

Entropy-based metrics for the analysis of partial and total occlusion in video object tracking

E. Loutas, I. Pitas and C. Nikou

Abstract: Metrics measuring tracking reliability under occlusion that are based on mutual information and do not resort to ground truth data are proposed. Metrics for both the initialisation of the region to be tracked as well as for measuring the performance of the tracking algorithm are presented. The metrics variations may be interpreted as a quantitative estimate of changes in the tracking region due to occlusion, sudden movement or deformation of the tracked object. Performance metrics based on the Kullback–Leibler distance and normalised correlation were also added for comparison purposes. The proposed approach was tested on an object tracking scheme using multiple feature point correspondences. Experimental results have shown that mutual information can effectively characterise object appearance and reappearance in many computer vision applications.

1 Introduction

Partial or full occlusion is an important issue in an object tracking process. A variety of algorithms handling occlusion exist [1–4], however, they do not handle total occlusion properly. Tracking is performed in [1] using sum of squared differences (SSD). Tracking does not rely on feature point sets. Partial occlusion and illumination changes are handled, nevertheless, the proposed algorithm does not handle full occlusion. A contour tracking algorithm is proposed in [2]. The resulting scheme is reliable in image clutter and partial occlusion, but it is not reliable for large amounts of occlusion. The algorithm presented in [3] relies on deformable templates and can handle moderate amounts of partial occlusion. In [4] the role of geometric invariants in tracking is examined. Feature point tracking verification using geometric invariants is presented. The aim of the method is to compute the target feature point set using geometric invariants. An algorithm insensitive to the disappearance and reappearance of feature points is described in [5]. Although the above mentioned methods handle partial occlusion, few of them behave well under total occlusion.

A model-based tracking scheme performing object tracking using edge information and capable of handling partial and total occlusion is proposed in [6]. The proposed method can handle partial and total occlusion events but is computationally expensive. A new approach to occlusion resistant object tracking using Kalman filtering and robust statistics has been proposed [7] that can handle full occlusion for short time periods. The way the tracking system recovers after total occlusion requires that the position of the disoccluded object lies within the tracker's search range. Another approach to tracking multiple

articulated objects in the presence of occlusion based on a Kalman filtering mechanism is presented in [8]. This system was tested in a surveillance scheme used to track moving people. The algorithm has shown good results in severe partial occlusion caused by inter-object and object–environment interference. Finally, a probabilistic multiple object tracking approach working under inter-object occlusions is presented in [9].

The performance measure of a tracking algorithm is also an open issue. Although most of the proposed techniques apply subjective evaluation methods, some of them use quantitative approaches based on ground truth [10]. Therefore, implementation of reliability measures not resorting to ground truth data is particularly important. Several metrics for performance evaluation of tracking algorithms without ground truth, based on colour and motion were introduced in [11]. More recent work on these metrics provides their incorporation in a tracking scheme in order to perform better tracking [12]. A variety of confidence measures for the analysis of optical flow techniques was presented in [13], however, the confidence measures analysed are only used for the evaluation of the velocity field and are application oriented.

The use of mutual information in object tracking as a tool for extracting information concerning the condition of a tracking object is assessed in this paper. The proposed scheme is efficient in extracting information under partial and total occlusion. Mutual information was first introduced in computer vision in [14] for medical image registration applications. In [15, 16] it was applied to combine the outputs of multiple tracking algorithms in order to improve the overall tracker performance.

In the proposed method, the tracking process is modelled as a communication task between a transmitter and a receiver through a channel. Information theory-based metrics are introduced. The mutual information is used as a quantitative measure of the tracking process. Its variations can improve understanding of the tracked region characteristics and are closely related to changes in the tracked region. These changes are caused by partial or total occlusion, movement of the occluding object and abrupt movements or deformations of the occluded (tracked)

© IEE, 2004

IEE Proceedings online no. 20040738

doi: 10.1049/ip-vis:20040738

Paper first received 9th September 2003 and in revised form 6th April 2004

The authors are with the Department of Informatics, University of Thessaloniki, Box 451, Thessaloniki 54124, Greece

object. Determining and understanding these changes may improve tracker performance and assist an event detection scheme. Measures based on the Kullback–Leibler distance and the normalised correlation are also implemented for comparison purposes. The entropy is used as a measure of the initialisation efficiency of the tracking process and is closely related to the first metric. The proposed metrics were tested on a feature point based tracking algorithm [17]. The algorithm is enhanced with an occlusion handling scheme, while an object reappearance verification scheme is also designed to allow tracking continuation after object reappearance. It relies on a mutual information-based metric measuring the similarity between a reference and a target region. The modified tracking algorithm performs better than [17] in partial and total occlusion situations.

The main contribution of the current work is the introduction of information theory based metrics as measures of tracking reliability. The use of the metrics does not impose utilisation of ground truth data and is extended to the analysis of partial and total occlusion in object tracking. Moreover, occlusion is processed without resorting to multiple camera systems fusing the outputs of different tracking cues.

2 Feature point generation and tracking

Object tracking is performed by minimising the sum of squared differences of a large set of feature points generated in the tracking region. The algorithm presented in [17] is used for feature point tracking. Kalman filtering motion prediction is employed to estimate the tracked region position during occlusion. The tracked region in the subsequent video frame is specified as the bounding rectangle of all the tracked feature points. Robustness to partial occlusion is achieved by estimating the motion of the lost feature points, using the estimated motion of the bounding box of the tracked object.

The displacement $\mathbf{d} = [d_x, d_y]^T$ between two feature point windows on images J_2 and J_1 is obtained by minimising

$$\epsilon = \int \int_W \left[J_2 \left(\mathbf{x} + \frac{\mathbf{d}}{2} \right) - J_1 \left(\mathbf{x} - \frac{\mathbf{d}}{2} \right) \right]^2 w(\mathbf{x}) d\mathbf{x} \quad (1)$$

where $\mathbf{x} = [x, y]^T$, W is the region of the intergration window and $w(\mathbf{x})$ is a weighting function that can be set to 1 for simplicity. Equation (1) uses

$$\left[J_2 \left(\mathbf{x} + \frac{\mathbf{d}}{2} \right) - J_1 \left(\mathbf{x} - \frac{\mathbf{d}}{2} \right) \right]$$

instead of $[J_2(\mathbf{x}) - J_1(\mathbf{x} - \mathbf{d})]$ used in [17], because of its symmetry with respect to both images [18]. In order to perform one iteration of the minimisation procedure of (1), the equation $\mathbf{Z}\mathbf{d} = \mathbf{e}$ must be solved where

$$\mathbf{Z} = \int \int_W \mathbf{g}(\mathbf{x}) \mathbf{g}^T(\mathbf{x}) w(\mathbf{x}) d\mathbf{x} \quad (2)$$

$$\mathbf{e} = 2 \int \int_W [J_1(\mathbf{x}) - J_2(\mathbf{x})] \mathbf{g}(\mathbf{x}) w(\mathbf{x}) d\mathbf{x} \quad (3)$$

and

$$\mathbf{g} = \begin{bmatrix} \frac{\partial(J_1 + J_2)}{\partial x} \\ \frac{\partial(J_1 + J_2)}{\partial y} \end{bmatrix} \quad (4)$$

Feature point occlusion is determined using the process described in [17] and is essentially controlled by the

residue ϵ . Large values of ϵ when compared to a predefined threshold imply that the feature point of interest should be rejected. The object tracking occlusion handling is based on feature point occlusion handling as presented in [17]. That is, an object part is considered lost when the feature points ‘belonging’ to that part are lost. Other methods of occlusion handling involve the use of constraints based on articulation [19] and layer representation [20]. The approach used in the context of the present work is general and can be used in a variety of applications. The layer representation method is very useful in coding and compression, while the role of articulation constraints in determining self-occlusions in human body part tracking is vital.

In order to avoid tracking stationary or slowly moving background feature points, we have introduced a clustering procedure. The mean (μ_x, μ_y) and the variance (σ_x, σ_y) of the feature point co-ordinates are computed for the tracked region in each frame. Let $[x, y]^T$ be the co-ordinates of a feature point at video frame t and (μ_x, μ_y) , (σ_x, σ_y) their mean and variance. A feature point is retained in video frame $t + 1$, if $x \in [\mu_x - \sigma_x, \mu_x + \sigma_x]$ and $y \in [\mu_y - \sigma_y, \mu_y + \sigma_y]$, otherwise it is rejected. Assuming that the object feature points have similar motion patterns, we can reject stationary or slow moving background features, after a number of frames, while retaining the moving object feature points. This procedure is particularly useful, if the initialised region to be tracked contains some portions of background regions.

