

Tracking and Modeling of Human Activity Using Laser Rangefinders

Anand Panangadan · Maja Matarić ·
Gaurav S. Sukhatme

Accepted: 23 December 2009 / Published online: 6 January 2010
© Springer Science & Business Media BV 2009

Abstract We describe a system that uses laser rangefinders to track the positions of people in typical environments and then builds predictive models of the observed movement patterns and interactions between persons. We represent all human activity as detected by the laser rangefinder system as a probability distribution over the space of possible displacements. The assumption is that different activities will map to distinct probability distributions. Position tracks are first segmented and clustered into short sequences representing different activities. The sequence of activity clusters is then used to build a stochastic model of the observed movement patterns and the typical frequency of occurrence. Interactions are assumed to occur between persons whose corresponding probability distributions exhibit a high degree of similarity. We describe the performance of the system on data recorded from unscripted activities that occurred in different environments: open layout laboratory, corridor, and an outdoor courtyard. In the laboratory environment, the system was able to detect interactions between people (ping-pong players) without utilizing a pre-defined model of specific interactions. In the courtyard environment, the system was able to flag a sudden increase in the number of people

in the courtyard as an anomalous occurrence without any pre-defined concept of occupancy of the environment.

Keywords Laser rangefinder · Tracking · Interaction · Markov models · Activity recognition

1 Introduction

As robots begin to operate in environments where human beings perform everyday activities, it will be important for the robots to detect, model and predict the movement of humans. In particular, the ability to detect interactions between people will be useful in order that the robot can modify its own behavior as necessary. Laser rangefinders are currently used by robots mainly for obstacle avoidance, localization, and mapping. In this work, we present a system that uses laser rangefinders for tracking people in both indoor and outdoor environments and then uses this information to build predictive models of the observed movement patterns.

In our work, we have modeled all human activity detected by the laser rangefinder system as arising from a probability distribution over the space of possible movements over a short period of time. In particular, we have assumed that different activities arise from different probability distributions. We then use this assumption to segment long movement tracks into sequences representing distinct activities and build a stochastic model of the frequency of occurrence of these activities. We measure the difference between probability distributions using a measure of entropy known as the Jensen-Shannon divergence. We use a recursive segmenting algorithm that maximizes this divergence between the activity sub-sequences contained within a movement track. We demonstrate our laser rangefinder tracking system in different environments: an open layout laboratory, corridor, and outdoor courtyard.

A. Panangadan (✉)
Saban Research Institute, Childrens Hospital Los Angeles,
4650 W Sunset Blvd., Los Angeles, CA 90027, USA
e-mail: APanangadan@chla.usc.edu

M. Matarić · G.S. Sukhatme
Computer Science Department, University of Southern
California, 942 West 37th Place, Los Angeles, CA 90089, USA

M. Matarić
e-mail: mataric@usc.edu

G.S. Sukhatme
e-mail: gaurav@usc.edu

We are particularly interested in detecting and modeling the situations where two persons interact with each other even though they are separated by a significant and changing distance. In such cases, applying a threshold to the distance between every pair of persons will not be sufficient to determine if a pair is interacting. For example, people following each other along a corridor may not be interacting, while ping-pong players who may be distant from each other and whose relative positions keep changing are actively interacting. Most work on detecting interactions therefore relies on a pre-defined model of the interactions to be detected. In contrast, our hypothesis is that these interactions can be detected by comparing the similarity of the probability distributions of the displacements of the two persons. We apply an entropy-based metric (Kullback-Leibler divergence) to the probability distributions representing the observed activities in order to determine the amount of similarity between the activities. As the range of human interactions is very large, we do not expect this method to detect all types of interaction. In this paper, we evaluate the feasibility of this approach in selected environments.

In the courtyard environment, we focus on modeling the frequency of occurrence of the observed activities and interactions. We assume that activities are generated by a Poisson process and model the number of occurrences of each activity type as a Poisson distribution. This enables us to compute the probability of detecting a given number of activities per unit time and flagging anomalous observations based on low probability of observation. An earlier version of this work appears in [34, 35].

This paper is structured as follows. Section 2 reviews some of the related work in modeling activities and interactions of people. Section 3 describes our laser rangefinder-based tracking setup. Section 4 describes how the output of the tracker is segmented and clustered into distinct activity segments. Section 5 describes the methods used to model activity and interactions. Section 6 describes modeling frequency of occurrence of activities as Poisson distributions for anomaly detection. Section 7 contains the results of applying our techniques to detecting a subset of the activities and interactions in real-world environments. Section 8 concludes with a summary of our findings.

2 Related Work

The problem of detecting the presence of humans has traditionally been studied using cameras as sensors. Schiele et al. [37] gave an overview of some of the vision-based techniques that are used for people detection. In order to recognize longer movement sequences, the tracks of foreground objects in an image are first classified into activity primitives such as hand gestures, walking, and paths between known

landmarks [5, 16, 30]. The recognition of these low level actions is usually performed by defining a hidden Markov model (HMM) for each action [28, 32]. More complex activities are defined as sequences of low level activities, i.e., the low level activities may be considered as symbols of a higher level representation. Kohlmorgen and Lemm [20] describe how a hidden Markov model with dynamically varying number of states can be used for segmentation. Kulić et al. [21] also use a stochastic approach for segmentation of data from a motion capture setup. As in these works, we also consider the probability distribution of the data rather than the raw values.

Laser rangefinders have also been used for people tracking. Schulz et al. [38] adapted Bayesian filtering algorithms to track multiple people using laser rangefinders mounted on a robot. Bruce and Gordon [7] used particle filters for motion tracking but use a more sophisticated motion model than the one used in our system. Glas et al. [13] used multiple laser rangefinders to estimate the orientation of the tracked person in addition to the position. Mozos et al. [29] used multiple laser rangefinders targeted at different portions of a person's body and merge the resulting scans to better estimate the location of persons. More recently, sensory input from both cameras and laser rangefinders have been used simultaneously in order to improve the detection accuracy [37, 40]. Luber et al. [26] describe an unsupervised learning approach to build models of objects observed by laser range-finders.

These works concentrate on the problem of tracking people reliably and obtaining the sequence of positions of the persons being tracked. In our work, we are interested in the problem of building models of higher-level activities from such tracks and predicting these activities over longer time-scales (on the order of minutes). As in vision-based systems, the tracks from the laser rangefinder system have been used to build activity models. Bennewitz et al. [3] used the Expectation-Maximization algorithm to build models of tracks in an indoor environment. A model was a sequence of positions with an associated Gaussian probability distribution for each position. The positions that comprise a learned motion track were then used as states in a HMM which can be used to estimate the positions of people in that environment [8]. The ability to model the movement of humans is useful for planning robot paths for better tracking. Bandyopadhyay et al. [2] used POMDPs to model the movement for this purpose, though their technique was demonstrated in a simulated environment. In these works, there is no *a priori* assumption of a model of human activity. This is in contrast to the work of Crowley et al. [9] and Brdiczka et al. [6] where a model from cognitive theory (Situation model) has been adapted to monitor human activity in smart environments. Frintrop et al. [12] also used a cognitive model but for image segmentation.

