

AUTOMATIC KERNEL WIDTH SELECTION FOR NEURAL NETWORK BASED VIDEO OBJECT SEGMENTATION

Dubravko Culibrk

*Department of Computer Science and Engineering, Florida Atlantic University
Boca Raton FL 33431, USA
dculibrk@fau.edu*

Daniel Socek

*Department of Computer Science and Engineering, Florida Atlantic University
Boca Raton FL 33431, USA
dsocek@fau.edu*

Oge Marques

*Department of Computer Science and Engineering, Florida Atlantic University
Boca Raton FL 33431, USA
omarques@fau.edu*

Borko Furht

*Department of Computer Science and Engineering, Florida Atlantic University
Boca Raton FL 33431, USA
borko@cse.fau.edu*

Keywords: Video processing, Object segmentation, Background modeling, Bayesian modeling, Neural Networks.

Abstract: Background modelling Neural Networks (BNNs) represent an approach to motion based object segmentation in video sequences. BNNs are probabilistic classifiers with non-parametric, kernel-based estimation of the underlying probability density functions. The paper presents an enhancement of the methodology, introducing automatic estimation and adaptation of the kernel width. The proposed enhancement eliminates the need to determine kernel width empirically. The selection of a kernel-width appropriate for the features used for segmentation is critical to achieving good segmentation results. Thus, the improvement results in a segmentation methodology truly general in terms of features used for segmentation.

1 INTRODUCTION

Object segmentation is a basic task in the domain of digital video processing. Diverse applications, such as scene understanding, object-based video encoding, surveillance applications and 2D-to-pseudo-3D video conversion, rely on the ability to extract objects from video sequences.

The research into object segmentation for video sequences grabbed from a stationary camera has yielded a number of approaches based on the detection of the motion of objects(). The approaches of this class scrutinize the changes observed between the consecutive frames of the sequence to detect pixels which corre-

spond to moving objects. The task is particularly difficult when the segmentation is to be done for natural scenes where the background contains shadows, moving objects, and undergoes illumination changes.

For purposes of automated surveillance and scene understanding it often of more interest not only to detect the objects moving in the scene, but to distinguish between the two classes of objects:

- *Background objects* corresponding to all objects that are present in the scene, during the whole sequence or longer than a predefined period of time.
- *Foreground objects* representing all other objects appearing in the scene.

The goal of the video object segmentation is to separate pixels corresponding to foreground from those corresponding to background.

Background Modeling Neural Network (BNN) represents a probabilistic approach to foreground detection(?). The network is a Bayes rule based unsupervised classifier designed to enable the classification of a single pixel as pertinent to foreground or background. A set of networks is used to classify all the pixels within the frame of the sequence.

The networks incorporate kernel-based estimators (Specht, 1991) for probability density functions used to model the background. While the networks are general in terms of the features of a pixel used to perform the classification, the accuracy of the process depends on the ability to determine the appropriate width for the estimator kernels. The process relies on empirical data and involves tedious experimentation. In addition the BNNs are unable to adapt to conditions occurring in specific sequences. Rather, a single value is typically used whenever the same features are used as basis for segmentation.

In this paper, an extension of the BNN methodology is proposed to incorporate automatic selection of the appropriate kernel width. The proposed enhanced BNNs do not suffer from the above mentioned shortcomings of the fixed-kernel-width BNNs and achieve comparable segmentation performance.

The rest of the paper is organized as follows: Section 2 provides a survey of related published work. Section 3 contains a short description of BNNs. The proposed enhancement of the BNNs is described in Section ???. Section 5 is dedicated to the presentation of experimental evaluation results. Section 6 contains the conclusions.

2 RELATED WORK

Early foreground segmentation methods dealing with non-stationary background are based on a model of background created by applying some kind of low-pass filter on the background frames. The high-frequency changes in intensity or color of a pixel are filtered out using different filtering techniques such as Kalman filters (Karmann and von Brandt, 1990) to create an approximation of background in the form of an image (reference frame). The reference frame is updated with each new frame in the input sequence and used to segment the foreground objects by subtracting the reference frame from the observed frame (?). These methods are based on the most restrictive assumption that observe pixel changes due to the background are much slower than those due to the objects to be segmented. Therefore they are not particularly effective for sequences with high-frequency

background changes, such as natural scene and outdoor sequences.