2.1 Initialisation

The region bounding rectangle is used to specify the region to be tracked. A large number of feature points is generated inside the tracked region using the process described in [17, 18, 21]. A good feature point is defined as one whose matrix \mathbf{Z} has two large eigenvalues that do not differ by several orders of magnitude [21]. In order to avoid loss of target, caused by too many lost feature points, the feature point set is periodically regenerated. Different strategies for the periodic feature point regeneration can be applied. It can be thorough (the entire feature point set is regenerated), periodic (it occurs after a fixed number of frames) or asynchronous (its occurrence is based on the tracking process metric value). Feature point generation and tracking are transparent to the observer.

The number of the generated feature points is essentially user controlled. The user controls the number of feature points by selecting their number and the minimum allowed distance between the feature points. Let N_s be the desired number of feature points selected by the user. The number N_k of feature points generated in the region to be tracked depends essentially on the minimal allowed distance between the feature points ($N_k \leq N_s$). Therefore, a set of the possible configurations of the ensemble of the possible feature point sets can be defined. Large minimum allowed distances between feature points may lead to a small N_k and poor tracker performance.

3 Robustness to partial and total occlusion

The previously described tracking process can be modelled as a communication between a transmitter (reference frame) and a receiver (target frame) with an N_{\max} symbol alphabet (the maximal number of greyscale levels). The tracking process is characterised by loss of information caused by feature point rejection and wrong feature point correspondences. Mutual information is a well known measure of the amount of information transmitted through the communication channel [22, 23]. Therefore, it can be used as a quantitative measure of tracking performance.

3.1 Mutual information as tracker evaluation metric

Let \mathbf{x}_i^r and \mathbf{x}_i^c represent the co-ordinate vectors of feature point i in the reference and current frame, respectively. During the tracking process, a feature point set of the initial video frame

$$\mathbf{S}_1 = [\mathbf{x}_1^r, \dots, \mathbf{x}_{N_k}^r]^T \quad (5)$$

is tracked to a feature point set

$$\mathbf{S}_2 = [\mathbf{x}_1^c, \dots, \mathbf{x}_N^c]^T \quad (6)$$

of the target video frame, with $N \leq N_k$, $N_k \leq N_s$, where N_s is the initial user preference for the number of the feature points.

Let U , V be two random variables with marginal probability mass functions $p(u)$, $p(v)$ and $u_i = J_1(\mathbf{x}_i^r)$, $v_j = J_2(\mathbf{x}_j^c)$ their possible outcomes, where J_1 and J_2 are the reference and target image respectively and $\mathbf{x}_i^r \in \mathbf{S}_1$, $\mathbf{x}_j^c \in \mathbf{S}_2$. The mutual information of the two random variables U , V with a joint probability mass function $p(u, v)$ is defined as

$$I(U, V) = \sum_{i=1}^{N_{\max}} \sum_{j=1}^{N_{\max}} p(u_i, v_j) \log_2 \frac{p(u_i, v_j)}{p(u_i)p(v_j)} \quad (7)$$

where N_{\max} is the maximum number of the available greyscale levels. In order to take into account the lost feature points during the tracking process, a cost function E_m is defined:

$$E_m(U, V, N, N_k) = c_1 \left(\frac{I(U, V)}{I_{\max}(U, V)} - \lambda_1 \frac{N_k - N}{N_k} + c_2 \right) \quad (8)$$

The term $I(U, V)/I_{\max}(U, V)$ is the mutual information part of the cost function. The maximum mutual information $I_{\max}(U, V)$ is [24]:

$$I_{\max}(U, V) = - \sum_{i=1}^{N_{\max}} p(u_i) \log_2 p(u_i) \quad (9)$$

The term $(N_k - N)/N_k$ is a penalising quantity depending on the number of the lost feature points during the tracking process. The use of the penalising term is necessary, because the mutual information part of the metric measures only the matching efficiency between the feature points that have not been lost. In the context of present work $c_1 = 0.5$, $\lambda_1 = 1$, $c_2 = 1$. The constants c_1 , c_2 , λ_1 are chosen to satisfy

$$0 \leq E_m \leq 1 \quad (10)$$

In the case of total occlusion:

$$\frac{I(U, V)}{I_{\max}(U, V)} = 0 \text{ and } \frac{N_k - N}{N_k} = 1 \quad (11)$$

leading to the minimum value of E_m . The maximum value of E_m occurs when

$$I(U, V) = I_{\max}(U, V) \text{ and } N = N_k \quad (12)$$

The metric E_m is a measure of the information flow during the tracking process. Large values of E_m represent large amounts of information carried from the reference region to the target output region. In this case, the similarity between the reference region and the target region and, consequently, the reliability of the tracker output, are high. Small values of E_m are an indication that the tracking process is unreliable.

3.2 Kullback–Leibler distance-based tracking metric

The Kullback–Leibler distance is defined as [25]

$$D(p(u)||p(v)) = \sum_{i=1}^{N_{\max}} p(u_i) \log_2 \frac{p(u_i)}{p(v_i)} \quad (13)$$

and measures the similarity between $p(u_i)$ and $p(v_i)$. It is not symmetric, i.e. in general $D(p(u)||p(v)) \neq D(p(v)||p(u))$. An upper bound of the Kullback–Leibler distance can be easily found as follows, since

$$\begin{aligned} D(p(u)||p(v)) &= \sum_{i=1}^{N_{\max}} p(u_i) \log_2 \frac{p(u_i)}{p(v_i)} \\ &= \sum_{i=1}^{N_{\max}} p(u_i) \log_2 p(u_i) - \sum_{i=1}^{N_{\max}} p(u_i) \log_2 p(v_i) \end{aligned} \quad (14)$$

The first term is negative or zero, while the second is positive. Therefore, an upper bound of the Kullback–Leibler distance is

$$D(p(u)||p(v)) \leq - \sum_{i=1}^{N_{\max}} p(u_i) \log_2 p(v_i) \quad (15)$$

A similar metric to $E_m(U, V, N, N_k)$ based on the Kullback–Leibler distance can be defined as

$$E_K(U, V, N, N_k) = c_1 \left(1 - \frac{D(p(u)||p(v))}{D_{\max}(p(u)||p(v))} - \lambda_1 \frac{N_k - N}{N_k} + c_2 \right) \quad (16)$$

and by construction is expected to behave similarly to E_m . Large values of E_K imply better matching between the reference and the target region. Both mutual information and Kullback–Leibler tracking metrics are expected to perform best when we have planar object motion with partial and total occlusions.

3.3 Normalised correlation-based metric

The normalised correlation between the reference and the target feature point sets can be defined as [26]

$$C_n = \frac{\sum_{i=1}^N J_1(\mathbf{x}_i^r) J_2(\mathbf{x}_i^c)}{\sqrt{\sum_{i=1}^N J_1^2(\mathbf{x}_i^r) \sum_{i=1}^N J_2^2(\mathbf{x}_i^c)}} \quad (17)$$

since a one by one correspondence exists between the feature point sets. Equation (17) expresses the similarity between J_1 and J_2 and can be used to construct a metric similar to those already presented in the context of present work (see (8) and (16)). The metric constructed is of the form

$$Corr = c_1 \left(C_n - \lambda_1 \frac{N_k - N}{N_k} + c_2 \right) \quad (18)$$

and was also tested under similar tracking conditions to the other two. It holds that

$$0 \leq Corr \leq 1 \quad (19)$$

The values of the constants c_1 , c_2 , λ_1 are the same as in (8) and (16).

