# Optical Flow Estimation and Moving Object Segmentation Based on Median Radial Basis Function Network

Adrian G. Borş and Ioannis Pitas

*Abstract*— **Various approaches have been proposed for simultaneous optical flow estimation and segmentation in image sequences. In this study, the moving scene is decomposed into different regions with respect to their motion, by means of a pattern recognition scheme. The inputs of the proposed scheme are the feature vectors representing still image and motion information. Each class corresponds to a moving object. The classifier employed is the median radial basis function (MRBF) neural network. An error criterion function derived from the probability estimation theory and expressed as a function of the moving scene model is used as the cost function. Each basis function is activated by a certain image region. Marginal median and median of the absolute deviations from the median (MAD) estimators are employed for estimating the basis function parameters. The image regions associated with the basis functions are merged by the output units in order to identify moving objects.**

*Index Terms*—**Energy minimization, median radial basis function neural network, moving object segmentation, optical flow estimation, robust training.**

## I. INTRODUCTION

**M**OTION representation and modeling are important steps toward the understanding of dynamic images. From the point of view of image analysis, the algorithms used for extracting the optical flow can be classified as gradient-based and feature-based. With respect to the method employed, they can be either deterministic or stochastic, local or global. Gradient-based techniques rely on the spatio-temporal differential equation describing the motion [1]. The methods based on feature matching detect the representative features from the image sequence and try to match them [2]. Block matching algorithms, widely used in video coding [3], estimate the motion field based on the correlation between each block from one frame and a block from the reference frame.

Scene segmentation formulation from optical flow, based on *a posteriori* criterion, was introduced in [4]. A more general framework for motion and line field estimation based on Markov random fields (MRF's) and Gibbs distributions [5] was used in [6]. Simulated annealing was employed in these approaches for estimating the optical flow and line field.

Iterated conditional modes (ICM's) were proposed for the restoration of color images in [8] and was applied to motion field smoothing in [9], [10]. The ICM is a deterministic local optimization approach used to minimize an energy function by employing constraints and can be viewed as a decision relaxation labeling algorithm [7]. This method was extended to take into account the moving region borders, which should not be smeared by the smoother [11], [12]. ICM relies on a piecewise MRF-based model for the optical flow, associating alternatively various motion vectors and segmentation labels with a given block site. The motion vector and the segmentation label are chosen in each iteration so that the displaced frame difference decreases and a smoothness constraint is fulfilled at the same time. An averaging operation is performed inside a given neighborhood, only in the regions considered to make up a moving object. ICM was modified to include the edge geometry and a gradient-based equation [13]. ICM performance depends on good initialization [12], [13]. The $k$-means clustering algorithm was employed for this purpose in [15].

Estimation of the optical flow based on clustering analysis was used in [16] and [17]. Image intensity information was incorporated in the clustering analysis algorithm in [18]. A robust regression technique where still image segmentation and motion estimation are performed separately was analyzed in [19]. Algorithms that perform simultaneous motion estimation and scene segmentation, in an iterative way, based on the minimization of an energy function, were provided in [12] and [20]. The approach suggested in [12] is based on local optimization methods like ICM and "highest confidence first," and provides improvements over the segmentation on the basis of a separate estimation. A gradient-based algorithm derived from the Mean field theory using global optimization was introduced in [20]. Some of these algorithms have been embedded in multiresolution structures [6], [12], [13], [20].

In this study, a cost function associated with the minimization of a global criterion is proposed for simultaneous estimation of the optical flow and segmentation of moving objects. The image is partitioned in block sites situated on a rectangular lattice. Each block site is associated with a five-dimensional (5-D) feature vector describing the position, the gray level, and the local motion information. The proposed method is based on the unsupervised classification of the feature vectors by considering the displaced frame difference as well. The classification is done according to a decision

criterion derived from the Bayesian theory and representing a distance in the parameter space [21]. The moving scene is split accordingly in moving regions. We consider the MRBF network for modeling the optical flow and moving object segmentation in the image sequence. This structure is embedded in a two-layer feedforward neural network, where each output is assigned to a moving object.

A radial basis function (RBF) decomposition is known to be a good functional approximator and has been used in many applications [22]–[27]. The first-layer units implement Gaussian functions. The classification criterion connects the Gaussian parameters to the set of feature vectors drawn from the image sequence. In the second layer, the moving regions associated to the basis functions are merged in order to model the moving objects. The mixture of basis functions approximates the probabilities associated to the optical flow estimation and segmentation of the moving objects.

The learning algorithm employed for estimating the set of parameters associated to the image sequence has two stages [26], [27]. In the first stage, the basis function center and variance are found based on a clustering approach, similar to learning vector quantization (LVQ) [28]. A robust statistics-based algorithm employing the marginal median and median of the absolute deviations from the median (MAD) [29], [30] and called MRBF was proposed for training the network in [27]. The number of moving objects does not need to be specified *a priori*. It is found according to a compactness measure. The backpropagation algorithm, applied for calculating the output weights of the RBF network [22], [31], is employed in the second learning stage.

In Section II the derivation of the classification criterion is presented. Section III describes the learning algorithm. In Section IV simulation results showing the effectiveness of the proposed method are provided. The conclusions of the present study are drawn in Section V.

## II. THE CLASSIFICATION CRITERION

### A. Feature Extraction

Let us consider an image sequence $f_t$. Each frame is partitioned in block sites $B_{IJ}$, for $I = 0, \cdots, n_x - 1$ and $J = 0, \cdots, n_y - 1$, situated on a rectangular grid. We associate a feature vector denoted as $\mathbf{u}_{IJ}$, describing the local image sequence properties, to each block site $B_{IJ}$. This vector contains a still image feature vector $\mathbf{S}_{IJ}$ and a motion vector $\mathbf{M}_{IJ}$:

$$\mathbf{u}_{IJ} = [\mathbf{S}_{IJ}, \mathbf{M}_{IJ}]. \tag{1}$$

The still image feature vector includes the block site coordinates and its mean gray level

$$\mathbf{S}_{IJ} = [I, J, l_{IJ}] \tag{2}$$

where the mean gray level is given by

$$\hat{l}_{IJ} = \frac{1}{D_x D_y} \sum_{h=-\frac{D_y}{2}}^{\frac{D_y}{2}} \sum_{g=-\frac{D_x}{2}}^{\frac{D_x}{2}} f(I + h, J + g). \tag{3}$$

$f(I + h, J + g)$ denotes the gray level values of the image elements and $D_x \times D_y$ is the block size. In the case of color image sequences, the color components are included in the feature vector $\mathbf{S}_{IJ}$ as well. The motion vector contains two components, corresponding to the local motion of the block

$$\mathbf{M}_{IJ} = [m_{x,IJ}, m_{y,IJ}]. \tag{4}$$

They can be calculated by any optical flow algorithm, e.g., by block matching. Block matching methods use the correlation of a given image block of size $D_x \times D_y$ from the frame $f_{t-1}$ with blocks of the same size, situated inside a search region $S_x \times S_y$, on the next frame $f_t$ [3]. The motion vector that minimizes the displaced frame difference (DFD) is chosen. The DFD is given by

$$d_{IJ}(\hat{\mathbf{M}}_{IJ}) = \sum_{k=-\frac{D_x}{2}}^{\frac{D_x}{2}} \sum_{l=-\frac{D_y}{2}}^{\frac{D_y}{2}} |f_t(I + k, J + l) - f_{t-1}(I + k + \hat{m}_{x,IJ}, J + l + \hat{m}_{y,IJ})| \tag{5}$$

$\hat{m}_{x,IJ}$, $\hat{m}_{y,IJ}$ are the motion component estimates.