Probabilistic techniques achieve superior results in case of such complex-background sequences. These methods rely on an explicit probabilistic model of the background, and a decision framework allowing for foreground segmentation. A Gaussian-based statistical model whose parameters are recursively updated in order to follow gradual background changes within the video sequence is proposed in (Boult et al., 1999). More recently, Gaussian-based modelling was significantly improved by employing a Mixture of Gaussians (MoG) as a model for the probability density functions (PDFs) related to the distribution of pixel values. Multiple Gaussian distributions, usually 3-5, are used to approximate the PDFs (Ellis and Xu, 2001)(Stauffer and Grimson, 2000)(Ya et al., 2001). The parameters of each Gaussian curve are updated with each observed pixel value. If an observed pixel value is within the 2.5 standard deviations (σ) from the mean (μ) of a Gaussian, the pixel value matches the Gaussian (Stauffer and Grimson, 2000). The parameters are updated only for Gaussians matching the observed pixel value, based on the following Equations:

$$\mu_t = (1 - \rho) * \mu_{t-1} + \rho * X_t \quad (1)$$

$$\sigma_t^2 = (1 - \rho) * \sigma_{t-1}^2 + \rho * (X_t - \mu_t)^T * (X_t - \mu_t) \quad (2)$$

where

$$\rho = \aleph(X_t, \mu_{t-1}, \sigma_{t-1}) \quad (3)$$

and \aleph is a Gaussian function. Equations 7 - 3 express a causal low-pass filter applied to the mean and variance of the Gaussian.

Using a small number of Gaussians leads to a rough approximation of the PDFs involved. Due to this fact, MoG achieves weaker results for video sequences containing non-periodical background changes (e.g. due to waves and water surface illumination, cloud shadows, and similar phenomena), as was reported in (Li et al., 2003). The Gaussian-based models are parametric in the sense that they incorporate underlying assumptions about the probability density functions (PDFs) they are trying to estimate.

In 2003, Li et al. proposed a method for foreground object detection employing a Bayes decision framework (Li et al., 2003). The method has shown promising experimental object segmentation results even for the sequences containing complex variations and non-periodical movements in the background. The primary model of the background used by Li et al. is a background image obtained through low pass filtering. However, the authors use a probabilistic model for the pixel values detected as foreground through frame-differencing between the current frame and the reference background image. The probabilistic model is used to enhance the results of primary foreground

detection. The probabilistic model is non-parametric since it does not impose any specific shape to the PDFs learned. However, for reasons of efficiency and improving results the authors applied binning of the features and assigned single probability to each bin, leading to a discrete representation of PDFs. The representation is equivalent to a kernel-based estimate with quadratic kernel. The width of the kernel used was determined empirically and remained fixed in all the reported experiments (Li et al., 2004). The system achieved performance better than that of the mixture of 5 Gaussians in the results presented in the same publication. However, when larger number of Gaussians is used, MoG achieved better performance (?).

A nonparametric kernel density estimation framework for foreground segmentation and object tracking for visual surveillance has been proposed in (?). The authors present good qualitative results of the proposed system, but do not evaluate segmentation quantitatively nor do they compare their system with other methods. The framework is computationally intensive as the number of kernels corresponds to the number of observed pixel values. The width of the kernels is adaptive and they use the observed median of absolute differences between consecutive pixel values. The rationale for the use of median is the fact that its estimate is robust to small number of outliers. They assume Gaussian (normal) distribution for the differences and establish a relation between the estimated median and the width of the kernel:

$$\sigma = \frac{m}{0.68\sqrt{2}} \quad (4)$$

where m is the estimated median.

The approach based on background modelling neural networks was proposed in (?). The networks employ represent a biologically plausible implementation of Bayesian classifiers and nonparametric kernel based density estimators. The weights of the network serve store a model of background, which is continuously updated. The PDF estimates consist of fixed number of kernels, which have fixed width. The appropriate width of the kernels is determined empirically. The kernel width depends on the features used to achieve segmentation. The results superior to those of Li et al. and MoG with 30 Gaussians are reported in (?). The BNNs address the problem of computational complexity of the kernel based background models by exploiting the parallelism of neural networks.

