

The active recovery of 3D motion trajectories and their use in prediction

Kevin J Bradshaw, Ian D Reid (Member, IEEE) and David W Murray
(Member, IEEE)

The authors are with the Department of Engineering Science, University of Oxford, Parks Road, Oxford OX1 3PJ, UK. Email: dwm@robots.ox.ac.uk

June 16, 2003

DRAFT

Previous publications have demonstrated, *inter alia*, the generation of saccades where the gaze direction of the system is rapidly redirected to a target of interest, followed by periods of smooth pursuit where fixation is maintained upon a target [6], [3]. The gaze controller, implemented as a finite state machine, selects and reselects these behaviours to provide robust performance in natural dynamic scenes over extended periods.

One of the housekeeping operations of the servo-controller is to propagate odometry — in particular the axis positions or *gaze angles* θ_e and θ_v — around the system. For an extended pursuit sequence, the sequence of gaze angles $(\theta_e(t), \theta_v(t))$ for the platform's elevation axis and one vergence axis defines an observer-based trajectory for the target. Our first aim here is to recover an observer trajectory as the camera pursues an object.

Various algorithms have been implemented to drive pursuit behaviour using active cameras, such as correlation [20], [21], deformable templates [22], feature-based affine transfer [4], [5] and segmented optical flow [6]. The last is used here, though any of the variety of methods might be employed to equal effect.

Optical flow is recovered in a small central, or foveal, area of the image. Because of its small size ($\sim 6^\circ$ field of view in both x and y), it is assumed that only one moving object is included but, because background may be visible, a process of segmentation is required. The initial data computed at 25Hz are edge-normal components \mathbf{v} of the motion field $\dot{\mathbf{r}}$ derived from the image irradiance $E(x, y, t)$ using the motion constraint equation [23]

$$E_t + \dot{\mathbf{r}} \cdot \nabla E = 0,$$

whence the components are found as

$$\mathbf{v} = -E_t \nabla E / |\nabla E|^2 .$$

Of course, whilst the camera pursues a target, a rotational image motion field $\dot{\mathbf{r}}_{rot}$ is induced throughout the image. Fortunately, such motion is independent of the scene depth, and so can be computed directly from the joint positions and velocities which are supplied by the servo-controller. It is subtracted from the observed motion to give components \mathbf{v}_{ind} of the motion field $\dot{\mathbf{r}}_{ind} = \dot{\mathbf{r}} - \dot{\mathbf{r}}_{rot}$ arising from independently moving objects

$$\mathbf{v}_{ind} = \mathbf{v} - [\dot{\mathbf{r}}_{rot}(\theta_e, \theta_v, \dot{\theta}_e, \dot{\theta}_v) \cdot \hat{\mathbf{v}}] \hat{\mathbf{v}} ,$$

where $\hat{\mathbf{v}}$ is a unit vector in the direction of \mathbf{v} . Image regions which correspond to background will, within a noise tolerance, have $\mathbf{v}_{ind} \sim \mathbf{0}$, and are excluded from further consideration. Foreground motion regions are grown by spatially grouping the non-background vectors using a constant velocity model, and the motion vectors and their positions in a segmented region are combined to yield a mean position and mean velocity estimate, $\langle \mathbf{r}_{ind} \rangle$ and $\langle \dot{\mathbf{r}}_{ind} \rangle$. As is apparent in Figure 2, multiple regions may be detected but, as noted earlier, we assume that they arise from the same object. (Other work on saccadic redirections of gaze direction [3] is based on a similar motion process running at coarse scale across the entire image and distinguishes between different regions using magnitude of motion, direction of motion, and area of the moving region. However, the “interest value” of a region based on these parameters remains hand-coded for different applications.) Our use of foveal motion for tracking is further described in [6].

Our aim is to centre the camera’s viewpoint on the target, and to match its velocity. Because we have already subtracted the motion induced by camera rotation, $\langle \mathbf{r}_{ind} \rangle$ and $\langle \dot{\mathbf{r}}_{ind} \rangle$ are *demands* which could, after passing through the inverse kinematics, be sent to the joint servo-controller. However, these raw signals are delayed by the visual latency. As indicated in the introduction, one method we have used to reduce the effect of latency is filter the position and velocity using a constant velocity Kalman Filter, and use the product of filtered velocity and the known latency to evaluate a corrected, prompt, positional demand. Later in the paper we explore whether filtering using 3D motion is more effective.

Three examples of the output from the foveal optical flow process during tracking of a car, a bus, and a person are shown in Figures 2(a–c). The camera platform was mounted on the sixth floor of an office block. In each case, the stills are every eighth frame from a 25 Hz sequence.

The observer trajectories captured from several persons entering and leaving the University’s Computing Centre are shown in Figure 3(a), overlaid on an image taken from the rest frame. (Note that the rest-frame image is for graphical illustration only, and was neither captured during pursuit nor used in the analysis.) The actual observer trajectories from one of these pursuit episodes are shown in Figures 3(b,c). The person exits the building, proceeds along the pavement to their left for some 30m, turns around and doubles back before cutting across the grass to re-enter the building. The time axis unit is one frame. Each took 40ms, and so the pursuit sequence covers a period of some 56 seconds.

I. PROJECTIVE MAPPING TO GROUND PLANE CO-ORDINATES

We now turn to the calibration method used to convert the observer trajectories $\theta_e(t), \theta_v(t)$ of the sort shown in Figure 3, into a trajectory describe in planar Euclidean coordinates based in a scene plane. The basic method requires no knowledge of the camera's position relative to the scene plane.

Figure 4 shows that any pair of gaze angles can be mapped to a point \mathbf{x} in the *frontal-plane* — a plane perpendicular to the resting gaze direction, $\theta_e = \theta_v = 0$, and an arbitrary distance in front of the rotation centre of the camera. The point, given by the forward kinematics, is

$$\mathbf{x} = (\sec \theta_v \tan \theta_e, \tan \theta_v, 1)^\top$$

where we have chosen the arbitrary distance to be unity. The 3-vector \mathbf{x} is a homogeneous coordinate in the 2D frontal plane.

The point corresponding to \mathbf{x} in the scene is \mathbf{X} , where \mathbf{X} is again a 3-vector, a homogenous coordinate in a 2D scene plane. The pair of corresponding points on the two planes is related projectively by an homography

$$\mathbf{X} \equiv [\mathbf{M}]\mathbf{x}$$

where the 3×3 matrix $[\mathbf{M}]$ has only eight degrees of freedom because scale is arbitrary under projective equivalence. The homography $[\mathbf{M}]$ can thus be recovered by establishing the correspondence between at least four known points, or at least four known lines, in each of the two planes [24].

A. Recovery from a 4-point calibration

To establish the calibration using points, the active camera is directed to view in succession the calibration points whose scene coordinates \mathbf{X}_i are known. The fronto-parallel plane position \mathbf{x}_i is derived from the joints angles $(\theta_e, \theta_v)_i$. Each point correspondence provides an equation

$$\begin{pmatrix} \lambda_i X_i \\ \lambda_i Y_i \\ \lambda_i \end{pmatrix} = \begin{bmatrix} M_{11} & M_{12} & M_{13} \\ M_{21} & M_{22} & M_{23} \\ M_{31} & M_{32} & 1 \end{bmatrix} \begin{pmatrix} x_i \\ y_i \\ 1 \end{pmatrix}$$

where by making the scale explicit we can enforce equality rather than projective equality, and where to constrain the degrees of freedom of $[\mathbf{M}]$ we set $M_{33} = 1$. Eliminating λ_i leaves two

equations contributing to the system

$$\begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_i & y_i & 1 & 0 & 0 & 0 & -X_i x_i & -X_i y_i \\ 0 & 0 & 0 & x_i & y_i & 1 & -Y_i x_i & -Y_i y_i \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \begin{pmatrix} M_{11} \\ M_{12} \\ M_{13} \\ M_{21} \\ M_{22} \\ M_{23} \\ M_{31} \\ M_{32} \end{pmatrix} = \begin{pmatrix} \vdots \\ X_i \\ Y_i \\ \vdots \end{pmatrix}.$$

or abbreviated $[\mathbf{A}]\mathbf{m} = \mathbf{b}$. For four or more points, the system can be solved in the least-squares sense as $\mathbf{m} = [\mathbf{A}^\top \mathbf{A}]^{-1}[\mathbf{A}^\top]\mathbf{b}$.

A.1 Implementation example

The view obtained in the resting $\theta_e = \theta_v = 0$ direction from the window of the office block is given in Figure 5(a). Overlaid on the image are the four points used for calibration. In the scene these points \mathbf{X}_i are the corners of a nominally rectangular lawn measured as $(X, Y)_A = (0.0, 0.0)$, $(X, Y)_B = (21.1, 0.0)$, $(X, Y)_C = (21.1, 7.9)$, $(X, Y)_D = (0.0, 7.9)$ in the ground plane, where each unit is 1 metre. Centering the active camera on these points gave head joints angles of $(\theta_e, \theta_v)_A = (1.24^\circ, 21.51^\circ)$, $(\theta_e, \theta_v)_B = (-3.62^\circ, 4.66^\circ)$, $(\theta_e, \theta_v)_C = (-0.48^\circ, -0.67^\circ)$, $(\theta_e, \theta_v)_D = (3.53^\circ, 16.40^\circ)$, giving

$$[\mathbf{M}] = \begin{bmatrix} 0.228 & 0.601 & -73.281 \\ 0.135 & -0.479 & 51.757 \\ 0.002 & 0.013 & 1.000 \end{bmatrix}$$

By registering the image with the frontal plane it is possible approximately to overlay image features onto the calibrated plane. Figure 5(b) shows this on edge features computed from the image in (a). The recovered observer-based trajectories of people entering and exiting the building on the far side of the road and walking along the pavement are shown again for ease of comparison in Figure 5(c), and Figure 5(d) shows the trajectories converted to ground plane trajectories using the computed homography.

If the pursued object lies above the ground plane, the 3D trajectory will be incorrect: it will be positioned where the gaze direction actually strikes the ground plane, that is, further from the camera. This problem is evident in the tracks of Figure 5(d) — empirically it is found that it is the torso centre that is tracked rather than the feet. (The same problem is evident in static features also: note that the the outline of the street lamp, which is of course above the ground plane, is reconstructed as though it were an object “painted” on the ground plane.)

Strictly within the present method, the trajectory can be corrected only by recalibration from points at the correct height h , provided that h is constant. If no such points can be found, then the simplicity of the method is lost somewhat and the ground plane coordinates of the camera (X_c, Y_c) must be known, together with either the height Z_c of the camera, or a knowledge of the plane parallel to the ground plane and containing the camera’s rotation centre. Using the latter, if the current gaze direction makes an angle α with the parallel plane then the corrected ground position is

$$\begin{pmatrix} X' \\ Y' \end{pmatrix} = \begin{pmatrix} X \\ Y \end{pmatrix} + \frac{h \cot \alpha}{\sqrt{\Delta X^2 + \Delta Y^2}} \begin{pmatrix} \Delta X \\ \Delta Y \end{pmatrix}.$$

where $\Delta X = X_c - X$ and $\Delta Y = Y_c - Y$. Figure 5(e) shows trajectories corrected using the latter method.

