

# Multiple Human Tracking in Dynamic Environment

Aras Akbari, U. of Southern California

**Abstract**—A mobile robot needs to perceive the objects around it and track their moving in real time. Also it need to identify the property of the object, for example, it's human or not. Many researches have been done by using Laser-Based system or vision-based system. The key challenge is that though laser-scanner can get a good estimated position of object in the laser range, it may miss some objects that are blocked by other objects sometimes and is very difficult to identify it's human or not. For Vision-Based system, it is difficult to estimate the position of objects in the real world. In our paper, we use both laser-based system and vision-based system to track the objects in the laser range and identify its human or not when the objects are in the camera range. So we have a global information about the objects in the laser range and its human or not. A Simple Bayes Filter is used in the Laser-Based System, Particle Filter and Adaboost object detection algorithm are used in the Vision-Based System.

## I. INTRODUCTION

For the mobile robot, tracking the moving objects around is an important capability. In the recent development, more and more mobile robots are designed to become the service robots, they work in the environment with human, and provide many different kinds of services such as cleaning the house, sending the mails in the building, even some works of taking care of the children or patients. For this kind of robot, it must own some specific capabilities to identify and track the moving objects, especially the human. To attain this aim, one possible direction is firstly to develop the proper methods to compensate the tracking camera itself, for this work, the image from camera must be kept stable during motion tracking process, therefore the well-designed control algorithms are required. Censi, Fusiello, and Roberto [6] used features tracking method to implement image stabilization. Ke and Sukthankar [7] proposed a new representation algorithm for local images. Lowe [5] presented a method to perform reliable matching through extracting invariant features from images. Zoghلامي, Faugeras and Deriche [23] used automatic homographies warping functions to build 2D mosaic from a series of images. Harris and Stephens [5] used feature tracking algorithms to maintain the consistency of 3D interpretation of image sequences.

After receiving the images from the tracking camera, the following works are to extract the image part of the moving objects. In the image processing, this is called Region of Interest (ROI). Stauer and Grimson [18] presented a algorithm which can monitor the moving objects and identify their moving patterns. Irani and Anandan [4] proposed a new indexing method to access the ROI in the sequential image frames. Lucas and Kanade [9] presented a improved image registration

technique for better efficiency and more capabilities.

For solving the problem of tracking moving objects in ROI of the images, the Gaussian Filter has been tested in single motion tracking. Sclaro and Rosales [15] implemented motion tracking for 3D trajectory through Kalman Filter. However the performance of Kalman Filter is not satisfactory under the imperfect Gaussian distribution conditions. Without the limitation of Gaussian assumptions, the Particle Filter is usually used here to be a substitute for Kalman Filter. Hue, Le Cadre and Prez [6] proposed an extension Particle Filter for multiple object tracking. MacCormick and Blake [10] presented a new sampling method of observation density for multiple object tracking. Some other methods from different approach were also proposed, Luber, Tipaldi and Arras [8] proposed a method with the probability representation - spatial affordance map to track person's motion. Wren, Azarbayejani, Darrell and Pentland [22] studied the performance of a real-time tracking system pfinder. Cox and Hingorani [2] used the MHT (Multiple Hypothesis Tracking) algorithm with Kalman Filter to track moving objects in the image frames. In Rasmussen and Hager's paper [14], Joint Probability Data Association Filtering (JPDAF) was used to track objects of complex geometry and quality. Sullivan, Nillius, and Carlsson [20] tested Bayesian network inference techniques in the problem of surveillance tracking. Lau, Arras, and Burgard's paper [16] showed that the method with the combination of both visual and laser-based means can perform much more well in person tracking tasks. Spinello, Triebel and Siegwart [17] also made use of both visual and laser based data as the observation source, and applied Kalman Filter to track moving objects. In their work, the Conditional Random Fields (CRF) was used to classify the laser data, and Implicit Shape Model (ISM) was used in image detecting system.