### B. Moving Object Classification

A moving scene can be seen as made up of regions having different motion parameters. We assume that each frame can be segmented into $L$ subsets, forming moving regions, denoted as $X_1, \cdots, X_L$. The moving objects are considered as compact moving entities, consisting of one or more moving regions. Each moving object is assigned a class. The partition of the moving objects in moving regions depends on the type of motion they employ and their local properties.

Each subset $X_k$ is associated to a 5-D representative vector $\mu_k$, describing the optical flow and segmentation information associated with a certain moving region

$$\mu_k = [\mathcal{S}_k, \mathcal{M}_k]. \tag{6}$$

The still image feature vector $\mathcal{S}_k$ is directly related to the segmentation label of the moving region $k$.

Let us denote by $\hat{\mathcal{S}}_k$ and $\hat{\mathcal{M}}_k$ the estimates of the segmentation label and optical flow associated to the moving region $k$. A block site $B_{IJ}$ is considered as belonging to a moving region $k$, $B_{IJ} \in X_k$, if it maximizes the *a posteriori* probability of the optical flow $\hat{\mathcal{M}}_k$ and moving region segmentation $\hat{\mathcal{S}}_k$ denoted as $P(\hat{\mathcal{M}}_k, \hat{\mathcal{S}}_k \mid f_{t-1}, f_t)$, when compared with the probabilities associated to the other moving regions

$$P(\hat{\mathcal{M}}_k, \hat{\mathcal{S}}_k \mid f_{t-1}, f_t) > P(\hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j \mid f_{t-1}, f_t) \tag{7}$$

for $j = 1, \cdots, L$, $j \neq k$, where $L$ is the number of moving regions.

In the second stage, the moving regions are merged, in order to describe moving objects, based on a neighboring criterion. Let us denote by $\hat{\mathcal{T}}_k$ the estimate of the optical flow and of the segmentation associated with a moving object. We consider a neighboring measure $V(X_l, X_k)$ between two subsets representing two moving regions, $X_l$ and $X_k$

$$V(X_l, X_k) = V(X_k, X_l) = \sum_{B_{IJ} \in X_l} |\mathcal{N}_{IJ} \cap X_k| \tag{8}$$

where $|\mathcal{N}_{IJ} \cap X_k|$ represents the number of block sites of a moving region $X_k$, which are situated in a certain neighborhood $\mathcal{N}_{IJ}$ of the block site $B_{IJ}$, belonging to the moving region $X_l$ [32]. This measure represents the boundary length between two moving regions. If two moving regions $X_k$ and $X_l$ do not have any common boundary, then $V(X_l, X_k) = 0$. We define a moving object as a moving region which contains a compact area in the image. In this case, the probability of estimating the optical flow and moving object segmentation $P(\hat{\mathcal{T}}_k \mid f_{t-1}, f_t)$ is

If $V(X_k, X_k) = \max\limits_{i=1}^{L} V(X_k, X_i)$ then

$$P(\hat{\mathcal{T}}_k \mid f_{t-1}, f_t) = P(\hat{\mathcal{M}}_k, \hat{\mathcal{S}}_k \mid f_{t-1}, f_t). \qquad (9)$$

A moving object $k$ contains a moving region $X_l$

If $V(X_l, X_k) = \max\limits_{i=1}^{L} V(X_l, X_i)$ then

$$\begin{aligned} P(\hat{\mathcal{T}}_k \mid f_{t-1}, f_t) &= \lambda_k P(\hat{\mathcal{M}}_k, \hat{\mathcal{S}}_k \mid f_{t-1}, f_t) \\ &\quad + \lambda_l P(\hat{\mathcal{M}}_l, \hat{\mathcal{S}}_l \mid f_{t-1}, f_t) \end{aligned} \qquad (10)$$

where $\lambda_k, \lambda_l$ are the weights denoting the contribution of each moving region probability to the moving object probability. This condition can be extended for moving objects containing many moving regions.

Let us express the *a posteriori* probabilities from (7) with respect to the features extracted from the image sequence. After applying Bayes' rule, each of the *a posteriori* distributions in (7) can be factored as follows:

$$\begin{aligned} &P(\hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j \mid f_t, f_{t-1}) \\ &= \frac{P(f_t \mid f_{t-1}, \hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j) P(\hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j \mid f_{t-1})}{P(f_t \mid f_{t-1})} \\ &= \frac{P(f_t \mid f_{t-1}, \hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j) P(\hat{\mathcal{M}}_j \mid \hat{\mathcal{S}}_j, f_{t-1}) P(\hat{\mathcal{S}}_j \mid f_{t-1})}{P(f_t \mid f_{t-1})}. \end{aligned} \qquad (11)$$

where $P(\hat{\mathcal{S}}_j \mid f_{t-1})$ represents the *a priori* probability of the segmentation and $P(\hat{\mathcal{M}}_j \mid \hat{\mathcal{S}}_j, f_{t-1})$ is the probability of the optical flow estimation depending on the segmentation map and image [6]. The probability $P(f_t \mid f_{t-1})$ does not depend on $[\hat{\mathcal{S}}_j, \hat{\mathcal{M}}_j]$ and can be neglected.