In Figure 6 we show the recovered $X(t)$ and $Y(t)$ ground plane trajectories corresponding to the observer trajectories given in Figure 3.

B. Recovery from a 4-line calibration

The calibration can also be achieved using four or more lines in the ground plane where, as the working below shows, there is no need to establish “endpoint to endpoint” correspondence.

The homography for lines is obtained from the dual relationships between lines and points in both planes, $\mathbf{u}^\top \mathbf{x} = 0$ and $\mathbf{U}^\top \mathbf{X} = 0$, as

$$\mathbf{U} = [\mathbf{M}^\top]^{-1} \mathbf{u}.$$

The homogeneous representation of a line in the ground plane is given in terms of the normal to the line $\hat{\mathbf{N}} = (\cos \Phi, \sin \Phi)^\top$ and the distance D (which may be negative) along the normal from the Cartesian origin to the line,

$$\mathbf{U} = (U, V, W)^\top = (\cos \Phi, \sin \Phi, -D)^\top$$

Similarly for lines in the frontal plane,

$$\mathbf{u} = (u, v, w)^\top = (\cos \phi, \sin \phi, -d)^\top$$

To constrain the degrees of freedom we set $L_{33} = 1$. (The $L_{33} = 0$ condition occurs only if an observed line at the horizon of the ground plane is mapped to line through the origin of the frontal plane, an occurrence we can safely ignore.)

Each line correspondence is represented as

$$\begin{pmatrix} \lambda_i u_i \\ \lambda_i v_i \\ \lambda_i w_i \end{pmatrix} = \begin{bmatrix} L_{11} & L_{12} & L_{13} \\ L_{21} & L_{22} & L_{23} \\ L_{31} & L_{32} & 1 \end{bmatrix} \begin{pmatrix} U_i \\ V_i \\ W_i \end{pmatrix}$$

and after eliminating λ_i contributes two equations to the system

$$[\mathbf{A}] (L_{11} \ L_{12} \ \dots \ L_{31} \ L_{32})^\top = \begin{pmatrix} \vdots \\ W_i u_i \\ W_i v_i \\ \vdots \end{pmatrix}$$

where

$$[\mathbf{A}] = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ U_i w_i & V_i w_i & W_i w_i & 0 & 0 & 0 & -U_i u_i & -V_i u_i \\ 0 & 0 & 0 & U_i w_i & V_i w_i & W_i w_i & -U_i v_i & -V_i v_i \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

As in the earlier case, the unknown elements $\mathbf{l} = (L_{11} \dots L_{32})^\top$ can be recovered by least-squares.

B.1 Implementation and example results

The implementation of the line calibration is somewhat more involved than that for points. The calibration lines are “traced” actively by the camera. The driving process is a video-rate implementation in the fovea of the Canny edge detection and hysteresis linking algorithms [25]. For each image frame, after computing edgels and linking them together into strings, the edgel nearest the centre of the fovea is located and its parent string identified. The gaze controller then traverses the string in the current direction to find the n th-neighbour of the central edgel, and

sends its position as a demand to the servo-controller. Repeating this process as each image is received causes the camera to trace smoothly along an extended edge, even though no point correspondences are made between frames. Figure 7 shows alternate frames from a 25Hz sequence of edgemaps obtained as the camera traces around a string.

As line tracing proceeds, at each frame the gaze angles (θ_e, θ_v) obtained from the platform odometry are used to determine the intersection of the gaze direction with the frontal plane and the co-ordinates x are stored. When a point of high curvature is found, indicating the end of a line, the succession of stored x values are deemed to form a straight edge, a u is fitted. The fitted lines are stored and used for matching and calibration with the known world lines U to recover the homography.

Figure 8(a) shows the observer trajectories obtained in the frontal plane while pursuing a radio-controlled buggy around the lab floor. The frontal image, which is included here for illustration and plays no part in the analysis, also shows part of the white-lined pattern used for calibration. The 50° field of view in the frontal image is much less than the near 180° field of view accessible to the active camera. The wider the field of view, the more accurate the calibration. Part (b) shows the Cartesian reconstruction of the trajectories and calibration pattern.

II. TRAJECTORY FILTERING

We have developed filters embodying three simple canonical motions to describe the ground plane trajectories using the following scene models based on speed and direction, rather than individual speeds in the X and Y directions:

1. CSD Constant speed and direction,
 ie, zero acceleration and turning rate.
2. CS Constant speed
 ie, having possibly non-zero turning rate.
3. CD Constant direction
 ie, having possibly non-zero acceleration.

Each model is incorporated in an Extended Kalman Filter (EKF), and the set embedded within the framework of the Interacting Multiple Model (IMM) [19] which selects the most appropriate filter to represent the target dynamics at any time.

A. The individual filters

For convenience, each of the filters for constant speed and direction (CSD), constant speed with turning rate (CS) and constant direction with acceleration (CD) has the state vector

$$\mathbf{x} = \begin{bmatrix} X & Y & s & \theta & \dot{s} & \dot{\theta} \end{bmatrix}^\top.$$

At each frame we obtain an estimate of the target position, but also allow for independent information about velocity to be included, giving an observation vector of

$$\mathbf{z} = \begin{bmatrix} X & Y & \dot{X} & \dot{Y} \end{bmatrix}^\top.$$

The observations are related to the state vector by

$$\mathbf{z} = \mathbf{h}(\mathbf{x}_k) + \mathbf{d}_k$$

where \mathbf{h} is the vector-valued observation model and \mathbf{d}_k is an uncorrelated, zero-mean, Gaussian noise sequence defined by $E[\mathbf{d}_k] = \mathbf{0}$ and $E[\mathbf{d}_i \mathbf{d}_j^\top] = \delta_{ij} \mathbf{R}$. The visual processes and robot dynamics implicitly and explicitly involved in the selection of a position in the image and its transfer onto the ground plane make it almost certain that these assumptions are violated. However, as the absolute size of noise in the Cartesian tracks, eg Figure 5(e), is small compared with variations we wish to account for using the different filters, the competition between filters should not be unduly biased.

In an EKF the prediction of state and variance are:

$$\begin{aligned} \hat{\mathbf{x}}_{(k+1|k)} &= \mathbf{f}(\hat{\mathbf{x}}_{(k|k)}, \mathbf{u}_k) + \mathbf{e}_k \\ \mathbf{P}_{(k+1|k)} &= \nabla \mathbf{f} \mathbf{P}_{(k|k)} \nabla \mathbf{f}^\top + \mathbf{Q}_k \end{aligned}$$

where \mathbf{f} is the vector of non-linear state transition functions, \mathbf{u}_k is the control input which here is zero, and \mathbf{e}_k is the process noise, an uncorrelated zero-mean Gaussian noise sequence with expectations $E[\mathbf{e}_k] = \mathbf{0}$ and $E[\mathbf{e}_i \mathbf{e}_j^\top] = \delta_{ij} \mathbf{Q}$. The update of state and variance after a new datum arrives are

$$\begin{aligned} \hat{\mathbf{x}}_{(k+1|k+1)} &= \hat{\mathbf{x}}_{(k+1|k)} + \\ &\quad \mathbf{W}_{(k+1)} [\mathbf{z}_{(k+1)} - \mathbf{h}(\hat{\mathbf{x}}_{(k+1|k)})] \\ \mathbf{P}_{(k+1|k+1)} &= \mathbf{P}_{(k+1|k)} - \\ &\quad \mathbf{W}_{(k+1)} \mathbf{S}_{(k+1)} \mathbf{W}_{(k+1)}^\top \end{aligned}$$

where the gain matrix and innovation covariance are respectively

$$\begin{aligned}\mathbf{W} &= \mathbf{P}_{(k+1|k)} \nabla \mathbf{h}^\top \mathbf{S}_{(k+1)}^{-1} \\ \mathbf{S}_{(k+1)} &= \nabla \mathbf{h} \mathbf{P}_{(k+1|k)} \nabla \mathbf{h}^\top + \mathbf{R}_{(k+1)}\end{aligned}$$

In the following derivations the time-step for information input is T and in each case there is no control input \mathbf{u}_k . In the CS and CD filters we introduce process noise to allow the system to adjust to changing parameter values. It is observed that larger values of process noise variance reduce the response time of the system to changes in the system state, at the expense of introducing increased noise into the estimated state values. The process noise values also account for errors introduced by the linearization of the EKF.

A.1 All filters

First, for all the filters the vector of functions \mathbf{h} is

$$\mathbf{h} = (X, Y, s \sin \theta, s \cos \theta)^\top$$

so that

$$\nabla \mathbf{h} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & \sin \theta & -s \cos \theta & 0 & 0 \\ 0 & 0 & \sin \theta & s \cos \theta & 0 & 0 \end{bmatrix}.$$

A.2 The Constant Speed and Direction (CSD) filter.

Both s and θ are constant, so that the update equation over a time interval T is

$$\mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) = \begin{bmatrix} x_k + s_k \sin \theta_k T \\ y_k + s_k \cos \theta_k T \\ s_k \\ \theta_k \\ \dot{s}_k \\ \dot{\theta}_k \end{bmatrix}$$

and the associated Jacobian matrix is

$$\nabla \mathbf{f} = \begin{bmatrix} 1 & 0 & \sin\theta T & -s\sin\theta T & 0 & 0 \\ 0 & 1 & \sin\theta T & s\sin\theta T & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

For this filter, \dot{s}_k and $\dot{\theta}_k$ are not recovered and are included for computational uniformity in the IMM. Strictly, their values at any time step are noise terms drawn from zero mean Gauss-random sequences with variances σ_s^2 and σ_θ^2 :

$$\begin{aligned} \dot{s}_k &= e_s = N(0, \sigma_s^2) \\ \dot{\theta}_k &= e_\theta = N(0, \sigma_\theta^2) . \end{aligned}$$

To evaluate the process noise covariance, both e_s and e_θ be constant over a frame interval T , so that to first order the noise unmodelled by the update equation is

$$\mathbf{e}_k = \begin{pmatrix} (e_s \sin\theta_k - e_\theta s_k \sin\theta_k) T^2 / 2 \\ (e_s \sin\theta_k + e_\theta s_k \sin\theta_k) T^2 / 2 \\ e_s T \\ e_\theta T \\ e_s \\ e_\theta \end{pmatrix}.$$

The covariance $\mathbf{Q}_k = E[\mathbf{e}_k \mathbf{e}_k^\top]$ is evaluated using the expectation values $E[e_s^2] = \sigma_s^2$, $E[e_\theta^2] = \sigma_\theta^2$, and $E[e_s e_\theta] = 0$. Again we stress that for the CSD filter only the first four members of \mathbf{e} and the upper-left 4×4 sub-matrix of \mathbf{Q} are strictly required for the EKF.