The tracks produced by a laser rangefinder system can be used to build occupancy maps as a representation of observed spatial activity [1, 10, 27, 41, 42]. Yan and Mataric [42] accumulated tracks to determine which parts of an office environment were occupied the most. Arbuckle et al. [1] extended occupancy grids to take into account the differences in occupancy of a space over different time-scales. The spatial activity maps produced by Luber et al. [27] included a Poisson distribution model of expected activity in every region of space. This is comparable to the Poisson distribution models of the frequency of observed activities in outdoor environments in our work.

In outdoor environments, GPS data can be used to build activity models of people [23, 24, 36]. Liao et al. [23, 24] used hierarchical Markov models and conditional random fields to learn the movement patterns. These predictive models are then used to detect unusual behaviors. Patterson et al. [36] used this approach to model activity patterns at a city level. In our work in outdoor environments which can accommodate large numbers of people, we use the tracking system to extract paths of individuals. However, it is also possible to track groups as a whole [22].

The above-mentioned systems do not consider interactions between people. Khan et al. [19] and Song et al. [39] demonstrated tracking systems that can operate even when people are interacting but these systems do not explicitly model the interactions themselves. Our main contribution is a system that uses the information from a laser rangefinder-based tracking system and no *a priori* models of interaction (such as distance) to detect certain types of interaction between persons in a typical environment. Detecting real-world interactions has been studied using mainly vision for tracking [15, 18, 31]. The type of interactions that were detected include meetings between people in an open outdoor area [15, 31] and pick-ups and drop-offs at a car parking lot [18]. In these works, the tracks of individuals are segmented into pre-defined low-level actions such as walking, stopping, and entering the parking lot. Thus, all the interactions that can be detected by the system are pre-defined. Moreover, no model of the observed frequency of activities or interactions is defined. Anomaly detection using these methods reduces to detecting if the observed behavior corresponds to one of the pre-defined models. Consequently, a behavior which happens much more often than expected will not be seen as unusual. In contrast, in our activity models, we have explicitly modeled frequency of activities (by assuming a Poisson distribution).

3 Laser Tracking

In our approach, the movements of people in an environment are measured using laser rangefinders that are placed

along the edges of that environment. Multiple lasers are used to offset the effect of occlusions. The lasers are mounted at approximately adult waist height, since tracking must continue even when people are seated. Each laser rangefinder returns the distance and bearing to every obstacle within its range. The readings from different rangefinders are transformed into a common coordinate system using *Mesh Relaxation* [17]. The range scans are used to maintain a model of the objects in the room that gave rise to these readings. The laser readings are divided into background and foreground readings. Background readings arise from static objects such as walls that are not relevant to the tracking process. The background readings are used to update a background model [11]. Readings that are not explained by the background model are assumed to come from objects that are to be tracked. The foreground model consists of a set of particle filters, one for each object being tracked. The association of foreground readings to particle filters is done in a greedy manner. Foreground readings are first clustered into “blobs”. A blob is a grouping of adjacent foreground readings from the laser that appear to be on a continuous surface. In our experiments, we have assumed that measurements that are separated by less than 10 cm belong to the same blob [11]. Each blob is then assigned to the nearest particle filter. The position estimates are updated using a linear kinematic model which enables tracking the object even through temporary occlusions. Figure 1(b) shows the tracking setup.

4 Activity as a Probability Distribution

The output of the laser tracker is a sequence of x, y positions (in global coordinates) for each tracked object. The positions are not measured at regular intervals. Hence, the sequence is resampled at fixed time intervals (0.5 s) by fitting a least-squared error line between every three consecutive positions. In order to facilitate analyzing the tracking data as arising from discrete probability distributions, every sequence of positions is converted into a sequence of displacements, $((r_i, \theta_i))$, $i = 1, 2, \dots, n$ where r_i is the distance moved in the i -th time interval and θ_i the direction. These displacements are then discretized into one of nine canonical displacements (denoted by \mathcal{X}) as shown in Fig. 1(c). If $r_i < 0.2$ m, then it is discretized as displacement “0”. If $r_i \geq 0.2$ m, then it is discretized as one of displacements “1”–“8”, depending on the sector in which θ_i lies. The value of 0.2 m is the minimum distance moved by a mobile person in a time-step, and so displacements “1”–“8” correspond to only moving persons. Thus, if a person is stationary, the displacements are very small in magnitude, and they are all discretized as canonical displacement “0”. Note that this scheme of discretizing displacements ignores the magnitude of the velocity of a person (provided the displacement

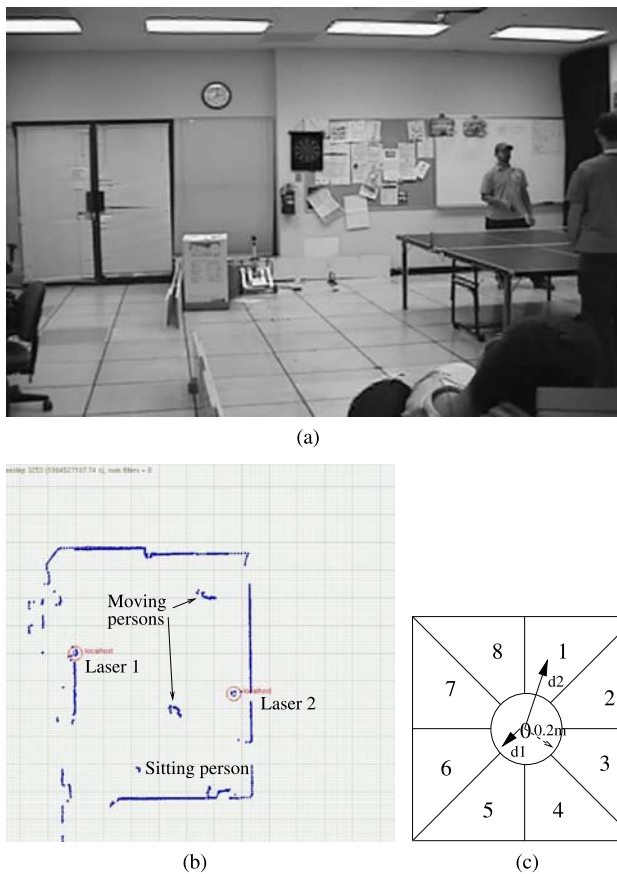


Fig. 1 (a) The open layout laboratory environment. Two people are playing ping-pong while a third is sitting at a desk (bottom right of the photograph). (b) Laser tracking the three people. (c) Discretizing displacements into one of nine canonical displacements. Displacement d_1 is discretized into bin “0”, while displacement d_2 is discretized into bin “1”

is above the threshold of 0.2 m used to distinguish stationary people). We found that the velocity of the people in our environment did not vary much while walking. This observation is similar to that made by [8]. In environments where objects display different velocity profiles, the displacement vectors can be further discretized according to magnitude (in addition to direction). This will likely only improve the discriminatory power of the approach as different objects will be mapped to different probability distributions.