## II. RELATED WORK

Several works have been done by using laser-scanner in multiple target tracking and object classification, generally related to the field of detecting and tracking moving objects, including application on localization and navigation [21], SLAM [4]. In order to classify objects; some authors have proposed a variety of approaches: i) voting schemes [11], ii) a method based on "heuristic" rules, [1], and iii) multi-hypotheses [19] are some examples of classification approaches using laser scanners. The first two methods have the disadvantage of not having a self-consistent mathematical framework in order to support their stability and consistency; nevertheless the results presented by the authors shows the feasibility.

Vision Systems are widely used for object detection [3] [12] and constitute a feasible option to be implemented separately, or a long with LRF [13]. In outdoor environments the performance of vision based systems is negatively affected by weather conditions, and the range information, when available, is in general not appropriate for long distance. Therefore laser-scanner in cooperation with camera are a useful solution to be exploited.

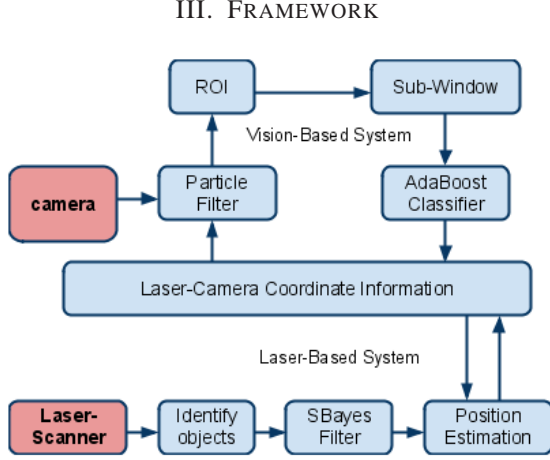


Figure 1. framework for multiple human tracking in dynamic system

Fig. 1 shows our framework of multiple human tracking in dynamic system, combining both laser information and vision information. Our framework consist of two components: laser-based system and vision based information. In the laser-base system, we use a simple Bayes Filter to track the position of all the objects in the range of the laser. The vision-based system uses Laser-Based System information and Particle Filter to locate the exact position of moving objects in the frame image. With having the region of interests (ROI) of the image, the human detection algorithm assigns a probability of humanities to each object. The probability of being human of each objects is transfer to Laser-Based System, by comparing the relative position of objects in the laser-based system and position in frame image.

#### IV. LASER-BASED SYSTEM

In this section, it is discussed the tasks performed in the laser-based system which are related to the processes of identifying Objects and SBayes Filter.

##### A. Identify objects

The sicklaser200 scanner computes for every half angle theta between 0-180 degrees a distance to the point of reection. Converting this data to Cartesian Coordinates gives an accurate depth map for the laser scan information. Fig.2(a) shows the laser depth map. As an ordered sequence of  $N$  measurements points  $P$ , where each scanned point can be defined in Cartesian  $(x_n, y_n)$  or in polar coordinates  $(r_n, \alpha_n)$ , that is

$$P = \left\{ P_n = \begin{pmatrix} r_n \\ \alpha_n \end{pmatrix} \right\} \quad n \in [1, N] \quad (1)$$

Objects in front of the sicklaser200 unit give rise to local points in the depth map which can be easily identified. The local points are shown in blue dot in Fig.2(b). Following the representation shown in Fig. 2(c), and in according to equation (1), a segment  $(S_i)$  can be expressed as

$$S_i = \{(r_i, \alpha_i); (x_i, y_i), : i = k : n, i \leq k \leq n \leq N\} \quad (2)$$

Applying, as input, a complete sequence of scanned points  $(P)$  the Identity algorithm output a pair of  $i^{th}$  values which represent the first and the final beam-point  $(r_n)$  in a detected segment. The algorithms concerned to this category have the general form:

**if**  $D(r_i, r_{i+1}) > D_{thd}$  **then**, segments are separated **else** segments are not separated.

Where  $D_{thd}$  is the threshold condition and  $D(r_i, r_{i+1})$  is the Euclidean distance between two consecutive scanned points

$$D(r_i, r_{i+1}) = \sqrt{r_i^2 + r_{i+1}^2 - 2r_i r_{i+1} \cos \Delta \alpha} \quad (3)$$

A well known approach uses the following threshold condition:

$$D_{thd} = C_0 + C_1 \min(r_i, r_{i+1}) \quad (4)$$

where  $C_1 = \sqrt{2(1 - \cos \Delta \alpha)} = D(r_i, r_{i+1})/r_i$ , and  $C_0$  is a constant parameter used for noise reduction. The Identifying algorithm us the information about the segments to find objects in the depth map as show in Fig.2(d).