A.3 The Constant Speed (CS) filter

The update equation is

$$\mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) = \begin{bmatrix} x_k + (s_k/\dot{\theta}_k)[\sin(\theta_k + \dot{\theta}_k T) - \sin\theta_k] \\ y_k + (s_k/\dot{\theta}_k)[\sin\theta_k - \sin(\theta_k + \dot{\theta}_k T)] \\ s_k \\ \theta_k + \dot{\theta}_k T \\ \dot{s}_k \\ \dot{\theta}_k \end{bmatrix}$$

from which the Jacobian is found by differentiation. Here, \dot{s} is not recovered but $\dot{\theta}_k$ is, so that we consider as unmodelled noise $e_{\dot{s}}$ and $e_{\dot{\theta}}$, with variances $\sigma_{\dot{s}}^2$ and $\sigma_{\dot{\theta}}^2$. In a similar manner to the CSD filter, we set these values to be constant over an update cycle to approximate the process noise vector \mathbf{e} and hence the covariance matrix.

A.4 The Constant Direction (CD) filter

Now θ is modelled as constant, but s may change, so that the update equation over a period T is

$$\mathbf{f}(\mathbf{x}_k, \mathbf{u}_k) = \begin{bmatrix} x_k + \sin\theta_k [s_k T + \dot{s}_k T^2/2] \\ y_k + \sin\theta_k [s_k T + \dot{s}_k T^2/2] \\ s_k + \dot{s}_k T \\ \theta_k \\ \dot{s}_k \\ \dot{\theta}_k \end{bmatrix}.$$

Again the Jacobian is found by differentiation and the process noise covariance matrix is found by considering noises $e_{\dot{s}}$ and $e_{\dot{\theta}}$ with variances $\sigma_{\dot{s}}^2$ and $\sigma_{\dot{\theta}}^2$.

B. Tests of individual filters

The individual filters were first tuned by applying them to simulated noisy data generated from models each was designed to filter. An example of tuning the process noise in the CS filter is shown in Figure 9. The synthesized motion has periods consisting of 50 observations of constant speed, zero turn rate motion (ie, straight line, no acceleration) followed by periods,

again of 50 observations, of constant speed, fixed turn rate (45°s^{-1}). The target therefore moves in a square with rounded corners. Figure 9(a), shows that with a process noise value of $\sigma_\theta^2 = 1 \text{ rad}^2.\text{s}^{-4}$, there is considerable overshoot in the estimated target direction but the recovered turning rate estimates are smooth. Increasing the process noise to $\sigma_\theta^2 = 100 \text{ rad}^2.\text{s}^{-4}$ provides much more rapid response to manoeuvres at the expense of increased noise in the recovered turning rate estimates. The chosen value of $30 \text{ rad}^2.\text{s}^{-4}$ was a compromise between promptness and smoothness.

Figure 10 shows the results of applying each filter *individually* to part of the ground plane trajectory data from a walking person recovered by active tracking. The person exited the Computing Centre, walked straight along the path, and then made a sharp left turn to walk along the pavement. We used measurement variance values of $\sigma_X^2 = \sigma_Y^2 = 0.1\text{m}^2$, $\sigma_{\dot{X}}^2 = \sigma_{\dot{Y}}^2 = 1.5\text{m}^2\text{s}^{-2}$. The initial covariances were taken as ten times the measurement noise values. Again we observe poor performance from the CSD and CD filters as the target manoeuvres, but the CS filter performs well as the speed remains constant throughout the turn. These results were obtained with tuned values for process noises of $\sigma_s^2 = 0.1\text{m}^2.\text{s}^{-4}$ and $\sigma_\theta^2 = 1 \text{ rad}^2.\text{s}^{-4}$.

These examples serve merely to highlight the obvious dilemma faced in the choice of *one* filter. Using a strong filter, one has a firm basis for prediction but one loses responsivity to manoeuvres. A weaker filter will work well all the time, but provides necessarily an uncertain basis for prediction.

This dilemma arises when observing any system which has dynamics which change abruptly and unexpectedly. Not surprisingly then, multiple-model-filtering has received considerable attention in work on radar tracking of aircraft. We have explored the Interacting Multiple Model (IMM) scheme for filter management, which from the literature appears to provide good state estimation in the presence of switching target dynamics with, importantly, a computational cost compatible with real-time operation [26], [27], [28], [29], [30], [19], [31].

III. THE INTERACTING MULTIPLE MODEL

The IMM algorithm operates with a set of N Kalman filters, each with a differing implicit motion model for the system dynamics. At each observation time-step, the outputs of the filters are first mixed according to a switching matrix \mathbf{J} generated randomly from a distribution. The various filters are updated using the normal Kalman predict-update step. The likelihood of each

filter is then evaluated, and sum of the state estimates of the filters, weighted by the respective filter likelihoods, formed as the output state of the system.

At the end of timestep k , let the i -th filter, $i = 1, \dots, N$, have state estimate $\hat{\mathbf{x}}_k(i)$, covariance $\mathbf{P}_k(i)$ and likelihood $p_k(i)$. The four steps in the IMM algorithm are:

1. For each filter model i , update the filter likelihoods to those appropriate to the start of the next step

$$p_{k+1}^*(i) = \sum_j J_k(ij) p_k(j) \quad .$$

Then mix the estimates of the state and covariance of each filter

$$\begin{aligned} \hat{\mathbf{x}}_k^*(i) &= \frac{1}{p_{k+1}^*(i)} \sum_j J_k(ij) p_k(j) \hat{\mathbf{x}}_k(j) \\ \mathbf{P}_k^*(i) &= \frac{1}{p_{k+1}^*(i)} \sum_j J_k(ij) p_k(j) \left(\mathbf{P}_j(j) \right. \\ &\quad \left. + [\hat{\mathbf{x}}_k(j) - \hat{\mathbf{x}}_k(i)][\hat{\mathbf{x}}_k(j) - \hat{\mathbf{x}}_k(i)]^\top \right) \quad . \end{aligned}$$

2. The values of $\hat{\mathbf{x}}_k^*(i)$ and $\mathbf{P}_k^*(i)$ are used in the usual way in the EKF cycle; the prediction phase gives $\hat{\mathbf{x}}_{(k+1|k)}(i)$ and $\mathbf{P}_{(k+1|k)}(i)$ and the update phase gives $\hat{\mathbf{x}}_{k+1}(i)$ and $\mathbf{P}_{k+1}(i)$.

3. The likelihoods of being in a filter state are updated as

$$p_{k+1}(i) = \kappa \frac{p_{k+1}^*(i) \exp[-\frac{1}{2} \mathbf{v}_{k+1}^\top(i) \mathbf{S}_{k+1}^{-1}(i) \mathbf{v}_{k+1}(i)]}{\|\mathbf{S}_{k+1}(i)\|^{1/2}}$$

where κ provides normalization $\sum_i p_{k+1}(i) = 1$. In this expression, the innovation vector is

$$\mathbf{v}_{k+1}(i) = \mathbf{z}_{k+1} - \mathbf{h}(\hat{\mathbf{x}}_{(k+1|k)}(i))$$

and, as earlier,

$$\mathbf{S}_{k+1}(i) = \nabla \mathbf{h} \mathbf{P}_{k+1|k}(i) \nabla \mathbf{h}^\top + \mathbf{R}_{k+1} \quad .$$

4. The useable output state and covariance $\hat{\mathbf{x}}_{k+1}$ and \mathbf{P}_{k+1} are generated by taking a linear sum of the updated state and covariance estimates for each filter weighted by the updated filter likelihoods:

$$\begin{aligned} \hat{\mathbf{x}}_{k+1} &= \sum_i p_{k+1}(i) \hat{\mathbf{x}}_{k+1}(i) \\ \mathbf{P}_{k+1} &= \sum_i p_{k+1}(i) \left(\mathbf{P}_{k+1}(i) \right. \\ &\quad \left. + [\hat{\mathbf{x}}_{k+1}(i) - \hat{\mathbf{x}}_{k+1}][\hat{\mathbf{x}}_{k+1}(i) - \hat{\mathbf{x}}_{k+1}]^\top \right) \end{aligned}$$

A. Multiple Model Tests

We now illustrate the benefits of the IMM approach, showing improved long-term filter response to motion trajectories, with the IMM detecting changes in the target motion state and selecting the appropriate filter accordingly.

The input mixing stage of the IMM performs filter mixing by resetting the state vector of each filter at each time-step to a linear combination of the output state vectors at the previous time-step, weighted by the filter likelihoods at the previous time-step and the model switching probabilities. The long-term performance of all filters is maintained as the filters with incorrect models at any time do not deviate significantly from the correct state vector, since mixing incorporates a large fraction of the correct state vector via the large filter likelihood $p_k(i)$ for the correct model i at step k .

The model switching probabilities are coded as \mathbf{J}_k , where $J_k(ij)$ represents the likelihood of switching from model i to model j at time-step k . As our expected targets are likely to continue with a single motion model for extended periods with only occasional model switching, we set the diagonal elements of the matrix as $J_k(mm) = J_d \approx 1$ and have small off-diagonal elements $J_k(mn) = (1 - J_d)/(N - 1)$, where $N = 3$ in our case. The filter likelihoods are initialized as $p_0(\text{CSD}) = p_0(\text{CS}) = p_0(\text{CD}) = 1/3$. Figure 11 illustrates the effect upon the model probabilities for the square trajectory of varying the switching probability from $J_d = 0.93$ to $J_d = 0.99$. The left hand figures (a1,b1) illustrate the model probability sequence for the trajectory, the right hand figures (a2,b2) the estimated turning rate values. In each figure, for the first 50 or so observations the filters exhibit similar performance to the above example of straight line motion. As the manoeuvre begins the likelihood of the CS filter rises and the CSD and CD filter likelihoods decrease. Around step 60 the CS filter becomes the most likely. During the manoeuvre the filter probabilities remain roughly constant, and when the manoeuvre completes the filter likelihoods adjust until the CSD filter again becomes most likely. The pattern continues as the state sequence repeats a further three times.

It can be seen that increasing J_d gives better discrimination between models as the filter probabilities are more widely separated. This is as we would expect, since increasing p reduces the amount of mixing between models, and the correct filter will obtain a more accurate estimate of the target state vector and attain a higher probability. However, increasing J_d reduces the

sharpness of response of the filter, as we see from the left hand graphs. The lag of the estimated turning rate is increased, since increasing J_d reduces the amount of mixing between filters, and the time taken after model switching for the new filter to affect the parameter estimates is thus increased. Empirically, a value of $J_d = 0.96$ proved a satisfactory compromise.

