameras. Moreover, the signal strength of mobile devices is employed to estimate the position of the person using particle filtering. The combination of all this information obtained by robots allows for accurate person tracking in more challenging situations. The system has been tested in an urban scenario, considering a camera network of 20 cameras, a WSN of 30 nodes and robots.

Very complex algorithms employing just one source of information are usually unable to cope with all the potential situations in these scenarios, affected by changes in illumination, clutter, and wide area coverage. The combination of complementary systems can be useful for this problem. Signal-based localization is less accurate than camera-based localization, but it is less affected by occlusions. Robots have usually narrower fields of view, but they can adapt to cover places occluded from the camera networks. However, camera networks can provide overall information on the scene though. Scalability is an issue in these systems, and thus decentralised algorithms are required. The system presented is a mixture between distributed or centralised subsystems that are linked through a decentralised data fusion scheme. The addition of new robots or sub-nets of cameras does not affect the rest of the perception system in terms of storage, as only local communication and local processing are used.

Future developments include the integration of active sensing behaviours in the system. The WSN can be actively controlled to save energy, activating those nodes more useful for tracking. The robots can also move maximizing the possibility of maintaining the person in the field of view. In some applications, the robots can use paths that are more informative as they are more covered by static cameras. Entropy-based information gain algorithms and Partially Observable Markov Decision Processes will be considered for these tasks.

Acknowledgements

This work is partially supported by URUS, Ubiquitous networking Robotics in Urban Settings, funded by the European Commission (EC) under FP6 with contract number FP6-EU-IST-045062. In addition the authors would like to thank to the rest of the partners of the URUS project for their help and support. Luis Merino is also funded by the EC through the project FROG (FP7-288235). Jesus Capitan is also funded by Fundação para a Ciência e a Tecnologia (ISR/IST pluriannual funding) through the PIDDAC Program funds and projects PEst-OE/EEI/LA0009/2011 and CMU-PT/SIA/0023/2009.

References

1. T. Bailey and H. Durrant-Whyte. Decentralised data fusion with delayed states for consistent inference in mobile ad hoc networks. Technical report, Australian Centre for Field Robotics, University of Sydney, 2007.
2. G.R. Bradski. Computer Vision Face Tracking as a Component of Perceptual User Interface. In *Proc. of Workshop on Applications of Computer Vision*, pages 214–219, 1998.
3. F. Caballero, L. Merino, P. Gil, I. Maza, and A. Ollero. A probabilistic framework for entire wsn localization using a mobile robot. *Journal of Robotics and Autonomous Systems*, 56(10):798–806, 2008.
4. Q. Cai and J.F. Aggarwal. Automatic Tracking of Human Motion in Indoor Scenes across Multiple Synchronized Sideo Streams. In *Proc. of IEEE International Conference on Computer Vision (ICCV'98)*, 1998.
5. J.P. Cánovas, K. LeBlanc, and A. Saffiotti. Robust multi-robot object localization using fuzzy logic. In *Proc. of the International Robocup Symposium*, 2004.
6. J. Capitán, L. Merino, F. Caballero, and A. Ollero. Delayed-State Information Filter for Cooperative Decentralized Tracking. In *Proceedings of the International Conference on Robotics and Automation, ICRA*, 2009.