Given a sequence of discrete displacements, we compute the corresponding Maximum Likelihood (ML) probability distribution p over the discrete space \mathcal{X} by counting each type of displacement and dividing by the total number of displacements in the sequence. The ML distribution assigns zero probabilities to those displacements that are not seen in the sequence D . We use Witten-Bell smoothing to assign non-zero prior probabilities to unobserved events. This ensures $p(i) > 0, \forall i \in \mathcal{X}$. In this technique, the probability mass is distributed to the unobserved events in proportion to the number of unique events (relative to the total number

of events). The intuition is that the probability that an event type is unobserved increases if there are many event types in the finite training set.

4.1 Segmenting Tracks

A person’s track may span more than one activity. For instance, a person may sit at a desk for a while, walk to another desk and then sit at that desk. We define an activity to be a set of canonical displacements drawn from a fixed distribution. We assume that different activities give rise to different probability distributions of displacements. For instance, the probability distribution corresponding to a person sitting at a desk would have a high probability for displacement “0”, and low probability for the other eight canonical displacements. Note that the probability distribution is obtained by considering the activity as a set, not a sequence, i.e., displacements are assumed to be independent of each other. The task of segmentation is to divide a track into individual activity segments.

We split each activity sequence into a number of consecutive sub-sequences such that these sub-sequences are distinct from each other in a probabilistic sense. We measure the difference between sub-sequences using a measure of entropy known as the Jensen-Shannon (JS) divergence [25]. Let p_1, p_2 be two probability distributions defined over the discrete space of canonical displacements, \mathcal{X} . Define the weighted average distribution p_{12} as

$$p_{12}(x) = w_1 p_1(x) + w_2 p_2(x) \quad \forall x \in \mathcal{X}$$

where the weights $w_1, w_2 \geq 0$ and $w_1 + w_2 = 1$. The JS divergence between the two probability distributions p_1, p_2 is defined as

$$JS(p_1, p_2) = H(p_{12}) - (w_1 H(p_1) + w_2 H(p_2))$$

where $H(p)$ is the Shannon entropy of distribution. The JS divergence is a modification of the Kullback-Leibler (KL) divergence that gives the JS divergence the property that it is always positive, symmetric, and equals 0 only when the two distributions are equal. The KL divergence corresponds to the likelihood that the newly observed samples (p_2) were generated from the same distribution that generated the earlier samples (p_1). An advantage of the Jensen-Shannon divergence is that it can weight the two probability distributions differently. This is useful in this segmentation algorithm as the number of samples in the left and right sub-sequences will be different for every candidate split point. We utilize the lengths of the sequences giving rise to the probability distributions p_l, p_r to weight the divergence. Thus, $w_1 = k/N$ and $w_2 = (N - k)/N$. Here N is the length of the track and k is the index of the position where the JS divergence is calculated. (This choice of weighting the left

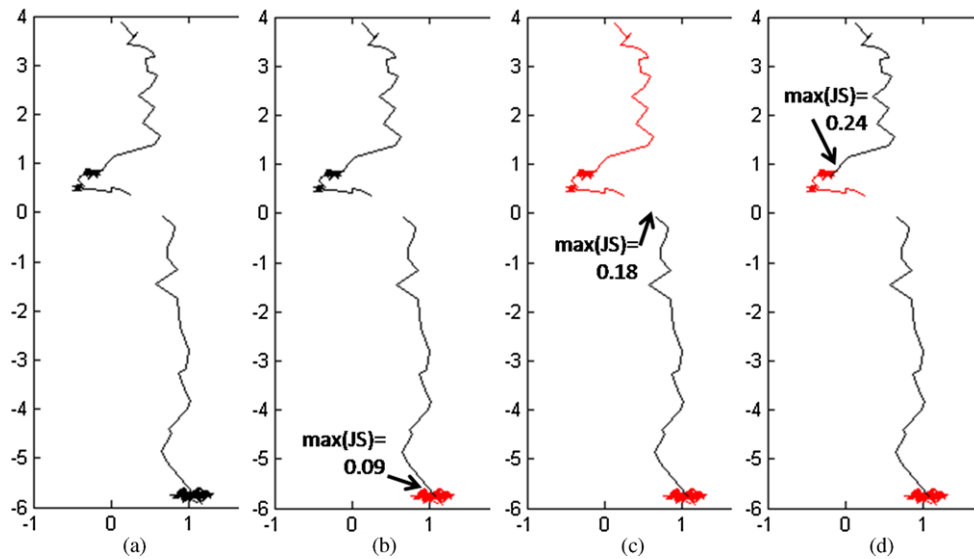


Fig. 2 Splitting a position track into four segments using the recursive segmentation algorithm. The track represents a person walking into a room and stopping at one location, and then walking to another location and stopping there. The two instances of walking and two places of stopping are detected as four distinct segments. (a) The complete

track. (b) The track is first split at the beginning of the second stop. The location and value of the JS divergence at the split point is shown. (c) The first segment is further split at the beginning of the second walking phase. (d) The track is split again at the end of the first walking phase

and right sub-sequences in the calculation of the JS divergence has also been used for segmenting DNA sequences by Bernaola-Galvan et al. [4].

We obtain all activity sub-sequences from a track of displacements by recursively splitting the track at the point that gives rise to maximum JS divergence of the two component tracks. Let $D = \{d_1, d_2, \dots, d_N\}$, $d_i \in \mathcal{X}$ be a sequence of displacements comprising a track. D is split into two sub-sequences D_l, D_r

$$D_l = \{d_1, d_2, \dots, d_k\}, \quad D_r = \{d_{k+1}, d_{k+2}, \dots, d_N\}$$

such that

$$k = \operatorname{argmax}_i (JS(p_{1,i}, p_{i+1,N}))$$

where $p_{i,j}$ is the probability distribution over \mathcal{X} obtained from the displacements d_i, d_{i+1}, \dots, d_j .

Even if all the displacements in a track were drawn from the same probability distribution, the ML probability distribution for the left and right sub-sequences could be different because of the inherent randomness involved in drawing from a probability distribution. The recursive segmenting procedure is therefore terminated when the confidence of seeing the JS divergence between the left and right sub-sequences due to a real difference in their underlying probability distributions falls below a certain threshold. Since we assume that each displacement is an independent identically distributed random variable, we can use the approximation given by Bernaola-Galvan [4] to estimate the minimum value of the JS divergence for a given confidence

value. This assumption is not valid in general but the approximation yields segments that correspond to manual segmentation (Sect. 7.1).

As an example, Fig. 2 shows a track recorded by our laser system (representing a person walking into a room and moving to two locations) split into four segments using the recursive segmentation algorithm. The recursive segmentation algorithm can be modified for online use with a sliding window technique [33]. In this method, the probability distributions are computed over a sliding window of a fixed length instead of the full sequence. Increasing the length of the sliding window makes the online segmentation perform similar to the offline method, but this also increases the latency of the system.

