

This is a repository copy of *Prediction and Tracking of Moving Objects in Image Sequences*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/942/>

Article:

Bors, A G orcid.org/0000-0001-7838-0021 and Pitas, I (2000) Prediction and Tracking of Moving Objects in Image Sequences. IEEE Transactions on Image Processing. pp. 1441-1445. ISSN 1057-7149

<https://doi.org/10.1109/83.855440>

Reuse

Unless indicated otherwise, fulltext items are protected by copyright with all rights reserved. The copyright exception in section 29 of the Copyright, Designs and Patents Act 1988 allows the making of a single copy solely for the purpose of non-commercial research or private study within the limits of fair dealing. The publisher or other rights-holder may allow further reproduction and re-use of this version - refer to the White Rose Research Online record for this item. Where records identify the publisher as the copyright holder, users can verify any specific terms of use on the publisher's website.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

Prediction and Tracking of Moving Objects in Image Sequences

Adrian G. Borş and Ioannis Pitas

Abstract—We employ a prediction model for moving object velocity and location estimation derived from Bayesian theory. The optical flow of a certain moving object depends on the history of its previous values. A joint optical flow estimation and moving object segmentation algorithm is used for the initialization of the tracking algorithm. The segmentation of the moving objects is determined by appropriately classifying the unlabeled and the occluding regions. Segmentation and optical flow tracking is used for predicting future frames.

Index Terms—Bayes procedures, image sequence analysis, tracking.

I. INTRODUCTION

Tracking of moving objects is important for video surveillance while future frame prediction is used in video coding. A Bayesian approach shows that we can estimate the location of moving objects and their associated velocity based on a set of initial estimates. Occluding and unlabeled regions are identified and classified in the context of a tracking algorithm. A few approaches have been adopted for solving these problems. In [1] an occlusion adaptive mesh is used for tracking moving features over several frames. In other approaches, features are extracted from a set of frames and afterwards they are tracked over the sequence. Kalman filters have been used for tracking in [2]–[4]. Objects are segmented based on clustering in [3] and [5]. Simultaneous optical flow estimation and moving object segmentation has been employed in [6]. In this approach, the moving scene is modeled based on the median radial basis function (MRBF) network [8]. Each output unit of the neural network corresponds to a moving object. The results provided by the MRBF modeling are used for the initialization of a tracking algorithm. The unlabeled regions in each frame are identified and classified appropriately based on the MRBF model. When new objects enter in the scene or when some objects leave the scene, retraining is necessary. In between two MRBF retraining stages, tracking is employed for following object movement. Using tracking we predict the moving object optical flow and segmentation. A future frame is represented as the union of its predicted moving objects.

The Bayesian model for motion and segmentation estimation over the entire image sequence is provided in Section II. Tracking the moving objects over a set of frames is described in Section III and frame reconstruction based on estimating the moving object location and optical flow is described in Section IV. Simulation results are presented in Section V and the conclusions are drawn in Section VI.

II. MOTION AND SEGMENTATION ESTIMATION

Let us consider that each frame of an image sequence $f(t)$, $t = 1, \dots, K$ is made up of a set of moving regions $\{X_i(t), i = 1, \dots, N\}$

Manuscript received October 20, 1998; revised March 1, 2000. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Steven D. Blostein.

A. G. Borş was with the Department of Informatics, University of Thessaloniki, Thessaloniki 540 06, Greece. He is now with the Department of Computer Science, University of York, York YO10 5DD, U.K. (e-mail: adrian.bors@cs.york.ac.uk).

I. Pitas is with the Department of Informatics, University of Thessaloniki, Thessaloniki 540 06, Greece (e-mail: pitas@zeus.csd.auth.gr).