7. T.H. Chang, S. Gong, and E. Ong. Tracking Multiple People under Occlusion using Multiple Cameras. In *Proc. of BMVA British Machine Vision Conference (BMVC'00)*, pages 566–575, 2000.
8. A. Dick and M. Brooks. A Stochastic Approach to Tracking Objects Across Multiple Cameras. In *Proc. of Australian Conference on Artificial Intelligence*, pages 160–170, 2004.
9. T.J. Ellis, D. Makris, and J.K. Black. Learning a Multi-Camera Topology. In *Proc. of Joint IEEE Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance (VS-PETS)*, pages 165–171, 2003.
10. A. Ess, B. Leibe, K. Schindler, and L. Van Gool. A Mobile Vision System for Robust Multi-Person Tracking. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2008.
11. M. Farenzena, L. Bazzani, M. Perina, A. Cristani, and V. Murino. Person re-identification by symmetry-driven accumulation of local features. In *Proc. of IEEE International Conference on Computer Vision and Pattern Recognition (CVPR'10)*, 2010.
12. B. Ferris, D. Hohnel, and D. Fox. Gaussian processes for signal strength-based location estimation. In *In Proc. of Robotics Science and Systems*, 2006.
13. S. Fintrop, A. Königs, F. Hoeller, and D. Schulz. Visual Person Tracking Using a Cognitive Observation Model. In *International Conference on Robotics and Automation (ICRA), Workshop in People Detection and Tracking*, 2009.
14. A. Gilbert and R. Bowden. Incremental, Scalable Tracking of Objects Inter Camera. In *Computer Vision and Image Understanding (CVIU)*, 3:43–58, 2008.
15. A. Gilbert, J. Capitán, R. Bowden, and L. Merino. Accurate Fusion of Robot, Camera and Wireless Sensors for Surveillance Applications. In *In Proc. Ninth IEEE International Workshop on Visual Surveillance (ICCV09)*, Kyoto, Japan, 2009.
16. A. Gilbert, J. Illingworth, and R. Bowden. Scale Invariant Action Recognition Using Compound Features Mined from Dense Spatio-Temporal Corners. In *Proc. of European Conference on Computer Vision (ECCV'08)*, 1:222–233, 2008.
17. D. Gray and H. Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *Proc. of European Conference on Computer Vision (ECCV'08)*, pages 262–275, 2008.
18. S. Grime and H. F. Durrant-Whyte. Data fusion in decentralized sensor networks. *Control Engineering Practice*, 2(5):849–863, Oct. 1994.
19. B. Grocholsky, A. Makarenko, T. Kaupp, and H. F. Durrant-Whyte. *Lecture notes in Computer Science*, volume 2634, chapter Scalable Control of Decentralised Sensor Platforms. Springer, 2003.
20. F. Gustafsson and F. Gunnarsson. Mobile Positioning using Wireless Networks. *IEEE Signal Processing Magazine*, pages 41–53, 2005.
21. G. Hollinger, J. Djughash, and S. Singh. Tracking a moving target in cluttered environments with ranging radios. In *IEEE International Conference on Robotics and Automation*, 2008.
22. T. Huang and S. Russell. Object Identification in a Bayesian Context. In *Proc. of International Joint Conference on Artificial Intelligence (IJCAI-97)*, pages 1276–1283, 1997.
23. O. Javed, K. Shafique, Z. Rasheed, and M. Shah. Modeling inter-camera space-time and appearance relationships for tracking across non-overlapping views. *Comput. Vis. Image Underst.*, 109(2):146–162, February 2008.
24. S.J. Julier and J.K. Uhlmann. A non-divergent estimation algorithm in the presence of unknown correlations. In *Proc. of the American Control Conference*, volume 4, pages 2369–2373, 1997.
25. P. KaewTrakulPong and R. Bowden. A Real-time Adaptive Visual Surveillance System for Tracking Low Resolution Colour Targets in Dynamically Changing Scenes. In *Journal of Image and Vision Computing*, 21(10):913–929, 2003.
26. P. KaewTrakulPong and R. Bowden. Towards automated wide area visual surveillance: Tracking objects between spatially separated, uncalibrated views. In *IEE Proc. Vision, Image and Signal Processing*, 152(10):213–224, 2005.
27. P. Kelly, A. Katkere, D. Kuramura, S. Moezzi, and S. Chatterjee. An Architecture for Multiple Perspective Interactive Video. In *Proc. of the 3rd ACE International Conference on Multimedia*, pages 201–212, 1995.
28. V. Kettner and R. Zabih. Bayesian Multi-Camera Surveillance. In *Proc. of International Conference on Computer Vision and Pattern Recognition*, pages 253–259, 1999.
29. B. Kusý, A. Ledeczi, and X. Koutsoukos. Tracking mobile nodes using RF Doppler shifts. In *Proceedings of SenSys*, pages 29–42, 2007.