4.2 Clustering Track Segments

Creating and updating a model of activities requires a method to check if two tracks correspond to the same kind of activity. We cluster activity segments to generate a smaller number of representative segments (denoted by C): two tracks arise from similar activities if they lie in the same cluster. We use different distance metrics between tracks from indoor and outdoor environments in the clustering step. The (symmetric) KL distance between two probability distributions is used as the distance metric in indoor environments. In outdoor environments, the mean distance between points on the tracks is used as the tracks mainly represent moving in a relatively straight line.

We use hierarchical clustering to generate the activity clusters. At every step in the clustering algorithm, the two tracks with the smallest divergence measure between them are replaced by a new track that is obtained by taking the mean of the corresponding points in the two tracks. This step is repeated until the minimum divergence among all pairs of tracks exceeds a pre-fixed threshold. A track is said to belong to a cluster if the divergence measure between the track and that cluster is minimal relative to all other clusters. In this work, the number of clusters was set by the authors based on the number of activities observed in each environment.

5 Predictive Model of Activities

We represent the activities performed by a person as Markov models. A Markov model is a triplet $M = \langle S, T, s_0 \rangle$, where S is a set of states, $T : S \times S \rightarrow (0, 1)$ is a transition function that assigns a probability $t(s_t, s_{t+1})$ to every transition between states, and $s_0 \in S$ is the start state. There is a one-to-one correspondence between the states of the Markov model and the representative activity segments discovered during the clustering step, $L : S \rightarrow C \cup \{C_0, C_\perp\}$ such that

$$L(s_i) = C_k, \quad s_i \in S, C_k \in C \cup \{C_0, C_\perp\}$$

and $L(s_i) \neq L(s_j)$ if $i \neq j$, and $L(s_0) = C_0$. C_0 and C_\perp denote special *start* and *end* activities that are generated when a person enters or exits the environment, respectively. The transition probabilities between two states represent how likely the person will perform the corresponding activities one after the other. Note that the duration of activities is not explicitly modeled in this model. The segmenting procedure divides a person's track into activities of different durations.

A Markov model is created for every new activity sequence. The transition probabilities are learned by counting transitions between every pair of representative activities in the sequence. We do not assume an *a priori* model of how likely each activity sequence is possible in the environment, i.e., every activity may be performed with equal probability. Given a sequence of probability distributions $P = \{p_0, p_1, \dots, p_n\}$ representing activities, the probability that it is generated by a Markov model $m = \langle S, T, s_0 \rangle$ is

$$p(P, m) = \prod_{i=1}^n t(s_i, s_{i-1})$$

where s_i is the state corresponding to p_i 's cluster.

In our approach, the states in the Markov models correspond to track segments that are obtained after the segmentation (and clustering) step. Segmentation is not performed at fixed time intervals but to maximize the difference between track segments. The states are then represented by

a probability distribution that is independent of the duration of the original segment. Thus, different states can correspond to different durations. This is one of the advantages of using observable state models compared to Hidden Markov Models (HMM). In the HMM approach all the variable-length segments (symbols output by a state) could potentially be observed at any state and hence time difference between segments has to be explicitly handled. Thus, these approaches interpolate segments to a common length [8], use fine-grained observations such as instantaneous velocity [36], or include duration as a model parameter as in Augmented Markov Models [14]. The reliability of this approach is related to the amount of data used to build the underlying distributions. We expect that using features with more discriminative power (instead of only positions) will improve the robustness of this technique.

5.1 Detecting Interactions

We now describe our procedure for detecting which pairs of activity segments represent interactions. In this work, we assume that two interacting people generate displacements from similar distributions. The similarity between two distributions p_1, p_2 is compared using the symmetric version of the Kullback-Leibler (KL) divergence between probability distributions:

$$\text{KL}(p_1, p_2) = D(p_1, p_2) + D(p_2, p_1)$$

where

$$D(p_1, p_2) = \sum_{x=1}^9 p_1(x) \log \frac{p_1(x)}{p_2(x)}$$

If the KL distance between the distributions of a pair of activity segments falls below a certain threshold, the segments are assumed to come from interacting persons. Currently, this threshold is set by the authors such that the results of automatic segmentation agree most closely with manual segmentation.

Using only the similarity measure to compare distributions means that two stationary objects will always be shown to be interacting. Hence, entropy of the distributions is used as a measure of *confidence* in that the similarity measure indicates a real interaction, i.e., two activity segments are more likely to represent interactions if their corresponding probability distributions have higher entropy. For instance, two people following each other along a zigzag path are more likely to be interacting than two people following each other along a straight line. Using the entropy of the distributions enables us to rank candidate interacting pairs based on our confidence in the interaction. Thus, our system would be more certain that two people following each other are interacting than two people sitting at their desks.

In environments where the range of interactions can be pre-defined, it may be possible to define more specific measures for detecting interactions (for instance, based on relative distance and velocity). Our similarity-based measure can be combined with such measures in order to capture a wider variety of interactions. For instance, Oliver et al. [31] use multiple metrics including position and degree of alignment to detect following, approaching, and walking together in an outdoor environment. However, this will require the estimation of additional parameters that specify the relative contribution of the different interaction measures. In this work, we have studied the similarity measure in isolation in order to determine the type of interactions that can be detected by this method.

Given that the activities of two people are generated by two Markov models, the probability of interaction between these two people can be estimated from the probability that co-occurring states of the two models correspond to the same activity cluster (since clustering in indoor environments is performed with the same metric used for detecting interactions). Predicting future interactions between two people who have generated partial activity sequences P_1 and P_2 now reduces to identifying the most probable Markov models m_1, m_2 representing the sequences:

$$m_1 = \operatorname{argmax}_{m \in M} (p(P_1, m))$$

$$m_2 = \operatorname{argmax}_{m \in M} (p(P_2, m))$$

where M is the set of learned Markov models, followed by calculating the probability that their respective next states represent the same activity type. The likelihood of the models that do not have the same partial initial sequence of activities will be low (due to the low transition probabilities corresponding to these states). Note that the interaction detection steps are performed after the segmentation step. In the current work, the segmentation algorithm is run off-line and hence the interaction detection is also performed off-line.

Figure 3 illustrates the steps of the interaction detection method for the case of segments from two persons playing ping-pong. The Markov models corresponding to the two tracks are shown in Fig. 3(b). The states in the Markov models correspond to the probability distribution clusters calculated in the clustering step (Sect. 4.2). The 15 probability distribution clusters are shown in Fig. 4.