Publisher Item Identifier S 1057-7149(00)06130-3.

with the properties

$$f(t) = \cup_{i=1}^N X_i(t) \quad (1)$$

$$X_j(t) \cap X_k(t) = 0, \quad \forall j \neq k. \quad (2)$$

A subset $X_k(t)$ is associated to a five-dimensional representative vector $\mu_k = [S_k, \mathcal{M}_k]$, describing the optical flow \mathcal{M}_k and segmentation information S_k of a certain moving region [6]. The still image feature vector S_k contains the location and the characteristic graylevel of the moving region. S_k is directly related to the segmentation label of the moving region k , while $\mathcal{M}_k = [\mathcal{M}_{k,x}, \mathcal{M}_{k,y}]$ is the velocity vector of the respective moving region. The classification of the image sequence in moving objects is done according to the maximization of the *a posteriori* probability

$$P(\hat{\mu}_k(t), t = 1, \dots, K - 1 | f(t), t = 1, \dots, K) > P(\hat{\mu}_j(t), t = 1, \dots, K - 1 | f(t), t = 1, \dots, K) \quad (3)$$

where each probability corresponds to the segmentation of a moving object and its optical flow in the entire image sequence. After repeatedly applying the Bayes theorem and after expressing the probabilities from one frame with respect to those corresponding to the previous frames, we obtain

$$\begin{aligned} &P(\hat{\mu}_k(t), t = 1, \dots, K - 1 | f(t), t = 1, \dots, K) \\ &= \prod_{t=p}^{K-1} [P(f(t+1) | \hat{\mu}_k(j), f(j), j = 1, \dots, t)] \\ &\quad \times \prod_{t=p}^{K-1} [P(\hat{\mu}_k(t) | \hat{\mu}_k(j), f(j), j = 1, \dots, t-1, f(t))] \\ &\quad \times \frac{P(\hat{\mu}_k(j), j = 1, \dots, p-1 | f(j), j = 1, \dots, p)}{\prod_{t=p}^{K-1} P(f(t+1) | f(j), j = 1, \dots, t)} \end{aligned} \quad (4)$$

where K is the number of frames and p is a given frame $p < K$. A component of the first probability product from the right side of this relationship is associated to the reconstruction of a frame from the previous frames using the moving object feature vectors. A component of the second probability product corresponds to the feature vector tracking over several frames. The third probability factor models the moving object characteristics derived from the first p frames. The probabilities from the denominator denote the dependence of a frame on the previous ones and it can be neglected in the following considerations.

In the following, we show how to initialize the algorithm which estimates the probabilities from (4). The first two frames are split in blocks and a feature vector denoted as $\mathbf{u}_{IJ} = [I, J, l]$ containing the site location $[I, J]$, the graylevel l and the motion vector is associated with each block. For $p = 2$, after using the Bayesian theorem, the third probability factor in (4), can be further described as

$$\begin{aligned} &P(\hat{\mathcal{M}}_j, \hat{S}_j | f(2), f(1)) = P(f(2) | \hat{\mathcal{M}}_j, \hat{S}_j, f(1)) \\ &\quad \times \frac{P(\hat{\mathcal{M}}_j | \hat{S}_j, f(1)) P(\hat{S}_j | f(1))}{P(f(2) | f(1))} \end{aligned} \quad (5)$$

where $P(\hat{S}_j | f(1))$ represents the *a priori* probability of the segmentation and $P(\hat{\mathcal{M}}_j | \hat{S}_j, f(1))$ is the probability of the optical flow estimation depending on the segmentation map and image [7]. After expressing each probability as an energy function, we model them with Gaussian functions. The Gaussian function associated with a moving region and implemented by a hidden unit of the MRBF network is given by

$$\begin{aligned} \phi_j(\mathbf{u}_{IJ}) = \exp \left[-(\mathbf{u}_{IJ} - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1} (\mathbf{u}_{IJ} - \hat{\mu}_j) \right. \\ \left. - \text{WDFD}(\hat{\mathcal{M}}_j) \right] \end{aligned} \quad (6)$$

where $\hat{\mu}_j$ and $\hat{\Sigma}_j$ are the center vector and covariance matrix estimates and WDFD($\hat{\mathcal{M}}_j$) represents the weighted displaced frame difference (a measure of confidence in the motion estimation algorithm) [6]. An unsupervised training algorithm provides the estimates of the MRBF network parameters while modeling the probabilities from (5) [6], [8].