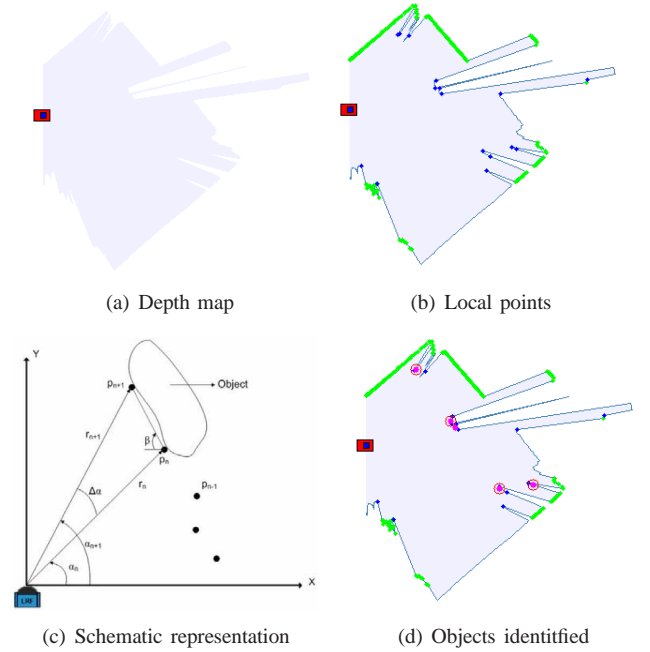


Figure 2. Identify objects in laser system

##### B. A Simple Bayes Filter

A simple Bayes Filter is used for tracking objects which has been identified in the depth map by Identifying object algorithm. The Information we have about objects at identifying objects stage are not stable and reliable. A simple

scenario is that when two objects have the same angle to the laser-scanner, then, only one object will show up in the laser depth map. That's because, the object closer to laser blocks the laser. Therefore, we miss some objects' information. In order to track all the objects in the laser range, we use a simple Bayes Filter to track the state of each object. For each tracking objects at time  $t$ , we define state  $s_i^t$ ,  $i = 0, 1, 2 \dots N$ . The state defined as:

$$s_i^t = [x, y] \quad (5)$$

We define the each object found in the laser information at time  $t$  as  $o_j^t$ ,  $j = 0, 1, 2 \dots M$ . And also each object has component of  $x$  and  $y$ .

1) *Prediction*: The frequency of laser-scanner scanning is very high,  $32Hz$ . So we assume that the velocity at time  $t$  and velocity at time  $t - 1$  has very little change. Therefore, we predict the the state at time  $t$  as:

$$s_i^t = s_i^{t-1} + v_i^{t-1} \Delta t \quad (6)$$

where  $v_i^{t-1}$  is the velocity at time  $t - 1$ ,

$$v_i^{t-1} = (s_i^{t-1} - s_i^{t-2}) / \Delta t \quad (7)$$

2) *Matching*: At time  $t$ , we also have information that, laser-scanner finds  $M$  objects in the laser range,  $o_j^t, i = 0, 1 \dots M$ . Notice that  $M$  is not necessary equal to  $N$ . As mentioned before, if some object stands front of one object, laser-scanner will lose the information of that object. In this case,  $M$  will smaller than  $N$ . So we need to choose a subset of tracking objects to match the  $M$  observations. In that scenario, new object may come into the laser range at time  $t$ . In that case  $M$  is larger than  $N$ . We need to create new tracking object to track the new object. One special case is that, these two things happen at the same time, some objects disappear in the laser-scanner, and some new objects come to the laser range. So the matching should consider all these three scenarios. First, we define the matching function between each tracking object and observation object.