Figure 12 shows results for the square trajectory when $J_d = 0.96$. In the CSD state the mean values of \dot{s} and $\dot{\theta}$ are small, confirming that the target is not manoeuvring, whilst in the CS state we see mean values of $\dot{\theta} \approx 0.6 \text{ rad s}^{-1}$, the actual value being 0.785 rad s^{-1} . It can be seen from Table I that the IMM detects changes in the motion model reliably, the correct filter being selected within ± 10 observations. This error is not so important as far as the filtering required for improved target pursuit is concerned, as the position and velocity estimates which generate controller demands are consistent throughout the switching periods because of mixing. The important point for a higher level motion classification scheme is a reliable estimate of the state duration between switches and the mean acceleration and turning rate values.

<i>Most likely filter</i>	<i>Duration in observations</i>		<i>Averaged Filter likelihoods</i>			<i>Averaged \dot{s} $\dot{\theta}$</i>	
	<i>Actual</i>	<i>Estimated</i>	\bar{p}_{CSD}	\bar{p}_{CS}	\bar{p}_{CD}	(m s^{-2})	(rad.s^{-1})
CSD	50	61	0.446	0.350	0.204	-0.006	0.025
CS	50	54	0.230	0.632	0.139	-0.002	0.616
CSD	50	46	0.441	0.360	0.199	-0.010	0.032
CS	50	57	0.201	0.667	0.132	-0.014	0.476
CSD	50	42	0.417	0.376	0.207	-0.008	0.046
CS	50	56	0.204	0.648	0.148	0.002	0.548
CSD	50	43	0.411	0.354	0.235	-0.008	0.023
CS	50	41	0.204	0.648	0.148	0.002	0.548

TABLE I

ACTUAL AND ESTIMATED STATE DURATIONS, STATE ESTIMATES AND MEAN MODEL LIKELIHOODS FOR THE “ROUNDED SQUARE” TRAJECTORY.

Figure 13 illustrates IMM results for the tracking of a walking person following the trajectory given in Figure 10(a) (also one of those shown in Figure 3(a)). The filter probabilities illustrate a period of straight line motion with occasional periods of slight acceleration, followed by a turning manoeuvre where the CS filter has highest likelihood, followed by further straightline motion. Table II illustrates mean model probabilities, acceleration and turning rate for each

motion period.

<i>Most likely filter</i>	<i>Estimated Duration (frames)</i>	<i>Averaged likelihoods</i>			<i>Average values</i>	
		\bar{p}_{CSD}	\bar{p}_{CS}	\bar{p}_{CD}	$\bar{\dot{s}}$ (m s ⁻²)	$\bar{\dot{\theta}}$ (rad s ⁻¹)
CSD	11	0.358	0.310	0.333	0.002	-0.001
CD	2	0.367	0.248	0.385	1.196	0.001
CSD	1	0.400	0.243	0.357	0.006	0.001
CD	1	0.383	0.229	0.388	0.957	0.001
CSD	77	0.636	0.182	0.182	0.000	0.000
CS	72	0.259	0.605	0.136	0.000	0.441
CSD	170	0.547	0.293	0.159	0.000	0.000

TABLE II

ESTIMATED STATE DURATIONS, AVERAGED FILTER LIKELIHOODS AND STATE ESTIMATES FOR THE WALKING PERSON.

The above tests have shown the advantages afforded by embedding filters within the IMM, namely much improved long term performance when trajectories incorporate periods of model switching. We now describe real-time implementation of the filter scheme and illustrate improved tracking performance of the system when filtering is incorporated.

IV. REAL-TIME FILTER IMPLEMENTATION

The IMM filter has been implemented in real-time to evaluate its potential for driving tracking from 3D rather 2D prediction. The first step of the filter loop is to evaluate the current state of the target predicted by each of the incorporated filters. Motion estimates received by the high-level gaze controller from the foveal motion process are then transformed to the world co-ordinate frame and the individual filter states updated according to the normal Kalman update step. The updated filter likelihoods are evaluated and the mixed output state vector calculated. The mixed state estimate of position and velocity are then converted to an angular position and velocity demands which are sent to the servo controller to drive platform motion.

The time taken to perform the prediction and update steps for each of the individual filters is 26 ms, of which 19 ms is a common overhead due to communication between processes and the storage and relay of results at each frame. Without mixing the total time taken to run predict and

update steps for the three CSD, CS and CD filters is 39 ms, barely under the 40 ms maximum permitted for operation at 25 Hz. With mixing incorporated into the filter scheme the time taken rises to 53 ms, and we are no longer able to run the complete filter cycle with mixing at 25 Hz. To overcome this limitation we run the prediction step for each filter at each frame, but only update the filter states estimate with every second frame of received data, allowing the system to maintain 25 Hz operation. Although such timings are processor dependent, this example illustrates that filter framework allows prediction not only over delays arising from latency, but also over delays due to missing data. Although direct timing comparisons between the IMM's three filters and a single image-based filter have not been performed, the above figures indicate that the computing time is increased less than four-fold.

Figure 14 shows results from the IMM while the head platform tracks a toy engine which follows an oval made of two straights and two semi-circles (as recovered in Figure 15). The filter probabilities show that around the corners the CS filter (filter 2) is selected and on straights periods of CSD and CD are mixed as the engine accelerates and decelerates. The periods of acceleration and four periods of turn (corresponding to two revolutions) are apparent in Figures 14(c) and (d). The total period of revolution was some 14 seconds, with some 4 seconds spent on the two sections of straight of total length 1.8 m, and 10 seconds spent on the circular sections of total length 3.3 m. The recovered value of $\dot{\theta} \sim 0.6 \text{ rad.s}^{-1}$ is close to the expected value of $2\pi/10$. The period of acceleration indicated in (c) occur along the straights, and indicate a change in speed from 0.4 ms^{-1} to 0.5 ms^{-1} along each straight section.

One of our key aims in this work was to demonstrate that prediction in 3D is more effective than that performed in 2D. Figure 15 illustrates the improvement when tracking the toy engine. The left sides of Figures 15(b) and (c) shows 2D tracking, where saccades have had to be resorted to in order to recapture the target when 2D pursuit fails. The right sides show the results when IMM filtering and 3D prediction is used. The system maintains pursuit without the use of corrective saccades.

V. CONCLUSIONS AND DISCUSSION

This paper has demonstrated first the tracking of targets moving in a plane in the scene using an active monocular camera and the recovery of “observer trajectories”, the intersection of the instantaneous direction of gaze with the frontal plane. By calibrating the homography between

the frontal plane and scene plane using plane to plane correspondences of at least four points or at least four lines, it was possible to reconstruct the scene trajectory of the target during tracking.

The paper then described the implementation of three Kalman filters based on simple canonical motions to filter the trajectory, and their embedding in an Interacting Multiple Model which allowed the most appropriate filter to describe the motion. An improvement in tracking performance was demonstrated in real time using the 3D trajectory and filter to provide prediction over the latency in the visual feedback loop, rather than using filtered 2D observer trajectories. The IMM was able to provide strong filtering and good prediction together with a robustness not available when using a single filter. In principle, there is no reason why the IMM method should not be extended to discriminate between fully 3D motions (recovered, say, from stereo tracking) rather than planar ones, although the addition of further state variables must increase uncertainty and hence make the distinction between filters more difficult within the IMM. In situations where the motion is not well described by any of the individual filters, the IMM is at best neutral, that is the latest datum is believed in preference to the filtered state. To what extent increasing the repertoire of simple motions and hence the number of filters is useful is open to question, but the early caveat about the competition between filters becoming more difficult must surely apply.

Although we have emphasized the use of filtering in 3D for prediction, an important issue in an active vision system is to what in the scene the system should devote its visual resources or attention. From this perspective, the detection of manoeuvres, or the lack of them, is more relevant than prediction. An interesting target, and one which demands more resources, is one which manoeuvres, whereas a boring one is one which does not. Our empirical evidence suggests that the Interacting Multiple Model provides not only strong filtering providing good prediction, but also a reliable — though inevitably delayed — indication of manoeuvre as the most likely model changes. The probability of changing model from moment to moment, which can be analyzed as a Markov chain [32], provides a characteristic motion signature for a target.

ACKNOWLEDGEMENT

This work was supported by EPSRC Grants GR/G30003 and GR/J65372, by an EPSRC Research Studentship to KJB and by the Violette and Samuel and Glasstone Fellowship in Science from the University of Oxford to IDR.

REFERENCES

- [1] P. M. Sharkey, D. W. Murray, S. Vandevelde, I. D. Reid, and P. F. McLauchlan, "A modular head/eye platform for real-time reactive vision," *Mechatronics*, vol. 3, no. 4, pp. 517–535, 1993.
- [2] D. W. Murray, P. F. McLauchlan, I. D. Reid, and P. M. Sharkey, "Reactions to peripheral image motion using a head/eye platform," *Proc. 4th Int. Conf. on Computer Vision, Berlin, 1993*, pp. 403–411, Los Alamitos, CA: IEEE Computer Society Press, 1993.
- [3] D. W. Murray, K. J. Bradshaw, P. F. McLauchlan, I. D. Reid, and P. M. Sharkey, "Driving saccade to pursuit using image motion," *International Journal of Computer Vision*, vol. 16, no. 3, pp. 205–228, 1995.
- [4] I. D. Reid and D. W. Murray, "Tracking foveated corner clusters using affine structure", *Proc. 4th Int. Conf. on Computer Vision, Berlin, 1993*, pp. 76–83, Los Alamitos, CA: IEEE Computer Society Press, 1993.
- [5] I. D. Reid and D. W. Murray, "Active tracking of foveated feature clusters using affine structure," *International Journal of Computer Vision*, vol. 18, no. 1, pp. 1–20, 1996.
- [6] K. J. Bradshaw, P. F. McLauchlan, I. D. Reid, and D. W. Murray, "Saccade and pursuit on an active head/eye platform," *Image and Vision Computing*, vol. 12, no. 3, pp. 155–163, 1994.
- [7] C. M. Brown, "Gaze control with interactions and delays," *IEEE Trans. Systems, Man and Cybernetics*, vol. 63, pp. 61–70, 1990.
- [8] C. M. Brown, "Prediction and cooperation in gaze control," *Biological Cybernetics*, vol. 63, pp. 61–70, 1990.
- [9] J. J. Clark and N. J. Ferrier, "Modal control of an attentive vision system," *Proc. 2nd Int. Conf. on Computer Vision, Tampa FL, 1988*, pp. 514–523, IEEE Computer Society Press.
- [10] P. M. Sharkey and D. W. Murray, "Coping with delays for real-time gaze control," *Proc. SPIE Sensor Fusion VI, Boston MA, September 1993*
- [11] R. H. S. Carpenter, *Movements of the Eyes*. London: Pion, 1988.
- [12] S. M. Fairley, I. D. Reid, and D. W. Murray, "Transfer of fixation for an active stereo platform via affine structure recovery," *Proc. 5th Int. Conf. on Computer Vision, Cambridge MA, 1995*, pp. 1100–1105, 1995. Los Alamitos, CA: IEEE Computer Society Press, 1995.
- [13] T. N. Tan, G. D. Sullivan, and K. D. Baker, "Structure from motion using ground plane constraint", *Proc. 2nd European Conf. on Computer Vision, Santa Margherita, Italy, 1992*, pp. 277–281, Heidelberg: Springer-Verlag.
- [14] G. D. Sullivan, "Visual interpretation of known objects in constrained scenes", *Phil. Trans. R. Soc. Lond. B.*, vol. B337, pp. 109–118, 1992.
- [15] M. Mohnhaupt and B. Neumann, "Understanding object motion: Recognition, learning and spatiotemporal reasoning", *Robotics and Autonomous Systems*, vol. 8, pp. 65–91, 1991.
- [16] D. Koller, K. Danilidis, T. Thørhallson, and H.-H. Nagel, "Model-based object tracking in traffic scenes", *Proc. 2nd European Conf. on Computer Vision, Santa Margherita, Italy, 1992*, pp. 437–452, Heidelberg: Springer-Verlag.
- [17] D. Koller, D. Danilidis, and H.-H. Nagel, "Model-based object tracking in monocular image sequences of road traffic scenes", *International Journal of Computer Vision*, vol. 10, no. 3, pp. 257–281, 1993.
- [18] S. S. Intille and A. F. Bobick. "Closed-world tracking", *Proc. 5th Int. Conf. on Computer Vision, Cambridge MA, 1995*, pp. 672–678, Los Alamitos, CA: IEEE Computer Society Press, 1995.
- [19] H. A. K. Blom, "An efficient filter for abruptly changing systems", *Proc. 23rd IEEE Conf. Decision Control, 1984*, pp. 656–658.
- [20] K. Pahlavan and J.-O. Eklundh, "A Head-Eye System - Analysis and Design," *CVGIP: Image Understanding*, vol. 56, no. 1, pp. 41–56, 1992.

