

Laser Tracking and Classification of Multiple Objects

Ajo Fod, Maja Mataric and Gaurav Sukhatme

Abstract—

Tracking people is a popular problem in machine vision. In this paper, we describe a method for real-time tracking of moving people with a laser range finder - a sensor rapidly gaining popularity in robotics. We show that with planar laser data we can track and distinguish between animate and inanimate objects, by using heuristic quantities. We also show how our method can be extended to estimate the path of objects behind occlusions. We present initial results on object classification based on their movement signatures and a Kalman Filter formulation tailored for laser tracking.

Keywords— **Kalman filter, Laser, Human Tracking, Occlusion Reasoning, Map Building, Human Interactions.**

I. INTRODUCTION

Tracking people, for surveillance, crowd control, activity recognition and characterization is a well studied problem in machine vision [21]. Our attempt to track people is inspired in part by the need in humanoid robotics to perceive, imitate and react to human beings. It is also motivated in part by the potential for effective tracking using a non-vision sensor, specifically a laser rangefinder.

In this paper, we discuss a real-time method of laser-based tracking of objects as a step towards the larger goal of identifying people and their activities. We demonstrate the use of a standard off-the-shelf laser rangefinder as an efficient sensor for tracking multiple humans in a workspace. We are not aware of other work in tracking people using planar scans. Existing commercial laser-based trackers servo on the location of the object being tracked using a single laser beam [1]. They are usually expensive, cannot handle occlusions, and impose limitations on the acceleration and velocity of the object being tracked.

Laser sensing differs significantly from vision in ways that can be exploited for tracking. In vision, variables like color, intensity, depth (if binocular) are available. In contrast, lasers are restricted by what they can observe: a 1D projection of events in one plane of the observable space. Most of the useful information for tracking is in just one parameter - the range to the nearest obstacle on small intervals over an arc. The range measurements are, however, of high accuracy, especially in comparison with other range sensors, such as ultrasound or infra red. Thus, lasers are rapidly gaining popularity for mobile robotic applications such as collision avoidance, navigation and localization[9], [23], and map building[27], [13].

In this paper we show that lasers are highly effective in a problem domain where vision has so far dominated. Compared to vision, laser data are sparse but more accurate

and efficient, as much less processing is required. Since lasers are based on measurements on laser light reflected from target objects, they are not very sensitive to noise from natural sources like ambient light in the case of vision and require almost no calibration. Furthermore, they have significant range, capable of covering entire rooms or corridors, and are thus ideal candidates to address the tracking problem. As with other line-of-sight devices, lasers are sensitive to occlusions, but their impact on tracking can be reduced by introducing mirrors or additional lasers.

We study the problem of tracking multiple moving people in a workspace. In this paper we discuss work using a single laser. We draw from vision-based approaches to estimate motion recursively from subsequent scans of the laser. We show what is involved in adapting the general ideas behind these methods to real-time laser data. We develop constraints for the Kalman filter formulation in [22], [25] specially tailored for our application. We also modify the nature of the constraints for computational efficiency.

II. APPROACH

A. Overview

Fig 1 illustrates the overall structure of the algorithm. The laser returns range measurements to the nearest obstacle in polar coordinates. From this we construct two models, the *Background model* and the *Foreground model*. The former is used to filter background information that is not of immediate relevance to tracking. The latter is used to predict range measurements and determine the most likely match in consecutive scans. This process, repeated over subsequent scans, enables us to estimate the trajectories of moving objects. Next, we verbally describe the approach, then go into more formal details.

The Background Model consists of a set of statistics that identify the background, namely the mean and the variance of the range to the stationary background. Thus each range measurement can be classified, within bounds of certainty, as part of the foreground or background. We stop updating background statistics whenever the current scan is known to not belong to the background.

Measurements that do not fit the background model are passed to the foreground model. We aggregate continuous measurements in the foreground into “blobs.” At each step, the foreground model consists of two sets of blobs - one consisting of the blobs from previous scan and the other from the current scan. Each blob in the current scan is matched or associated with one or more blobs from the previous scan. Blobs are modeled as rigid bodies moving in a plane. Their velocities are stored as a state vector.

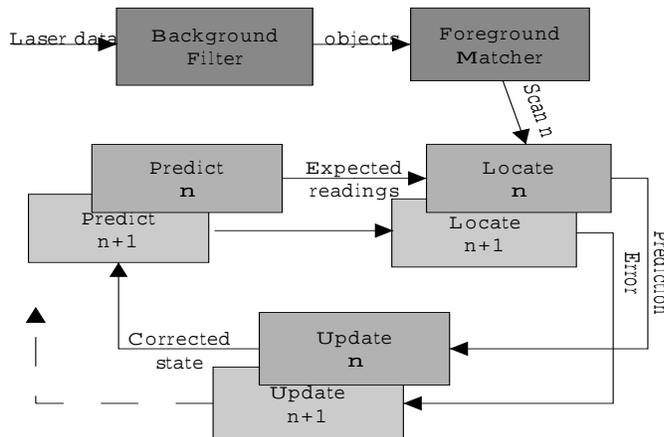


Fig. 1. **Background Filter:** Laser data is first filtered to remove background readings. **Foreground Match and Locate:** blobs in each scan are matched with blobs from the previous scan. **Update:** The error in match is converted to a correction in state estimate. **Predict:** Position of the blob in the next scan is predicted. **Note:** There are as many instances of the Predict-Locate-Update loop as there are blobs, but there is only one instance of the Background Filter and the Foreground Matcher

The state vector associated with a blob in a scan is a linear combination of the states from all matching blobs from the previous scan. After the completion of a scan, the error in prediction of velocity is used to re-estimate the state of the blob.