III. MOVING OBJECT TRACKING

Let us neglect the dependence on all the frames excepting the previous. In this case we can express each probability in the first product of (4) as an energy function measuring the accuracy of reconstructing the frame $f(t+1)$ from the displaced moving objects which had been segmented in the frame $f(t)$

$$P(f(t+1) | \hat{\mu}(t), f(t)) \approx \frac{1}{Z} \exp \left\{ -\delta \left[\cup_i^N (X_i(t) \oplus \hat{\mathcal{M}}_i(t)), f(t+1) \right] \right\} \quad (7)$$

where Z is a normalizing constant, $X_i(t) \oplus \hat{\mathcal{M}}_i(t)$ represents the translation of the moving region $X_i(t)$ resulted from the segmentation of the frame $f(t)$ with its corresponding motion vector $\hat{\mathcal{M}}_i(t)$ and $\delta[f(t), g(t)]$ represents a function which counts in how many locations $f(t)$ and $g(t)$ have a different segmentation level. The maximization of this probability represents the minimization of the difference between the given frame and its prediction based on the previous frame segmentation and its estimated optical flow. It can be observed that by displacing the set of pixels $X_i(t)$ representing a moving region in the frame $f(t)$, certain pixels from $X_i(t) \oplus \hat{\mathcal{M}}_i(t)$ have uncertain assignment. When regions from one frame do not have a correspondent in the next frame (uncovered regions), (1) is not respected any more. When two or more different objects project in the same region of the next frame (occluding regions), (2) is not valid. Both situations occur in regions located at the margins of the moving objects and can be easily identified as providing a probability equal or smaller than $\exp(-1)/Z$ in (7). If we have a one-to-one correspondence between the frames $f(t)$ and $g(t)$ based on the given model then the probability from (7) is equal to $1/Z$. After detecting the unlabeled regions, we estimate their feature vectors \mathbf{u}_{IJ} considering only the likely correlations given by the motion vectors of the neighboring moving objects. The trained MRBF network, can be applied in a multiresolution approach where the network parameters obtained from the initial block-based segmentation are used for image segmentation at pixel resolution [6]. We apply the already trained MRBF network only in the regions decided as uncertain according to (7).

The components of the second product from the expression (4) representing the dependency of a feature vector on the values of the same feature vector in the previous frames, can be expressed as in maximum-likelihood regression estimation [9]

$$P(\hat{\mu}_k(t) | \hat{\mu}_k(j), f(j), j = 1, \dots, t-1, f(t)) = \frac{1}{Z} \exp \left[- \left| \hat{\mu}_k(t) - \sum_{i=1}^M W_i \psi_i(\hat{\mu}_k) \right| \right] \quad (8)$$

where $\psi(\hat{\mu}_k)$ are a set of functions modeling the variation of the k th object feature vector in time, W_i are their associated weights, M is the number of previous frames used for feature estimation, and Z is a normalizing constant. However, in most of the cases, moving objects have slow changing motion, which can be modeled by a linear system. Under this assumption, the model (8) can be simplified

$$P(\hat{\mu}_k(t) | \hat{\mu}_k(j), f(j), j = 1, \dots, t-1, f(t)) = \frac{1}{Z} \exp \left[- \left| \hat{\mu}_k(t) - \mathbf{W}_k \Phi_k^T \right| \right] \quad (9)$$

where Φ_k consists of the feature vectors from the last M frames

$$\Phi_k = [\hat{\mu}_k(t-1) \quad \hat{\mu}_k(t-2) \quad \dots \quad \hat{\mu}_k(t-M)] \quad (10)$$

and \mathbf{W}_k is a matrix of size $5M \times 5$ whose entries represent the dependency of a feature vector component at time t with respect to all feature entries in the previous M frames. The features that are tracked over time correspond to object location, graylevel changes and optical flow. The components of the matrix \mathbf{W}_k can be found by using the least mean squares (LMS) algorithm [10]. LMS is a fast on-line algorithm which can ensure feature tracking over several frames based on minimizing the prediction mean square error. Kalman filters can be seen as an extension of the LMS algorithm which however requires a much larger computational complexity. Changes in the moving object representative vectors are reflected in the moving object segmentation. In order to maximize the probability in (4), we should maximize its components from (5), (7), and (9). The relationship (5) provides the initial estimate, while (9) gives an estimate of the moving object feature vector from its previous values. This estimate must be consistent with an accurate frame reconstruction as given by (7).