$$p_{ij} = \begin{cases} \frac{0.25}{\|s_i^t - o_j^t\|^2} & \text{if } \|s_i^t - o_j^t\|^2 < Max \\ 0 & \text{if } \|s_i^t - o_j^t\|^2 > Max \end{cases} \quad (8)$$

From the matching function, we can see that the smaller the distance is, that bigger matching score is for the two objects. And also, if the distance is too large, for example, the distance is larger than half of width or length of the laser-scanner range, the matching score will be 0. That's true, because one object will not jump to a very large distance during a very small time  $\Delta t$ . We choose matching pair base on the matching score (for pairs with positive matching score), and of course, this is one-by-one matching, which means one object can not match more than one other objects.

After this matching, we will assign the state of the observation object  $o_j^t$  to the state of matched tracking objects  $s_i^t$ . For those tracking objects that has no matching object, we will keep the state. And for the observation objects that do not found matching, which mean they are new objects, we

will create the new tracking object for that object. Fig.3(a) shows the objects observed by laser-scanner. Some objects are not shown because they are behind of other objects. However, Fig.3(b) shows that our model still can track these objects.

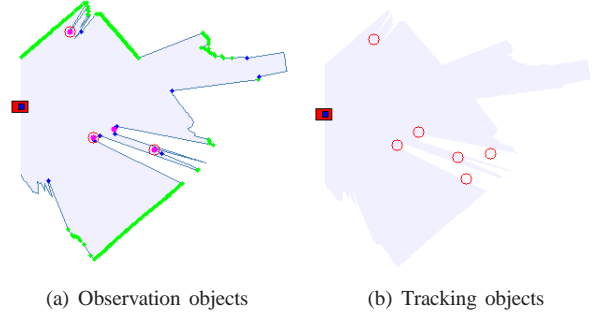


Figure 3. Simplified Bayes Filter for Tracking

## V. VISION-BASED SYSTEM

The vision-based system uses *particle filter* to track the object moving in the image frames. By clustering, we get the area of moving objects. Then we use *Haar Like* features to extract the information and detect the objects are human or not.

### A. Particle Filter

The *Particle Filter* is a simple but effective algorithm to estimate the posterior probability distribution recursively, which is appropriate for real-time application. In addition, its ability to perform multi-model tracking is *attractive* for unsegmented objects detection and tracking from camera images.

1) *Frame Differencing*: We represent our image sequence as  $I^0, I^1 \dots I^{t-1}, I^t, I^{t+1}$ . And all the images have been converted to grayscale images. From the image sequence, we can get the difference between adjacent images, whose pixel values represent the amount of motion occurred in the position. defined as:

$$I_d(x, y) = |I^t(x, y) - I^{t-1}(x, y)| \quad (9)$$

For the image sequence, we can generate the *Frame Difference Sequence*  $I_d^0, I_d^1 \dots I_d^{t-1}, I_d^t$ . Fig. 4(a) shows one example of frame difference.

2) *Probability model*: To deal with errors, a probability approach is adopted. The normalized pixel values in the difference images can be interpreted as the probability of the existence of moving objects in that position, and the position and size of moving objects. This estimation process can be modeled using a Baye formulation. Let  $X^t$  represent the state of moving object.

$$\mathbf{x} = [x, y]^T \quad (10)$$

The posterior probability distribution  $P_m(x^t)$  of the state is derived as follows.

$$\begin{aligned}
P_m(x^t) &= P(x^t | I_d^0, I_d^1) \\
&= \eta^t P(I_d^t | x^t, I_d^0, \dots, I_d^{t-1}) P(x^t | I_d^0, \dots, I_d^{t-1}) \\
&= \eta^t P(I_d^t | x^t) P(x^t | I_d^0, \dots, I_d^{t-1}) \\
&= \eta^t P(I_d^t | x^t) \int P(x^t | x^{t-1}) P_m(x^{t-1}) dx^{t-1}
\end{aligned}$$

Now the posterior probability distribution can be updated recursively by applying a motion model  $P(x^t | x^{t-1})$  and sensor model  $P(I_d^t | x^t)$  over time.