Every object manifests itself as a set of blobs. Consequently, groups of blobs representing the same object (such as for example the arms and the torso of a moving human) tend to stay together. Such groups are parsed into a higher level entity called an “object” which is tracked by an *ObjectTracker*. *ObjectTrackers* have a long term existence. They are used to store information and discriminate between blobs and people.

B. Background Model

Background models used in most vision applications utilize pixel-level statistics [17], [21]. Each new pixel value is classified as the background or a moving object. A similar function could be achieved using lasers by keeping the statistics of the laser reading in a background map and processing outliers as objects.

To apply this general idea to laser data, we assume that the farthest known stationary object can be assumed to be a part of the background. In order to be robust to small errors, we also maintain the mean and the variance of the measurements that are classified as background. We use the variance and the deviation of any given reading to test the hypothesis that the reading could be from the background.

Any measurements that cannot be classified as the back-

ground are treated as follows:

- Objects that are in front of the current background are processed as part of the foreground.
- Any reading greater than 8m is assumed to be an error c.f. [II-H]; readings above this threshold are ignored.
- If the range measurement is greater than the current background but below the error threshold, the measurement is considered part of a new background.

C. Tracking

The tracking process is divided into two subtasks: 1) tracking continuous surfaces and 2) grouping surfaces into identities. A *BlobTracker* tracks the continuous surfaces, and an *ObjectTracker* groups blobs into objects and estimates paths of people in the workspace. The *ObjectTracker* is discussed further in section II-G.

At a high level, tracking involves matching across different scans. Consequently, we need to group measurements that have similar properties together into blobs. Formally, we define *Blobs* as groupings of consecutive foreground readings from the laser that appear to be on a continuous surface. The set of measurements that form this continuous surface constitute their *signature*. In our experiments, we assume that measurements that are spatially separated by less than 10cm belong to the same blob.

The Foreground model is a sequence of groups of blobs. It contains two generations of blobs at any given instant: 1) **gen-x**, the old blobs and 2) **gen-y**, the new blobs. To correct the state estimate of every blob, we need to match blobs in gen-x with those in gen-y. To reduce the search space in computing matching blobs, we draw a rectangular bounding box around any blob. Based on a bound on average acceleration of $10m/s^2$ (about 1g), we estimate the maximum error in the prediction of the state to be about 50mm. We expand the bounding box in all directions by this amount and search for matches between the signatures of a blob in gen-x and a blob in gen-y only if an overlap exists between the expanded bounding box of gen-x and the simple bounding box of gen-y. Within a blob, a match is said to exist between two points if the distance between them is less than 50mm. We attempt to pair every point in gen-y with two points in gen-x that match and have the minimum distance to the new point. We use these two points for a linear approximation of the surface of the blob. This concept is along the lines of the chamfer method suggested in [4] and is illustrated in figure 2.

We use the Extended Kalman filter (EKF) to predict the state of the blob. Newly discovered gen-y blobs inherit their state vector as a linear combination of the predicted state vectors from the matched gen-x blobs. The weighting of each of the gen-x blobs in determining the state of the gen-y blobs is given by the number of point matches between the blobs:

$$q_y = \frac{\sum_x m_{x,y} q_x}{\sum_x m_{x,y}} \quad (1)$$

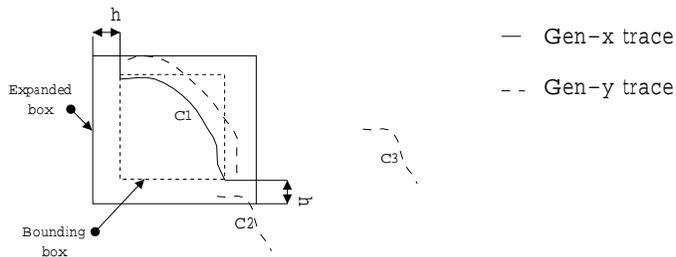


Fig. 2. Consider a single blob belonging to gen-x. The bounding box is first expanded to correct for possible errors. Blobs C1, C2, and C3 are gen-y blobs. Only blobs C1 and C2 are considered for a match. Of these, C1 is strongly linked to the gen-x blob, and C2 is weakly linked to it.

where $m_{x,y}$ is the number of point matches between blob x of gen-x and blob y of gen-y.