Each of the above conditional probabilities can be expressed as an energy function $E(\mathbf{X})$

$$P(\mathbf{X}) = \frac{1}{Z} \exp\left[ -\frac{E(\mathbf{X})}{\beta} \right] \qquad (12)$$

where $Z$ is a normalizing constant and $\beta$ is a constant controlling the properties of $E(\mathbf{X})$ [5]. Thus, the probability estimation problem (7) is converted into the minimization of an energy function, derived from (11) and (12):

$$\begin{aligned} E_j &= E(f_t \mid f_{t-1}, \hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j) + E(\hat{\mathcal{M}}_j \mid \hat{\mathcal{S}}_j, f_{t-1}) \\ &\quad + E(\hat{\mathcal{S}}_j \mid f_{t-1}). \end{aligned} \qquad (13)$$

The first component of the energy function is related to the estimation of the frame $f_t$ based on the previous frame $f_{t-1}$, the optical flow and the segmentation. The second component represents the energy associated to the estimated

optical flow and the third the energy associated to the estimated segmentation. In order to minimize $E_j$, all three components should be simultaneously minimized. This corresponds to the simultaneous optical flow and frame segmentation.

### C. The Cost Function

The performance criterion is related to the total squared error minimization in the feature space [33]. The energy function in (13) is expressed as a clustering metric in the feature space. This metric relates the moving region feature vectors (6) to the block site feature vectors (1).

The energy $E(f_t \mid f_{t-1}, \hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j)$ in (13) is represented as a weighted function of the displaced frame difference corresponding to the moving region $j$ and denoted as $\mathrm{WDFD}(\hat{\mathcal{M}}_j)$:

$$\begin{aligned} E(f_t \mid f_{t-1}, \hat{\mathcal{M}}_j, \hat{\mathcal{S}}_j) &= \mathrm{WDFD}(\hat{\mathcal{M}}_j) \\ &= \sum_{\substack{I=0 \\ B_{IJ} \in X_j}}^{n_x-1} \sum_{J=0}^{n_y-1} [w_{IJ}(\hat{\mathcal{M}}_j) d_{IJ}(\hat{\mathcal{M}}_j)]^2 \end{aligned}$$
$$(14)$$

where $d_{IJ}(\hat{\mathcal{M}}_j)$ is the DFD estimate (5) for the motion vector $\hat{\mathcal{M}}_j$, and $w_{IJ}(\hat{\mathcal{M}}_j)$ is a weighting factor corresponding to the block $B_{IJ}$ and depending on the motion vector $\hat{\mathcal{M}}_j$. We consider as the weighting factor a reliability coefficient in the output of the block matching:

$$w_{IJ}(\hat{\mathcal{M}}_j) = \frac{d_{IJ}(\hat{\mathcal{M}}_j)}{\sum_{k=-\frac{S_x}{2}}^{\frac{S_x}{2}} \sum_{l=-\frac{S_y}{2}}^{\frac{S_y}{2}} d_{I+k, J+l}(\hat{\mathcal{M}}_j)} \qquad (15)$$

where $S_x \times S_y$ is the search region for the block matching algorithm. This coefficient is small when we have good matching and large in the case of poor matching. For regions with constant gray level in the search area, this coefficient becomes

$$w_{IJ}(\hat{\mathcal{M}}_j) = \frac{D_x D_y}{S_x S_y}. \qquad (16)$$

The energy function $E(\hat{\mathcal{M}}_j \mid \hat{\mathcal{S}}_j, f_{t-1})$ in (13) is associated with motion vector clustering

$$E(\hat{\mathcal{M}}_j \mid \hat{\mathcal{S}}_j, f_{t-1}) = \sum_{\substack{I=0 \\ B_{IJ} \in X_j}}^{n_x-1} \sum_{J=0}^{n_y-1} (\mathbf{M}_{IJ} - \hat{\mathcal{M}}_j)^T (\mathbf{M}_{IJ} - \hat{\mathcal{M}}_j)$$
$$(17)$$

where $\hat{\mathcal{M}}_j$ is the motion vector estimate of the moving region $j$ and $\mathbf{M}_{IJ}$ is the motion vector associated with a block site (4) belonging to the region $j$.

The cost function associated to the moving region segmentation $E(\hat{\mathcal{S}}_j \mid f_{t-1})$ is related to vector clustering with respect to their gray level, and their geometrical proximity

$$E(\hat{\mathcal{S}}_j \mid f_{t-1}) = \sum_{\substack{I=0 \\ B_{IJ} \in X_j}}^{n_x-1} \sum_{J=0}^{n_y-1} (\mathbf{S}_{IJ} - \hat{\mathcal{S}}_j)^T (\mathbf{S}_{IJ} - \hat{\mathcal{S}}_j) \qquad (18)$$

where $\hat{\mathcal{S}}_j$ and $\mathbf{S}_{IJ}$ are still image feature vectors of the moving region $j$ and its component block sites (2), respectively. This cost function component entails that block sites assigned to a certain moving region are situated in the same image area and have similar average gray levels.

By replacing (14), (17) and (18) in (13) we obtain the $j$th moving object energy

$$E_j(\mathbf{u}_{IJ}) = \sum_{\substack{I=0 \\ B_{IJ} \in X_j}}^{n_x-1} \sum_{J=0}^{n_y-1} (\mathbf{u}_{IJ} - \hat{\mu}_j)^T (\mathbf{u}_{IJ} - \hat{\mu}_j)$$
$$+ \mathrm{WDFD}(\hat{\mathcal{M}}_j) \qquad (19)$$

where $\hat{\mu}_j$ is the $j$th moving region feature estimate (6), $\mathbf{u}_{IJ}$ is a block site feature vector (1) and $\mathrm{WDFD}(\hat{\mathcal{M}}_j)$ is provided in (14). The first part of this expression is associated to clustering in the feature space.

## III. MEDIAN RADIAL BASIS FUNCTION LEARNING ALGORITHM

### A. The Estimation of the Hidden Unit Parameters

The cost function (19) corresponds to image partition in moving regions. If we take into account the covariance matrix and we express (19) as an unnormalized probability (12), we obtain a so called radial basis function

$$\phi_j(\mathbf{u}_{IJ}) = \exp\left[-(\mathbf{u}_{IJ} - \hat{\mu}_j)^T \hat{\Sigma}_j^{-1}(\mathbf{u}_{IJ} - \hat{\mu}_j)\right.$$
$$\left. - \mathrm{WDFD}(\hat{\mathcal{M}}_j)\right] \qquad (20)$$

where $\hat{\mu}_j$ is the estimated center vector and $\hat{\Sigma}_j$ is the estimate of the covariance matrix. Each basis function must be defined such that it maximizes the probability of the optical flow estimation and segmentation of a certain moving region, $P(\hat{\mathcal{M}}_k, \hat{\mathcal{S}}_k \mid f_t, f_{t-1})$. The RBF network consists of a two-layer topology where each hidden unit implements a radial-activated region function (20).