IV. FRAME RECONSTRUCTION FROM MOVING OBJECT PREDICTION

A prediction function provides an estimate of the moving object segmentation and its corresponding optical flow in a future frame based on the data extracted from the previous frames. Let us denote by $\pi_t(X_k(t+1))$ and $\pi_t(\hat{\mathcal{M}}_k(t+1))$ the prediction of the location for the moving region k and the prediction of its optical flow respectively, from the frame t into the frame $t+1$. The prediction function for the velocity uses the matrix \mathbf{W}_k , derived from the maximization of the probability from (9). The optical flow for a certain moving object is predicted for each consecutive frame by using the dependency on its previous values

$$\pi_t(\hat{\mathcal{M}}_{k,x}(t+1)) = \mathbf{W}_{xx} \hat{\mathcal{M}}_{k,x} + \mathbf{W}_{yx} \hat{\mathcal{M}}_{k,y} \quad (11)$$

$$\pi_t(\hat{\mathcal{M}}_{k,y}(t+1)) = \mathbf{W}_{xy} \hat{\mathcal{M}}_{k,x} + \mathbf{W}_{yy} \hat{\mathcal{M}}_{k,y} \quad (12)$$

where $\hat{\mathcal{M}}_{k,x}$, $\hat{\mathcal{M}}_{k,y}$ represent the motion vector components on x and y directions associated with the k -th moving object for the last M frames and \mathbf{W}_{xy} , \mathbf{W}_{xx} , \mathbf{W}_{yx} , \mathbf{W}_{yy} are their corresponding weighting vectors found by the LMS algorithm [10] as in (9). This prediction function can easily model complex movements such as rotation and acceleration. The number of frames M to be taken into account for the prediction system must be larger when the motion is smooth and smaller when the motion is fast changing. Similarly to (11) or (12), we can derive a prediction system for the luminance by tracking the change in the average graylevel of a certain moving object.

The location of a moving object in a future frame is given by the segmentation in the actual frame and the prediction of its associated optical flow

$$\pi_t(X_k(t+1)) = X_k(t) \oplus \pi_t(\hat{\mathcal{M}}_k(t+1)) \quad (13)$$

where we consider the displacement for all the pixels composing the moving object k , and where $\pi_t(\hat{\mathcal{M}}_k(t+1))$ components are derived in (11) and (12). Given a prediction function for the optical flow associated with the moving object k , we can predict the frame $t+1$ considering the segmentation of the individual objects

$$\hat{f}(t+1) = \cup_{k=1}^N \pi_t(X_k(t+1)) \quad (14)$$

where $\hat{f}(t+1)$ is the predicted image. As it was shown in the previous section, certain regions do not have a clear assignment. Such regions are classified based on an overlapping priority assumption. For example,



Fig. 1. First frame of the "Hamburg taxi" image sequence.

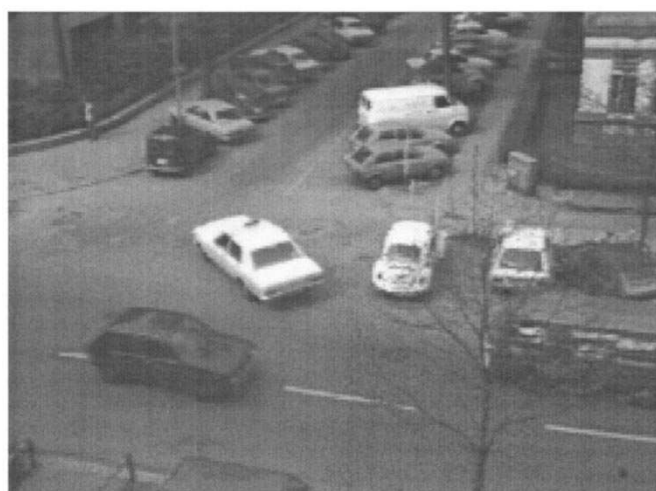


Fig. 2. Twentieth frame of the "Hamburg taxi" image sequence.

if the background is known, it will get the lowest priority and it will be covered in the case of moving objects pointing to the same region, or it will fill in the regions which remained uncovered. The values to be used in the unlabeled regions of the predicted frame are taken from one of the previous frames by considering the optical flow.