At the end of the scan, every gen-y blob's state vector is a linear combination of gen-x state vectors. Blobs that do not have any parents are initiated with a state vector of 0. The gen-y blobs now predict their state at the end of the next scan. The gen-x blobs are removed and the current gen-y blobs become the gen-x blobs for the next scan.

D. Motion Model

Many methods exist for determining structure independent of motion [20], structure from motion [15], [19], and motion of a known structure [26], [10], [11]. Koller et al [16] have described a method of tracking objects using a template of the object. We use a method of recursive estimation of motion using an EKF aided by the constraint that the resulting images lie on the previously determined surface. Effectively, our template of the object for any given scan is given by the previous scan. Here, we use a version of the subspace filter proposed by Soatto et al [22], [25] modified extensively to exploit the properties of laser data and to perform in real time.

We assume that the scene is composed of moving opaque (to the laser) objects. Each reading from the laser corresponds to the nearest obstacle along a straight line. Since the laser is practically oblivious to the component of motion perpendicular to the plane of scan of the laser, we do not attempt to model motion or changes along this direction. In other words, we only consider objects that move and rotate in a plane. Within reasonable margins, human beings tend to stay in a plane parallel to the ground in a planar workspace.

For simplicity, we model every laser range reading as a point in the Cartesian space \mathbf{x}_i . The change in any point δx is then given by:

$$\delta \mathbf{x}_i = \mathbf{C} \mathbf{q} \delta t \quad (2)$$

where

$$\mathbf{C} = \begin{pmatrix} 1 & 0 & -r_y \\ 0 & 1 & r_x \end{pmatrix} \quad (3)$$

$$\mathbf{q} = (v_x \quad v_y \quad \omega)^T \quad (4)$$

where v_x and v_y are the x and y components of velocity, ω is the angular velocity and δt is a small change in time.

E. Rigidity Constraints

We use a rigid body model of motion and the signature from the previous scan to predict the next scan. The matched blob's signature is also used to update the rotational and rectilinear velocities of the object. For a blob with N points, the space which represents the blob would have $2N$ dimensions. However, a rigid body moving in a plane has only 3 degrees of freedom corresponding to our state vector. The extra dimensions represent the plasticity of the object and noise in laser ranging. Consequently, the space in which the points representing a rigid body move forms a 3-D subspace of the $2N$ -D space. Here we formulate the constraints that help us evaluate our predictions in each step.

As presented in [22], the method of subspace filtering requires one-to-one correspondences between features. As a consequence of their inability to distinguish between objects at the same depth, laser measurements are sparse in features. In the laser sensor space, our strength is in the ability to determine the range reading accurately.

We model the surface of objects as a sequence of lines connecting adjacent estimates of the signature. Each new j^{th} measurement $\mathbf{x}_{j,t+1|t+1} = (x_j, y_j)$ is paired with the two closest adjacent i and $i+1$ predictions $\mathbf{x}_{i,t+1|t}$ and $\mathbf{x}_{i+1,t+1|t}$ based on the previous observations. We measure the distance from a line joining $\mathbf{x}_{i,t+1|t}$ and $\mathbf{x}_{i+1,t+1|t}$ and the point $\mathbf{x}_{j,t+1|t+1}$ as a measure of the error in our prediction of the point. We thus attempt to fit the new point to a linear approximation to the previous surface. This error is summed over all points in the new surface to obtain a metric that measures the performance of the prediction.

A line can be represented by the equation:

$$y \cos \theta - x \sin \theta - c = 0 \quad (5)$$

Substituting another point on the LHS gives us the perpendicular distance of the point from the line.

As mentioned earlier, each new point (x_i, y_i) is paired with exactly two predicted points - here they are represented by $\mathbf{x}_{i,t+1|t} = (a_x, a_y)$ and $\mathbf{x}_{i+1,t+1|t} = (b_x, b_y)$. For small errors, the distance E_i of a generic point (x_i, y_i) from

the line passing through and \mathbf{a} and \mathbf{b} .

$$\Delta x_i = a_x - b_x \quad (6)$$

$$\Delta y_i = a_y - b_y \quad (7)$$

$$\tan \theta_i = \frac{\Delta y_i}{\Delta x_i} \quad (8)$$

$$\bar{x}_i = \frac{a_x + b_x}{2} \quad (9)$$

$$\bar{y}_i = \frac{a_y + b_y}{2} \quad (10)$$

$$E_i = (y_i - \bar{y}_i) \cos \theta_i - (x_i - \bar{x}_i) \sin \theta_i \quad (11)$$

To determine the impact of errors in the estimate of the velocity of the blob, we determine this distance as a function of these small errors in the velocities Δv_x , Δv_y and $\Delta \omega$. The predicted surface is shown here as a function of the original surface and the parameters.