3.4 Tracker initialisation evaluation metric

Since the feature point set \mathbf{S}_1 generated on the initial frame belongs to the power set of the possible feature point set

configurations, a metric measuring the reliability of S_1 can be defined. It can characterise the efficiency of the initially selected region for tracking. Each feature point set S_k is characterised by its entropy:

$$H_{S_k} = - \sum_{i=1}^{N_{\max}} p_k(u_i) \log_2 p_k(u_i) \quad (20)$$

where $u = J(x)$ are the image luminances at feature point locations on the initial frame. Let N_k be the number of feature points generated in the tracked region. In general, $N_k \leq N_s$. The maximal value of H_{S_k} depends on N_k , if $N_k \leq N_{\max}$, since in that case the number of greyscale levels, belonging to the feature point set, cannot reach N_{\max} . Then the distribution $p_k(u_i) = 1/N_k$ can create an upper bound H_{S_k} if $N_k \leq N_{\max}$. Therefore

$$p_k(u_i) = \begin{cases} \frac{1}{N_k} & N_k \leq N_{\max} \\ \frac{1}{N_{\max}} & N_k > N_{\max} \end{cases} \quad (21)$$

Clearly H_{S_k} is maximised when $N_k \geq N_{\max}$ and $p_k(u_i) = 1/N_{\max}$. The maximal symbol value of the communication alphabet is N_{\max} (maximum number of greyscale levels). In order to handle degenerative cases, where the number of the generated feature points N_k is much smaller than the initial user preference N_s , a penalising term depending on the number of not generated feature points is added. Such cases occur when the minimum allowed distance between feature points is large, compared to the region size. Therefore, the metric, measuring the efficiency of the feature point sets produced during the initialisation step, is defined as

$$E_i(H_{S_k}, N_k, N_s) = \begin{cases} \frac{H_{S_k}}{\log_2 N_k} & N_T \leq N_k < N_{\max} \\ \lambda_H \frac{H_{S_k}}{\log_2 N_k} + \lambda_F \frac{N_k}{N_s} & N_k < N_{\max}, \quad N_k < N_T \\ \frac{H_{S_k}}{\log_2 N_{\max}} & N_T \leq N_k, \quad N_{\max} \leq N_k \\ \lambda_H \frac{H_{S_k}}{\log_2 N_{\max}} + \lambda_F \frac{N_k}{N_s} & N_{\max} \leq N_k, \quad N_k < N_T \end{cases} \quad (22)$$

Threshold N_T is usually a fraction of the user specified feature point number N_s . In the context of present work we have chosen $N_s = 180$, $N_T = N_s/4$, $\lambda_H = 0.5$, $\lambda_F = 0.5$. The penalising term is introduced only when $N_k < N_T$. In such cases the number of the feature points N_k is small and the penalising term of (22) has to be added. The metric E_i is a measure of efficiency of the initial feature point set configuration. It imposes feature point selection based on feature point set entropy. The initialisation metric imposes large feature point set luminance variation by using entropy maximisation. The most effective way of controlling the feature point set configuration is by changing the minimum distance between the feature points. Small distances lead to a feature point concentration in certain parts of the object being tracked. Larger distances usually help to provide feature point sets with better coverage of the object being tracked and to attain better tracking results. Ideally, the average feature point distance should be greater than the texture cell or grain size.

The entropy based selection criterion aims at imposing a large feature point intensity dispersion in order to provide better tracking results. The penalising term is introduced to prevent a feature point set choice with too large distances between feature points that contains a small number of feature points. The initialisation criterion can also be

applied to untextured objects with limited success. The choice of the initial feature point set configuration is important for the success of the object tracking process.

4 Tracking algorithm enhancement

The tracking algorithm presented in Section 2 is enhanced by using an occlusion handling scheme. It is capable of handling partial and total occlusion in a variety of cases. The occlusion handling scheme is assisted by an object verification scheme, applied to total occlusion situations. The object verification scheme is based on the metric E_m , in the context of the present work. Nevertheless, other techniques, such as elastic graph matching, can also be used.

4.1 Occlusion handling

In order to cope with partial occlusion, a prediction scheme is applied. The lost features are not tracked. However, their co-ordinates are updated using the estimated movement of the upper left and the lower right corner of the bounding rectangle of the tracked object. The procedure is stopped if the occlusion is total, that is, when none of the feature points comprising the feature point set can be further tracked correctly due to occlusion. In order to handle large variations of the bounding box size, caused by feature point loss, the area of the tracked region is introduced as a reliability measure of the update of the upper left and lower right bounding box co-ordinates. The feature points, whose co-ordinates are updated, are considered lost if the bounding box area exceeds a threshold T_{\max} or is smaller than a threshold T_{\min} . Periodical regeneration of the lost feature points during the tracking process using the procedure presented in Section 2 is also a useful tool in order to handle partial occlusion and allow tracking for long time periods. The feature points not lost in the tracking process are not regenerated. In order to cope with total occlusion, the position of the occluded region is updated using the velocity estimates of the region corners obtained from the measurements before total occlusion with the help of a Kalman filtering scheme.

The Kalman filtering prediction process is applied on the upper left corner and the lower right corner of the region bounding rectangle before total occlusion. A constant acceleration model is used [27]. Let $\mathbf{d}(k)$, $\mathbf{u}(k)$, and $\mathbf{a}(k)$ denote the displacement velocity and acceleration for each corner of the bounding box at time k respectively. The state-transition equation for each corner is [27]

$$\mathbf{s}(k) = \mathbf{C}\mathbf{s}(k-1) + \mathbf{w}(k) \quad k = 1, \dots, N \quad (23)$$

where $\mathbf{w}(k)$ is a zero-mean, white random sequence and \mathbf{s} is a 6×1 vector containing the co-ordinates of displacement velocity and acceleration, for each corner of the bounding box:

$$\mathbf{s} = [d_x \quad d_y \quad u_x \quad u_y \quad a_x \quad a_y]^T \quad (24)$$

The measurements $\mathbf{d}(k)$ are related to the state variables $\mathbf{s}(k)$ with

$$\mathbf{d}(k) = \mathbf{H}\mathbf{s}(k) + \mathbf{v}(k) \quad k = 1, \dots, N \quad (25)$$

where $\mathbf{v}(k)$ denotes a zero-mean, white observation noise sequence. The matrices describing the model are given below. The 2×1 observation vector and the 2×6 measurement matrix are given by:

$$\mathbf{d} = \begin{bmatrix} d_x \\ d_y \end{bmatrix} \quad (26)$$

$$H = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix} \quad (27)$$

The observation equation states that the noisy displacement co-ordinates of each bounding box corner can be observed.

The 6×6 state transition matrix describing the model is [27]

$$C = \begin{bmatrix} 1 & 0 & 1 & 0 & 0.5 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0.5 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \quad (28)$$

4.2 Object reappearance prediction and verification

Reappearance prediction is obtained by estimating and tracking the occluding region. To estimate the occluding region bounding box, a simple region growing segmentation algorithm is used. A seed is determined by the last position of the occluded region before total occlusion. The occluded object is considered to be entirely disoccluded when

$$A_1 \cap A_2 = \emptyset \quad (29)$$

where A_1 is the occluding region and A_2 is the predicted occluded region. The occluded region reappears provided that the above condition is satisfied. Reappearance is associated with the regeneration of a set of feature points, as in the initialisation step. Again, the selected feature points for object reappearance are those that have large eigenvalues of matrix Z . The feature point regeneration is thorough after total occlusion, that is the entire feature point set is regenerated inside the bounding rectangle specified by A_2 . Object tracking continues after the feature point set regeneration.

When the tracker predicts that reappearance has taken place, it has to decide if the reappearing region is similar enough to the tracked region before occlusion. This can be achieved by using the mutual information metric E_m (8). A feature point set is generated in the tracked region belonging to the frame before total occlusion. The position of this feature point set on the current frame is predicted. The metric E_m is calculated using the feature point sets of the reference and target frames, while the predicted tracked region is allowed to change slightly. The maximum of the E_m value is compared with a threshold. The threshold value

can be chosen according to the value of E_m before total occlusion.

Graph matching is an alternative technique that can be used for object reappearance verification. Nevertheless, the use of E_m as previously described is preferred, in the context of present work for simplicity and uniformity.

5 Experimental results

The proposed tracking algorithm was tested on both real and artificially generated image sequences. In order to evaluate the efficiency of the proposed scheme, image sequences containing total occlusion and partial occlusion were used. Curves showing the variations of the metrics E_m , E_K and $Corr$ during the tracking process were calculated for different occlusion cases. The metric E_i of the tracking algorithm initialisation efficiency, was tested on both artificial and real image sequences.