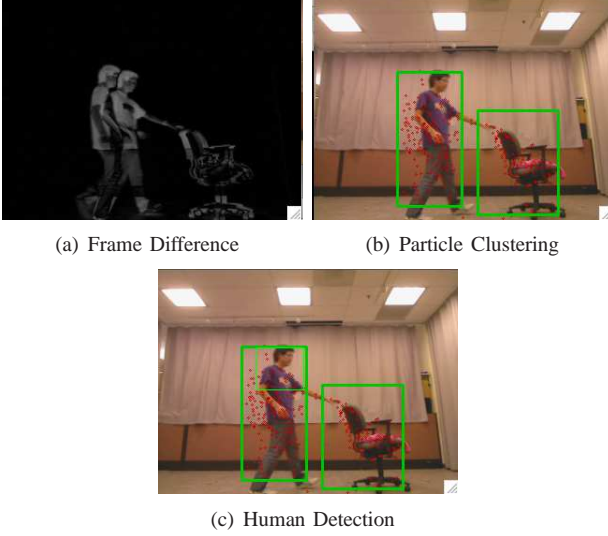


Figure 4. Vision-Based System

3) *Motion Model*: The motion model  $P(x^t | x^{t-1})$  captures the best guess about the motion of moving object.

$$\mu = [x^{t-1} + vx^{t-1}\Delta t, y^{t-1} + vy^{t-1}\Delta t] \quad (11)$$

where  $vx$  and  $vy$  are the velocities in  $x$  and  $y$  direction. The position changes a bit from time  $t-1$  to  $t$ , but we assume that the velocity changes only very little. So we define the  $vx^{t-1}$  and  $vy^{t-1}$  as the  $vx^{t-2}$  and  $vy^{t-2}$ , which is the displacement from time  $t-2$  to time  $t-1$ .

$$[vx, vy] = [x^{t-1} - x^{t-2}, y^{t-1} - y^{t-2}] \quad (12)$$

4) *Sensor model*: The sensor model  $P(I_d^t | x^t)$  is used to evaluate a particle and compute its weight. This captures the idea that if there is a motion at position  $x^t$ , then the difference values of the pixel in the position and its neighbors should be big. The neighborhood can be modeled using multi-variate Gaussian, and the Sensor Model can be define as

$$P(I_d^t | x^t) = \int I_d^t(x) \frac{1}{\sqrt{(2\pi)^d |\Sigma_s|}} e^{-\frac{1}{2}(x-x^t)^T \Sigma_s^{-1}(x-x^t)} dx \quad (13)$$

where  $d$  is the dimension of the state  $\mathbf{x}$ , and the covariance matrix  $\Sigma_x$  controls the range of effective neighborhood. However, for efficiency, we simplify the Sensor Model. A

step function is used instead of a multi-variate Gaussian, and evaluation range is also limited to  $m \times m$  fixed area. The weight  $w_i^t$  of the  $i_{th}$  particle ( $s_i^t = [x_i^t, y_i^t]$ ) is computed by

$$w_i^t = \frac{1}{m^2} \sum_{j=-m/2}^{m/2} \sum_{k=-m/2}^{m/2} I_d(x_i^t - j, y_i^t - k) \quad (14)$$

So for each particle, we only need to evaluate  $m * m$  pixels.

5) *Particles Clustering*: Base on the Laser-Based system, we can know how many object moving in the camera range, which means how many moving objects will show up in the frame image. We assume that there are  $K$  moving objects. So we can divide all the particles in the frame into  $K$  clusters. For each cluster, we can calculate the mean value of the particles and the area of the cluster. Therefore, these areas will become our Region of Interest(ROI) for human detection. Fig. 4(b) shows the two cluster in the image. One is a person, and another one is a chair.