- [21] K. Pahlavan, T. Uhlin, and J.-O. Eklundh, "Dynamic fixation," *Proc. 4th Int. Conf. on Computer Vision, Berlin, 1993*, pp. 412–419, Los Alamitos CA: IEEE Computer Society Press, 1993.
- [22] A. Blake, R. Curwen, and A. Zisserman, "A framework for spatiotemporal control in the tracking of visual contours," *International Journal of Computer Vision*, vol. 11, no. 2, pp. 127–146, 1993.
- [23] B. K. P. Horn and B. G. Schunck, "Determining Optical Flow," *Artificial Intelligence*, vol 17, pp. 185–203, 1981.
- [24] J. L. Mundy and A. P. Zisserman, editors, *Geometric Invariance in Computer Vision*. Cambridge MA: MIT Press, 1992.
- [25] J. Canny, "A Computational Approach to Edge Detection", *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 8, no. 6, pp. 679–698, 1986.
- [26] G. A. Watson and W. D. Blair, "IMM algorithm for tracking targets that maneuver through coordinated turns", *SPIE Proc. Signal and Data Processing of Small Targets*, SPIE vol. 1698, pp. 236–247, 1992.
- [27] G. A. Watson and W. D. Blair, "The IMM algorithm and aperiodic data," *SPIE Proc. Acquisition, Tracking and Pointing VI*, SPIE vol. 1697, pp. 83–91, 1992.
- [28] Y. Bar-Shalom, K. C. Chang, and H. A. Blom, "Tracking a Maneuvering Target Using Input Estimation Versus the Interacting Multiple Model Algorithm," *IEEE Trans. Aerospace and Electronic Systems*, vol. 25, no. 2, pp. 296–300, 1989.
- [29] C-B. Chang and J. A. Tabaczynski, "Application of State Estimation to Target Tracking," *IEEE Trans. Automatic Control*, vol. 29, no. 2, pp 98–109, 1984.
- [30] P. Andersson, "Adaptive Forgetting in recursive identification through multiple models," *International Journal of Control*, vol. 42, no. 5, pp. 1175–1193, 1985.
- [31] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, San Mateo CA: Morgan Kauffman, 1988.
- [32] K. J. Bradshaw, *Active Visual Surveillance*, DPhil Thesis, University of Oxford, 1995.

Kevin Bradshaw graduated with 1st class honours in Physics in 1991 from Heriot-Watt University, Edinburgh and received a DPhil in Engineering Science in 1995 from the University of Oxford. His doctoral work was concerned with the use of a closed-loop active vision system to perform automated tracking within dynamic scenes, part of which is summarized in this paper. Since 1995 he has worked with Smith System Engineering, a UK-based science and technology consultancy, as a consultant in the Environment and Science Policy sector.

Ian Reid received the BSc with 1st class honours in Computer Science from the University of Western Australia in 1987 and was awarded the Western Australian Rhodes Scholarship for that year. He completed a D.Phil. in Engineering Science at the University of Oxford in 1991 on the subject of recognising parametric models in range data. Since then he has been a postdoctoral researcher in Engineering Science investigating the theory and practice of uncalibrated active vision for tracking and navigation and has published over thirty papers on these and other topics. He currently holds a Junior Research Fellowship (The Queen's College) and a Violette and Samuel Glasstone Research Fellowship in Science at the University of Oxford.

David Murray received a BA with 1st class honours in Physics in 1977 and a DPhil in Physics in 1980, both from the University of Oxford. After a periods as a Research Fellow at California Institute of Technology and as a staff scientist at the GEC Research Laboratories in London he returned to Oxford in 1989 as a University Lecturer in the Department of Engineering Science. Since 1983 his research interests have centred on the computation of image motion and structure from motion, more recently applying these to active and teleoperated visual systems. He has published some ninety papers in physics and machine vision, and co-authored with Bernard Buxton a book (Experiments in the Computation of Visual Motion, MIT Press, 1990). He holds a Tutorial Fellowship in Engineering Science at St Anne's College, Oxford.

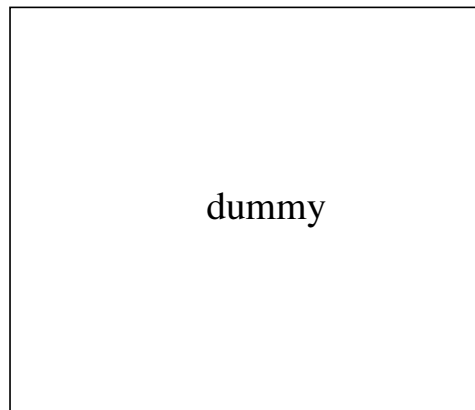


Fig. 1. The architecture of the visuo-control loop in our active camera system. Images from the platform cameras are processed by a number of parallel low-level vision algorithms. Vision results are sent to the high-level gaze controller which selects the visual feedback appropriate to the current task and generates a synchronous angular demand to the servo-controller.

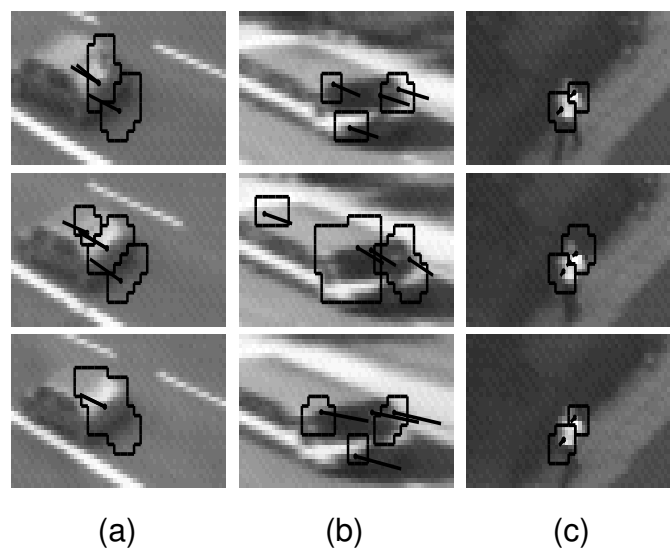


Fig. 2. Optical flow detected, segmented and fitted to the foveal imagery while tracking (a) a car, (b) a bus and (c) a person. The number of independent motion regions detected at each frame can change, but that in this application all are assumed to be from the same object.

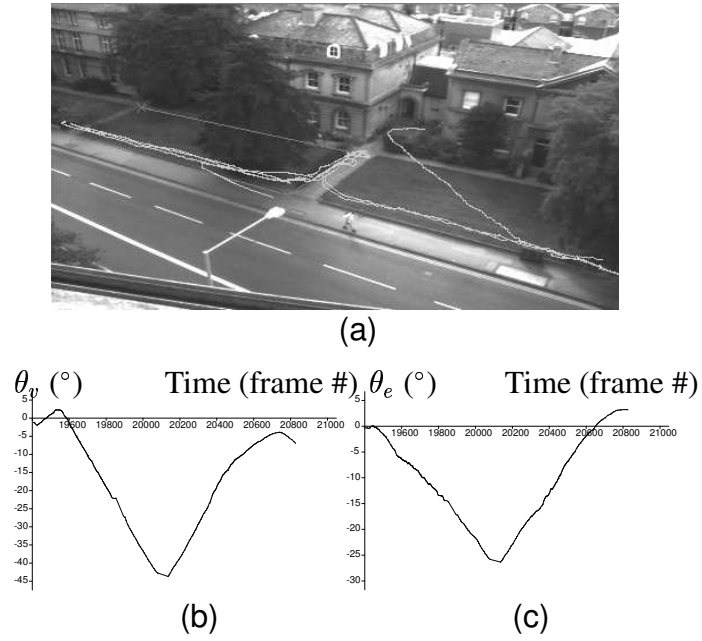


Fig. 3. Observer-based trajectories: (a) overlaid on a rest-frame image; and (b) and (c) expressed for one trajectory as vergence and elevation joint angles θ_v, θ_e (degrees) versus time (in frame units of 40ms).

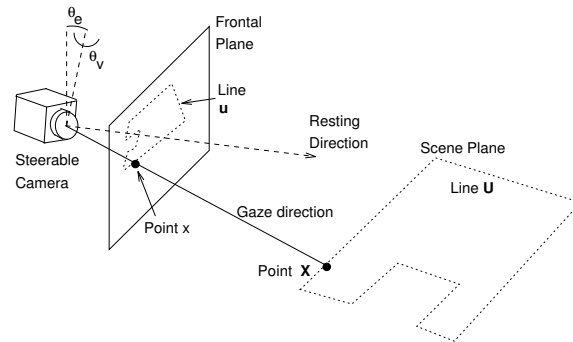
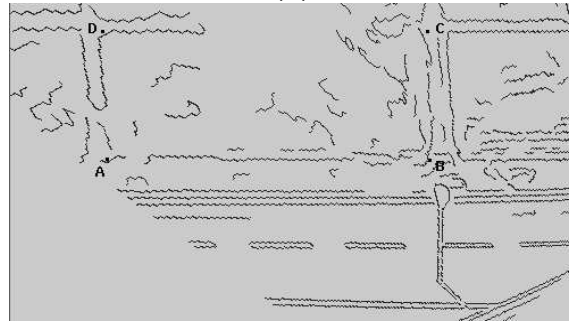


Fig. 4. Scene elements are projected into the frontal plane.