The PSNR between the predicted frame $f_p(t + 1)$ and the real one $f(t + 1)$, when it is available, is considered for checking the validation of the assumed model

$$PSNR = 20 \log_{10} \left(\frac{255R \times S}{\sqrt{\sum_{i=1, j=1}^{R, S} (\hat{f}_{ij}(t + 1) - f_{ij}(t + 1))^2}} \right) \quad (15)$$

where $R \times S$ is the size of the image. If the PSNR between the two images is below a certain threshold, then the model is not valid at the respective frame. Usually, this is caused because a moving object enters or leaves the scene. In such a case, the MRBF network is retrained in order to obtain the appropriate moving object segmentation and optical flow (5) [6]. The new model is tracked over the following frames as described in the previous section.

V. SIMULATION RESULTS

We provide simulation results when the proposed algorithm is applied in the "Hamburg taxi" image sequence. The first and the 20th frames are displayed in Figs. 1 and 2. In the center of a frame from this image sequence a white taxi turns around the corner, a black car moves from left to right while a van moves from right to left. The moving object segmentation as provided by the MRBF network for the first frame is shown in Fig. 3. Its corresponding optical flow is provided in Fig. 4. The segmentation and optical flow parameters are used for the initialization of the tracking algorithm. The occluding and unlabeled regions for the first frame are shown in Fig. 5. They are located at the moving object boundaries according to a small local frame reconstruction probability in (7). The pixels of these regions are classified using the MRBF network parameters. The moving object segmentation resulted after this classification is displayed in Fig. 6. After tracking the moving objects as described in Section III, we obtain the segmentation of the 20th frame, as provided in Fig. 7. Six past frames ($M = 6$) have been used for tracking. It can be observed that the segmentation of the white taxi in the center of the frame is quite good despite the fact that, due to the three-dimensional perspective view, its projection changes

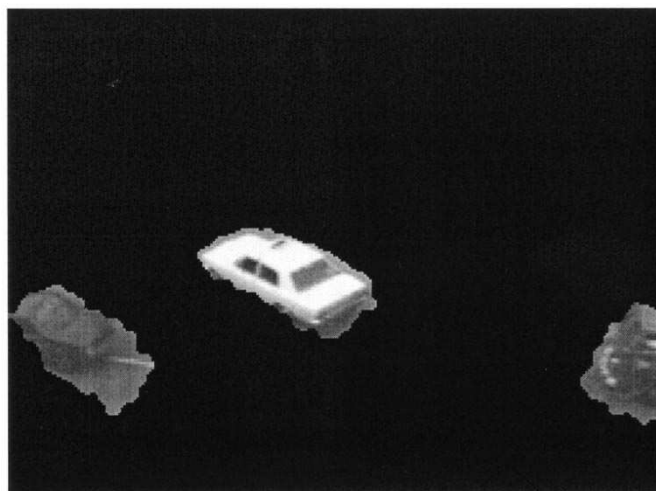


Fig. 3. Moving object segmentation.

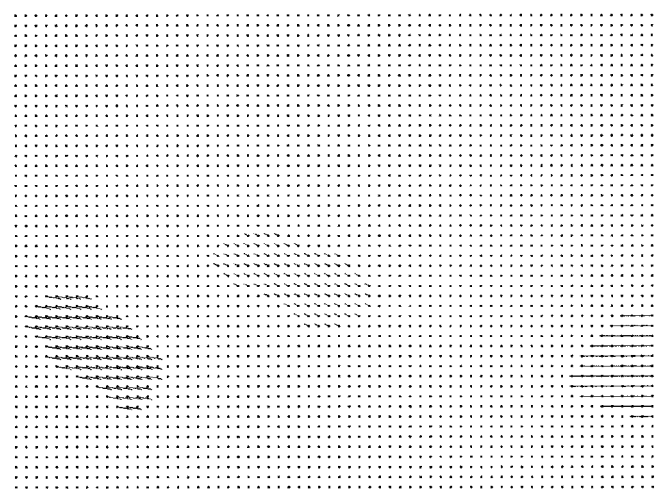


Fig. 4. The optical flow of the first frame from the "Hamburg taxi" image sequence.

while turning around the corner. The optical flow corresponding to the tracked objects in the 20th frame is represented in Fig. 8. The predicted 20th frame, reconstructed from the predicted segmentation and moving