3 BACKGROUND MODELING NEURAL NETWORKS (BNNs)

Background Modeling Neural Network (BNN) is a neural network designed specifically to serve as a statistical model of the background at each pixel position

in the video sequences and highly-parallelized background subtraction algorithm. The network is an unsupervised learner. It collects statistics related to the dynamic processes of pixel feature values changes. The learnt statistics are used to classify a pixel as pertinent to a foreground or background object in each frame of the sequence.

Note that a video sequence can be viewed as a set of pixel feature values changing in time. Probabilistic motion (change) based background subtraction methods rely on a commonsensical supposition:

Pixel feature values corresponding to background objects will occur most of the time, i.e. more often than those pertinent to the foreground.

Thus, if a classifier is able to effectively distinguish between the values occurring more frequently than others it should be able to achieve accurate segmentation.

The structure of BNN, shown in Figure 1, has three distinct subnets. The classification subnet is a central part of BNN concerned with approximating the Probability Density Function (PDF) of pixel feature values belonging to background/foreground. It is a neural network implementation of a Parzen (kernel based) estimator (Parzen, 1962). This class of PDF estimators asymptotically approaches the underlying parent density, provided that it is smooth and continuous.

The classification subnet contains three layers of neurons. Input neurons of this network simply map the inputs of the network, which are the values of the features for a specific pixel, to the pattern neurons. The output of the pattern neurons is a nonlinear function of Euclidean distance between the input of the network and the stored pattern for that specific neuron:

$$p_i = \exp\left[-\frac{(w_i - x_t)^T(w_i - x_t)}{2\sigma^2}\right] \quad (5)$$

where w_i is the vector of weights between the input neurons and the i -th pattern neuron, x_t is the pixel feature value vector and p_i is the output of the i -th pattern neuron. The only parameter of this subnet is a so-called smoothing parameter (σ) used to determine the shape of the nonlinear function.

The output of the summation units of the classification subnet is the sum of their inputs. The subnet has two summation neurons: one to calculate the probability of the observed pixel value corresponding to background and the other to calculate the probability of the value pertaining to foreground.

Weights between the pattern and summation neurons are used to store the confidence with which a pattern belongs to the background/foreground. The weights of these connections are updated with each new value of a pixel at a certain position received (i.e. with each frame), according to the following recursive

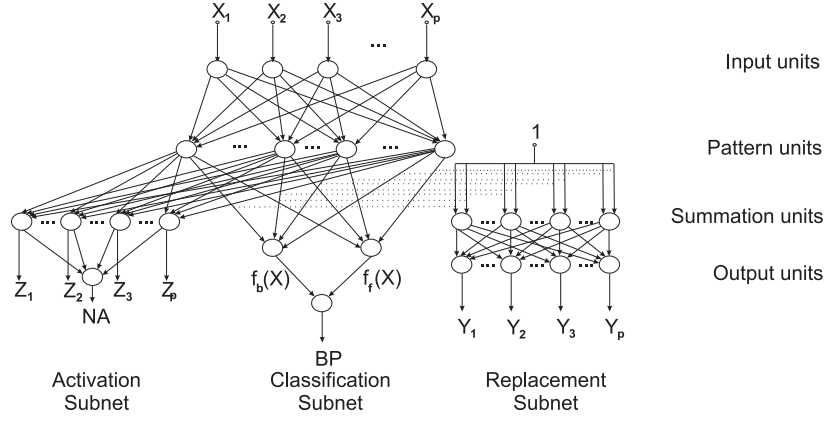


Figure 1: Structure of Background Modeling Neural Network.

