

REAL-TIME VISUAL SERVOING

Peter K. Allen Billibon Yoshimi Aleksandar Timcenko

Department of Computer Science
Columbia University
New York, NY 10027

ABSTRACT

This paper describes a new real-time tracking algorithm in conjunction with a predictive filter to allow real-time visual servoing of a robotic arm that is tracking a moving object. The system consists of two calibrated (but unregistered) cameras that provide images to a real-time, pipelined-parallel optic-flow algorithm that can robustly compute optic-flow and calculate the 3-D position of a moving object at approximately 5 Hz rates. These 3-D positions of the moving object serve as input to a predictive kinematic control algorithm that uses an $\alpha - \beta - \gamma$ filter to update the position of a robotic arm tracking the moving object. Experimental results are presented for the tracking of a moving model train in a variety of different trajectories.

1. INTRODUCTION

Tracking the three-dimensional movement of objects by a vision system in real-time is an important problem. It has been addressed by researchers in a number of different fields including target tracking, surveillance, automated guidance systems, inspection, and monitoring. The focus of our work is to achieve a high level of interaction between a real-time vision system that is capable of tracking moving objects in 3-D and a robot arm that contains a dexterous hand that can be used to intercept, grasp and pick up a moving object. We are interested in exploring the interplay of hand-eye coordination for dynamic grasping tasks such as grasping of parts on a moving conveyor system, assembly of articulated parts or for grasping from a mobile robotic system. However, the algorithms we have developed are quite general, and applicable to a wider range of domains including the ones described above.

The real-time tracking system we have built encompasses many of the components that are necessary for the creation of intelligent robotic systems. Specifically, the system is able to operate in real-time, at approximately human arm movement rates, using visual feedback in an active sense for a dynamically changing (as opposed to static) environment.

Previous efforts in the areas of motion tracking and real-time control are too numerous to exhaustively list here. We instead list some notable efforts that have inspired us or use similar approaches: Burt et al. [7], Lee and Wahn [15] Corke, Paul and Wahn [8], Goldenberg et al [9], Safadi [18], Brown [6], and an earlier work of ours [2] that performed tracking in the X-Y plane using a single calibrated camera that was mounted on an arm.

The system we have built addresses three distinct problems in real-time motion tracking: fast computation of 3-D motion parameters, predictive control of a moving robotic arm to track a moving object, and grasp planning which entails coordination with a dexterous robotic hand. The focus of this paper is on the first two problems.

This work is notable for the following reasons: First, the vision algorithm is robust since it is based upon detecting optic-flow in real-time as opposed to simple differencing or thresholding methods. Second, the system is capable of tracking objects in three dimensions, in real-time, using unregistered (but calibrated) cameras. Third, the robotic control system is able to accurately predict kinematic parameters for trajectory following of moving objects in real-time. We have tested the system on circular, oval and arbitrary trajectories with good results.

2. COMPUTING OPTIC-FLOW

In a visual tracking problem, motion in the imaging system has to be translated into 3-D scene motion. Our approach is to initially compute local optic-flow fields that measure image velocity at each pixel in the image. A variety of techniques for computing optic-flow fields have been used with varying results including matching based techniques [3], gradient based techniques [12] and spatio-temporal energy methods [1]. Optic-flow was chosen as the primitive upon which to base the tracking algorithm for the following reasons:

- The ability to track an object in three dimensions implies that there will be motion across the retinas (image planes) that are imaging the scene. By identifying this motion in each camera, we can begin to find the actual 3-D motion.
- The principal constraint in the imaging process is high computational speed to satisfy the update process for the robotic arm parameters. Hence, we needed to be able to compute image motion quickly and robustly. The Horn-Schunck optic-flow algorithm (described below) is well suited for real-time computation on our PIPE image processing engine.
- We have developed a new framework for computing optic-flow robustly using an estimation-theoretic framework [19]. While this work does not specifically use these ideas, we have future plans to try to adapt this algorithm to such a framework.