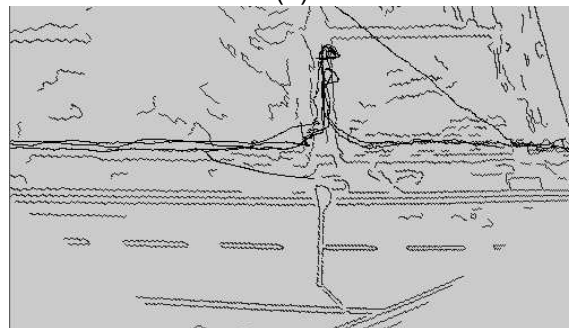
(a)



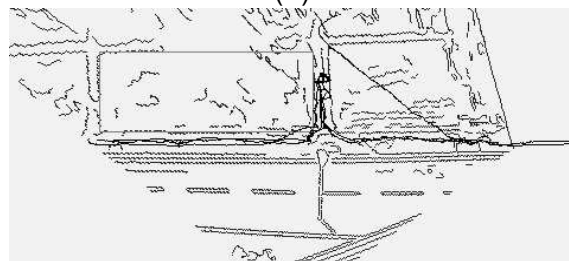
(b)



(c)



(d)



(e)

June 16, 2003

DRAFT

Fig. 5. Calibration of the outdoor scene: (a) calibration points in rest-frame image; (b) image features lines projected onto the ground plane; (c) trajectories in the observer space; (d) trajectories converted to the ground plane; (e) corrected ground plane trajectories.

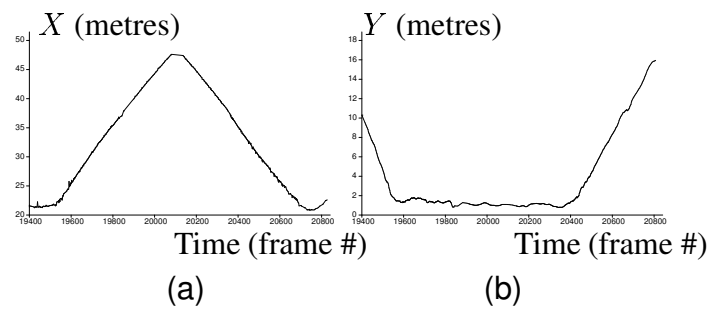


Fig. 6. The ground plane positions X and Y of the tracked person over some 56 seconds. Each frame lasts 40 ms.

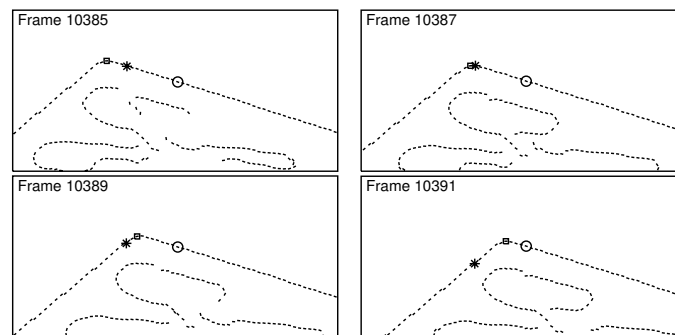


Fig. 7. Alternate frames from a 25 Hz sequence of edgemaps obtained as the camera traces around around a string. The near-centre string is marked with a circle, the n th-neighbour edgel which is creating the positional demand is marked with a star. A point of high curvature on the traced string is marked with a square.

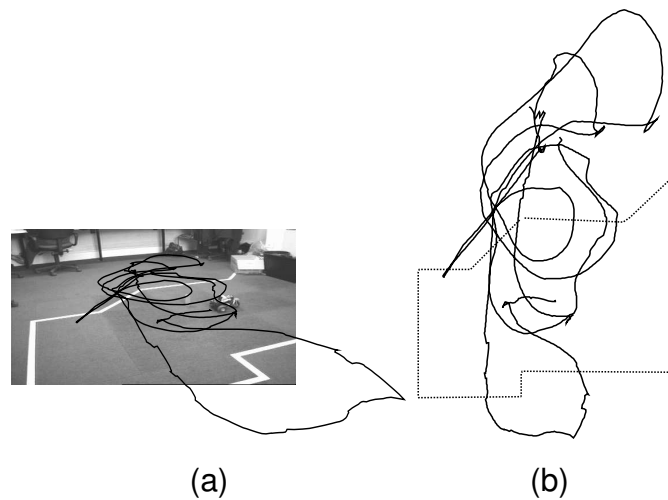


Fig. 8. (a) Buggy and observer track (b) Ground plane trajectory.

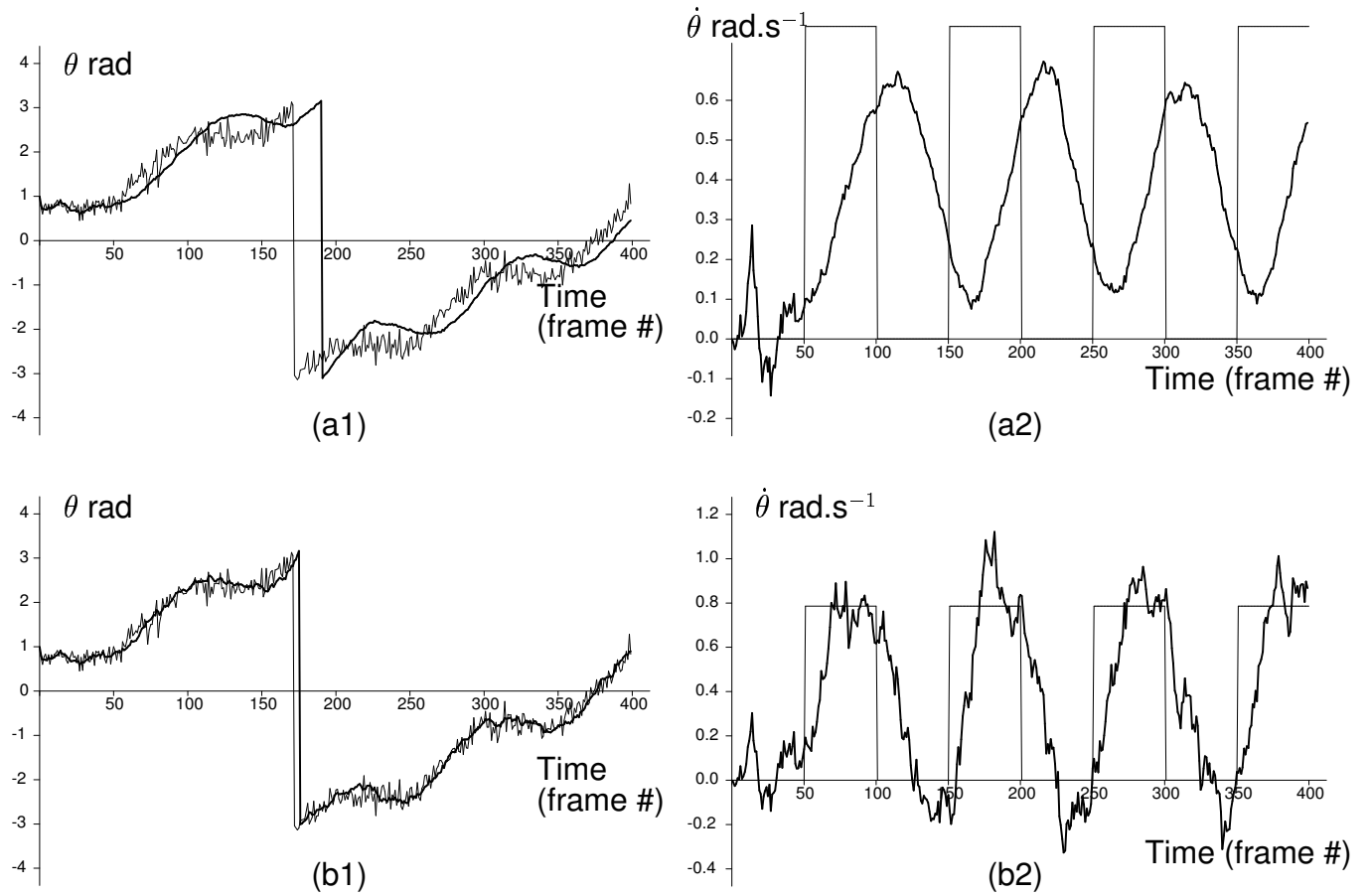


Fig. 9. The effects of increasing the process noise on estimated direction θ and turning rate $\dot{\theta}$ parameters for the CS filter. (a) shows the temporal response when $\sigma_{\theta}^2 = 1 \text{ rad}^2 \cdot \text{s}^{-4}$ and (b) shows the response when $\sigma_{\theta}^2 = 100 \text{ rad}^2 \cdot \text{s}^{-4}$. With higher process noise, the response is prompter, at the expense of increased noise in the turning rate estimates. The thinner lines represent data, the thicker the filter output. The discontinuity in direction at $\pm\pi$ is merely one of winding, and does not affect the filter performance.