equations:

$$W_{ib}^{t+1} = f_c\left(\left(1 - \frac{\beta}{N_{pn}}\right) * W_{ib}^t + MA^t \beta\right) \quad (6)$$

$$W_{if}^{t+1} = 1 - W_{ib}^{t+1} \quad (7)$$

where W_{ib}^t is the value of the weight between the i -th pattern neuron and the background summation neuron at time t , W_{if}^t is the value of the weight between the i -th pattern neuron and the foreground summation neuron at time t , β is the learning rate, N_{pn} is the number of the pattern neurons of BNN, f_c is a clipping function defined by (8) and MA^t indicates the neuron with the maximum response (activation potential) at frame t , according to (9).

$$f_c(x) = \begin{cases} 1, & x > 1 \\ x, & x \leq 1 \end{cases} \quad (8)$$

$$MA^t = \begin{cases} 1, & \text{for neuron with maximum response;} \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

Equations 1-4 express the notion that whenever an instance pertinent to a pattern neuron is encountered, the probability that that pattern neuron is activated by a feature value vector belonging to the background is increased. Naturally, if that is the case, the probability that the pattern neuron is excited by a pattern belonging to foreground is decreased. Vice versa, the more seldom a feature vector value corresponding to a pattern neuron is encountered the more likely it is that the patterns represented by it belong to foreground objects. By adjusting the learning rates, it is possible to control the speed of the learning process.

The output of the classification subnet indicates whether the output of the background summation neuron is higher than that of the foreground summation

neuron, i.e. that it is more probable that the input feature value is due to a background object rather than a foreground object. More formally, the following classification criterion is evaluated:

$$p(b|v)p(v) - p(f|v)p(v) > 0 \quad (10)$$

where $p(f|v)$ and $p(b|v)$ represent estimated conditional PDFs of an observed pixel value v indicating a foreground and background object, respectively. $p(v)$ is the estimated PDF of a feature value occurring.

If the criterion 10 is satisfied then the pixel is classified as background, otherwise it is classified as foreground.

The activation and replacement subnets are Winner-Take-All (WTA) neural networks. The activation subnet performs a dual function: it determines which of the neurons of the network has maximum activation (output) and whether that value exceeds a threshold provided as a parameter to the algorithm. If it does not, the BNN is considered inactive and replacement of a pattern neuron's weights with the values of the current input vector is required. If this is the case, the feature is considered to belong to a foreground object.

The first layer of this network has the structure of a 1LF-MAXNET network and a single neuron is used to indicate whether the network is active. The output of the neurons of the first layer of the network can be expressed in the form of the following equation:

$$Y_j = X_j \times \prod_{i=1}^P \{F(X_j - X_i | i \neq j)\} \quad (11)$$

where:

$$F(z) = \begin{cases} 1, & \text{if } z \geq 0; \\ 0, & \text{if } z < 0; \end{cases} \quad (12)$$

The output of the first layer of the activation subnet will differ from 0 only for the neurons with maximum activation and will be equal to the maximum activation. In Figure 1 these outputs are indicated with Z_1, \dots, Z_P . Figure 2 shows the inner structure of a neuron in the first layer of the subnet. A single neuron in

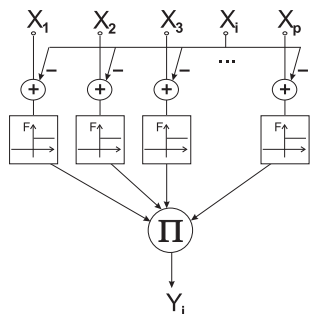


Figure 2: Structure of processing neurons of the activation subnet.

the second layer of the activation subnet is concerned with detecting whether the BNN is active or not and its function can be expressed in the form of the following equations:

$$NA = F\left(\sum_{i=1}^P Z_i - \theta\right) \quad (13)$$

where F is given by Equation 12 and θ is the activation threshold, which is provided to the network as a parameter. Finally, the replacement subnet in Figure 1 can be viewed as a separate neural net with the unit input. However, it is inextricably related to the classification subnet since each of the replacement subnet first-layer neurons is connected with the input via synapses that have the same weight as the two output synapses between the pattern and summation neurons of the classification subnet. Each pattern neuron has a corresponding neuron in the replacement net. The function of the replacement net is to determine the pattern neuron that minimizes the criterion for replacement, expressed by the following equation:

$$\text{replacement_criterion} = W_{if}^t + |W_{ib}^t - W_{if}^T| \quad (14)$$

The criterion is a mathematical expression of the idea that those patterns that are least likely to belong to the background and those that provide least confidence to make the decision should be eliminated from the model.