Our method begins with an implementation of the Horn-Schunck method of computing optic-flow [11]. The underlying assumption of this method is the optic-flow constraint equation, which assumes image irradiance at time t and $t+\delta t$ will be the same:

$$I(x+\delta x, y+\delta y, t+\delta t) = I(x, y, t).$$

If we expand this constraint via a Taylor series expansion, and drop second and higher order terms, we obtain the form of the constraint we need to compute normal velocity

$$I_x u + I_y v + I_t = 0$$

where u and v are the velocities in image-space, and I_x , I_y , and I_t are the spatial and temporal derivatives in the image. This constraint limits the velocity field in an image to lie on a straight line in velocity space. The actual velocity cannot be determined directly from this constraint due to the aperture problem, but one can recover the component of velocity normal to this constraint line as:

$$V_n = -\frac{I_t}{\sqrt{I_x^2 + I_y^2}}$$

While computationally appealing, this method of determining optic-flow has some inherent problems. First, the computation is done on a pixel by pixel basis, creating a large computational demand. Second, the information on optic flow is only available in areas where the gradients defined above exist. A second, iterative process is usually employed to propagate velocities in image neighborhoods, based upon a variety of smoothness and heuristic constraints.

We have overcome the first of these problems by using the PIPE image processor [14]. The PIPE is a pipe-lined computer capable of processing 256x256x8 bit images at frame rate speeds, and it supports the operations necessary for optic-flow computation in a pixel-parallel method (a typical image operation such as convolution, warping, addition/subtraction of images can be done in one cycle - 1/60 second). The second problem is alleviated by our not needing to know the actual velocities in the image. What we need is the ability to locate and quantify gross image motion robustly. This rules out simple differencing methods which are too prone to noise and will make location of image movement difficult. Hence, a set of normal velocities at strong gradients is adequate for our task, precluding the smoothing step of the algorithm.

3. A REAL-TIME OPTIC-FLOW ALGORITHM

Our goal is to track a single moving object in real-time. We are using 2 static cameras that image the scene and need to report motion in 3-D to a robotic arm control program. Each camera is calibrated with the 3-D scene, but there is no explicit need to use registered (i.e scan-line coherence) cameras. Our method computes optic-flow fields in each camera and then use a triangulation to intersect the flow fields in areas of image motion in each camera. To implement our algorithm in the PIPE, we used 4 processors on the PIPE. The processors are assigned as 2 per camera - one each for the calculation of X and Y motion energy centroids in each image. We also use a

special processor board (ISMAP) to perform real-time histogramming. The steps below correspond to the numbers in Figure 1 (single camera):

1. The camera images the scene and the image is sent to processing stages in the PIPE.
2. The image is smoothed by convolution with a Gaussian mask. The convolution operator is a built in operation in the PIPE and it can be performed in one frame cycle.
- 3-4. In the next 2 cycles, two more images are read in, smoothed and buffered, yielding smoothed images I_0 and I_1 and I_2 . The ability to buffer and pipeline images allows temporal operations on images, albeit at the cost of processing delays (lags) on output. There are now 3 smoothed images in the PIPE, with the oldest image lagging by 3/60 second.
5. Images I_0 and I_2 are subtracted yielding the temporal derivative I_t .
6. In parallel with step 5, Image I_1 is convolved with a 3x3 horizontal spatial gradient operator, returning the discrete form of I_x . In parallel, the vertical spatial gradient is calculated yielding I_y (not shown).
- 7-8. The results from steps 5 and 6 are held in buffers and then are input to a look-up table that divides the temporal gradient at each pixel by the absolute value of the summed horizontal and vertical spatial gradients. This yields the normal velocity in the image at each pixel.
- 9-10. In order to get the centroid of the motion information, we need the X and Y coordinates of the motion energy. For simplicity sake we show only the situation for the X coordinate. The gray-value ramp in Figure 1 encodes the horizontal coordinate value (0-255) for each point in the image. If we threshold the computed normal velocities, and then AND the above threshold velocities with the positional ramp, we have an image which encodes high velocity with its positional coordinates in the image. In our experiments, we thresholded all velocities below 10 pixels per 60 msec. to zero velocity.
11. By taking this result and histogramming it, via a special stage of the PIPE which performs histograms at frame rate speeds, we can find the centroid of the moving object by finding the mean of the resulting histogram. Histogramming the high velocity position encoded images yields 256 16-bit values (a result for each intensity in the image). These 256 values can be read off the PIPE via a parallel interface in about 10 ms. This operation is performed in parallel to find the moving objects Y centroid (and in parallel for X and Y centroids for camera 2). The total associated delay time for finding the centroid of a moving object becomes 15 cycles or 0.25 seconds.