$$\bar{x}_i \Rightarrow \bar{x}_i + \Delta v_x t - y \Delta \omega t \quad (12)$$

$$\bar{y}_i \Rightarrow \bar{y}_i + \Delta v_y t + x \Delta \omega t \quad (13)$$

$$\theta_i \Rightarrow \theta_i + \Delta \omega t \quad (14)$$

We substitute the above in the distance metric and use the following to linearize:

$$\cos(\theta + \Delta \omega t) = \cos \theta - \Delta \omega t \sin \theta \quad (15)$$

$$\sin(\theta + \Delta \omega t) = \sin \theta + \Delta \omega t \cos \theta \quad (16)$$

Substituting these into equation 11, linearizing the Δ terms and rearranging we get:

$$\begin{aligned} E_i &= t \sin \theta_i (\Delta v_x) - t \cos \theta_i (\Delta v_y) \\ &\quad - t(x_i \cos \theta_i + y_i \sin \theta_i) (\Delta \omega) \\ &\quad + (y_i - \bar{y}_i) \cos \theta_i - (x_i - \bar{x}_i) \sin \theta_i \end{aligned} \quad (17)$$

This equation is in the form:

$$E_i = \frac{\partial E_i}{\partial v_x} \Delta v_x + \frac{\partial E_i}{\partial v_y} \Delta v_y + \frac{\partial E_i}{\partial \omega} \Delta \omega + c_i \quad (18)$$

$$= \sum_j \frac{\partial E_i}{\partial q_j} \Delta q_j + c_i \quad (19)$$

We now formulate a sum of squares error metric $M = \sum_i E_i^2$. We determine the best correction that can be applied to our estimate of the state. Since the error is quadratic, the exact solution Δq_j^* which minimizes the error in the q_j direction is given by:

$$\Delta q_j^* = - \frac{\sum_i \frac{\partial E_i}{\partial q_j} c_i}{\sum_i \frac{\partial E_i}{\partial q_j}^2} \quad (20)$$

In the absence of confidence in the pre-existing state estimate, the correction term applied to the estimate of the velocity would be exactly as required by the above equation. Since the blob trackers information about the state from previous scans, they use an EKF to incorporate the state estimate from the current scan to update their state estimate using the above information.

In their implementation of an EKF for vision, Soatto and Perona [25] define an operator $C(\mathbf{x}, \mathbf{q})^\perp$ and a locally linear constraint of the form $C(\mathbf{x}, \mathbf{q})^\perp \dot{\mathbf{x}} = \mathbf{0}$. Each constraint in their formulation measures the deviation of a point from the corresponding point in the observation. Consequently, they have as many constraints as points. The computational time complexity of the update steps of their EKF formulation is at least $O(n^2)$ in the number of constraints used. Their formulation also requires the computation of the differential of the SVD of the constraint matrix and the computation of the inverse of a matrix of the size of the constraints. Our formulation replaces these constraints with a set of three constraints of the form $\Delta q_j^* = 0$, one for each state variable, which express our desire that the correction used by the EKF be as close as possible to our optimal estimate derived above.

Our modified version of C^\perp is the transform on the points (x_i, y_i) that gives us the optimal state correction. We compute the derivatives of the above constraints analytically and use them in the update equations of the EKF formulated in [22].

The matrices $D_+(t+1)$ and $C(t+1)$ are:

$$\frac{\partial \Delta q_j^*}{\partial q_k} = - \frac{\sum_i \frac{\partial E_i}{\partial q_j} \frac{\partial E_i}{\partial q_k}}{\sum_i \frac{\partial E_i}{\partial q_j}^2} \quad (21)$$

$$\frac{\partial \Delta q_i^*}{\partial x_j} = - \frac{\frac{\partial E_i}{\partial q_i} \frac{\partial E_i}{\partial x_j}}{\sum_k \frac{\partial E_k}{\partial q_i}^2} \quad (22)$$

$$C_{ij} = \frac{\partial \Delta q_i^*}{\partial q_j} \quad (23)$$

$$D_{+ij} = \frac{\partial \Delta q_i^*}{\partial x_j} \quad (24)$$

The above are used in the EKF described in Eqn.25 through Eqn.38.

F. Kalman Filtering

We use a constrained version of the Kalman filter given by Soatto et.al in [22], [25]. The Kalman filter implementation is summarized below.

Initialization:

$$q(0) = q_0 \quad (25)$$

$$P(0|0) = P_0 \quad (26)$$

Prediction step:

$$q(t+1|t) = q(t|t) \quad (27)$$

$$P(t+1|t) = P(t|t) + R_q(t) \quad (28)$$

Update step:

$$\hat{q}(t+1|t+1) = \hat{q}(t+1|t) \quad (29)$$

$$+ L(t+1)C^\perp(y(t), \hat{q}(t+1|t))y' \quad (30)$$

$$P(t+1|t+1) = \Gamma(t+1)P(t+1|t)\Gamma^T(t+1) \quad (31)$$

$$+ L(t+1)D_+(t)R_n(t+1)D_+^T(t)L^T(t+1) \quad (32)$$

Where:

$$D_+(t+1) = \frac{\partial C^\perp \dot{x}}{\partial x(t)}_{y(t),q(t)} \quad (33)$$

$$C(t+1) = \frac{\partial C^\perp \dot{x}}{\partial q(t)}_{y(t),q(t)} \quad (34)$$

$$\Lambda(t+1) = C(t+1)P(t+1|t)C^T(t+1) \quad (35)$$

$$+ D_+(t+1)R_n(t+1)D_+^T(t+1) \quad (36)$$

$$L(t+1) = P(t+1|t)C^T(t+1)\Lambda^\dagger(t+1) \quad (37)$$

$$\Gamma(t+1) = I - L(t+1)C(t+1) \quad (38)$$

R_n is the measurement estimate error and R_q is the state estimation error. The initial estimate of the velocity q is set to the vector 0 and the estimate P_0 is set to the identity matrix I when a new object is detected.