### B. Human Detection

For human detection, we use *Haar-Like* features. Each *Haar-Like* features is represented by a template, its coordinate relative to the search window original and the size of the features. A subset of the features prototypes used is shown in Fig 5. The *Haar-Like* features values are calculated

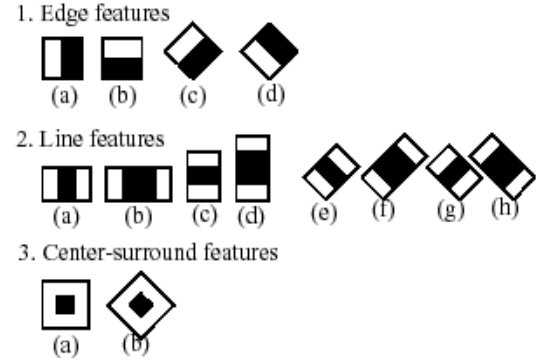


Figure 5. Haar Like Features

as weighted sum of two components: the pixel gray level values sum over the black rectangle and sum over the whole features area. First, a classifier is trained with a few hundreds of samples views of particular object(i.e, a face or a body), called positive examples, that are scaled to the same size(say, 20X20), and negative examples- arbitrary images of the same size. After a classifier is trained, it can be applied to a region of interest(of the same size as used during the training) in an input image. The classifier classify find object or not. To search for the object in the whole image one can move the search window across the image and check every location using the classifier. The classifier is designed so that it can be easily "resized" in order to be able to find the objects of interest at different sizes, which is more efficient than resizing the image itself. So, to find an object of an unknown size in the image the scan procedure should be done several times at different scales.



To detect the object is human or not, we do several kinds of object detection at the same time, face detection and upper-body detection in the Region of Interest(ROI), and combine the detection results together and give a overall evaluation on it's human or not. In the frame image by the robot camera, most of time, we can see the upper body of human, we do face detection and upper body detection. In the Fig. 4(c) show the detection result. In the ROI of person, we detect a face(the small rectangle in ROI) and no body are not detected in the ROI of chair.

## VI. EXPERIMENTS

### A. Experiment Setup

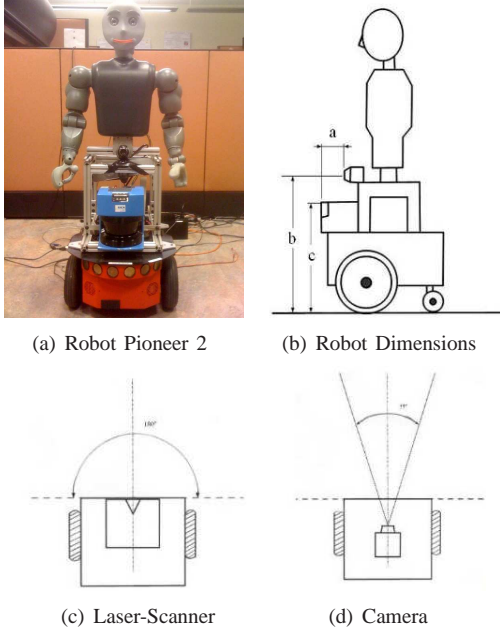


Figure 6. Vision-Based System

In our experiments, the mobile robot platform is *Pioneer2*, with laser and camera installed at the center line of the robot as shown in Fig.6(a). The dimensions of the robot platform is shown in Fig.6(b), where  $a = 5inch$ ,  $b = 24inch$  and  $c = 18inch$ . The camera is *Logitech* camera with 1.3 mega-pixels, the vision angles is 55 degrees. The laser scanner has a 180 degrees scan range as shown in Fig.6(c), the scan frequency is  $32Hz$ . Our framework consist of two systems: Laser-Based system and Vision-Based system. In Laser-Based system, we use Particle Filter used to tracking moving objects in the image frame and do human detection. In the Laser-Based system, we use a probability model to tracking the moving objects in the laser range. We do several experiments on different scenaro and see how each system works and how the whole framework works.