The same algorithm is run in parallel on the PIPE for the second camera. Once the motion centroids are known for each camera, they are back-projected into the scene using the camera calibration matrices and triangulated to find the actual 3-D location of the movement. Because of the pipelined nature of the PIPE, a new X or Y coordinate is produced every 1/60

second with this delay.

The system exhibits an interesting mix of local and global computations, separated by processors and update rates. The PIPE is able to perform the local optic-flow computation at video rates (but with delay), the ISMAP board gathers global histogram statistics, and these are then shipped via a high-speed interface to a host where the stereo triangulation and kinematic control algorithms reside, updating the arm parameters every 30 msec.

4. ROBOTIC ARM CONTROL

The second part of the system is the arm control. The robotic arm has to be controlled in real-time to follow the motion of the object, using the output of the vision system. The vision system output is not sufficient as a control parameter since its output is both noisy as well as delayed in time. The control system needs to do the following:

- Filter out the noise with a digital filter
- Predict the position to cope with delays introduced by both vision subsystem and the digital filter
- Perform the kinematic transformations which will map the desired manipulator's tip position from a Cartesian coordinate frame into joint coordinates, and actually perform the movement

We have adopted a simple, modular solution for the control subsystem where each task performs its job independently from the others. This approach, although suboptimal, exhibits a high degree of modularity and, nonetheless, simplicity; characteristics which make a complex robotic system more robust and extensible.

4.1. DIGITAL FILTER

The 3-D values of the moving object's position determined by the vision system are both noisy and out of date due to the processing delays. The need for signal smoothing emerged after our first experiments were conducted. Since the sampling period of the control subsystem is 30ms and the noise components were about 100Hz, we have adopted simple second-order filter with a transfer function in the s-domain:

$$G(s) = \frac{Y(s)}{U(s)} = \frac{1}{(1+as)^2}$$

where $Y(s)$ is the Laplace picture of filter's output, $U(s)$ Laplace transform of filter's input with s being Laplace operator and a time constant selected experimentally. Since we have second-order filter, according to Bode's theory we have a 40 dB/decade decay of the filter's amplitude characteristics with a breaking frequency of $1/a$ Hz. The time constant a is experimentally set to 16ms. Upon discretization we get a discrete filter which has a computational expense of 6 multiplications and 4 additions.

4.2. PREDICTIVE FILTERING

In choosing a predictive filter for tracking, the major constraints are knowledge of the nature of the motion to be tracked, knowledge of the noise characteristics of the sensor, and computational cost of the filter. Two common approaches are auto-regressive models (ARM) and Kalman filters. Kalman filtering generates time-variable tracking coefficients that are determined by *a priori* models for the statistics of measurement noise and target dynamics. Because of its dynamic nature, it imposes a computational cost which can be difficult to incur in certain real-time applications. Fixed-coefficient filters have the advantage of simple implementation using fixed parameters for the filter gains [5]. In a few particularly simple target tracking problems, it is possible to derive closed-form steady-state solutions for the associated Kalman filter covariance matrix and the corresponding filter gains. Such solutions can be used to avoid real-time computation of the complete filter equations [13], and can approach the Kalman filter in steady-state performance.

We have implemented an $\alpha - \beta - \gamma$ filter [4, 18] which includes an acceleration estimate as well as a position and velocity estimate.

The gains for the fixed coefficient filters represent a compromise between noise reduction and maneuver-following capability in the steady-state. Kalata [13] has shown that the optimal steady-state gains for a $\alpha - \beta - \gamma$ filter are given by:

$$\frac{\gamma^2}{4(1-\alpha)} = \lambda^2$$

$$\beta = 2(2-\alpha) - 4\sqrt{1-\alpha} \equiv \alpha = \sqrt{2\beta} - \frac{1}{2}\beta$$

$$\gamma = \frac{\beta^2}{\alpha}$$

where λ is the *target maneuvering index* given as

$$\lambda = \frac{T^2 \sigma_a}{\sigma_n}$$

where T is the sample period, σ_a is the position uncertainty due to variance of the acceleration, and σ_n is measurement noise variance. In practice, it is very difficult to accurately estimate the variances above. We treat this ratio of variances as a free parameter in calculating the filter gains, and have found a suitable value empirically.

The updated information from the cameras is used to predict the next set point for the robotic arm that is tracking the moving object. We are using RCCL [10] to control the robotic arm (a PUMA 560). RCCL (Robot Control Library) allows the use of C programming constructs to control the robot as well as defining transformation equations (as described in [17]). The transformation equations permit dynamic updating of arm position by generating the 4x4 transform of the moving object's position from the vision system and sending this information to the arm control algorithm.