The output layer function consists of a weighted sum of hidden-unit outputs, scaled to the interval $(0, 1)$ by a sigmoidal function

$$Y_k(\mathbf{u}_{IJ}) = \frac{1}{1 + \exp\left[-\sum_{j=1}^{L} \lambda_{kj}\phi_j(\mathbf{u}_{IJ})\right]} \qquad (21)$$

for $k = 1, \cdots, N$, where $\lambda_{kj}$ is the parameter associated to the connection between the hidden unit $j$ and the output unit $k$, $L$ is the number of basis functions, $N$ is the number of moving objects and $\phi_j(\mathbf{u}_{IJ})$ is provided in (20). The *a posteriori* probability of optical flow estimation and segmentation associated to each moving object, is modeled by the output unit

$$P(\hat{\mathcal{T}}_k \mid f_t, f_{t-1}) = \sum_{I=0}^{n_x-1} \sum_{J=0}^{n_y-1} \sum_{j=1}^{L} \lambda_{kj}\phi_j(\mathbf{u}_{IJ}). \qquad (22)$$

RBF networks have good functional approximation capabilities and they have been employed in many applications [22]–[27]. The diagram of the proposed system, as applied to moving object segmentation, is shown in Fig. 1. The RBF network entries correspond to still image features (2) and
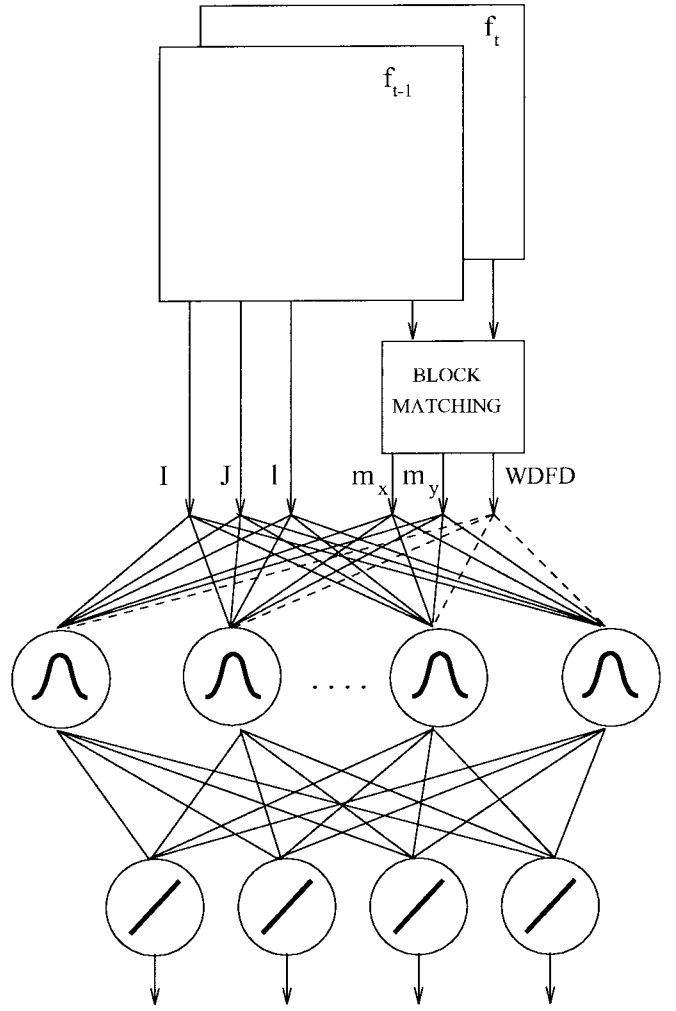


Fig. 1. MRBF structure for optical flow estimation and segmentation. $I$ and $J$ are the location coordinate features assigned to a block site, $l$ is the gray level, and $m_x$, $m_y$ are the motion vector components provided by the block matching algorithm.

block matching motion estimates (4). Each basis function (20) consists of a Gaussian function and a constraint representing the weighted DFD (14). In the context of regularization theory, the RBF networks were shown to approximate smooth functions arbitrary well [25].

The network parameters are estimated based on a learning algorithm. A well-known RBF training algorithm [26] employs in the first training stage the LVQ algorithm [28]. For a given feature vector $\mathbf{u}_{IJ}$, we update only the parameters of that basis function that provides a minimal energy in (19):

$$\text{If } \min_{k=1}^{L} E_k(\mathbf{u}_{IJ}) = E_j(\mathbf{u}_{IJ}), \quad \text{then update } \hat{\mu}_j, \hat{\Sigma}_j. \qquad (23)$$

A robust statistics-based training algorithm called median radial basis function (MRBF) is proposed in [27] for estimating the RBF network parameters. According to the results provided by the theoretical analysis, the algorithms based on robust statistics give a better accuracy in estimating the parameters of overlapping Gaussian mixtures [27]. The basis function center is randomly initialized in the input space range. Its updating is based on the marginal median LVQ algorithm [34].

Fig. 2. Frame from the Hamburg taxi sequence.



Fig. 3. Optical flow produced by the full search block matching algorithm.

This algorithm orders the data samples, on each dimension separately, and takes their median:

$$\hat{\mu}_k = \mathrm{med}\{\mathbf{u}_0, \mathbf{u}_1, \cdots, \mathbf{u}_{p-1}\} \tag{24}$$

where $\mathbf{u}_{p-1}$ is the last data sample assigned to that basis function. The MAD estimator is used as a robust dispersion estimator:

$$\hat{\mathbf{r}}_k = \frac{\mathrm{med}\{|\mathbf{u}_0 - \hat{\mu}_k|, \ldots, |\mathbf{u}_{p-1} - \hat{\mu}_k|\}}{0.6745} \tag{25}$$

where $\hat{\mathbf{r}}_k$ denotes the diagonal vector of the covariance matrix $\hat{\boldsymbol{\Sigma}}_k$ (20) and 0.6745 is the scale parameter chosen in order to make the estimator Fisher consistent for the normal distribution [29]. The model parameter estimation based on robust algorithms like marginal median and MAD enables the rejection of a certain amount of data not belonging to the given model. A fast implementation algorithm for (24) and (25), based on data sample histogram modeling, is provided in [27].