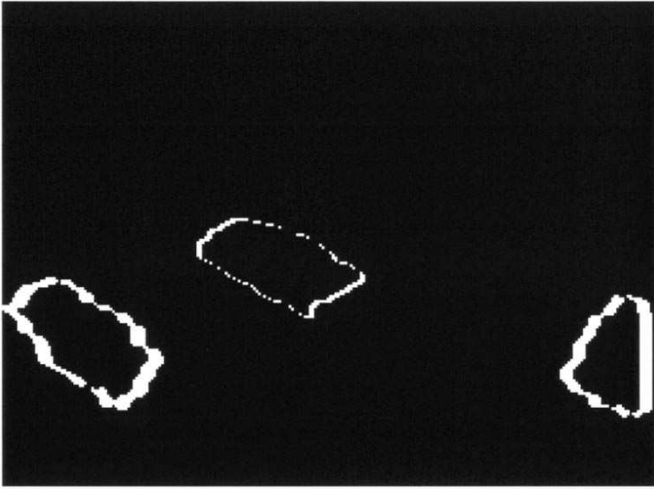


Fig. 5. Occluding and unlabeled regions.

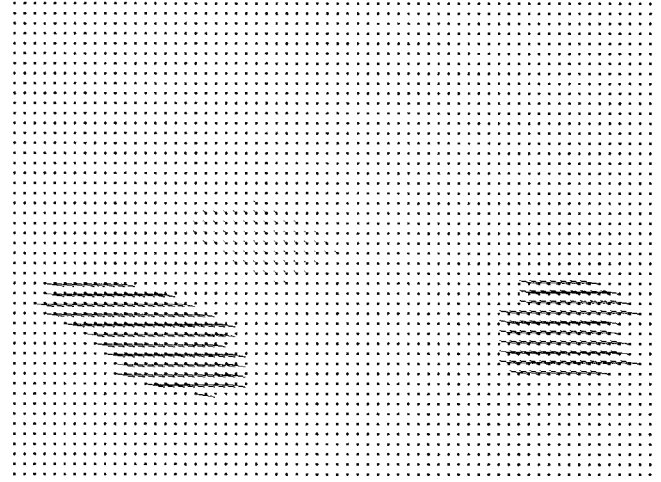


Fig. 8. Estimated optical flow of the 20th frame.

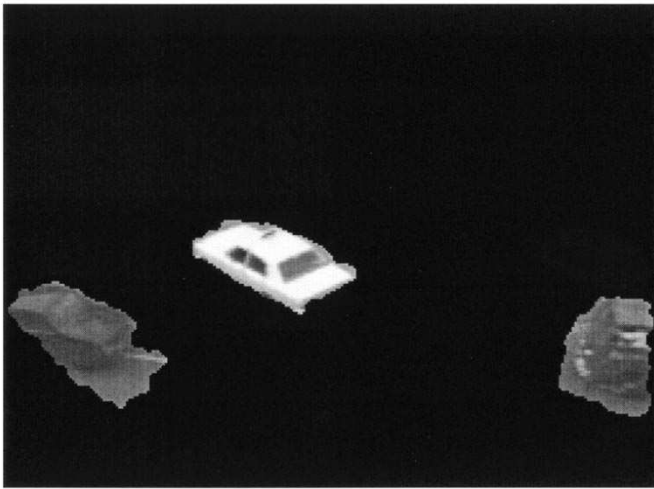


Fig. 6. Segmentation of the moving objects after appropriately classifying the occluding regions.

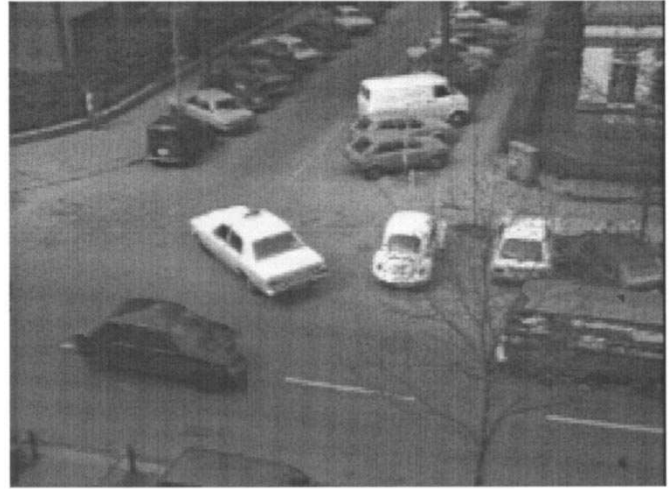


Fig. 9. Predicted 20th frame.