G. Heuristics

Intille et al [24] describe heuristic methods to explain the scenes described by blob trackers in a nine-step procedure. Although our algorithm is based on some of these principles, the accuracy of range measurements changes our approach. For example, we do not need to distinguish between blobs that occupy the same space, since we can identify them easily by using our 2D positioning accuracy.

Our ObjectTrackers are associated with groups of blobs that are likely to map on to a single entity. Each blob inherits its ObjectTracker from its parent in the previous generation. When two blobs form the parent of a blob, the closest ObjectTracker is associated with the new blob. Blobs that appear spontaneously close enough to the ObjectTracker are associated with the closest ObjectTracker. Whenever no known ObjectTracker is sufficiently close to a new blob, a new ObjectTracker is constructed to describe the blob. Once instantiated, ObjectTrackers persist until there is evidence to indicate that the object they represent has disappeared, perhaps by exiting the workspace. We eliminate ObjectTrackers that are not supported by blobs for over 5 seconds, an empirically determined threshold.

Problems arise when objects move out of or into the plane of the laser. When an object disappears from the plane of the laser, the tracker continues scanning for objects that can be described by the tracker with less error than by other active trackers. In this way, we can predict the position of the object after occlusion. If the occluded object does not reappear during a fixed time interval, its tracker is eliminated.

The ObjectTracker also maintains statistics about the blobs that help us to classify objects into inanimate or animate categories. Such statistics can potentially be used to identify human beings by gross properties of their signature, such as their diameter, mobility, etc. Below we show have shown the use of such statistics to distinguish between inanimate and animate objects by their plasticity,

given by:

$$p = \frac{1}{N} \sum_i \left\{ \sum_j \frac{\partial E_i}{\partial q_j} \Delta q_j + c_i \right\}^2 * 10000 \quad (39)$$

This value is almost always below 1 for solid objects like chairs and tables that have been moved by a person. The rationale is that inanimate objects are typically stationary and rigid, and thus tend to conform to the rigid body constraints stated earlier.

This measure could be used in conjunction with spatial distribution of objects to further qualify our statements. For example, furniture would move only in the presence of animate objects in its vicinity. Interestingly, robots would be considered animate by this method. Another heuristic that could be used is the diameter of a human, which is approximately 0.3 m on average.

The tracker we implemented was designed to work in real time on an IBM 600e laptop. To do this we made small modifications to the above algorithm. To extract the rigid body motion component of the movement, most of the state information can be obtained from a few laser readings from the blob. We use at most five readings from each blob.

We address occlusions by creating bounding boxes around known objects. During an occlusion, parts of the bounding box are beyond the reach of the laser. When this happens, the visible part of the signature is used to update the state and the ObjectTracker continues to attempt to match new blobs within the bounding box. We withhold the EKF measurement updates during the occlusion as suggested in [12], [21]. In most cases, when the object reappears, it can be matched with its estimated position.

H. Errors

While lasers are very accurate in general, they do return erroneous readings occasionally. These errors can be better understood by considering the physical principles involved. The laser scanner uses the time between the transmission and the reception of the reflection of the pulse to determine the range to the observed object. Some surfaces, including very dark objects, do not reflect enough laser light back to the receiver to be visible. Sometimes, the angle of incidence of the laser on a surface is high enough for the beam to be reflected away in a specular fashion. Finally, if nothing is perceived within the range of a laser beam, the laser times out. Consequently, the laser cannot distinguish among the above cases and the resulting range reading always exceeds 8m for our sensor. We handle such errors by ignoring measurements above 8m.

To ascertain the effect of small errors in laser range measurements on our predictions, we find the transformation from measurements in polar space of the form (r, q) into the Cartesian space:

$$\begin{pmatrix} \partial x \\ \partial y \end{pmatrix} = \begin{pmatrix} \cos \alpha & -r \sin \alpha \\ \sin \alpha & r \cos \alpha \end{pmatrix} \begin{pmatrix} \partial r \\ \partial \alpha \end{pmatrix} \quad (40)$$

$$\Delta \mathbf{x} = H \Delta \mathbf{r} \quad (41)$$

The advertised standard deviation of the error of range is 5mm. Since the accuracy is 0.5° , we assume that the standard deviation is half this amount. We also assume that the co-variance matrix is diagonal or that the error in α is independent of the error in r . The covariance of the observation vector \mathbf{x} is then given by:

$$R_n = HR_l H^T \quad (42)$$

Other sources of error arise from our modeling of the process. The laser scans provide an estimate of the surface of the object being tracked. Since the object is not rigid (people are not), we can expect the predictions of the signature of the object to be inaccurate. For simplicity, we threshold the distance from a predicted scan to the actual scan to match objects in subsequent scans. Increasing this threshold would result in more frequent matches which might be inaccurate. On the other hand, decreasing the estimate might lead to lost objects.