### B. The Estimation of the Output Layer Parameter

By estimating the component radial basis functions, the image is split in moving regions. A given block site is assigned to that moving region which corresponds to the most activated radial basis function:

$$\text{If } \phi_k(\mathbf{u}_{IJ}) = \max_{j=1}^{L} \phi_j(\mathbf{u}_{IJ}) \quad \text{then } B_{IJ} \in X_k. \tag{26}$$

If an image region $X_k$ has not been assigned sufficient data samples after the first learning stage, according to (26), the respective hidden unit is discarded. Evidently, that hidden unit does not contribute significantly to the optical flow estimation and segmentation probability. Afterward, the training is reiterated with a smaller number of hidden units. The number of hidden units has to be increased in the case when new moving regions appear in scene. For each two moving regions, we evaluate their boundary measure (8), assuming a four-block site neighborhood system. The moving regions that have a high interconnectivity among their component block sites, are considered as moving objects (9), and the other regions are merged with the existing moving objects according to
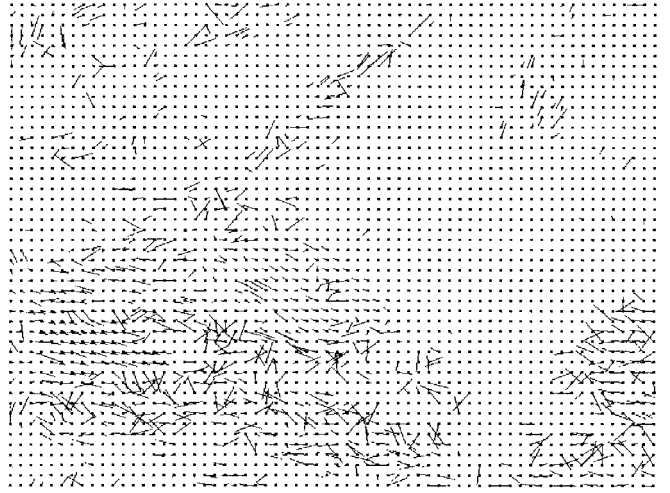
(10), (22). An output unit is assigned to each moving object. The block sites assigned to the moving regions are labeled according to the decision (26) and (9) or (10):

$$\text{If } B_{IJ} \in X_k \quad \text{then } F_k(\mathbf{u}_{IJ}) = 1, \ F_l(\mathbf{u}_{IJ}) = 0$$
$$\text{for } l = 1, \ldots, N, \quad l \neq k. \tag{27}$$

The second learning stage estimates the parameters $\lambda_{kj}$ associated with the connections between hidden and output layers (22). In order to find the weights $\lambda_{kj}$ we use a gradient based learning algorithm [22]:

$$\lambda_{kj} = \sum_{I=1}^{n_x} \sum_{J=1}^{n_y} [F_k(\mathbf{u}_{IJ}) - Y_k(\mathbf{u}_{IJ})] Y_k(\mathbf{u}_{IJ})$$
$$\times [1 - Y_k(\mathbf{u}_{IJ})] \phi_j(\mathbf{u}_{IJ}) \tag{28}$$

where $Y_k(\mathbf{u}_{IJ})$ are the outputs of the network (21) and $F_k(\mathbf{u}_{IJ})$ are the target labels assigned to the block sites according to (27).

After the training stage is completed, when entering a feature vector (1), the maximum activated output unit will show the corresponding moving object. This procedure leads to the division of the image sequence into moving objects. The trained network can be applied in other frames of the same image sequence, if their optical flow estimation and segmentation probabilities are consistent with those of the frames used in the training stage. The consistency is given by the optical flow and moving object similarity measures and can be expressed as the energy function from (13).

The MRBF network records in its weights the motion, luminance, and localization of the moving objects. The network can be used in a multiresolution (hierarchical) representation of the image. Let us consider that the MRBF network was trained on a certain image partition and afterward we input feature vectors from a different block size image partition. If the block size is large, then the segmentation will provide rough boundaries but the training time will be short. The network can be trained with features corresponding to big blocks and, afterwards be applied on an image partition in blocks of smaller size. The location features must be scaled according to the ratio between
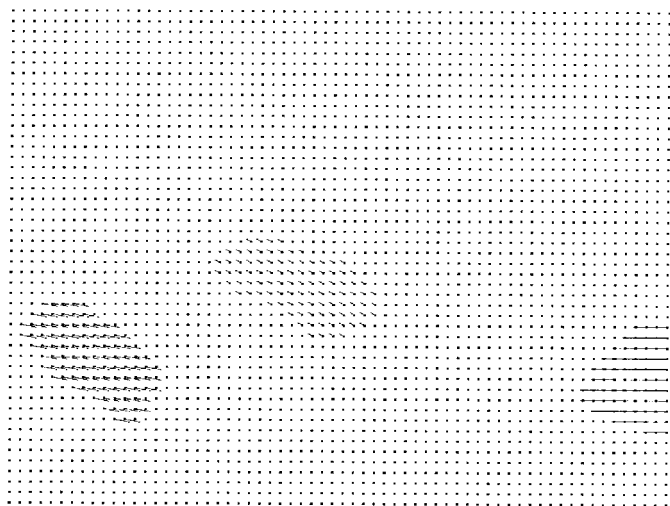
Fig. 4.   Optical flow smoothed by the MRBF network.



Fig. 6.   Optical flow smoothed by ICM.



Fig. 5.   Segmentation of the moving objects based on the MRBF network.



Fig. 7.   Segmentation of the moving objects based on ICM.

the block sizes used in training and testing. In order to save time, the network can be applied at a higher resolution only for the block sites situated at the moving object boundaries.

## IV. SIMULATION RESULTS

The MRBF network was tested on various image sequences. In this study we provide two representative results. In Fig. 2 a frame from the "Hamburg taxi" sequence is shown. This sequence contains three important moving objects: 1) a taxi near the center turning around the corner; 2) a car in the lower left, driving from left to right; and 3) a van in the lower right, driving from right to left. The optical flow provided by the full search block matching algorithm [3], is shown in Fig. 3. The feature vectors are drawn from the first and third frames of the Hamburg taxi sequence. The learning algorithm described in Section III is applied to the given data and the optical flow smoothed by the MRBF network is provided in Fig. 4. The moving objects segmented by means of the MRBF network are shown in Fig. 5.

The optical flow and moving object segmentation provided by ICM [8]–[15] for Hamburg taxi sequence are shown in
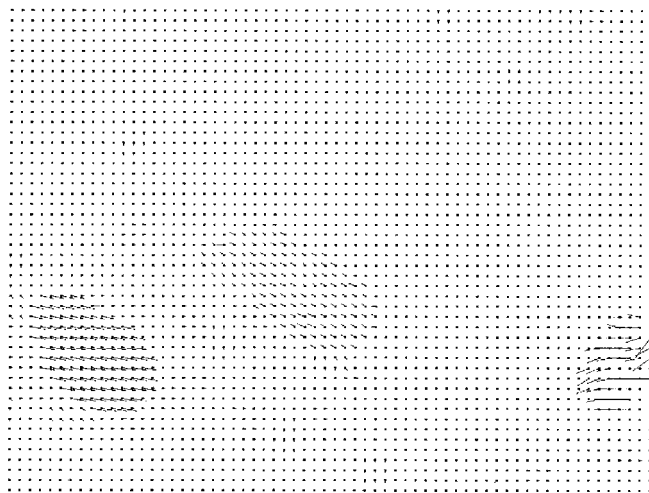
Figs. 6 and 7, respectively. We have used the full search block matching and $k$-means clustering algorithm for ICM initialization [15]. Each block site has associated a label and a motion vector, which are iteratively modified based on an optimization method, employing local constraints. This algorithm takes into account the borders of the moving objects, but smooths the optical flow only locally, inside a given neighborhood. Certain discontinuities can be observed in the optical flow smoothed by the ICM algorithm as shown in Fig. 6.