### B. Experimental Processing and Results

In the firs part of the experiment, the only moving object in the scene is humane, firstly, we do detection for only one person, then two persons. In the second part of the experiment,

other moving objects are added into the scene. Firstly a chair moves in the scene, then two chairs move in the scene. In the third part, human and other moving objects move in the scene together.

In our experiment, human is represented by blue circle in the window, non-human object is represented by red circle. The experiment results are in Fig.7. The human head is the feature to differentiate the object, if head is detected in ROI, the object is identified to be human, otherwise the object is not human.

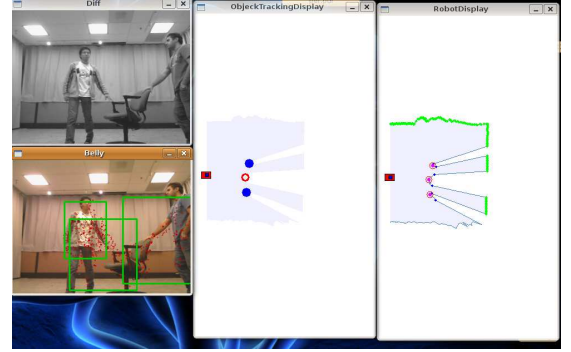


Figure 7. Experiment 1

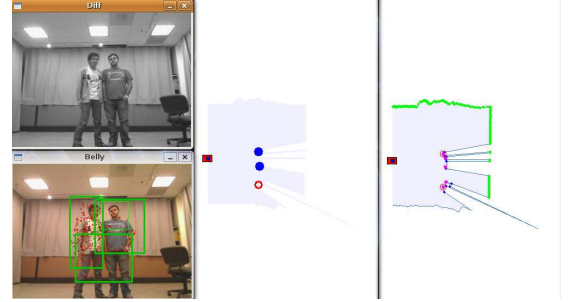


Figure 8. Experiment 2

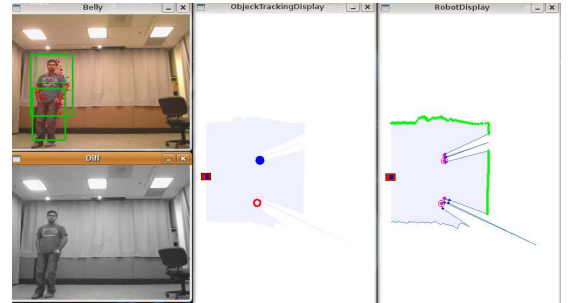


Figure 9. Experiment 3

## VII. CONCLUSION

In this project, we implemented the human tracking function through the combination of laser-based system, vision-based system and image processing techniques.

Laser-based system has excellent capability to trace the exact position of the moving objects, it even provides the detailed locations of the different parts of the object. In our experiment, we found that laser system can locate the exact position of hands and legs, however when the traced object is behind other objects, it is totally invisible from laser's point of view, Bayes filter is suitable for providing the missing information in this problem, and becomes a information source to vision-based system. Particle filter is good in analyzing the information of visual features, giving the exact location of human in the image, then the local new image is generated for human identification. We demonstrate the efficiency and the effectiveness of the integration of these three systems.

In the experiment, our robot is static in the environment, for future work, we plan to test the system in more dynamic environment, and study the possible dynamic factors which may influence its performance, such as the velocity, acceleration and deceleration of the robot, the vibration happening in the robot motion. The different design of the filter for laser-based system is another issue, such as Extended Kalman Filter. Because our real environment is full of the non-linear conditions, the EKF is a better choice in this situation.

#### ACKNOWLEDGMENT

Thanks for Prof. Maja Mataric' and Prof. Gaurav S.Sukhatme's instruction and suggestions.