The algorithm involves the choice of system parameters in order to work. The parameters' values are kept constant during the experiments. The choice of N_s is left to the user and depends mainly on the tracked object size. N_T is a fraction of N_s acquired by experience. The choice of c_1 , c_2 , λ_1 , λ_F and λ_H is imposed by the requirement $0 \leq E_m \leq 1$ and $0 \leq E_i \leq 1$. Their value is kept constant throughout the entirety of experiments. The choice of the minimum distance between feature points is crucial to the tracking process and is obtained by using the initialisation metric E_i .

Results on an artificial image sequence are presented in Fig. 1. A small circular object (Fig. 1a) moves slowly from right to left and is fully occluded by a faster moving elliptical object that moves in the opposite direction (Fig. 1b). The tracked region bounding rectangle is recalculated after total disocclusion. The algorithm performs well, even when the occluding object reappears suddenly without previous appearance in the image sequence. The cost function E_m for various image sequence frames is shown in Fig. 2. A decrease in E_m begins after frame 15, marking the beginning of partial occlusion. The minimal value of $E_m = 0$ marks the beginning of total occlusion. Object reappearance is marked by an abrupt increase of E_m . The frames corresponding to the start of partial occlusion and the start of total occlusion are shown in Fig. 3.

In Fig. 4, the tracked region (head of the football player) is occluded by the foot of another football player. The cost function E_m for each frame of the image sequence is shown in Fig. 5. E_m drops to its minimum value $E_m = 0$ during total occlusion. The results on real and artificial image sequences

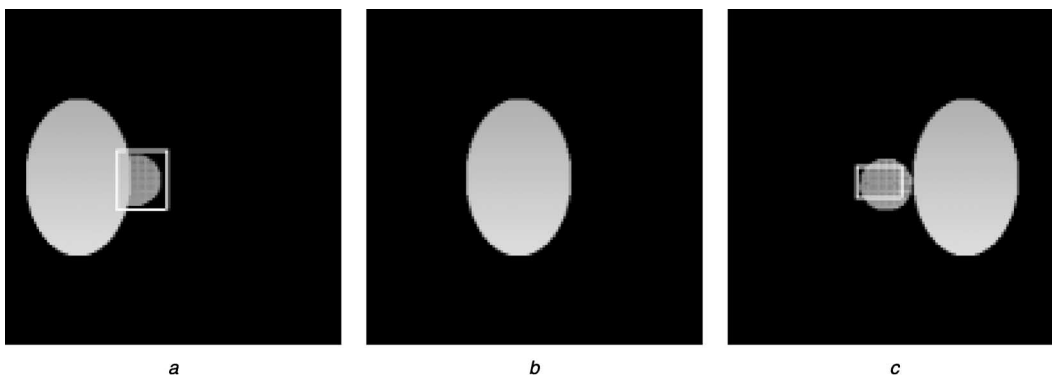


Fig. 1 Artificial image sequence

- a Before total occlusion
- b During total occlusion
- c Region reappearance

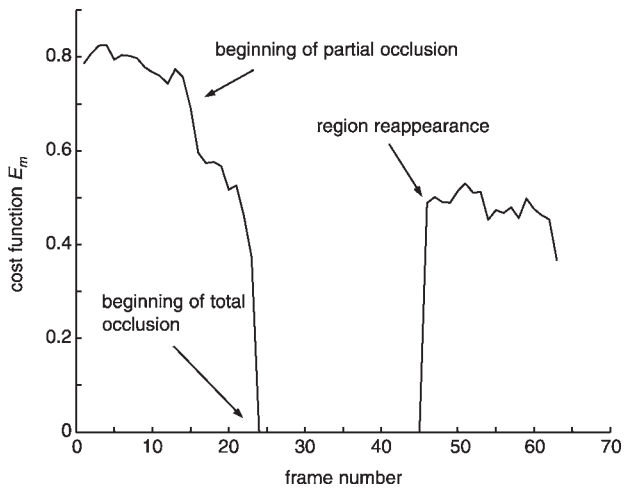


Fig. 2 Cost function E_m against frame number for artificial image sequence of Fig. 1

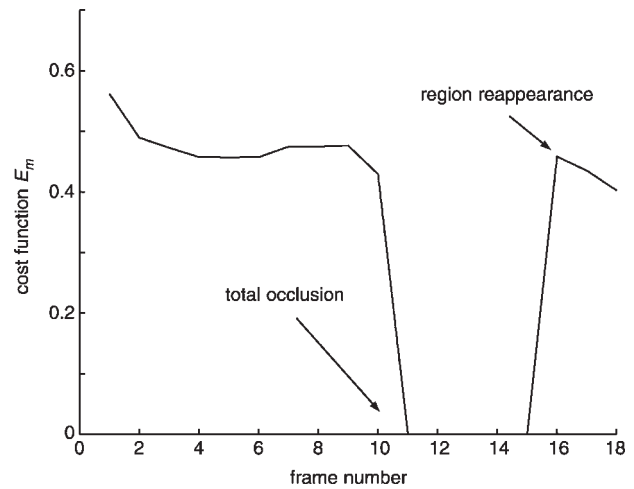


Fig. 5 Value of cost function E_m against frame number for part of 'Football' image sequence (Fig. 4)

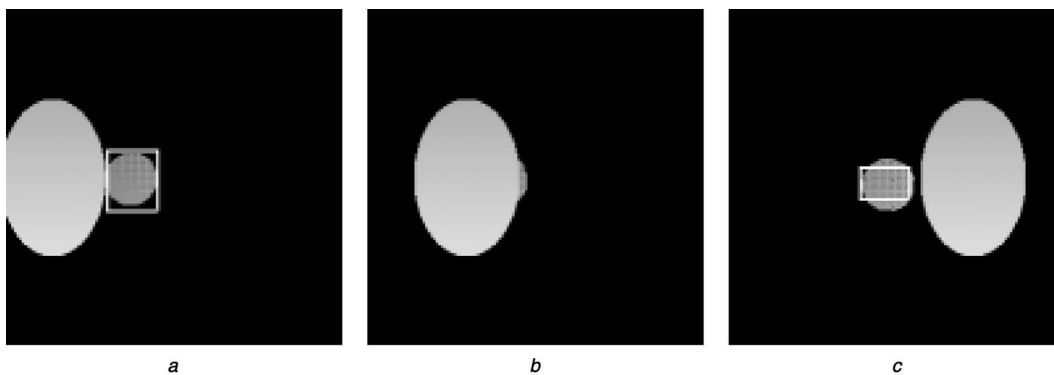


Fig. 3 Artificial image sequence frames characterised by mutual information cost function

- a Start of partial occlusion frame (frame no. 15)
- b Start of total occlusion frame (frame no. 23)
- c First frame after total object reappearance (frame no. 45)

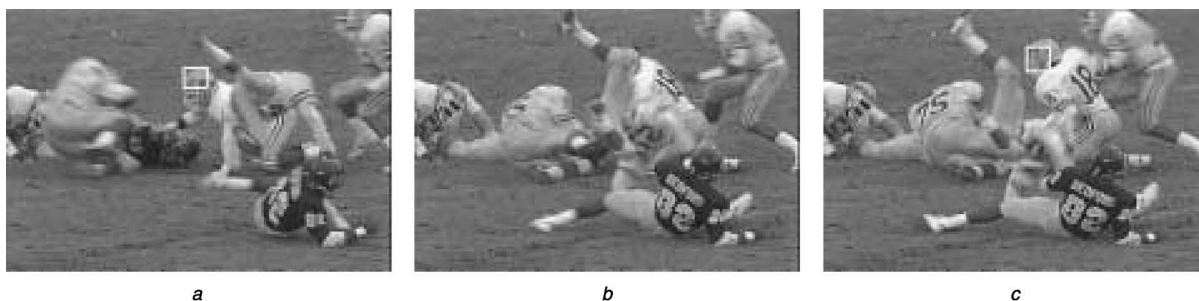


Fig. 4 'Football' image sequence

- a Before total occlusion
- b During total occlusion
- c Region reappearance

show that E_m can be useful in the analysis of partial and total occlusion in object tracking. Partial occlusion is accompanied by a drop of E_m , while total occlusion is characterised by a zero E_m value. The sudden increase of cost function E_m after object reappearance in the example of Fig. 4 is caused by generation of a feature point set during object reappearance, as described in Section 4.2. An increase of E_m is possible whenever a feature point set regeneration occurs.