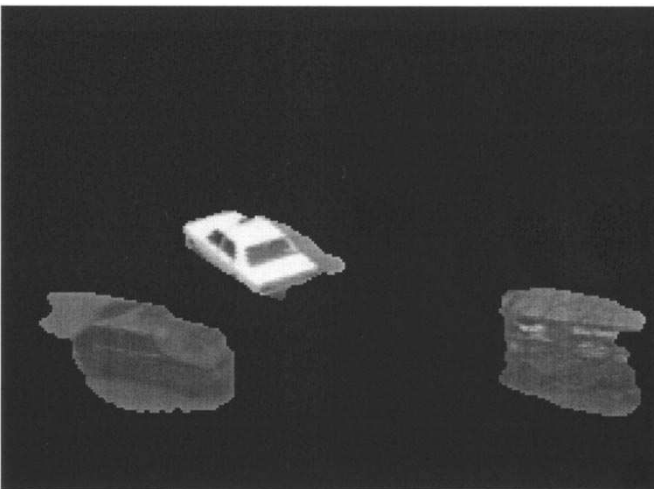


Fig. 7. Moving object segmentation after tracking them 20 frames.



Fig. 10. Difference between the predicted and the real 20th frame from the "Hamburg taxi" image sequence.

object velocities, is provided in Fig. 9. The difference between the predicted and the real 20th frame is shown in Fig. 10. We can see from this Figure that many errors in the prediction of the 20th frame are due to

changes in illumination. In Fig. 11, the PSNR of the predicted image when tracking the moving objects is plotted for a set of frames. The MRBF network training took 33.3 s when using a Silicon Graphics

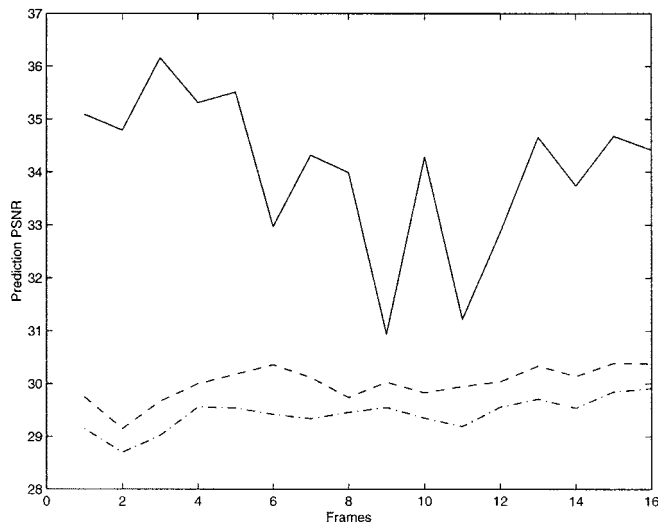


Fig. 11. PSNR of the predicted frame in the "Hamburg taxi" image sequence. "-" denotes the PSNR of the proposed tracking algorithm. "- -" represents the PSNR prediction considering the initial MRBF model on 4×4 pixel blocks. "..." denotes the PSNR between the actual frame and that used for prediction.

Indy Workstation. The trained network, can be used for those successive frames which match the model according to a criterion [6]. In this case, 95 s are required for segmenting the moving objects and the optical flow for 20 frames when using 4×4 pixel blocks. When employing tracking as described in this study, only 68 s are necessary for the same frames using pixel resolution segmentation. In the first case only 3040 vectors had been processed while in the second case their number was 48 640. The segmentation provided by the tracking algorithm is quite good as it can be observed from the experimental results and provides a good basis for prediction-based frame reconstruction. The prediction PSNR of the tracking algorithm is better than when considering the initial MRBF model for segmenting all the frames and assuming just the previous moving object features for reconstruction, as it can be observed from Fig. 11.

VI. CONCLUSION

We propose a moving object tracking algorithm derived from the Bayesian theory. The optical flow and the segmentation features are jointly modeled by the MRBF network in the initial stage. The occluding and unlabeled regions are detected and classified appropriately. The proposed algorithm provides good moving object tracking capabilities. Such capabilities are used for segmenting and estimating the moving object velocity and segmentation in a future frame. The proposed algorithm is employed for frame prediction.