The MRBF network processes entire image regions assigned to the same moving object and exploits the interdependency among their block sites. Starting with various numbers of hidden units, the network always detected four moving objects in the Hamburg taxi sequence and provided accurate optical flows. The moving object segmentation is quite accurate despite the poor quality of the object edges. The optimal optical flow and moving object segmentation, evaluated semi-automatically by averaging the displacements of various clear features and manually segmenting the image, are shown in Figs. 8 and 9. The results provided by the MRBF network are better than those given by the ICM and they are quite close

TABLE I
COMPARISON BETWEEN MRBF NETWORK AND ICM WHEN APPLIED IN THE HAMBURG TAXI IMAGE SEQUENCE

| Algorithm | Training Frame | | | | | Image Sequence (9 Frames) | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Classification Error (%) | MAE | MSE | Number of Parameters | Number of Iterations | Classification Error (%) | Total Processing Time (s) |
| MRBF | 3.02 | 0.17 | 0.85 | 210 | 13 | 3.55 | 74.2 |
| ICM | 4.07 | 0.33 | 1.15 | 9216 | 23 | 8.70 | 57.4 |



Fig. 8.   Reference optical flow for the Hamburg taxi sequence.



Fig. 10.   MRBF network-based segmentation when applied at pixel level partition in the Hamburg taxi sequence.
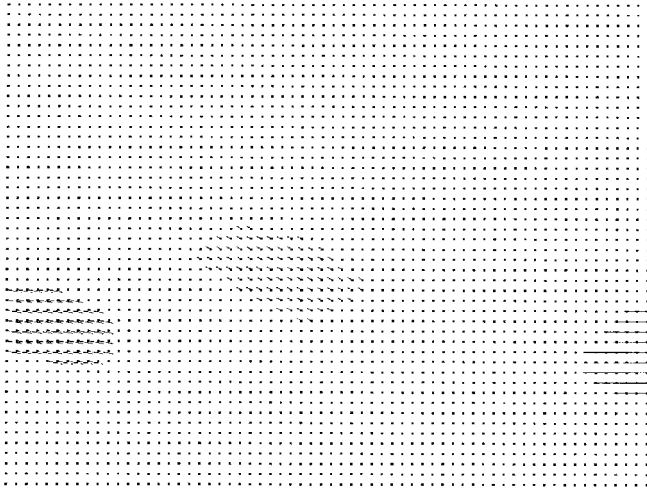


Fig. 9.   The reference moving object segmentation for the Hamburg taxi sequence.



Fig. 11.   Moving object segmentation in the eighth frame obtained after applying the network trained with data samples drawn from the first and third frames.

to the optimal results.

The MRBF network trained on a 4 × 4 pixel block partition, whose results are presented in Figs. 4 and 5, is applied on a pixel-level partition of the same frame, in a hierarchical approach. The segmentation of the moving objects is provided in Fig. 10. The "taxi" and "van" moving objects are better segmented, as can be seen from Figs. 5 and 10. The segmentation results obtained when applying the MRBF network trained using the first and third frames, on the eighth frame are shown in Fig. 11. This result shows the network's capability to embed in its weights the parameters associated with the moving objects.

A frame from Trevor White image sequence and its corresponding optical flow provided by the block matching algorithm are shown in Figs. 12 and 13. The optical flow smoothed by the MRBF network and by ICM are displayed in Figs. 14 and 15. The movement segmentation when using the MRBF network and ICM algorithm is illustrated in Figs. 16 and 17. These results show that MRBF network provides a good moving object segmentation and smooth optical flow even in the case when the movement is complex as it is with the Trevor White sequence.

Numerical comparisons when applying the MRBF network and ICM algorithm in Hamburg taxi sequence are provided in Table I. The comparison criteria when using the training frames are the optical flow mean absolute error (MAE), mean squared error (MSE), the misclassification error, the number

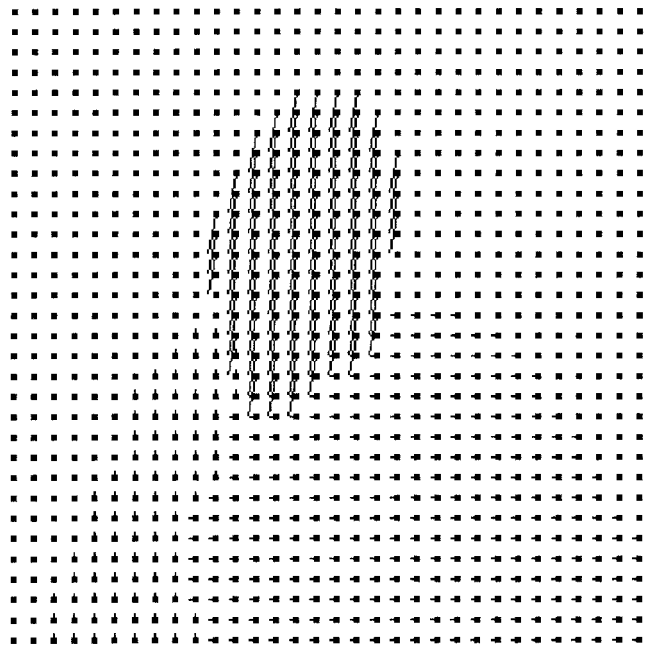Fig. 12. Frame from the Trevor White image sequence.



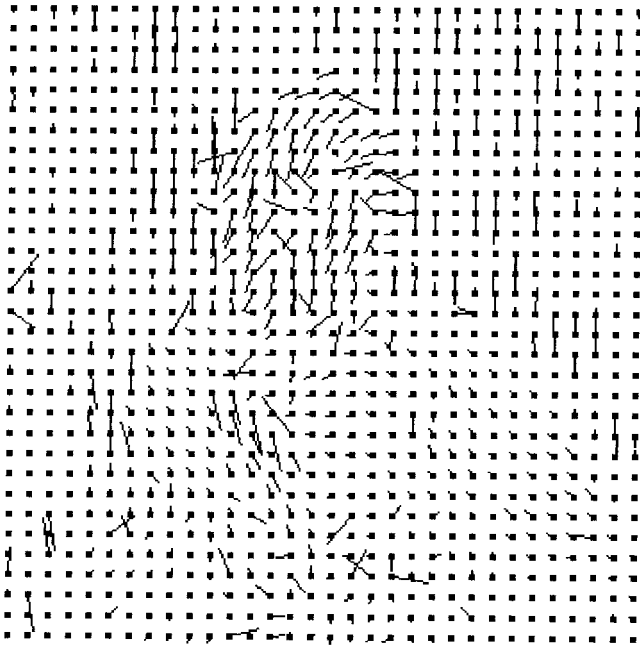Fig. 14. Optical flow smoothed by the MRBF network.