#### REFERENCES

- [1] Q. Cai, A. Mitiche, and JK Aggarwal. Tracking human motion in an indoor environment. In *Proceedings of the 2nd International Conference on Image Processing (ICIP95)*, pages 215–218, 1995.
- [2] I.J. Cox, S.L. Hingorani, et al. An efficient implementation of Reid's multiple hypothesis tracking algorithm and its evaluation for the purpose of visual tracking. *IEEE Transactions on pattern analysis and machine intelligence*, 18(2):138–150, 1996.
- [3] U. Franke and S. Heinrich. Fast obstacle detection for urban traffic situations. *IEEE Transactions on Intelligent Transportation Systems*, 3(3):173–181, 2002.
- [4] J. Guivant, E. Nebot, and H. Durrant-Whyte. Simultaneous localization and map building using natural features in outdoor environments. In *Intelligent Autonomous Systems*, volume 6, pages 581–586, 2000.
- [5] C. Harris and M. Stephens. A combined corner and edge detector. In *Alvey vision conference*, volume 15, page 50. Manchester, UK, 1988.
- [6] C. Hue, J.P. Le Cadre, and P. Pérez. A particle filter to track multiple objects. In *IEEE Workshop on Multi-Object Tracking*, pages 61–68, 2001.
- [7] B. Jung and G.S. Sukhatme. Real-time Motion Tracking from a Mobile Robot, 2005.
- [8] M. Luber, G.D. Tipaldi, and K.O. Arras. Spatially Grounded Multi-Hypothesis Tracking of People.
- [9] B.D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *International joint conference on artificial intelligence*, volume 3, page 3. Citeseer, 1981.
- [10] J. MacCormick and A. Blake. A probabilistic exclusion principle for tracking multiple objects. *International Journal of Computer Vision*, 39(1):57–71, 2000.
- [11] A. Mendes, L.C. Bento, and U. Nunes. Multi-target detection and tracking with a laserscanner. In *Proc. of 2004 IEEE Intelligent Vehicles Symposium*.
- [12] G. Monteiro, P. Peixoto, and U. Nunes. Vision-based pedestrian detection using Haar-like features. In *Proc. 6th National Festival of Robotics, Scientific Meeting (ROBOTICA)*, 2006.
- [13] J. Neira, J.D. Tardos, J. Horn, and G. Schmidt. Fusing range and intensity images for mobile robot localization. *IEEE Transactions on Robotics and Automation*, 15(1):76–84, 1999.
- [14] C. Rasmussen and G.D. Hager. Probabilistic data association methods for tracking complex visual objects. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(6):560–576, 2001.
- [15] R. Rosales and S. Sclaroff. 3D trajectory recovery for tracking multiple objects and trajectory guided recognition of actions. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2, pages 117–123. IEEE, LOS ALAMITOS, CA,(USA), 1999.
- [16] B. Schiele, M. Andriluka, N. Majer, S. Roth, and C. Wojek. Visual People Detection–Different Models, Comparison and Discussion. In *Proc. of the IEEE ICRA2009 Workshop on People Detection and Tracking*, pages 1–8, 2009.
- [17] L. Spinello, R. Triebel, and R. Siegwart. A Trained System for Multimodal Perception in Urban Environments.
- [18] C. Stauffer and W.E.L. Grimson. Learning patterns of activity using real-time tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(8):747–757, 2000.
- [19] D. Streller and K. Dietmayer. Object tracking and classification using a multiple hypothesis approach. In *IEEE Intelligent Vehicles Symposium*, 2004.
- [20] J. Sullivan, P. Nillius, and S. Carlsson. Multi-target Tracking on a Large Scale: Experiences from Football Player Tracking.
- [21] J. Vandonpe, H. Van Brussel, and H. Xu. Exact dynamic map building for a mobile robot using geometricalprimitives produced by a 2D range finder. In *1996 IEEE International Conference on Robotics and Automation, 1996. Proceedings.*, volume 1, 1996.
- [22] C.R. Wren, A. Azarbayejani, T. Darrell, and A.P. Pentland. Pfnder: Real-time tracking of the human body. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):780–785, 1997.
- [23] I. Zoghlami, O. Faugeras, and R. Deriche. Using geometric corners to build a 2D mosaic from a set of images. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 420–425. Citeseer, 1997.