The neurons of the first layer calculate the negated value of the replacement criterion for the pattern neuron they correspond to. The second layer is a 1LF-MAXNET that yields non-zero output corresponding to the pattern neuron to be replaced.

To form a complete background-subtraction solution a single instance of a BNN is used to model the features at each pixel of the image.

4 AUTOMATIC KERNEL WIDTH ESTIMATION

BNNs employ a Parzen-estimator-based (Parzen, 1962) representation of the PDFs needed to achieve classification. This class of estimators is nowadays also known as kernel-based density estimators and was used in the approach presented in (?), as discussed in Section 2. A Parzen estimator of a PDF based on a set of measurements used within BNNs has the following analytical form:

$$p(v) = \frac{1}{(2\pi)^{n/2}\sigma^n} \frac{1}{T_o} \sum_{t=0}^{T_o} \exp\left[-\frac{(v - v_t)^T(v - v_t)}{2\sigma^2}\right] \quad (15)$$

where n is the dimension of the feature vector, T_o is the number of patterns used to estimate the PDF (observed pixel values), v_t are the pixel values observed up to the frame T_o , σ is a smoothing parameter.

The Parzen estimator defined by (15) is a sum of multivariate Gaussian distributions centered at the observed pixel values. As the number of observed values approaches infinity, the Parzen estimator converges to its underlying parent density, provided that it is smooth and continuous. To reduce the memory and computational requirements of the segmentation, the BNNs employ a relatively small number of kernels (up to 30 kernels showed good results in our experiments), but the kernels are used to represent more than one observation and assigned weights in a manner similar to that discussed in (Specht, 1991).

The smoothing parameter controls the width of the Gaussians. Fig. ?? shows the plot of a Parzen estimator for three stored points with values in two-dimensional plane (e.g. if only R and G values for a pixel are considered). The horizontal planes in Fig. ?? represent the threshold values used to decide which feature values are covered by a single Gaussian. The threshold was set to 0.5 in the plots. All features within the circle defined by the cross-section of the Parzen estimate and the threshold plane are deemed close enough to the center of the peak to be within the cluster pertinent to the Gaussian. The selection of smoothing parameter value and the threshold controls the size of the circle and the cluster. Larger values of σ lead to less pronounced peaks in the estimation, i.e. make the estimation "smoother".

The value of smoothing parameter has profound impact on the quality of segmentation and requires tedious experimentation to determine for a particular

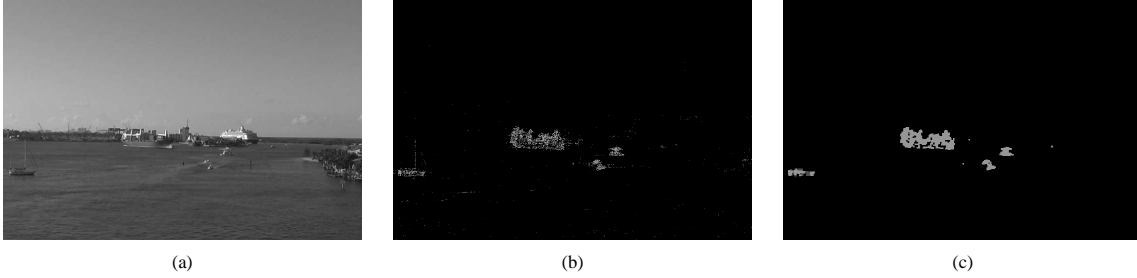


Figure 3: Results obtained for representative frames: (a) frame from the original sequence, (b) segmentation results obtained for the frame shown, (c) segmentation result with morphological enhancement.

pixel feature used for segmentation. To alleviate this deficiency of the BNN approach, a automatic procedure is proposed for learning and adaptation of the smoothing parameter based on the properties of the segmented sequence.