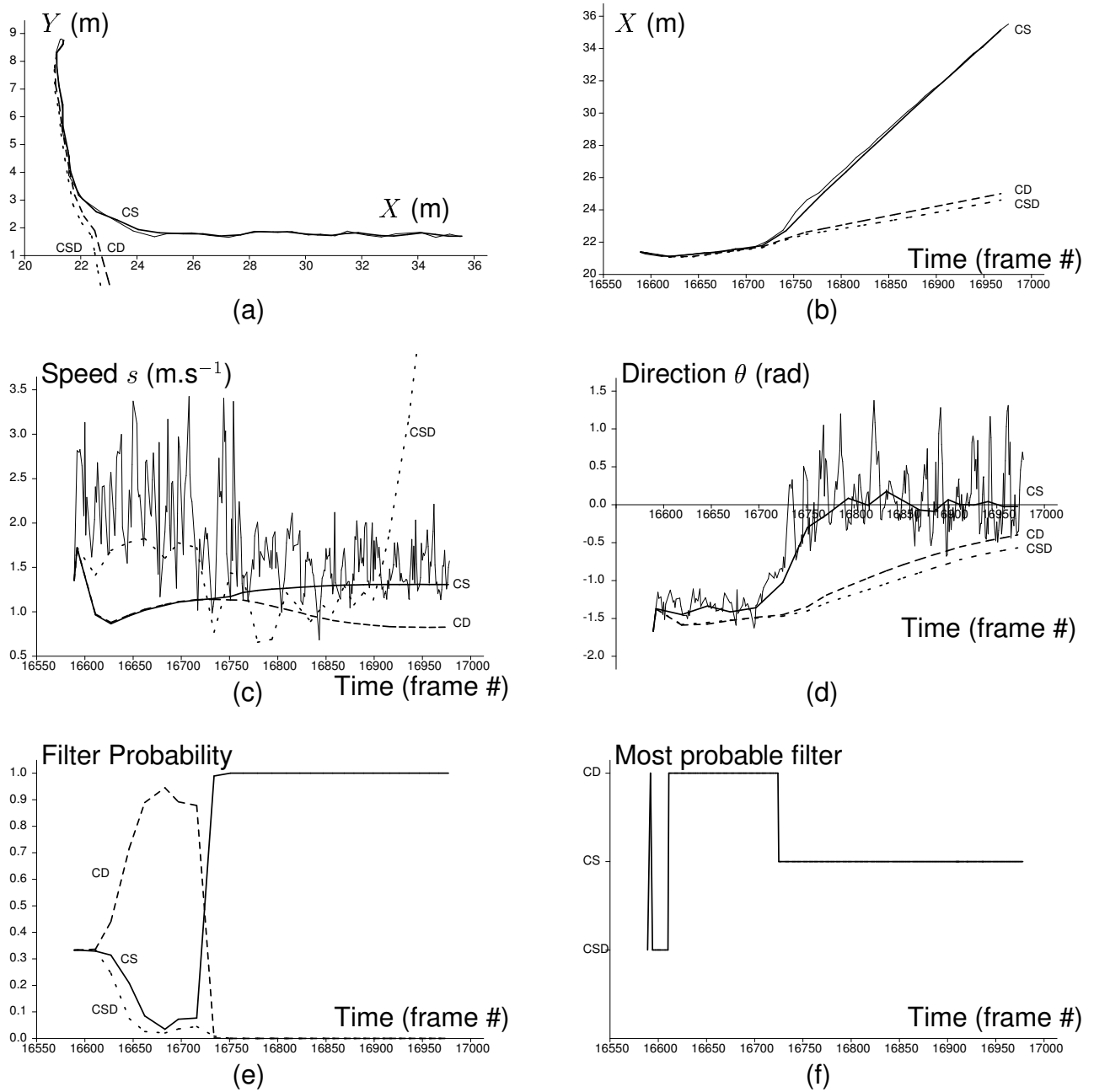


Fig. 10. Results for individual (ie, non-interacting) EKF filters for a walking person trajectory. (a) Shows the Cartesian trajectory (in metres) and (b)-(f) the filters' performance over time measured in frames of 40ms duration. In (a)-(e) the thin solid line records the measurements, and the thick solid, dashed and dotted lines are outputs from the CS, CD and CSD filters respectively.

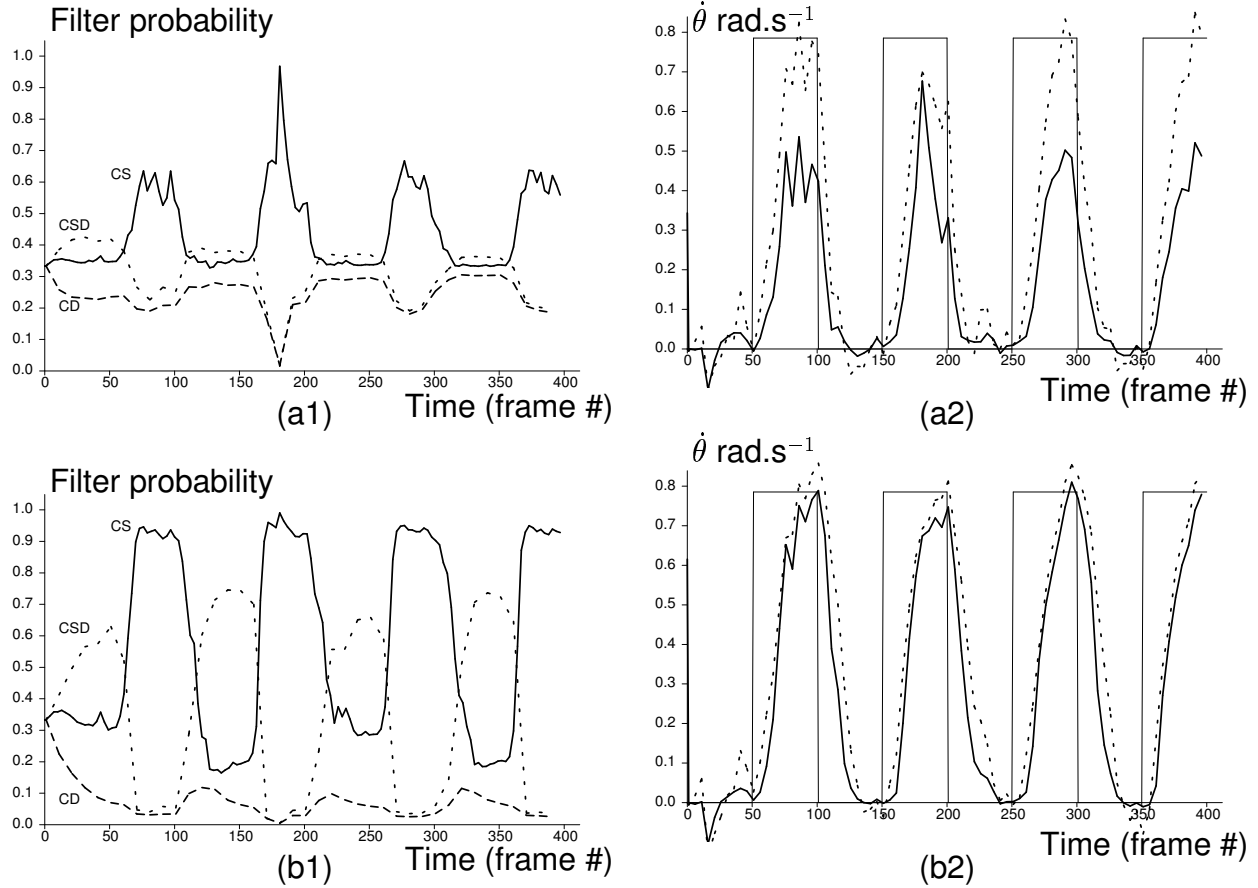


Fig. 11. The effect of the switching parameter on recovered turning rate and filter probability. (a) shows the response when $J_d = 0.93$ and (b) shows that when $J_d = 0.99$. The discrimination between filter models is increased with a higher J_d , but response to change is lessened. The thick solid, dashed and dotted lines are outputs from the CS, CD and CSD filters respectively. The thin line in the turning rate diagrams is the veridical value.

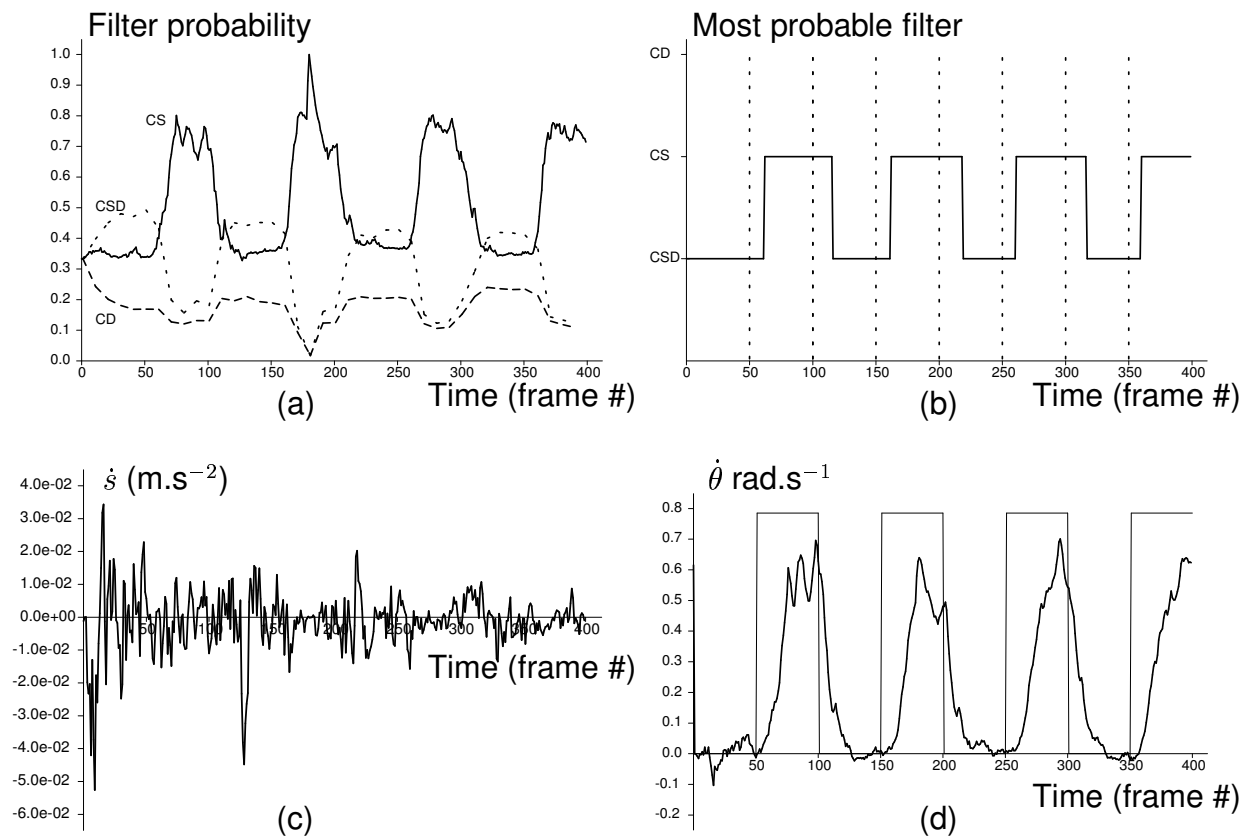


Fig. 12. IMM output for "rounded-square" trajectory data: (a) the filter likelihoods (where solid, dashed and dotted lines are outputs from the CS, CD and CSD filters) (b) the most likely filter, where the dashed lines show when the actual manoeuvres took place; (c) and (d), the overall IMM estimates of acceleration \dot{s} in m.s^{-2} and turn rate $\dot{\theta}$ in rad.s^{-1} . The thin line in (d) indicates the true value of $\dot{\theta}$.