30. B. Kusý, J. Sallai, G. Balogh, A. Ledeczki, V. Protopopescu, J. Tolliver, F. DeNap, and M. Parang. Radio interferometric tracking mobile wireless nodes. In *Proceedings of MobySys*, pages 139–151, 2007.
31. A. Makarenko, A. Brooks, S. Williams, H. Durrant-Whyte, and B. Grocholsky. A decentralized architecture for active sensor networks. In *Proceedings IEEE International Conference on Robotics and Automation, ICRA*, volume 2, pages 1097–1102, 2004.
32. O. Martinez-Mozos, R. Kurazume, and T. Hasegawa. Multi-part people detection using 2D range data. *International Journal of Social Robotics*, 2010.
33. V.I Morariu and O.I Camps. Modeling Correspondences for Multi-Camera Tracking using Nonlinear Manifold Learning and Target Dynamics. In *Proc. of IEEE International Conference on Computer Vision and Pattern Recognition (CVPR'06)*, I:545–552, 2006.
34. C. Morelli, M. Nicoli, V. Rampa, U. Spagnolini, and C. Alippi. Particle filters for rssi-based localization in wireless sensor networks: An experimental study. In *International Conference on Acoustics, Speech and Signal Processing*, volume 4, page IV, 2006.
35. L. E. Navarro-Serment, C. Mertz, and M. Hebert. Pedestrian detection and tracking using three-dimensional ladar data. In *Proc. of The 7th Int. Conf. on Field and Service Robotics*, July 2009.
36. E. Nettleton, H. Durrant-Whyte, and S. Sukkarieh. A robust architecture for decentralised data fusion. In *Proc. of the International Conference on Advanced Robotics (ICAR)*, 2003.
37. P.J. Norlund, F. Gustafsson, and F. Gunnarsson. Particle Filters for Positioning in Wireless Networks. In *Proceedings of EUSIPCO*, 2002.
38. R. Olfati-Saber, J. Fax, and R. Murray. Consensus and Cooperation in Networked Multi-Agent Systems. *Proceedings of the IEEE*, 95(1):215–233, 2007.
39. B. Prosser, S. Gong, and T. Xiang. Multi-camera Matching using Bi-Directional Cumulative Brightness Transfer Functions. In *BMVC'08*, pages –1–1, 2008.
40. V. Ramadurai and M. L. Sichitiu. Localization in wireless sensor networks: A probabilistic approach. In *Proceedings of the 2003 International Conference on Wireless Networks (ICWN 2003)*, pages 275–281, Las Vegas, NV, June 2003.
41. W. Ren, R. Beard, and E. Atkins. Information consensus in multivehicle cooperative control. *IEEE Control Systems*, 27(2):71–82, 2007.
42. M. Rosencrantz, G. Gordon, and S. Thrun. Decentralized sensor fusion with distributed particle filters. In *Proc. Conf. Uncertainty in Artificial Intelligence*, 2003.
43. A. Sanfeliu, J. Andrade-Cetto, M. Barbosa, R. Bowden, J. Capitán, A. Corominas, A. Gilbert, J. Illingworth, L. Merino, J. M. Mirats, P. Moreno, A. Ollero, J. Sequeira, and M. T. J. Spaan. Decentralized sensor fusion for ubiquitous networking robotics in urban areas. *Sensors*, 10(3):2274–2314, 2010.
44. V. Seshadri, G.V. Zaruba, and M. Huber. A bayesian sampling approach to in-door localization of wireless devices using received signal strength indication. In *International Conference on Pervasive Computing and Communications*, pages 75 – 84, 2005.
45. L. D. Stone, T. L. Corwin, and C. A. Barlow. *Bayesian Multiple Target Tracking*. Artech House, Inc., Norwood, MA, USA, 1999.
46. S. Sukkarieh, E. Nettleton, J.-H. Kim, M. Ridley, A. Goktogan, and H. Durrant-Whyte. The ANSER Project: Data Fusion Across Multiple Uninhabited Air Vehicles. *The International Journal of Robotics Research*, 22(7-8):505–539, 2003.
47. J. K. Uhlmann. Covariance consistency methods for fault-tolerant distributed data fusion. *Information Fusion*, (4):201–215, 2003.
48. N. Ukita. Probabilistic-topological calibration of widely distributed camera networks. *Mach. Vision Appl.*, 18(3):249–260, May 2007.
49. P. Viola and M. Jones. Robust Real-Time Face Detection. *International Journal of Computer Vision*, 57:137–154, 2004.