In Fig. 6, results showing robustness to partial occlusions are presented. A person's face is partially occluded and, at the end of partial occlusion, the tracked face reappears

completely. The beginning of partial occlusion in frame 34 (Fig. 7) is marked by a sudden drop in E_m (Fig. 12). The mutual information does not increase after face disocclusion, since many feature points were lost during partial occlusion that have not been regenerated after disocclusion. Two frames showing the feature point sets before and after partial occlusion are presented in Fig. 8. Notice the loss of feature points, which is caused by partial occlusion. The tracked region size is computed correctly with the help of the partial occlusion handling scheme.

The object tracking algorithm containing the occlusion handling scheme and the object reappearance prediction

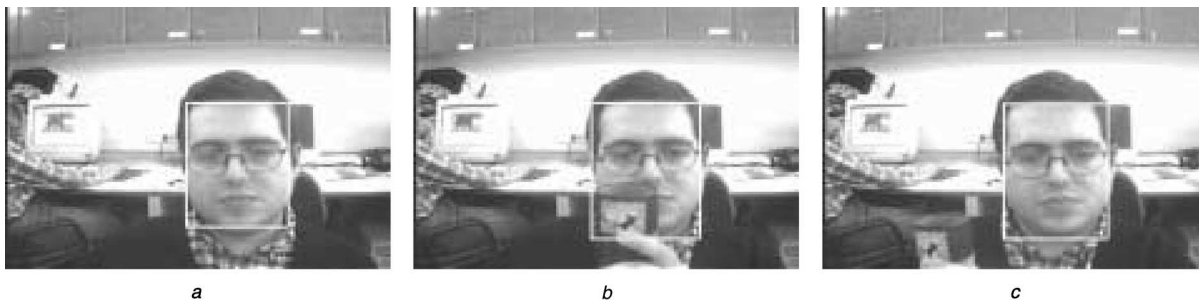


Fig. 6 *Lab image sequence*

- a* Tracked region
- b* Partial occlusion
- c* Region after occlusion

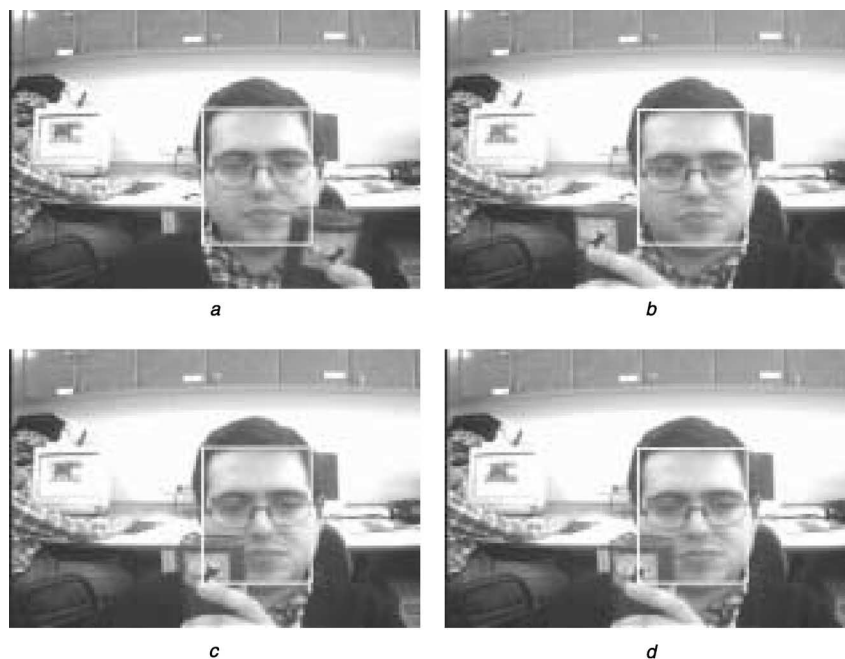


Fig. 7 *Lab image sequence*

- a* Beginning of partial occlusion (frame 34)
- b* Disocclusion (frame 101)
- c* Movement of the occluding region (frame 85)
- d* Movement of the occluding region (frame 86)



Fig. 8 *Lab image sequence*

- a* Feature point set before partial occlusion
- b* Feature point set after partial occlusion

and verification scheme performs better than an object tracking algorithm based, on [17] without these new additions. In Fig. 9 the results of [17] without the new additions on the 'Football' image sequence are presented. Notice the performance degradation before total occlusion

and the loss of target after total occlusion versus the results shown in Fig. 4. Similar results on the artificial image sequence are presented in Fig. 10. Performance degradation before total occlusion and loss of target after total occlusion is also noticed, when compared with the results shown in

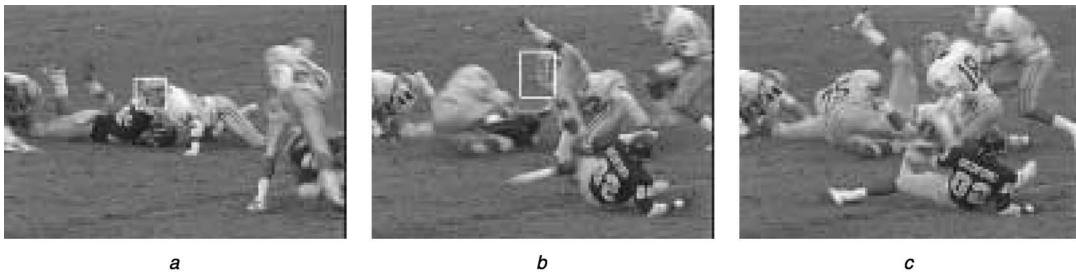


Fig. 9 'Football' image sequence: tracking without occlusion handling and object reappearance prediction and verification

- a Initial frame
- b Before total occlusion
- c After total occlusion

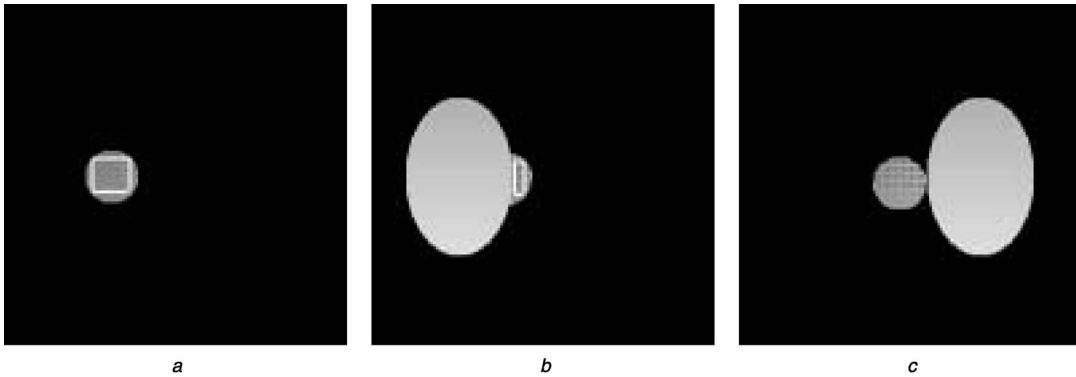


Fig. 10 Artificial image sequence: tracking without occlusion handling and object reappearance prediction and verification

- a Initial frame
 - b Before total occlusion
 - c After total occlusion
- Notice tracking degradation in b and in c

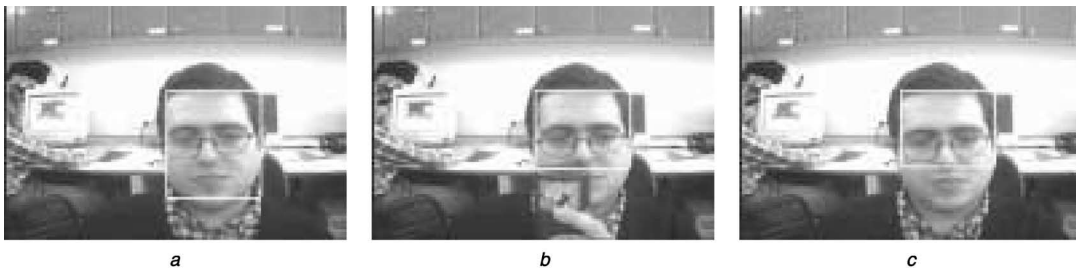


Fig. 11 Lab image sequence: tracking without occlusion handling and object reappearance prediction and verification