Fig. 13. Optical flow in the Trevor White image sequence provided by the full search block matching algorithm.

of parameters required by each algorithm in order to estimate the optical flow and segment the motion, and the necessary number of iterations in order to achieve the convergence. The optical flow MAE and MSE are

$$\text{MAE} = \frac{1}{n_x n_y} \sum_{I=0}^{n_x-1} \sum_{J=0}^{n_y-1} (|\hat{m}_{x,IJ} - m_{x,IJ}|$$
$$+ |\hat{m}_{y,IJ} - m_{y,IJ}|) \qquad (29)$$
$$\text{MSE} = \frac{1}{n_x n_y} \sum_{I=0}^{n_x-1} \sum_{J=0}^{n_y-1} (\hat{\mathbf{M}}_{IJ} - \mathbf{M}_{IJ})^T (\hat{\mathbf{M}}_{IJ} - \mathbf{M}_{IJ})$$
$$(30)$$

where $\hat{\mathbf{M}}_{IJ}$ is the block motion estimate. The misclassification error represents the percentage of erroneous decisions when assigning each block site to a moving object. The optical flow used as reference $\mathbf{M}_{IJ}$, and the optimal segmentation are provided in Figs. 8 and 9. The total number of necessary parameters for the MRBF network is $(10+N)L$, where $L$ is the number of hidden units and $N$ that of outputs. The total number of parameters required by the ICM algorithm is $3n_x n_y$. The MRBF network trained using data samples drawn from the first and third frames has been applied on the frames two to nine of the Hamburg taxi sequence. The misclassification error and the total processing time corresponding to the nine frames (including the feature computation and the training time for the MRBF) evaluated on a Silicon Graphics Indy Workstation are provided in Table I. The time spent in the MRBF training stage, for this example, is 27.9 seconds.

It can be seen from Table I, as well as from Figs. 4–17, that the MRBF network outperforms the ICM algorithm. The average time per frame for the MRBF network is larger when compared with the ICM algorithm, but the MRBF network requires fewer iterations in training. The reduction in the number of parameters, provided by the MRBF network is very high, when compared to the ICM.

## V. CONCLUSION

The MRBF neural network when used for optical flow estimation and moving object segmentation is analyzed in this study. The considered model consists of a mixture of kernel functions whose parameters are found by training. The criterion for segmenting the moving objects is derived from the *a posteriori* probability maximization criterion. Consequently, a cost function is obtained and used as a feature space metric in the learning stage. The cost function takes into account the
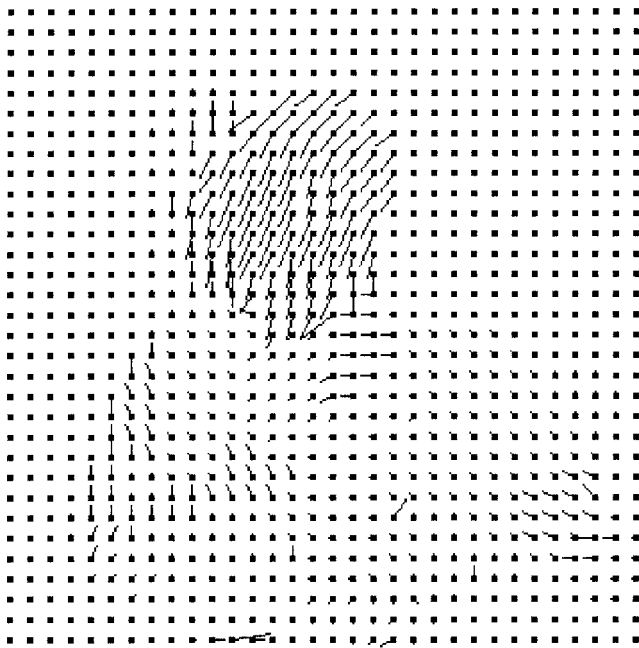
Fig. 15.   Optical flow smoothed by the ICM.



Fig. 17.   Movement segmentation provided by the ICM.



Fig. 16.   Movement segmentation provided by the MRBF network.

number of parameters required by the MRBF network to model the optical flow and moving object segmentation is small. The information contained in the network parameters can be further used for image sequence analysis and region-based coding.

## REFERENCES

[1] J. K. Kearney, W. B. Thompson, and D. L. Boley, "Optical flow estimation: An error analysis of gradient based methods with local optimization," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-9, pp. 229–244, Mar. 1987.

[2] J. Weng, N. Ahuja, and T. S. Huang, "Matching two perspective views," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 14, pp. 806–825, Aug. 1992.

[3] A. M. Tekalp, *Digital Video Processing*. Englewood Cliffs, NJ: Prentice-Hall, 1995.

[4] D. W. Murray and B. F. Buxton, "Scene segmentation from visual motion using global optimization," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-9, pp. 220–228, Mar. 1987.

[5] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-6, pp. 721–741, Nov. 1984.

[6] J. Konrad and E. Dubois, "Bayesian estimation of motion vector fields" *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 14, pp. 910–927, Sept. 1992.

[7] H. M. Kalayeh and G. Horgan, "Adaptive relaxation labeling," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. PAMI-6, pp. 369–372, 1984.

[8] J. Besag, "On the statistical analysis of dirty pictures," *J. R. Stat. Soc. B*, vol. 48, pp. 259–302, 1986.

[9] T. Aach, A. Kaup, and R. Mester, "Combined displacement estimation and segmentation of stereo image pairs based on Gibbs random fields," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing*, Albuquerque, NM, Apr. 1990, pp. 2301–2304.

[10] M. M. Chang, A. M. Tekalp, and M. I. Sezan, "Motion-field segmentation using an adaptive MAP criterion," in *Proc. IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, Minneapolis, MN, Apr. 1993, pp. V-33–V-36.

[11] C. Stiller, "Motion-estimation for coding of moving video at 8 kbit/s with Gibbs modeled vector field smoothing," in *Proc. SPIE Conf. Visual Communications and Image Processing*, 1990, vol. 1360, pp. 468–476.

[12] M. M. Chang, M. I. Sezan, and A. M. Tekalp, "An algorithm for simultaneous motion estimation and scene segmentation," in *Proc. IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, Adelaide, Australia, May 1994, pp. V-221–V-224.