5. EXPERIMENTAL RESULTS

Figure 2 shows the experimental hardware that was used for the initial experiments of the system. It consists of 2 CCD cameras with 25mm lenses mounted approximately 1.5 meters apart, with the moving object's full trajectory in the field of view. In the experiments described below, the lenses were increased to 50mm and the baseline increased, with the cameras placed approximately orthogonal to each other. The cameras were calibrated using the two-plane method of Martins et al [16] in which a planar test pattern is imaged twice in each camera. The calibration was able to determine the three-dimensional position of the test patterns to within 3-4 millimeters in X, Y and Z coordinates. The experimental results were obtained with the cameras mounted approximately 1 meter above the plane of the table pointing down at an angle of about 10 degrees from horizontal.

In the experiments, a model train was tracked by vision as it moved around a circular track and an oval track and the arm was commanded to follow each trajectory. The train was moving at a velocity of approximately 25 cm/sec (the algorithm was tested at faster speeds and works with increasing overshoot in the control algorithm). The results shown below were obtained without the robotic hand mounted on the system. Figure 3 shows the thresholded motion-energy for each camera found by the algorithm. Figure 4 is the actual arm trajectory obtained by following the train as it moves on an oval track. The trajectory is quite close to the actual trajectory except as the camera distance to the moving object increases. The algorithm is less accurate here for two reasons: first, the inherent stereo digitization error increases with distance from the cameras, and secondly, the profile of the moving object is quite different at this position, yielding two views whose motion centroids may not correspond to the same exact point in space, thus inducing triangulation error. The tracking begins with oscillations until the filter kicks in fairly rapidly as shown. The path is for 5 revolutions and is very repeatable due to the robustness of the vision algorithm. At the end of the path the arm is commanded to smoothly end its motion. Figure 5 is the trajectory followed for a circular path.

In addition to the trajectories above, we had the arm follow a semi-random trajectory created by a moving toy robot that changes its path direction arbitrarily when it hits an obstacle. In this experiment, the toy robot is placed inside the oval train track and it moves in a path until it hits the tracks, changing direction as it bounces off the tracks. While no ground truth is available for this trajectory, the arm follows the motion quite well, even though the robot moves in a way that is hard to model or predict. The results of these experiments are also available on video tape.

6. FUTURE WORK AND CONCLUSIONS

The algorithm described here has worked quite well in providing real-time visual servoing of a robotic arm, a fairly difficult task. We hope to continue to improve this algorithm, specifically by reducing the stereo error problems discussed above. One approach is to perform a better calibration of the camera system, and the other is try to better isolate centroids in

each separate camera to be the same physical point. Various heuristics are being investigated on this front.

We also are interested in using a Kalman filter approach with dynamic gains for prediction. There is an added computational cost with this approach, but it may yield more accurate trajectories, particularly with more complex trajectories. In addition, we are exploring tracking multiple objects using the PIPE's ability to work on a region by region basis.

Finally, we are working on using the attached robotic hand to intercept and pick up the moving objects. An interesting approach we are currently exploring is to use vision to servo the hand and object together. The parallel processing capabilities of the PIPE may allow us to servo on both moving objects (target and robotic hand) to allow precise contact and grasping.

Acknowledgement: This work was supported in part by DARPA contract N00039-84-C-0165, NSF grants DMC-86-05065, DCI-86-08845, CCR-86-12709, IRI-86-57151, IRI-88-1319, North American Philips Laboratories, Siemens Corporation and Rockwell Inc.