- a Initial frame
 - b During partial occlusion
 - c After partial occlusion
- Notice tracking degradation in b and in c

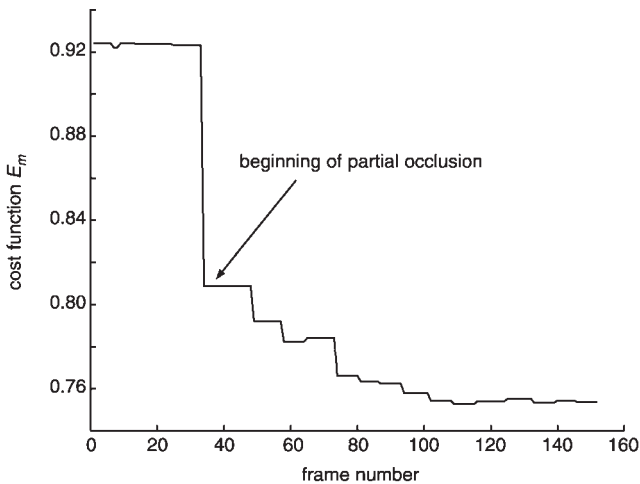


Fig. 12 Cost function E_m for lab image sequence (Fig. 7) against frame number

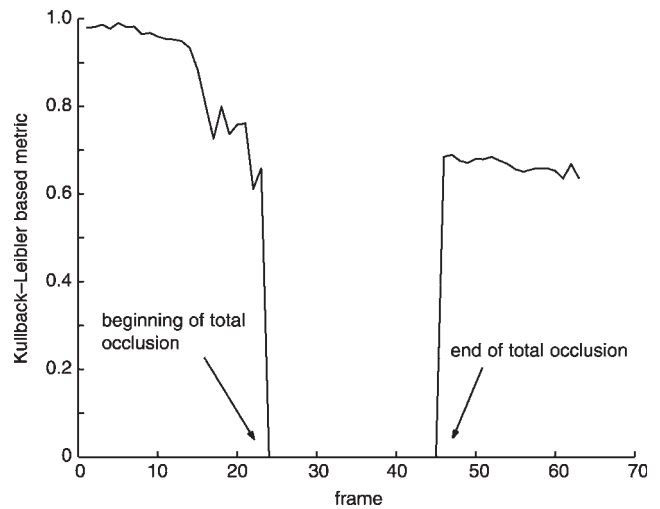


Fig. 13 Cost function E_K for artificial image sequence (Fig. 3) against frame number

Fig. 3. Results on the lab image sequence are presented in Fig. 11. Partial occlusion affects the tracking performance. One part of the tracked object is lost during and after partial occlusion, as can be seen in Fig. 11, in contrast to that shown in Fig. 6.

The variations of the metric E_m for the three sequences are presented in Figs. 2, 5 and 12 respectively, while the metric E_k variations are presented in Figs. 13, 14 and 15. Finally, the variations of the *Corr* metric are presented in Figs. 17, 18 and 19. As can be seen, metric E_k performs similarly to E_m . Further tests have shown that no significant change in E_k behaviour was caused by its asymmetry (Fig. 16). The normalised correlation-based metric *Corr* does not behave as well as the information theory based metrics in partial occlusion situations (Figs. 17, 18 and 19). The authors believe that the information theory based metrics should be preferred over the normalised correlation metric. Mutual information can be very useful as it provides spatial information and is symmetrical. The Kullback–Leibler distance can provide a variety of metrics with similar performance.

The variations of the initialisation performance metric E_i (22) with respect to the minimum allowed distance in pixels between feature points in the reference frame are presented in Figs. 20 and 21 for the artificial image sequence and the ‘Football’ image sequence, respectively. The cost function

values are generally bigger in the ‘Football’ image sequence than in the artificial image sequence case due to the fact that the initialised region in the artificial image sequence is uniformly textured. The value of E_i increases when the minimum allowed distance between features increases, provided that $N_k \cong N_s$. A rapid decrease in the E_i value is

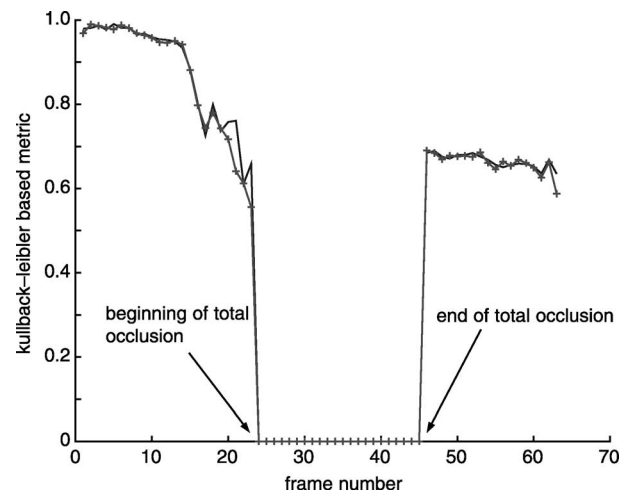


Fig. 16 Values of two forms of cost function E_K based on $D(p(u)||p(v))$ and $D(p(v)||p(u))$ for artificial image sequence (Fig. 3)

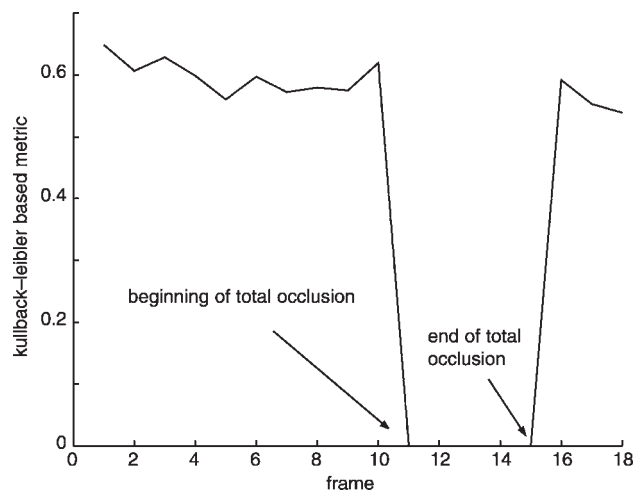


Fig. 14 Values of cost function E_K against frame number for part of ‘Football’ image sequence (Fig. 4)

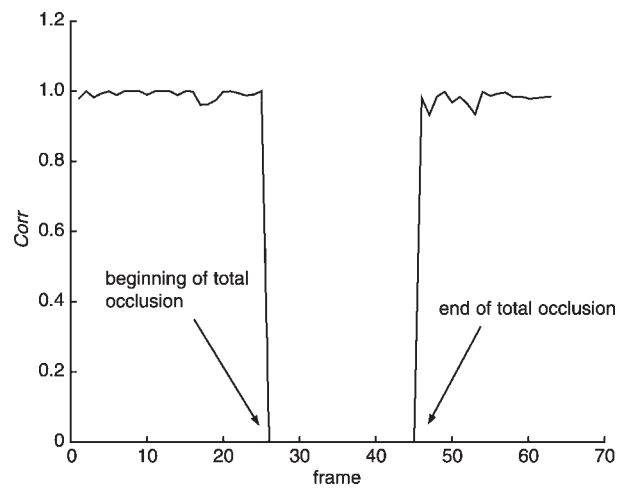


Fig. 17 Normalised correlation for artificial image sequence (Fig. 3) against frame number