6 Detecting Anomalous Activities

To identify anomalous behavior in the outdoor environment, we build a model of the frequency at which similar activities and interactions are observed. During the monitoring phase, we count the number of times different activities are

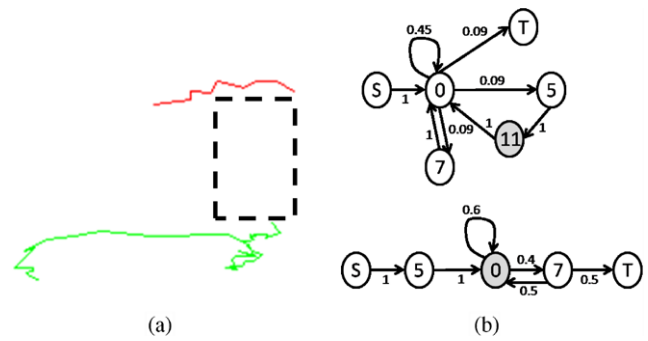


Fig. 3 Example illustrating the interaction detection method. (a) Track segments from two ping-pong players on either side of the table (dashed lines). (b) Markov models that matched the tracks corresponding to the two segments. The current state is shown shaded. The interaction is detected by comparing the most probable next state from the current state in each Markov model

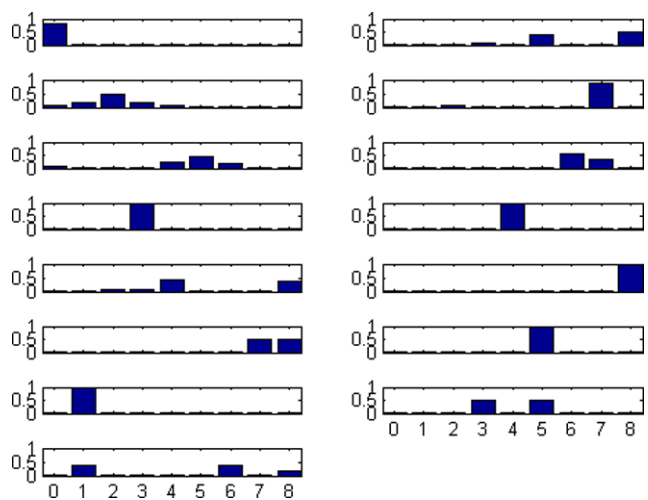


Fig. 4 15 discrete probability distribution clusters corresponding to the states in the Markov models. The discrete space is the range of possible displacements shown in Fig. 1(c)

observed in a small interval. If the probability of observing the detected number of activities in that interval falls below a certain threshold, an anomalous behavior is flagged.

We compute the mean distance between corresponding points on two tracks as a distance metric for use in the clustering procedure. As the number of points in any two tracks is generally different, the tracks are normalized before they are compared to each other: a cubic spline is fitted through all the points in a track. The fitted spline is then sampled at 20 equidistant points to generate the normalized track.

We assume that the number of similar activities performed in a given time period follow a Poisson distribution. The probability $p(n)$ of observing n activities in a given time interval is then given by

$$p(n) = \frac{\mu^n e^{-\mu}}{n!}$$

where μ is the mean number of activities expected to occur in that time interval. The Poisson distribution is completely specified by the parameter μ . This parameter is computed for each activity cluster by counting the total occurrences of activities from that cluster and dividing by the duration of the observation period.

7 Results

We monitored people in two indoor environments: an open layout laboratory environment measuring approximately 8 m × 8 m, and a long, narrow corridor measuring approximately 27 m × 1.5 m. The laboratory environment contained desks and chairs arranged along the sides and a ping-pong table near the middle. Typical activities in this room included sitting at desks, moving between desks, and playing ping-pong (Fig. 1a). Our interaction detection methods were evaluated by trying to detect the relatively complex interaction between ping-pong players. Typical activities in the corridor were less varied: walking along the corridor, or waiting by a door. Three SICK laser rangefinders were placed in the open layout laboratory environment and two laser rangefinders were placed in the corridor environment to capture position data. We conducted three motion capture sessions in each environment, each lasting three hours. Normal unscripted activity proceeded throughout the recording sessions.

7.1 Quality of Segmentation

The ability to detect interactions is related to the quality of segmentation. If two consecutive activities are not segmented into separate sub-sequences, then the resulting (combined) probability distribution will not accurately represent either of the activities. On the other hand, segmenting a single activity into two sub-sequences raises the chance that the ML distribution differs significantly from the true distribution. We stop the recursive segmentation algorithm when the confidence in the JS divergence falls below 90%.

Segmentation works well when a stationary person starts to move or vice-versa. In these cases, there is 100% agreement with manual segmentation. During ping-pong, segmentation should ideally occur only when the ball goes out of play (when one of the players moves away to pick it). However, not every such event is detected because the ball may be picked up close to the playing area (false negatives). There are also instances when the pace of the game changes significantly causing incorrect segmentation (false positives). We considered a playing segment of approximately 600 seconds and compared the activity sub-sequences that were segmented out with manual segmentation (by watching a video recording of the activities). The results are shown in Table 1.

Table 1 Comparison of JS recursive segmentation with manual segmentation for a ping-pong sequence

	Hand segmented	True positive	False positive	False negative
Player 1	13	9	13	4
Player 2	13	7	12	6

7.2 Interactions in the Laboratory Environment

The types of interactions that were identified are (in order of decreasing entropy):

1. *Two people playing ping-pong*: The main result from this experiment is that our similarity-based measure is able to determine that the laser tracks generated by the ping-pong players correspond to interacting persons without any *a priori* model of ping-pong playing. This interaction is particularly interesting because the two players are separated by a significant distance, are constantly moving, and movement can occur in many directions. Thus measures like distance between people, or direction of movement, will fail to pick up this interaction. Though the positions of the players are influenced by each other, the statistical correlation between the sequences is not significant since the players do not move in lock-step. Ping-pong is made up of a series of hitting the ball across the net, picking up the ball, and standing at one location (waiting for the other player to pick up the ball). Our entropy-based measure detects the interaction during the times when the players are hitting the ball. This is because the two players share a characteristic type of movement that is captured by the corresponding probability distributions. As the interaction criterion does not take into account distance, the person sitting close to one of the players is correctly not identified as interacting with the players (Fig. 1a). However, a player standing still cannot be distinguished from a sitting person (false positive).

2. *Two people following each other along the corridor*: Two people moving along the corridor in the same direction are detected as interacting. Since a narrow corridor effectively limits the range of interactions possible (interacting people may either walk beside or follow closely), there are few false negatives. However, not all detected interactions correspond to actual interactions since there is no way for the tracking system to capture the intent of the walkers (false positives). People walking in opposite directions are not detected as interacting.

3. *Two people standing or sitting*: Whenever a person is stationary, the probability distribution has a high value for canonical displacement “0” and very low values for all other displacements. Thus all stationary people are shown to be interacting albeit with low entropy (false positive).