REFERENCES

- [1] Y. Altunbasak and A. M. Tekalp, "Occlusion-adaptive, content-based mesh design and forward tracking," *IEEE Trans. Image Processing*, vol. 6, pp. 1270–1280, Sept. 1997.
- [2] K. J. Bradshaw, I. D. Reid, and D. W. Murray, "The active recovery of 3-D motion trajectories and their use in prediction," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 19, pp. 219–224, Mar. 1997.
- [3] S. M. Smith and J. M. Brady, "ASSET-2: Real-time motion segmentation and shape tracking," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 17, pp. 814–820, Aug. 1995.
- [4] J. Weber and J. Malik, "Rigid body segmentation and shape description from dense optical flow under weak perspective," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 19, pp. 139–143, Feb. 1997.

- [5] A. Kumar, Y. B.-Shalom, and E. Oron, "Precision tracking based on segmentation with optimal layering for imaging sensors," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 17, pp. 182–188, Feb. 1995.
- [6] A. G. Borş and I. Pitas, "Optical flow estimation and moving object segmentation based on median radial basis function network," *IEEE Trans. Image Processing*, vol. 7, pp. 693–702, May 1998.
- [7] J. Konrad and E. Dubois, "Bayesian estimation of motion vector fields," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 14, pp. 910–927, Sept. 1992.
- [8] A. G. Borş and I. Pitas, "Median radial basis function neural network," *IEEE Trans. Neural Networks*, vol. 7, pp. 1351–1364, Nov. 1996.
- [9] V. Vapnik, *Estimation of Dependences Based on Empirical Data*. New York: Springer-Verlag, 1982.
- [10] B. Widrow and S. D. Stearns, *Adaptive Signal Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1985.

Tomographic Reconstruction Using Nonseparable Wavelets

Stéphane Bonnet, Françoise Peyrin, Francis Turjman, and Rémy Prost

Abstract—In this paper, the use of nonseparable wavelets for tomographic reconstruction is investigated. Local tomography is also presented. The algorithm computes both the quincunx approximation and detail coefficients of a function from its projections. Simulation results showed that nonseparable wavelets provide a reconstruction improvement versus separable wavelets.

Index Terms—Local tomography, McClellan transformation, nonseparable wavelets.

I. INTRODUCTION

Computerized tomography (CT) consists of recovering a function from a set of its projections and relies on the inversion of the Radon transform. According to the nature of the data set, this problem may be ill-posed. The use of wavelets for inverse problems in general, and CT in particular, presents several interesting features to stabilize the inversion process [1]. As a matter of fact, wavelets may bring valuable solutions to the problem of local tomography [2]–[4].

The relationships between the continuous wavelet transform and the Radon transform have first been established in several independent works [5], [6]. Olson was the first to devise a reconstruction scheme from a customized sampling of the Radon transform [2]. Delaney [3] and Rashid-Farrokhi [4] proposed a multiresolution tomographic reconstruction algorithm to recover the two-dimensional (2-D) separable discrete wavelet transform (2-D DWT) of the image from its projections, and applied it to local tomography. Both algorithms are based on 2-D wavelets, constructed by tensor products of one-dimensional (1-D) wavelets. The 2-D separable wavelets impose a rectangular tiling of the

Manuscript received December 28, 1998; revised February 27, 2000. S. Bonnet was supported by a grant from Siemens, France. This work is in the scope of the scientific topics of the PRC-GDR ISIS research group of the French National Center for Scientific Research (CNRS). The associate editor coordinating the review of this manuscript and approving it for publication was Prof. William Clem Karl.

S. Bonnet and R. Prost are with CREATIS, CNRS Research Unit, 69621 Villeurbanne Cedex, France (e-mail: bonnet@creatis.insa-lyon.fr).

F. Peyrin is with CREATIS, CNRS Research Unit, 69621 Villeurbanne Cedex, France. She is also with ESRF, 38043 Grenoble Cedex, France.

F. Turjman is with CREATIS, CNRS Research Unit, 69621 Villeurbanne Cedex, France. He is also with Hôpital Neurologique, 69500 Bron, France.

Publisher Item Identifier S 1057-7149(00)06139-X.