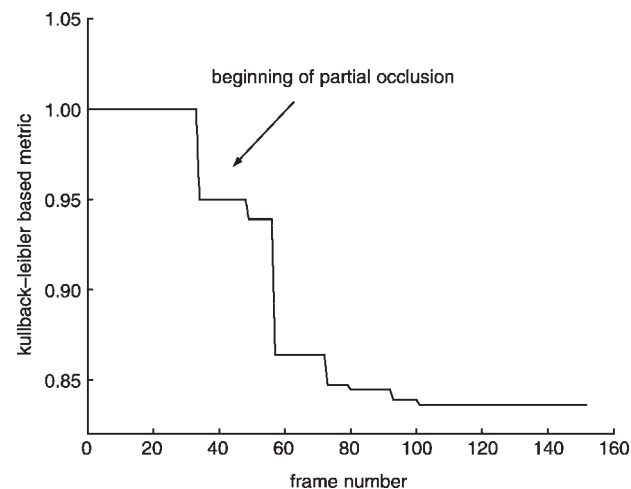


Fig. 15 Value of cost function E_K against frame number for lab image sequence (Fig. 7)

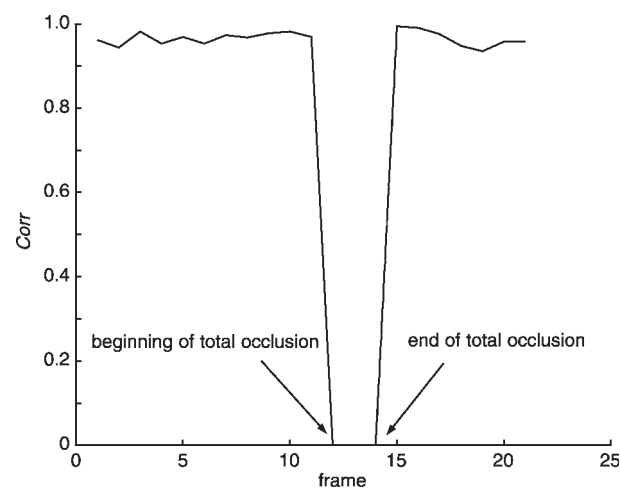


Fig. 18 Normalised correlation for ‘Football’ image sequence (Fig. 4) against frame number

noticed when the minimum allowed distance between feature points increase causes the number of feature points generated in the tracking region to be much smaller than the initial user-preferred feature point number ($N_k \ll N_s$).

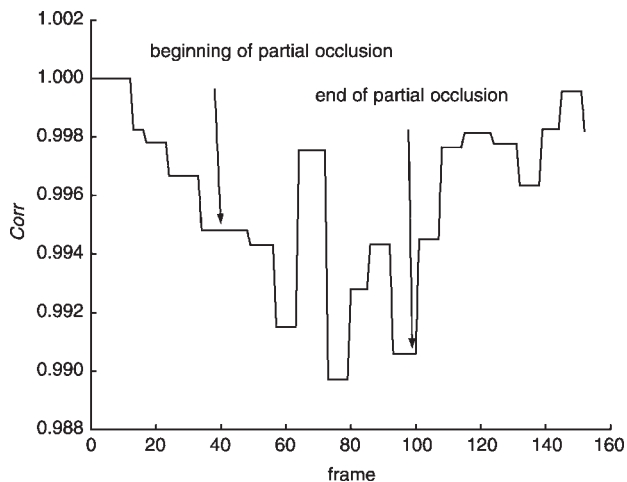


Fig. 19 Normalised correlation for lab image sequence (Fig. 7) against frame number

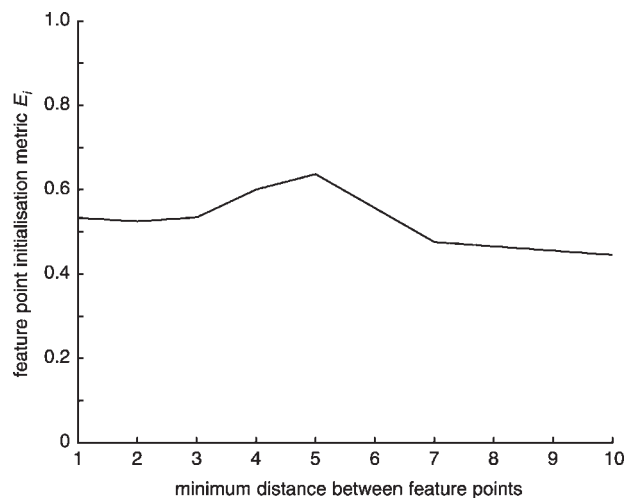


Fig. 20 Cost function E_i for algorithm initialisation of artificial image sequence

Notice that the texture grain size is 5 pixels (Figs. 1, 2)

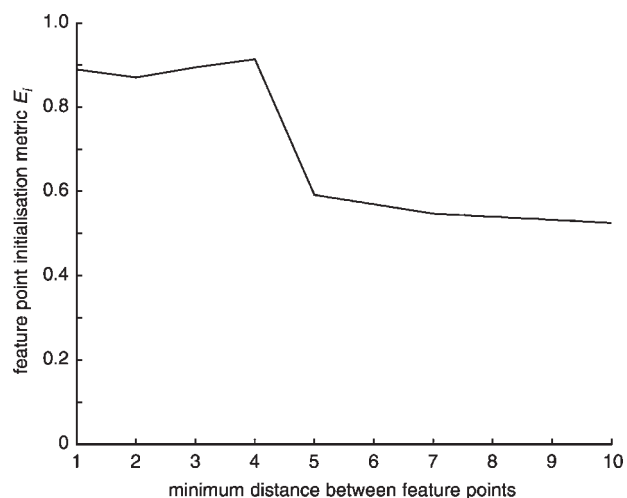


Fig. 21 Cost function E_i for initialisation process of 'Football' image sequence (Fig. 4)

The effectiveness of the proposed tracker initialisation metric was tested by performing object tracking of the 'Football' and artificial image sequences under different minimum feature points distances. The algorithm performs well when the minimum distance between feature point in the 'Football' image sequence case (Fig. 22) is less than five pixels. A rapid decrease in performance was noticed when the feature point distance increased above five pixels. Tests performed on the artificial image sequence case have shown no significant change in algorithm performance for feature point distances in the range $[3, \dots, 7]$ pixels. A decrease in the algorithm performance was noticed for a feature point distance around 10 pixels. It can be noticed in the artificial image sequence that the 'best' 5-pixel value is equal to the texture grain size. This shows a possible relationship between texture grain size and feature point distance.

The effectiveness of the partial occlusion handling scheme during the tracking process is shown in Figs. 23 and 24. In Fig. 23, the loss of feature points is caused by partial occlusion. Figure 24 demonstrates the usefulness of the updating procedure in a case not containing partial occlusion, since loss of feature points can be caused by

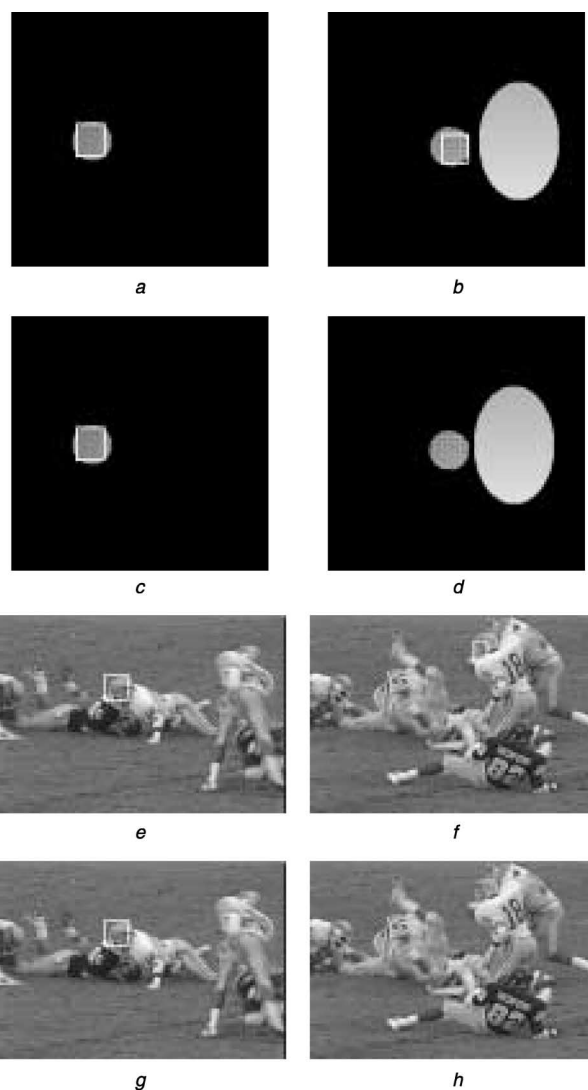


Fig. 22 Tracker outputs for artificial image sequence (a–d) and 'Football' image sequence I (e–h)

a and b minimum distance between feature points 5 pixels
 c and d minimum distance between feature points 10 pixels
 e and f minimum distance between feature points 4 pixels
 g and h minimum distance between feature points 5 pixels
 Notice the performance degradation at (c),(d) and (g),(h)