The spatial distribution of all the tracks in the laboratory environment is shown in Fig. 5(a). The tracks correspond to

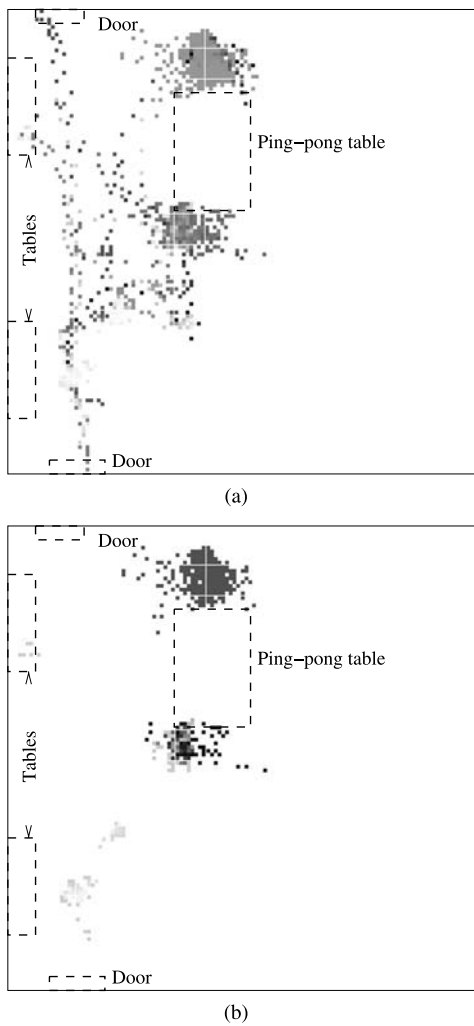


Fig. 5 (a) Spatial distribution of tracks in the laboratory environment. This is a plan view of the room with the main features marked by dashed lines. Darkness corresponds to the average entropy in that area. (b) Spatial distribution of interactions.

walking between the two doors, sitting/standing at the desks, and playing ping-pong. The spatial distribution of the interactions in the laboratory is shown in Fig. 5(b). The chief areas of interaction are at the ends of the ping-pong table and in front of the desks (lower entropy since sitting at a desk is a lower entropy activity than playing ping-pong).

7.3 Predicting Interactions

The number of people in the environment ranged from zero to seven. The resulting tracks were segmented and clustered into 15 representative activity segments using the hierarchical clustering algorithm. 81 Markov models of activities were learned during this session. We used these models of activity to infer the activities during a 30 minute session in the same environment.

To test the ability of the models to predict interactions we presented pairs of sequences of activities to the models

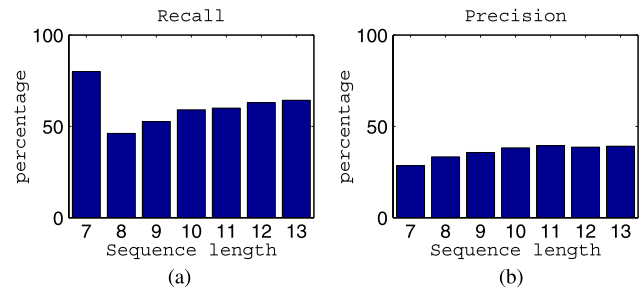


Fig. 6 (a) Recall (as percentage) while predicting if an interaction will take place between two people given a sequence of activities performed by them (length of sequence along x-axis). (b) Corresponding precision as percentage

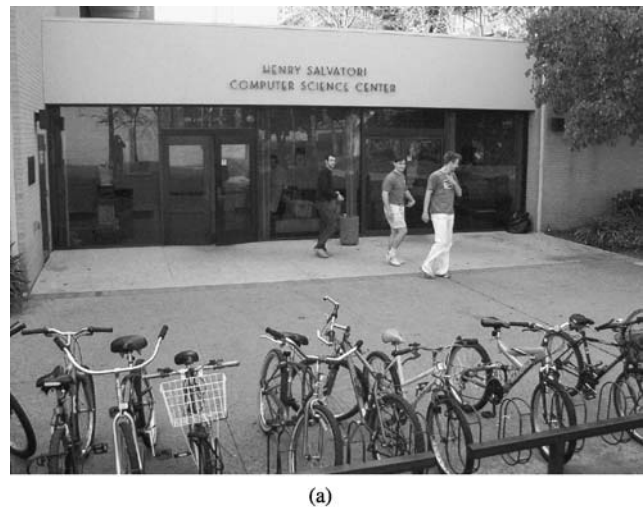
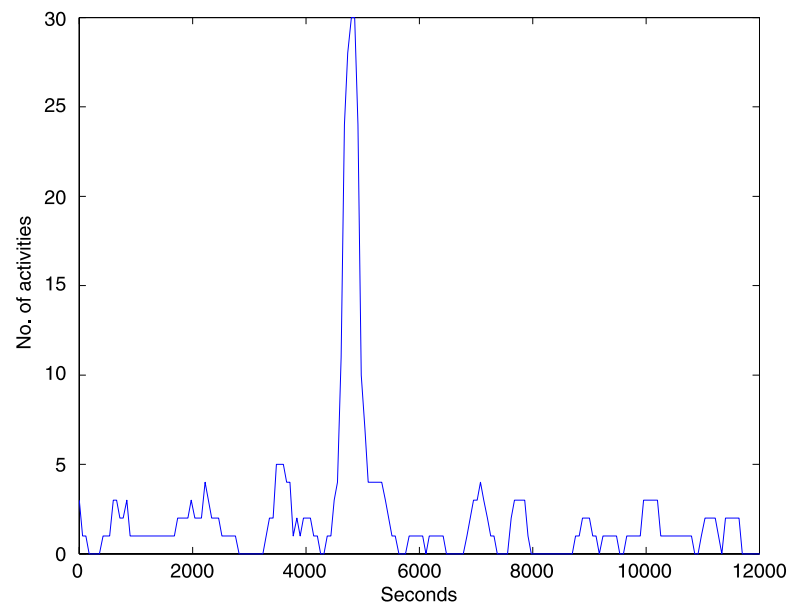


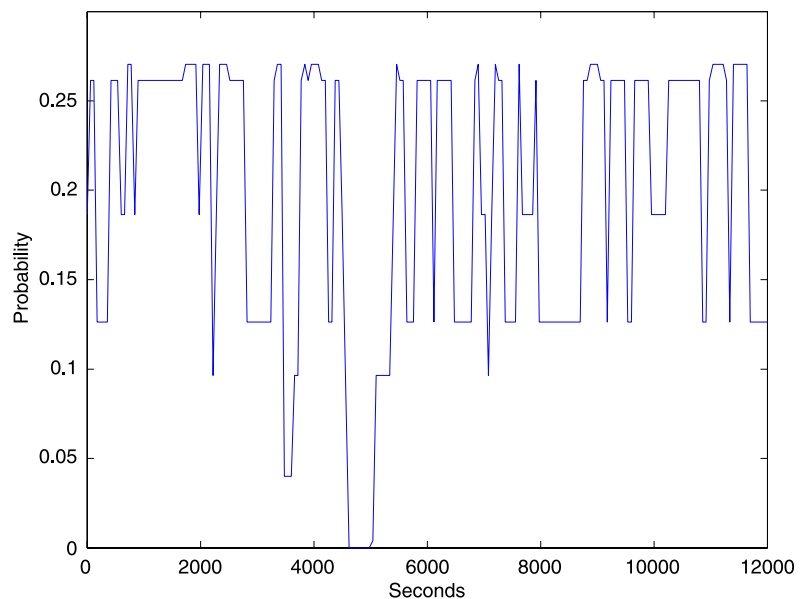
Fig. 7 (a) People crossing the courtyard where we monitored activity. (b) The lines indicate the 20 clusters of activity in the environment

and calculated the most probable Markov model and current state that fit these sequences. The next states of the selected models were used to predict an interaction during the next activity. The prediction was compared with the actual activities that were next performed by the peo-

Fig. 8 (a) Number of activities observed in 5-minute intervals. (b) Expected probability of seeing the number of activities in 5-minute intervals according to the learned Poisson distribution



(a)



(b)

ple. A true interaction is assumed to have occurred if the probability distributions corresponding to the activities that were performed by the two persons were clustered to the same state in the Markov model. The recall and precision curves for predicting interactions in the next activity change are shown in Fig. 6. As expected, a longer history gives better prediction. Though the recall value for a sequence of length 7 is higher than for longer sequences, this has a lower precision value. The unexpected higher recall value is due to the relatively small number of training instances. The interactions that took place during the longer activity sequences corresponded to ping-pong exchanges.