III. EXPERIMENTAL PROCEDURE

We have experimented with real scan sequences. In this section we report our observations in experiments conducted on data obtained in our lab.

We use Player [8], a server and protocol that connects robots, sensors and control programs across the network to communicate with our laser rangefinder. The latest version of this software can be found at [18].

A. Setup

Our laser is located at one corner of the workspace, on a table. It scans approximately an inch or so above the hip level. This height was chosen to minimize the distortions caused by hands and feet, and because people do not move out of this plane during normal workspace interactions.

Scans are performed at approximately 10Hz. In each scan, we obtain range measurements closest to the nearest millimeter every 0.5° of a 180° sweep. The laser has a guaranteed accuracy of 5mm (Standard Deviation). In practice however, we have obtained about half that standard deviation.

The system is designed to track human movement. We collected data over normal workspace interactions observed over the span of a few minutes. Our blob detectors track all moving objects, including chairs and doors that are moved by people.

B. Evaluation

In order to compare the performance of the EKF in estimating the velocity state of the object, we integrate the state predicted by the Kalman filter and compare it with the actual path of an object. For simplicity, we consider the case of a stationary object.

We define the center of gravity (c.g.) of an object $\bar{\mathbf{c}}_t$ to be the vector mean of all the laser observations attributed to an object. We define two quantities $\Delta\bar{\mathbf{c}}_t$, the observed change in c.g. and $\Delta\hat{\mathbf{c}}_t$ the predicted change in c.g. at time t . We now define two error metrics:

$$\epsilon_1 = \frac{1}{T} \sum_{t=1}^T |\Delta\bar{\mathbf{c}}_t - \Delta\hat{\mathbf{c}}_t|^2 \quad (43)$$

$$\epsilon_2 = \frac{1}{T} \left| \sum_{t=1}^T (\Delta\bar{\mathbf{c}}_t - \Delta\hat{\mathbf{c}}_t) \right| \quad (44)$$

ϵ_1 is a measure of the instantaneous uncertainty or variance in the error of the position estimate of the blob while ϵ_2 is a measure of the drift in the EKF estimates. Ideally, the drift would be 0 and the variance very small. We measured the above error metrics for a book placed on a table - a stationary object. The variance was in the vicinity of a $1\text{cm}^2/\text{s}$. The EKF's ungrounded estimate drifts about in the vicinity of the location of the actual object.

To measure the effectiveness of the ObjectTracker in estimating the number of people, we performed empirical tests with multiple people in the room. The tracker correctly identified the people in the scene and differentiated between people and inanimate objects. When tracking more than 4 people in a room, our tracker is slightly slower than real time. This is not a significant problem since trackers are independent of one another and the time complexity of the algorithm is $O(n)$ in the number of objects. One solution is to simply use a faster processor or, for dedicated trackers, one optimized for this task. Another solution to the problem is to time-slice the tracking of objects, but this would imply some performance degradation of the tracker. Assuming that the process is a Wiener process, the error measure ϵ_1 would increase proportional to the time between updates. Figure 3 shows the evolution of a scene and its interpretation by our sensor.

IV. DISCUSSION

In this paper we have introduced a novel laser-based tracking approach. Our method exploits the advantages of lasers while overcoming some of the shortcomings. Lasers provide a "gist" of the scene in a plane very accurately. Unlike vision, where multiple cameras are required to cover a scene, a single laser can effectively cover a reasonably large workspace. Owing to the sparse and accurate nature of laser data, it is easier and computationally cheaper to handle. However, it is difficult to obtain a set of features that uniquely distinguish one object from another in laser data. Occlusions are the greatest problem; while it is possible estimate the identities of objects that have been occluded, based on the velocity and position at the beginning of occlusion, it is difficult to distinguish between them positively.

In order to take advantage of the particular features of the laser sensor, we adapted and modified algorithms used in vision-based tracking. We designed a tracking algorithm to be computationally light-weight; it works in real time on an IBM 600e Thinkpad with 64kB RAM and a 300MHz P-II processor. We did not use any special code or hardware optimizations. We believe that with a few optimizations this method can be executed in real time with multiple lasers.