[13] F. Heitz and P. Bouthemy, "Multimodal estimation of discontinuous optical flow using Markov random fields," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 15, pp. 1217–1232, Dec. 1993.

local motion information, the gray level or color components, the geometrical proximity and considers the displaced frame difference as well. The weights of the network are updated by means of an unsupervised learning algorithm. A robust training algorithm is employed for estimating the hidden unit parameters. The mixture of hidden units is fed into the output units, each of them associated to a moving object. The moving region areas found in the first learning stage are merged, based on a compactness measure, forming moving objects. The optical flow provided by the proposed algorithm proved to be accurate and smooth. After the training stage, the network records in its weights the moving object information. The

[14] C. Stiller, "Object-oriented coding employing dense motion fields," in *Proc. IEEE Int. Conf. Acoustics, Speech, and Signal Processing*, Adelaide, Australia, May 1994, pp. V-273–V-276.

[15] T. N. Pappas, "An adaptive clustering algorithm for image segmentation," *IEEE Trans. Signal Processing*, vol. 40, pp. 901–914, Apr. 1992.

[16] B. G. Schunck, "Image flow segmentation and estimation by constraint line clustering," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 11, pp. 1010–1027, Oct. 1989.

[17] J. Y. A. Wang and E. H. Adelson, "Layered representation for motion analysis," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, New York, NY, June 1993, pp. 361–366.

[18] D. P. Kottke and Y. Sun, "Motion estimation via clustering matching," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 16, pp. 1128–1132, Nov. 1994.

[19] B. Duc, P. Schroeter, and J. Bigun, "Spatio-temporal robust motion estimation and segmentation," in *Proc. Conf. Computer Analysis of Images and Patterns*, Prague, Czech Republic, Sept. 6–8, 1995, pp. 238–245.

[20] T. Y. Tian and M. Shah, "Motion segmentation and estimation," in *Proc. IEEE Conf. Image Processing*, Austin, TX, Nov. 1994, vol II, pp. 785–789.

[21] A. G. Borş and I. Pitas, "Segmentation and estimation of the optical flow," in *Proc. Int. Conf. Computer Analysis of Images and Patterns*, Prague, Czech Republic, Sept. 1995, pp. 680–685.

[22] A. G. Borş and M. Gabbouj, "Minimal topology for a radial basis functions neural network for pattern classification," *Digit. Signal Process.: Rev. J.*, vol. 4, pp. 173–188, July 1994.

[23] E. J. Hartman, J. D. Keeler, and J. M. Kowalski, "Layered neural networks with Gaussian hidden units as universal approximations," *Neural Comput.*, vol. 2, pp. 210–215, 1990.

[24] J. Park and J. W. Sandberg, "Universal approximation using radial basis functions network," *Neural Comput.*, vol. 3, pp. 246–257, 1991.

[25] T. Poggio and F. Girosi, "Networks for approximation and learning," in *Proc. IEEE*, vol. 78, pp. 1481–1497, Sept. 1990.

[26] L. Xu, A. Krzyzak, and E. Oja, "Rival penalized competitive learning for clustering analysis, RBF net, and curve detection," *IEEE Trans. Neural Networks*, vol. 4, pp. 636–649, July 1993.

[27] A. G. Borş and I. Pitas, "Median radial basis function neural network," *IEEE Trans. Neural Networks*, vol. 7, pp. 1351–1364, Nov. 1996.

[28] T. K. Kohonen, *Self-Organization and Associative Memory*. New York: Springer-Verlag, 1989.

[29] G. Seber, *Multivariate Observations*. New York: Wiley, 1984.

[30] I. Pitas and A. N. Venetsanopoulos, *Nonlinear Digital Filters: Principles and Applications*. Boston, MA: Kluwer, 1990.

[31] S. Haykin, *Neural Networks: A Comprehensive Foundation*. New York: Macmillan, 1994.

[32] I. M. Elfadel and R. W. Picard, "Gibbs random fields, cooccurrences, and texture modeling," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 16, pp. 24–37, Jan. 1994.

[33] J. Marroquin, S. Mitter, and T. Poggio, "Probabilistic solution of ill-posed problems in computational vision," *J. Amer. Stat. Assoc.*, vol. 82, pp. 76–89, 1987.

[34] I. Pitas *et al.*, "Order statistics learning vector quantizer," *IEEE Trans. Image Processing*, vol. 5, pp. 1048–1053, June 1996.

**Adrian G. Borş** was born in Piatra Neamţ, Romania, in 1967. He received the M.S. degree in electronics engineering from the Polytechnic University of Bucharest, Romania, in 1992. He is pursuing the Ph.D. degree in the Department of Informatics at the University of Thessaloniki, Greece.

From September 1992 to August 1993, he was a Visiting Researcher at the Signal Processing Laboratory, Tampere University of Technology, Finland. Since 1993, he has been contributing to various European research projects while at the University of Thessaloniki. His research interests include neural networks, computer vision, pattern recognition, and nonlinear digital signal processing.

**Ioannis Pitas** received the Diploma of Electrical Engineering degree in 1980 and the Ph.D. degree in electrical engineering in 1985, both from the University of Thessaloniki, Greece.

Since 1994, he has been a Professor at the Department of Informatics, University of Thessaloniki. From 1980 to 1993, he served as Scientific Assistant, Lecturer, Assistant Professor, and Associate Professor in the Department of Electrical and Computer Engineering at the same university. He served as a Visiting Research Associate at the University of Toronto, Toronto, Ont., Canada, University of Erlangen-Nuernberg, Germany, Tampere University of Technology, Finland and as Visiting Assistant Professor at the University of Toronto. He was lecturer in short courses for continuing education. His current interests are in the areas digital image processing, multidimensional signal processing, and computer vision. He has published over 200 papers and contributed to eight books in his area of interest. He is the co-author of *Nonlinear Digital Filters: Principles and Applications* (Boston, MA: Kluwer, 1990). He is author of *Digital Image Processing Algorithms* (Englewood Cliffs, NJ: Prentice-Hall, 1993). He is editor of *Parallel Algorithms and Architectures for Digital Image Processing, Computer Vision and Neural Networks* (New York: Wiley, 1993).

Dr. Pitas has been member of the European Community ESPRIT Parallel Action Committee. He has also been invited speaker and/or member of the program committee of several scientific conferences and workshops. He is associate editor of IEEE TRANSACTIONS ON CIRCUITS AND SYSTEMS and co-editor of the *Journal of Multidimensional Systems and Signal Processing*. He was chair of the 1995 IEEE Workshop on Nonlinear Signal and Image Processing (NSIP'95).