7.4 Anomalous Event Detection

We monitored people in an outdoor courtyard measuring approximately $10\text{ m} \times 10\text{ m}$ (Fig. 7(a)). People can enter/exit the courtyard via the doors of a lecture hall or through a walkway passing through the courtyard. Most of the activities in this environment consisted of people crossing the courtyard from one of the entrances to one of the exits. There is a large increase in people exiting the lecture hall after a class. Such activities often happened in small groups. Occasionally, a small group of people stopped in the courtyard for a conversation before exiting the courtyard. We next describe our modeling methodology that considers only the

frequency of the observed activities. Our system flags the sudden exodus of students from the lecture hall as an anomalous occurrence without any pre-defined concept of the occupancy of an environment.

Two SICK laser rangefinders were placed at the corners of the courtyard to capture position data. The motion capture session lasted 3.5 hours. Normal unscripted activity proceeded throughout the tracking session. The number of occupants of the environment during the experiment varied widely from none to tens of people crossing the courtyard at the end of a class. The resulting tracks were segmented using the recursive segmentation method and clustered into 20 representative activity segments using the hierarchical clustering algorithm (Fig. 7(b)). The movement path of a person from the tracking setup may not be continuous especially when there is a large number of people in the environment. Moreover, objects such as large bicycles are not always distinguishable from people. Hence, not all the clustered segments correspond to complete tracks.

We computed the mean number of activities performed per unit time for each activity cluster by dividing the total number of such activities by the total duration of the experiment. The Poisson distribution then gave us the probability of seeing a particular number of activities performed in any 5-minute interval. Figure 8(a) shows the number of activities performed in 5-minute intervals for one of the 20 activity clusters shown in Fig. 7(b). Figure 8(b) shows the corresponding probability of seeing that number of activities. The probability drops sharply at $t = 4850$ s, where the number of activities in the 5-minute interval goes up to 30, well over the number of activities seen at other times. This time corresponds to the end of a class and the number of students exiting the building into the courtyard suddenly increases.

The Poisson distribution is used to model the typical behavior. In the courtyard environment we assume this to be persons entering and exiting independently. The large increase in the number of people at the end of a class is the unexpected event and it is identified as an anomalous event by our system as it does not fit the Poisson model (since the appearance of the persons in this case is not independent and the expected probability of this event according to the Poisson model is very low).

8 Conclusion

We presented a system that is capable of tracking the movement of people in both indoor and outdoor environments using only laser rangefinders. We then showed how this tracking data can be segmented into sequences arising from distinct activities by representing the tracks as probability distributions. We presented different techniques for modeling the occurrence of activities and classifying activities as interactions for use in indoor and outdoor environments. These

techniques attempt to identify interactions without requiring a pre-defined and detailed characterization of the interactions of interest. We evaluated the techniques designed for indoor environments in an open layout laboratory environment where occupants performed everyday unscripted activities. The system is able to detect that persons playing ping-pong are interacting, without having any pre-defined model of player movements or utilizing the distance between players. The techniques for outdoor environments were evaluated in a courtyard. The system was able to flag a sudden increase in the number of people in the courtyard (due to students leaving a lecture hall after a class) as an anomalous occurrence without any pre-defined concept of occupancy of the environment. The work thus demonstrates the feasibility of detecting and modeling certain kinds of interaction using only the position tracks from a laser rangefinder-based tracking system. The techniques have to be further validated in different environments and over longer durations.

Acknowledgements The authors thank Helen Yan for developing the laser tracking system. This work was supported in part by the US Department of Energy RIM grant DE-FG03-01ER45905, and in part by the US Office of Naval Research MURI grant SA3319.

References

1. Arbuckle D, Howard A, Mataric MJ (2002) Temporal Occupancy Grids: a method for classifying spatio-temporal properties of the environment. In: Proceedings of the IEEE/RSJ international conference on intelligent robots and systems (IROS), EPFL, Switzerland, pp 409–414
2. Bandyopadhyay T, Rong N, Ang M, Hsu D, Lee WS (2009) Motion planning for people tracking in uncertain and dynamic environments. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE International Conference on Robotics and Automation (ICRA), Kobe, Japan
3. Bennewitz M, Burgard W, Thrun S (2002) Learning motion patterns of persons for mobile service robots. In: Proceedings of the 2002 IEEE international conference on robotics and automation (ICRA), May 11–15, 2002, Washington, DC, USA, pp 3601–3606
4. Bernaola-Galvan P, Roman-Roldan R, Oliver J (1996) Compositional segmentation and long-range fractal correlations in DNA sequences. *Phys Rev E* 53(5):5181–5189
5. Bobick AF, Ivanov YA (1998) Action recognition using probabilistic parsing. In: 1998 conference on computer vision and pattern recognition (CVPR). IEEE Comput Soc, Santa Barbara, CA, USA, pp 196–202
6. Brdiczka O, Crowley J, Reignier P (2009) Learning situation models in a smart home. *IEEE Trans Syst Man Cybern, Part B, Cybern* 39(1):56–63
7. Bruce A, Gordon G (2004) Better motion prediction for people-tracking. In: Proceedings of the IEEE international conference on robotics and automation (ICRA), New Orleans, USA, vol 2
8. Cielniak G, Bennewitz M, Burgard W (2003) Where is . . . ? Learning and utilizing motion patterns of persons with mobile robots. In: Eighteenth international joint conference on artificial intelligence (IJCAI-03), Acapulco, Mexico, pp 909–914
9. Crowley JL, Reignier P, Barranquand R (2009) Situation models: A tool for observing and understanding activity. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking,