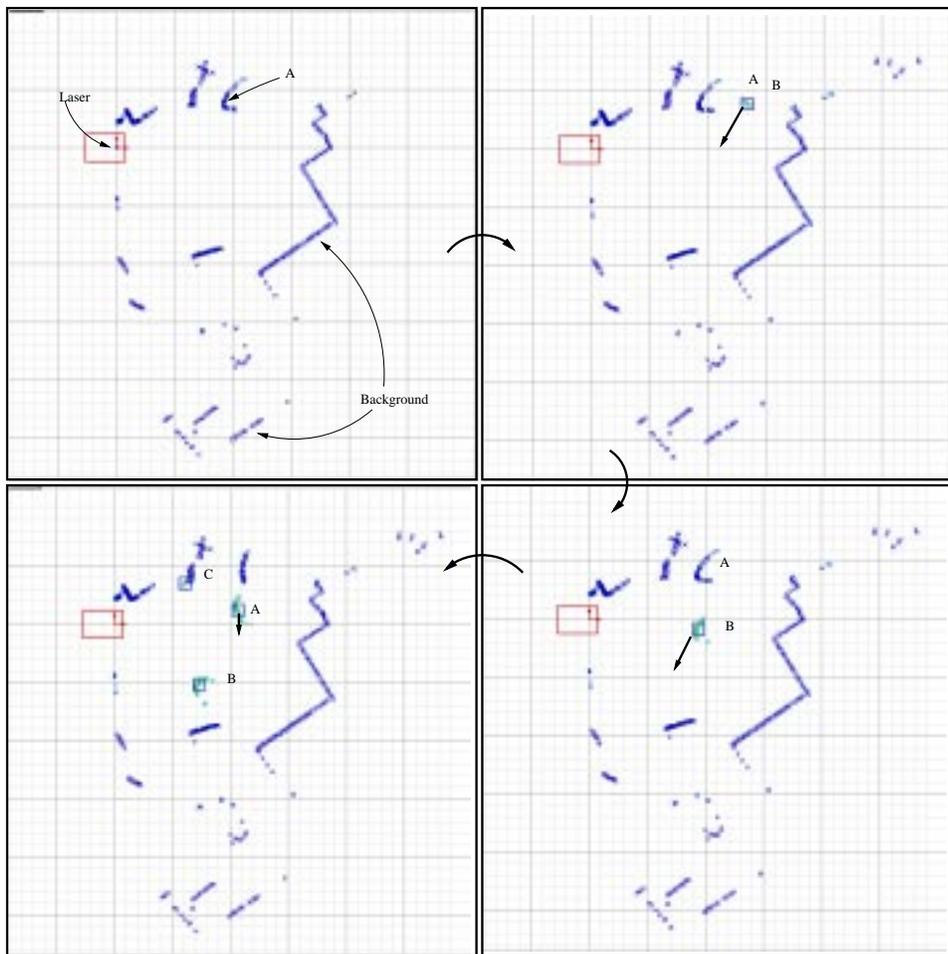


Fig. 3. **A scene:** This shows the evolution of a scene. The laser is the block located at the left of each of the images. Each cross in the figure represents a laser range reading. The figure (a) shows a stationary person A who has not yet moved being picked up by the sensor as part of the background. In figure (d), he has moved and is identified as part of the foreground and a new background (his stationary chair) is now the new background. The person B enters at the top right corner in (b) and is continuously tracked as part of the foreground.

Lasers are oblivious to objects outside their plane of scan, thus making it difficult to perceive details of activity. It is possible to use multiple lasers in different planes to build up a 3-D representation of the scene. Multiple lasers can also be used in large workspaces with a vast numbers of people, where occlusions prohibit effective tracking. Unfortunately, as a practical matter, lasers are currently expensive relative to cameras.

A. Applications

Gait recognition and action modeling are among the interesting applications of trackers that can store and process entity-specific information. Bregler and Brand [7], [5], [6] use motion capture to record joint orientation data for movement analysis. One could modify these to use laser data to determine patterns in the movement in a plane possibly to identify the task being performed or even to identify characteristics of the person.

Microsoft’s research on EasyLiving Technologies [14], [3] requires locating and tracking individuals in a closed workspace. They use this information to guess the intent

of the user and help in automating some of the mundane tasks like switching appliances on in the presence of a user, setting parameters like the temperature in a workspace automatically or even logging on automatically in the presence of a user.

Robotic security devices [2] need to identify moving objects and possibly determine any helpful information about them. One feature would be to find out the number of objects, if the object is an intruder, and if so to determine his physical characteristics. These features can potentially be determined by using a laser tracker.

In this tracking application, we have filtered out the background. However, we could instead filter out the foreground. Moving the laser would then provide us with an estimate of the velocity of the sensor. We have used the EKF to estimate a differential state. We could instead use the EKF to estimate the displacement directly. This would be applicable to laser-based map building, navigation and localization as suggested in [13], [9], [23].

V. CONCLUSION

We have argued that lasers are an effective sensor for object and people tracking. We have evaluated the performance of our EKF to show its effectiveness in tracking objects behind occlusions. We have also developed a computationally efficient algorithm for multiple object tracking which runs in real time on a P-II processor, and validated it in a laboratory workspace.