References

1. Adelson, E. H. and J. R. Bergen, "Spatio-temporal energy models for the perception of motion," *Journal of the Optical Society of America*, vol. 2, no. 2, pp. 284-299, 1985.
2. Allen, Peter, "Real-time motion tracking using spatio-temporal filters," *Proceedings of DARPA Image Understanding Workshop*, Morgan-Kaufman Publishers, Palo Alto, May 1989.
3. Anandan, P., *Measuring visual motion from image sequences*, Technical Report COINS-TR-87-21, COINS Department, University of Massachusetts, Amherst, 1987.
4. Bar-Shalom, Y. and T. Fortmann, *Tracking and data association*, Academic Press, 1988.
5. Blackman, Samuel, *Multiple-target tracking with radar applications*, Artech House, Dedham, MA, 1986.
6. Brown, Christopher, "Gaze controls with interaction delays," *Proc. DARPA Image Understanding Workshop*, pp. 200-218, Morgan-Kaufman, Palo Alto, May 23-26, 1989.
7. Burt, P. J., J. R. Bergen, R. Hingorani, R. Kolczynski, W. A. Lee, A. Leung, J. Lubin, and H. Shvayster, "Object tracking with a moving camera," *IEEE Workshop on Visual Motion*, pp. 2-12, Irvine, CA, March 20-22, 1988.
8. Corke, Peter, Richard Paul, and K. Wahn, *Video-rate visual servoing for sensory-based robotics*, Technical Report, Grasp Laboratory, Department of Computer and Information Science, University of Pennsylvania, Philadelphia.

9. Goldenberg, Richard, Wan Chi Lau, Alfred She, and Alan Waxman, "Progress on the prototype PIPE," *IEEE Conference on Robotics and Automation*, Raleigh, N. C., March 31 - April 3, 1987.
10. Hayward, Vincent and Richard Paul, "Robot manipulator control under UNIX," *Proc. of the 13th ISIR*, pp. 20:32-20:44, Chicago, April 17-21, 1983.
11. Horn, B. K. P., *Robot vision*, M.I.T. Press, 1986.
12. Horn, B. K. P. and B. Schunck, "Determining optical flow," *Artificial Intelligence*, vol. 17, pp. 185-203, 1983.
13. Kalata, Paul, "The tracking index: A generalized parameter for $\alpha-\beta$ and $\alpha-\beta-\gamma$ trackers," *IEEE Transactions on Aerospace and Electronic Systems*, vol. AES-20, no. 2, pp. 174-182, March 1984.
14. Kent, E. W., M. O. Shneier, and R. Lumia, "PIPE: Pipelined image processing engine," *Journal of Parallel and Distributed Computing*, no. 2, pp. 50-78, 1985.
15. Lee, Sang Wook and K. Wohn, *Tracking moving objects by a mobile camera*, Technical Report MS-CIS-88-97, University of Pennsylvania, Department of Computer and Information Science, Philadelphia, November 1988.
16. Martins, H. A., J. R. Birk, and R. B. Kelley, "Camera models based on data from two calibration planes," *Computer Graphics and Image Processing*, vol. 17, pp. 173-180, 1981.
17. Paul, Richard, *Robot Manipulators*, MIT Press, Cambridge, MA, 1981.
18. Safadi, Reem Bassam, "An adaptive algorithm for robotics and computer vision application," Tech. Report MS-CIS-88-05, Department of Computer and Information Science, University of Pennsylvania, January 1988.
19. Singh, Aji, "An estimation-theoretic framework for image-flow computation," *Proc. International Conference on Computer Vision (ICCV-90)*, Kyoto, Japan, December 4-7, 1990.