- 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
10. Emaduddin M, Shell DA (2009) Estimation of pedestrian distribution in indoor environments using multiple pedestrian tracking. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 11. Fod A, Howard A, Mataric MJ (2002) A laser-based people tracker. In: Proceedings of the 2002 IEEE international conference on robotics and automation (ICRA), Washington, DC, USA, pp 3024–3029
 12. Frintrop S, Königs A, Hoeller F, Schulz D (2009) Visual person tracking using a cognitive observation model. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 13. Glas DF, Miyashita T, Ishiguro H (2009) Laser-based tracking of human position and orientation using parametric shape modeling. *Adv Robot* 23(4):405–428
 14. Goldberg D, Mataric M (1999) Augmented Markov models. Tech. Rep. IRIS-99-367, University of Southern California Institute for Robotics and Intelligent Systems
 15. Hongeng S, Nevatia R (2001) Multi-agent event recognition. In: IEEE international conference on computer vision (ICCV), Vancouver, Canada, vol 2, pp 84–91
 16. Hongeng S, Bremond F, Nevatia R (2000) Representation and optimal recognition of human activities. In: 2000 conference on computer vision and pattern recognition (CVPR). IEEE Comput Soc, Hilton Head, pp 1818–1825
 17. Howard A, Mataric MJ, Sukhatme GS (2001) Relaxation on a mesh: a formalism for generalized localization. In: Proceedings of the IEEE/RSJ international conference on intelligent robots and systems (IROS), October 2001, Wailea, Hawaii, pp 1055–1060
 18. Ivanov Y, Stauffer C, Bobick A, Grimson W (1999) Video surveillance of interactions. In: 1999 computer vision and pattern recognition (CVPR) workshop on visual surveillance, Fort Collins, Colorado, USA
 19. Khan Z, Balch T, Dellaert F (2005) MCMC-based particle filtering for tracking a variable number of interacting targets. *IEEE Trans Pattern Anal Mach Intell* 27(11):1805–1819
 20. Kohlmorgen J, Lemm S (2002) A dynamic HMM for on-line segmentation of sequential data. In: Advances in neural information processing systems 14 (NIPS), pp 793–800
 21. Kulić D, Takano W, Nakamura Y (2009) On-line segmentation and clustering from continuous observation of whole body motions. *IEEE Trans Robot* 25(5):1158–1166
 22. Lau B, Arras KO, Burgard W, (2009) Multi-model hypothesis group tracking and, group size estimation. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 23. Liao L, Fox D, Kautz H (2007a) Extracting places and activities from GPS traces using hierarchical conditional random fields. *Int J Robot Res* 26(1):119–134
 24. Liao L, Patterson D, Fox D, Kautz H (2007b) Learning and inferring transportation routines. *Artif Intell* 171(5–6):311–331
 25. Lin J (1991) Divergence measures based on the Shannon entropy. *IEEE Trans Inf Theory* 37(1):145–151
 26. Luber M, Arras KO, Plagemann C, Burgard W (2009a) Classifying dynamic objects. *Auton Robots* 26(2–3):141–151
 27. Luber M, Tipaldi GD, Arras KO (2009b) Spatially grounded multi-hypothesis tracking of people. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 28. Medioni G, Cohen I, Bremond F, Hongeng S, Nevatia R (2001) Event detection and analysis from video streams. *IEEE Trans Pattern Anal Mach Intell* 23(8):873–889
 29. Mozos OM, Kurazume R, Hasegawa T (2009) Multi-layer people detection using 2d range data. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 30. Nguyen NT, Phung DQ, Venkatesh S, Bui H (2005) Learning and detecting activities from movement trajectories using the hierarchical hidden Markov models. In: Proceedings of the 2005 IEEE computer society conference on computer vision and pattern recognition (CVPR), IEEE Comput Soc, pp 955–960
 31. Oliver N, Rosario B, Pentland A (1998) Statistical modeling of human interactions. In: IEEE computer society conference on computer vision and pattern recognition, workshop on the interpretation of visual motion (CVPR-98), Santa Barbara, USA
 32. Oliver N, Horvitz E, Garg A (2002) Layered representations for recognizing office activity. In: Proceedings of the fourth IEEE international conference on multimodal interfaces (ICMI 2002), Pittsburgh, USA, pp 3–8
 33. Panangadan A, Sukhatme GS (2005) Data segmentation for region detection in a sensor network. Tech. Rep. CRES 05-005, University of Southern California
 34. Panangadan A, Mataric MJ, Sukhatme G (2004a) Detecting anomalous human interactions using laser range-finders. In: Proceedings of the IEEE/RSJ international conference on intelligent robots and systems (IROS). IEEE Press, New York, pp 2136–2141
 35. Panangadan A, Mataric MJ, Sukhatme G (2004b) Learning models of human interaction in indoor environments. In: Proceedings of the third international joint conference on autonomous agents and multi-agent systems (AAMAS), IEEE Comput Soc, pp 1308–1309
 36. Patterson D, Liao L, Fox D, Kautz H (2003) Inferring high-level behavior from low-level sensors. In: Proceedings of the fifth international conference on ubiquitous computing (UBICOMP-03), Seattle, USA, pp 73–89
 37. Schiele B, Andriluka M, Majer N, Roth S, Wojek C (2009) Visual people detection: Different models, comparison and discussion. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 38. Schulz D, Burgard W, Fox D, Cremers AB (2003) People tracking with mobile robots using sample-based joint probabilistic data association filters. *Int J Robot Res* 22(2):99–116
 39. Song X, Cui J, Wang X, Zhao H, Zha H (2008) Tracking interacting targets with laser scanner via on-line supervised learning. In: IEEE international conference on robotics and automation (ICRA), Pasadena, California, pp 2271–2276
 40. Spinello L, Triebel R, Siegwart R (2009) A trained system for multimodal perception in urban environments. In: Arras KO, Mozos OM (eds) Workshop on people detection and tracking, 2009 IEEE international conference on robotics and automation (ICRA), Kobe, Japan
 41. Wolf DF, Sukhatme GS (2007) Semantic mapping using mobile robots. *IEEE Trans Robot* 24(2):245–258
 42. Yan H, Mataric MJ (2002) General spatial features for analysis of multi-robot and human activities from raw position data. In: Proceedings of the 2002 IEEE/RSJ international conference on intelligent robots and systems (IROS-2002), September 30–October 4, 2002, EPFL, Switzerland, pp 770–775

Anand Panangadan is a Research Specialist at the Saban Research Institute of the Childrens Hospital Los Angeles and a research affiliate at the NASA Jet Propulsion Laboratory. He received his PhD in Computer

Science from the University of California Los Angeles in 2002, and B.Tech in Computer Science and Engineering from the Indian Institute of Technology, Bombay in 1996. He was a post-doctoral researcher at the University of Southern California's Robotics Research Lab in 2003–2004. His current research interests are in body sensor networks and integrating environment sensor networks with remote sensed data.

Maja Matarić is a professor of Computer Science and Neuroscience at the University of Southern California (USC), founding director of the USC Center for Robotics and Embedded Systems, co-director of the USC Robotics Research Lab and Senior Associate Dean for Research in the USC Viterbi School of Engineering. She received her PhD in Computer Science and Artificial Intelligence from MIT in 1994, MS in Computer Science from MIT in 1990, and BS in Computer Science from the University of Kansas in 1987. Her Interaction Lab's research

into socially assistive robotics is aimed at endowing robots with the ability to help people through individual assistance (for convalescence, rehabilitation, training, and education) and team cooperation (for habitat monitoring and emergency response).

Gaurav S. Sukhatme is a Professor of Computer Science (joint appointment in Electrical Engineering) at the University of Southern California (USC). He received his undergraduate education at IIT Bombay in Computer Science and Engineering, and M.S. and Ph.D. degrees in Computer Science from USC. He is the co-director of the USC Robotics Research Laboratory and the director of the USC Robotic Embedded Systems Laboratory which he founded in 2000. His research interests are in multi-robot systems and sensor/actuator networks. He has published extensively in these and related areas.