In our continuing work, we are scaling up the presented approach to track large groups of interacting humans in complex environments with multiple lasers. Our method of attributing qualities to objects to help distinguish between inanimate and animate objects can be extended to qualify objects in other ways, including distinguishing people by height and gait. With cues from spatial configuration of the objects or people in the scene, we can qualify interactions between objects.

ACKNOWLEDGEMENTS

The research described here was supported in part by the NSF Career Grant IRI-9624237 to M. Matarić, and in part by ONR DURIP Grant N0014-00-1-0140 for equipment support to M. Matarić and G. Sukhatme. We thank Andrew Howard for technical assistance.

REFERENCES

- [1] Laser tracking system. <http://members.aol.com/apisensor/lastrack.htm>. Last accessed March 2001.
- [2] Security robots. <http://www.nosc.mil/robots/land/robart/history.html>. Last accessed March 2001.
- [3] J. Krumm B. Brumitt, B. Meyers. Easyliving: Technologies for intelligent environments. In *Handheld and Ubiquitous Computing*, September 2000.
- [4] H. G. Barrow. Parametric correspondence and chamfer matching: Two new techniques for image matching. In *Proc. IJCAI*, volume 2, pages 650–663, 1977.
- [5] M. Brand. Understanding manipulation in video. In *2nd International Conference on Face and Gesture Recognition*, Killington, VT, 1996.
- [6] M. Brand. Style machines. In *SIGGRAPH 2000*, New Orleans, Louisiana, USA, July 23-28 2000.
- [7] C. Bregler. Learning and recognizing human dynamics in video sequences. In *CVPR97*, 1997.
- [8] Richard T. Vaughan, Brian Gerkey, Kasper Stoy. "player robot server". Institute for Robotics and Intelligent Systems Technical Report IRIS-00-391, University of Southern California, 2000.
- [9] S. Thrun D. Fox, W. Burgard. Markov localization for mobile robots in dynamic environments. *Artificial Intelligence*, 11:391–427, 1999.
- [10] D.Gennery. Tracking known 3-dimensional objects. In *Proc. AAAI 2nd Natl. Conf. Artif. Intell.*, pages 13–17, Pittsburg, PA, 1982.
- [11] D.Gennery. Visual tracking of known 3-dimensional object. In *Int. J. of Computer Vision*, volume 7, pages 243–270, 1992.
- [12] J.Malik D.Koller, J.Weber. Robust multiple car tracking with occlusion reasoning. In *Third European Conference on Computer Vision*, LNCS 800. Springer-Verlag, 1994.
- [13] A. Reina J. Gonzalez, A. Ollero. Map building for a mobile robot equipped with a laser range scanner. In *IEEE Int Conference on Robotics and Automation, San Diego, CA, USA*, San Deigo, CA, USA, May 1994.
- [14] B.Meyers J.Krumm, S.Harris. Multi-camera multi-person tracking for easyliving. In *IEEE Workshop on Visual Surveillance*, 2000.
- [15] Thomas J.Oliensis, J.Inigo. Recursive multi-frame structure from motion incorporating motion error. In *DARPA Image Understanding Workshop*, 1992.
- [16] D. Koller. Moving object recognition and classification based on recursive shape parameter estimation. In *12th Israel Conference on Artificial Intelligence, Computer Vision*, December 27-28 1993.
- [17] B. Brumitt K.Toyama, J.Krumm. Wallflower principles and practice of background maintenance. In *International Conference on Computer Vision*, pages 255–261, Corfu, Greece, Sep 1999.
- [18] Robotics Lab. Stage home page. <http://robotics.usc.edu/player/>.
- [19] T.Kanade L.Matthies, R.Szeliski. Kalman filter based algorithms for estimating depth from image sequences. *Int. J. of computer vision*, 1989.
- [20] G.Medioni M.Lee. Structure and motion from a sparse set of views. In *Proceedings of the IEEE Conference on Computer Vision*, Coral Gable, Florida, Nov 1995.
- [21] S.Sclaroff R.Rosales. Improved tracking of multiple humans with trajectory prediction and occlusion modeling. In *IEEE Conf. on Computer Vision and Pattern Recognition*, Workshop on the interpretation of Visual Motion, Santa Barbara, CA, 1998.
- [22] P.Perona S. Soatto, R.Frezza. Motion estimation via dynamic vision. In *Proc. of the 33rd IEEE Conf. on Decision and Control, CDC 94*, 1994.
- [23] W. Burgard S. Thrun, D. Fox. Robust monte carlo localization for mobile robots. *To appear in Artificial Intelligence*, 2000.
- [24] A.Bobick S.Intille, J.Davis. Real-time closed world tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 697–703, Puerto Rico, June 1997.
- [25] P.Perona S.Soatto. 3d visual mtion estimation using subspace constraints. In *Proc. of the 1st IEEE Conf. on Image Processing*, 1994.
- [26] R.Chellappa T.Broida. Estimation of object motion parameters from noisy images. In *IEEE trans. PAMI*, Jan 1986.
- [27] S. Thrun. Learning metric-topological maps for indoor mobile robot navigation. *Artificial Intelligence*, 1:21–71, 1999.