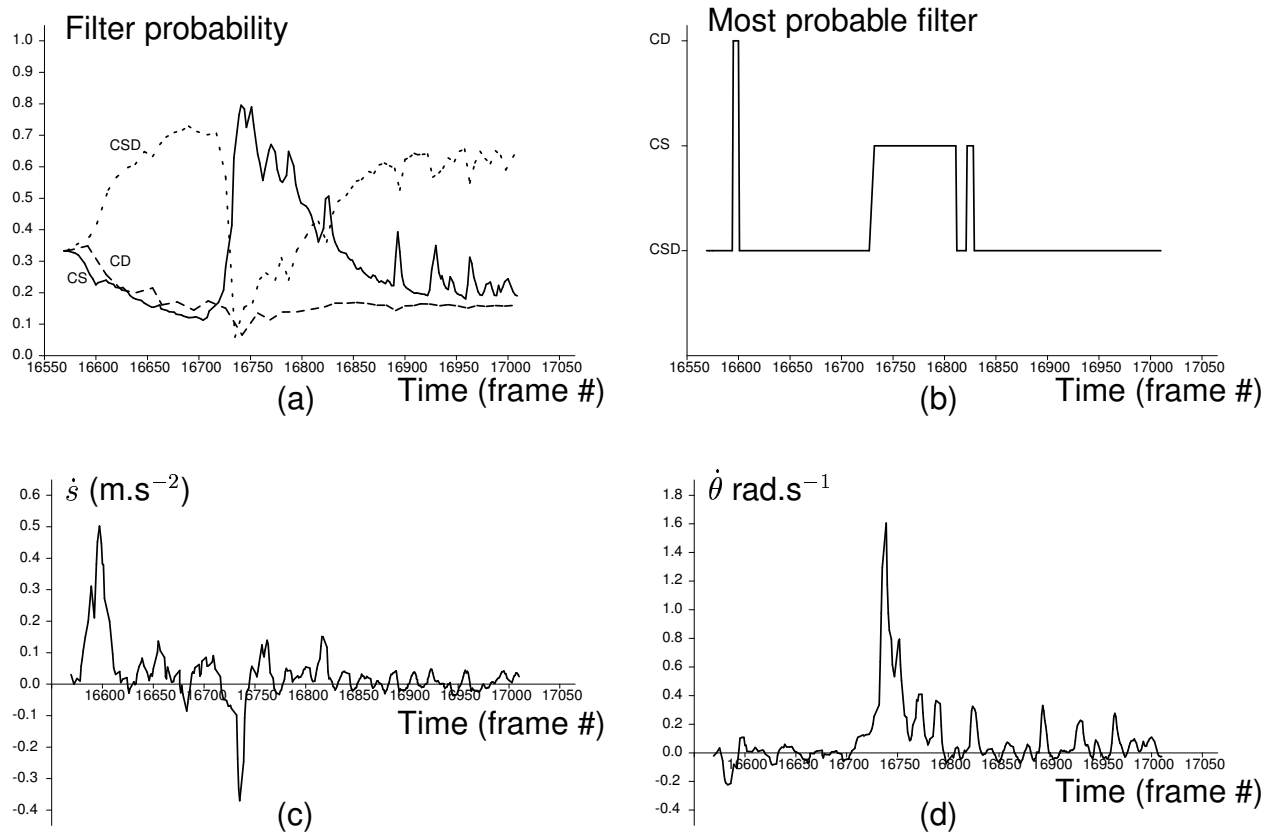


Fig. 13. IMM output for person trajectory. (a) shows the filter likelihoods (where solid, dashed and dotted lines are outputs from the CS, CD and CSD filters). (b) gives the most likely filter, where the dashed lines show when the actual manoeuvres took place; (c) and (d) show the overall IMM estimates of acceleration \dot{s} in ms^{-2} and turn rate $\dot{\theta}$ in rad.s^{-1} .

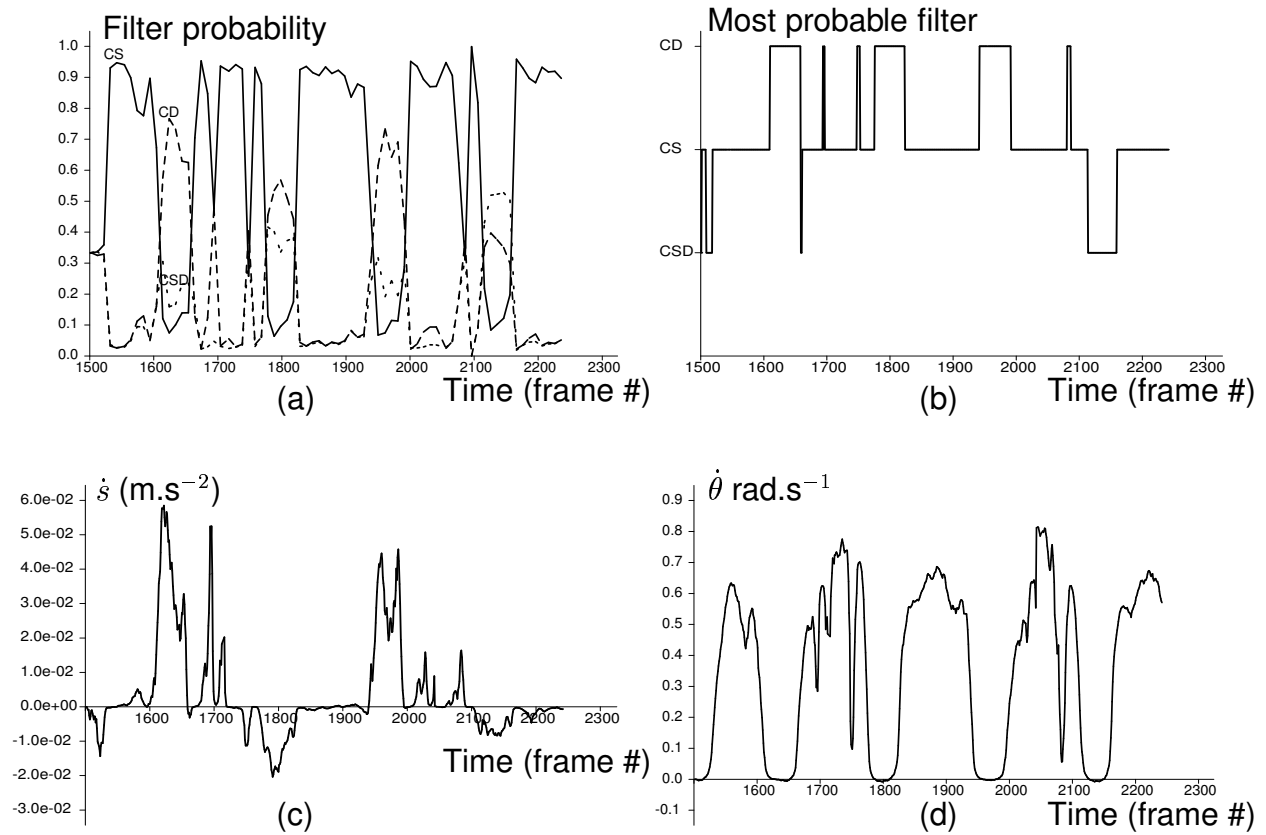


Fig. 14. IMM output for train in oval trajectory. (a) shows the filter likelihoods (where solid, dashed and dotted lines are outputs from the CS, CD and CSD filters). (b) gives the most likely filter, where the dashed lines show when the actual manoeuvres took place; (c) and (d) show the overall IMM estimates of acceleration \dot{s} in ms⁻² and turn rate $\dot{\theta}$ in rad.s⁻¹.

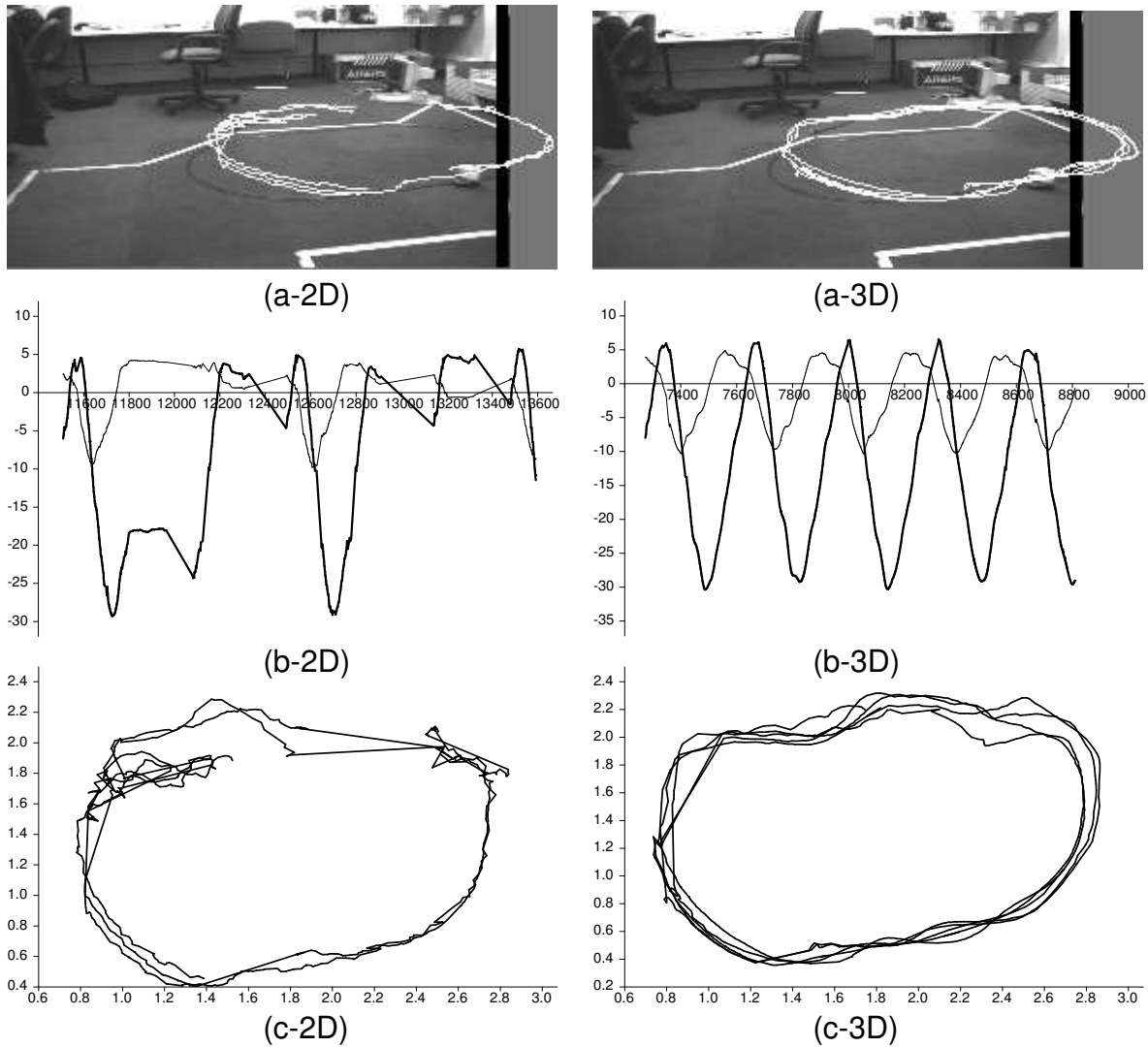


Fig. 15. The pursuit performance using 3D filtering and prediction contrasted with that from 2D filtering alone: (a) the trajectory mapped into reference image; (b) vergence (thick line) and elevation (thin line) gaze angles over sequence; (c) the recovered ground plane trajectory. Figures (b-2D) and (c-2D) show that, with 2D filtering alone, multiple saccades are required to recapture the target when pursuit breaks down.