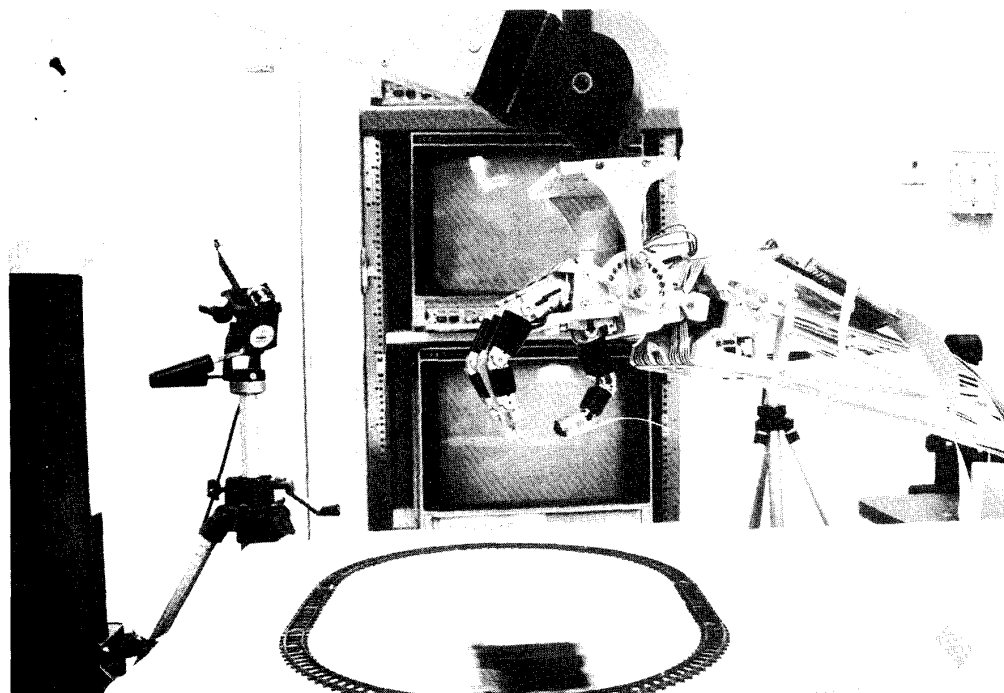


Figure 2. Experimental hardware.

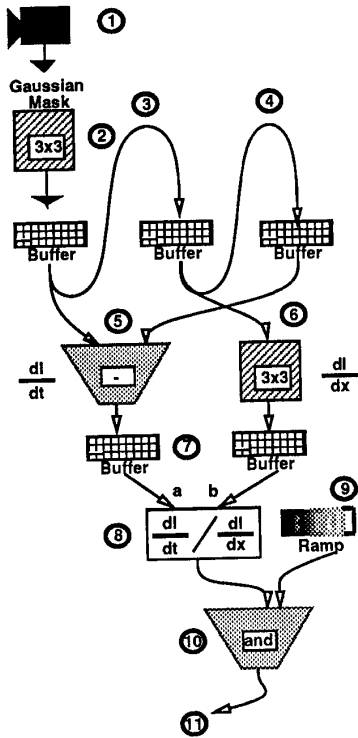


Figure 1. PIPE motion tracking algorithm.

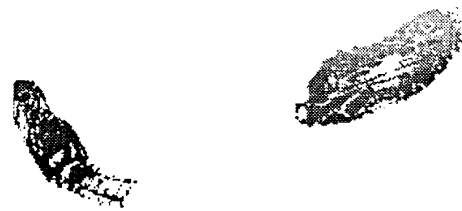


Figure 3. Motion Energy from Optic-Flow (left and right cameras).

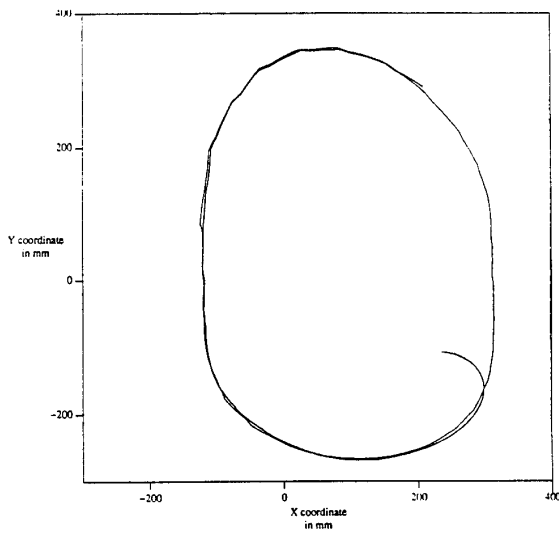


Figure 4. Tracked arm path, oval trajectory

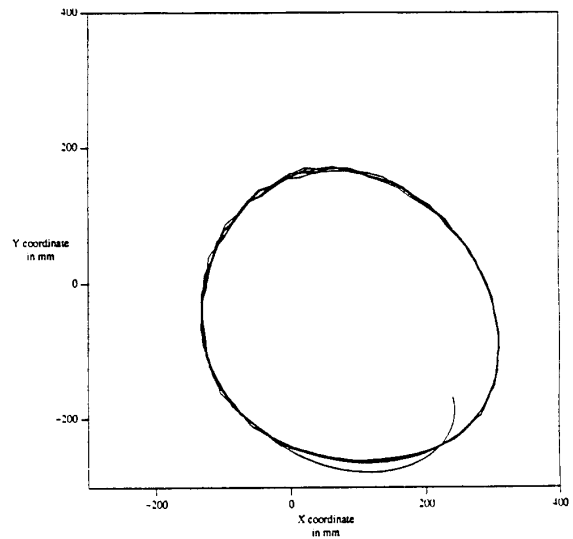


Figure 5. Tracked arm path, circular trajectory