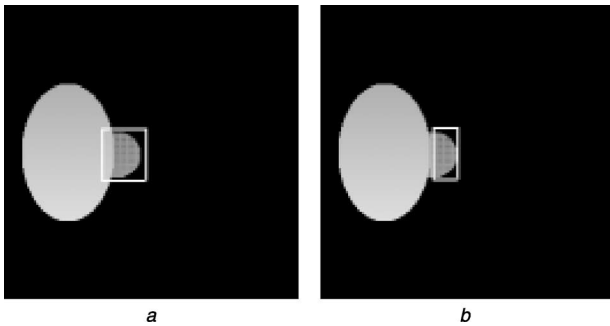


Fig. 23 Tracker outputs obtained with and without applying the partial occlusion handling scheme in a frame of the artificial image sequence

a With partial occlusion handling
b Without partial occlusion handling

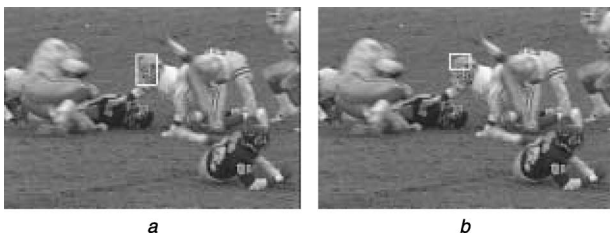


Fig. 24 Tracker outputs obtained with and without partial occlusion handling scheme in the 'Football' image sequence

a With partial occlusion handling
b Without partial occlusion handling

illumination changes, deformations of the tracked objects, abrupt motion or a combination of them.

6 Conclusions

This paper has presented an object tracking algorithm that is robust to partial and full occlusion. Information theory-based metrics were used as a reliability measure for algorithm initialisation and tracking procedures. The mutual information and Kullback–Leibler-based metrics provide the means to detect abrupt changes (partial occlusion, full occlusion or movement of the occluding object). Furthermore, motion detection of the tracked object is also possible in static scenes. Finally, an object verification process based on mutual information was also proposed and applied after object disocclusion. The use of information theory-based metrics combined with an occlusion handling scheme provide an object tracking algorithm that performs better than [17] in partial and total occlusion situations.

Experimental results have shown that the algorithm correctly detects and processes partial and total occlusion situations. The interpretation of variations of the proposed metrics may lead to a thorough understanding of the object tracking process in many computer vision applications.

The information theory-based metrics behave better in partial occlusion situations than the normalised correlation-based metric. A clear distinction in performance between the two information theory based metrics cannot easily be extracted. Nevertheless, mutual information, having the

advantage of being symmetrical and including spatial information, seems to be the preferred choice.

7 References

- Hager, G.D., and Belhumeur, P.N.: 'Efficient region tracking with parametric models of geometry and illumination', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1998, **20**, (10), pp. 1025–1039
- Peterfreund, N.: 'Robust tracking of position and velocity with kalman snakes', *IEEE Trans. Pattern Anal. Mach. Intell.*, 1999, **21**, (6), pp. 564–569
- Zhong, Yu., Jain, Anil K., and Dubuisson-Jolly, M.-P.: 'Object tracking using deformable templates', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2000, **22**, (5), pp. 544–549
- Uenohara, U., and Kanade, T.: 'Geometric invariants for verification in 3-d object tracking'. Proc. IEEE Int. Conf. Intelligent Robots and Systems, IROS 96, Osaka, Japan, 1996, Vol. 2, pp. 785–790
- Manku, S., Jain, P., Aggarwal, A., Kumar, L., and Banerjee, S.: 'Object tracking using affine structure for point correspondences'. Proc. IEEE Computer Society Conf. on Computer Vision and Pattern Recognition, San Juan, Puerto Rico, 1997, pp. 704–709
- Huttenlocher, D.P., Noh, J.J., and Rucklidge, W.J.: 'Tracking non rigid objects in complex scenes'. Proc. of the Int. Conf. on Computer Vision, Berlin, Germany, 1993, pp. 93–101
- Nguyen, H.T., Worring, M., and van den Boomgaard, R.: 'Occlusion robust adaptive template tracking'. Proc. of the Int. Conf. on Computer Vision, Vancouver, Canada, 2001, Vol. I, pp. 678–683
- Dockstader, S.L., and Tekalp, A.M.: 'Multiple camera tracking of interacting and occluded human motion', *Proc. IEEE*, 2001, **89**, (10), pp. 1441–1455
- Rasmussen, C., and Hager, G.D.: 'Probabilistic data association methods for tracking complex visual objects', *IEEE Trans. Pattern Anal. Mach. Intell.*, 2001, **23**, (6), pp. 560–576
- Schoepflin, T., Chalana, V., Haynor, D.R., and Kim, Y.: 'Video object tracking with a sequential hierarchy of template deformations', *IEEE Trans. Circuits Syst. Video Technol.*, 2001, **11**, (11), pp. 1171–1182
- Erdem, C., Tekalp, A.M., and Sankur, B.: 'Metrics for performance evaluation of video object segmentation and tracking without ground truth'. Proc. 2001 Int. Conf. on Image Processing, Thessaloniki, Greece, 2001, Vol. II, pp. 69–72
- Erdem, C.E., Tekalp, A.M., and Sankur, B.: 'Video object tracking with feedback of performance measures', *IEEE Trans. Circuits Syst. Video Technol.*, 2003, **13**, (4), pp. 310–324
- Fleet, D., Barron, J., and Beauchemin, S.: 'Performance of optical flow techniques', *Int. J. Comput. Vis.*, 1994, **12**, (1), pp. 43–77
- Viola, P., and Wells, W.M.: 'Alignment by maximization of mutual information', *Int. J. Comput. Vis.*, 1997, **24**, (2), pp. 137–154
- Kruppa, H., and Schiele, B.: 'Context-driven model switching for visual tracking'. Proc. 9th Int. Symp. on Intelligent Robotic Systems, Toulouse, France, 2001
- Kruppa, H., and Schiele, B.: 'Using mutual information to combine object models'. Proc. 8th Int. Symp. Intelligent Robotic Systems, Reading, UK, 2000
- Tomasi, C., and Kanade, T.: 'Shape and motion from image streams: a factorization method - Part 3 Detection and tracking of point features'. Technical report CMU-CS-91-132, Computer Science Department, Carnegie Mellon University, USA, 1991
- Birchfield, S.: 'Depth and motion discontinuities'. PhD thesis, Stanford University, USA, 1999
- Bregler, C., and Malik, J.: 'Tracking people with twists and exponential maps'. Proc. Int. Conf. on Computer Vision and Pattern Recognition, Santa Barbara, CA, USA, 1998
- Wang, J., and Adelson, E.: 'Representing moving images with layers', *IEEE Trans. Image Process.*, 1994, **3**, (5), pp. 625–638
- Shi, J., and Tomasi, C.: 'Good features to track'. Proc. Int. Conf. on Computer Vision and Pattern Recognition, Hawaii, USA, 2000, pp. 593–600
- Haykin, S.: 'Communication systems' (Wiley, New York, 1994, 3rd edn.)
- Reza, F.M.: 'An introduction to information theory' (Dover, New York, 1994)
- Skouson, M., Guo, Q., and Liang, Z.: 'A bound on mutual information for image registration', *IEEE Trans. Med. Imaging*, 2001, **20**, (8), pp. 843–846
- Do, M.N., and Vetterli, M.: 'Texture similarity measurement using kullback-leibler distance on wavelet subbands'. Proc. Int. Conf. on Image Processing, Vancouver, Canada, 2000
- Papoulis, A.: 'Probability, random variables, and stochastic processes' (McGraw-Hill, Inc, New York, 1991)
- Murat Tekalp, A.: 'Digital video processing' (Prentice Hall, New Jersey, 1995)