The number of kernels in a BNN is fixed, and determined by the available computational resources. The smoothing parameter should be selected so that the predetermined number of kernels is able to cover all the pixel values occurring due to background. In addition, a more accurate approximation of PDFs can, in general, be achieved by the kernel of smaller width. Thus, the goal of smoothing parameter estimation is to determine the minimum width of kernels needed to account for all the background pixel values, based on a predefined number of kernels and the BNN activation threshold. To achieve this goal, the kernel width is updated with each new pixel value observed.

Let σ_i and μ_i be estimates of the standard deviation and mean of the background pixel values in the observed part of the video sequence, along the i -th dimension of the pixel feature vector. In order not to increase the computational and memory requirements of the BNN, it is desirable to estimate σ_i and μ_i based on the information already contained in the BNN.

For an estimate of the mean μ_i we use the average value of the i -th coordinate of patterns stored in the network, which correspond to the weights between the i -th input neuron and each pattern neuron of the classification subnet of BNN :

$$\mu_i = \frac{\sum_{j=1}^N w_{ij}}{N} \quad (16)$$

To estimate the standard deviation σ_i , the average of the distance of each center form μ_i is calculated, weighted with the weights of the connections between each pattern neuron and summation neuron corresponding to the background. This way the contribution of patterns likely to correspond to background is exacerbated, while the influence of the patterns due to foreground is diminished. The formula for σ_i is given

by Equation 17.

$$\sigma_i = \sqrt{\left(\frac{\sum_{j=1}^N w_{bi} * (w_{ij} - \mu_i)^2}{N - 1} \right)} \quad (17)$$

Since the width of the BNN kernels is the same along each dimension of the feature vector, maximum standard deviation over all the dimensions is used to estimate the smoothing parameter:

$$\sigma = 0.5 * \sqrt{\frac{-2 * \max_{i \in 1..p} \sigma_i}{\log \theta}} \quad (18)$$

where θ corresponds to the BNN activation threshold. Equation 18 corresponds to the kernel that will be active for all patterns within two estimated maximum standard deviations. The factor of two is introduced, since the estimator based on 17 tends to underestimate the real deviation. Equation 17 would give a precise estimate based on stored patterns, if all of them corresponded to background and the confidence of of them pertaining to background was one. It is unlikely that all the stored patterns in the BNN will correspond to background. In addition, the confidence that these patterns corresponds to the background will usually be less than one.

5 EXPERIMENTS AND RESULTS

To evaluate the approach, a PC-based foreground segmentation application based on the BNN methodology and employing adaptive kernels, has been developed. The application is sequential in its execution and cannot provide a valid estimate of the speed the methodology could achieve in a parallelized hardware system. However, it can demonstrate the segmentation ability of the system.

A set of diverse sequences containing complex background conditions, provided by Li *et al.* (Li et al., 2004) and publicly available at <http://perception.i2r.a-star.edu.sg>, was used. The results of the segmentation were evaluated both qualitatively and quantitatively, using a set of ground truth

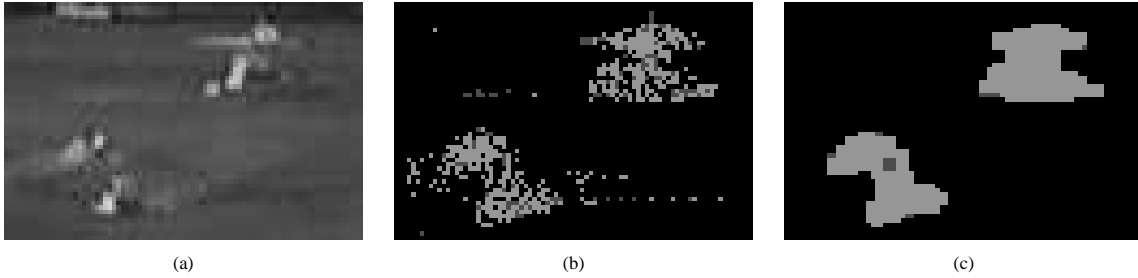


Figure 4: Details of: (a) a frame of the original sequence, (b) segmentation result and (c) segmentation result enhanced using morphological operations.

frames provided by the same authors for the different sequences.

A representative frame from the sequence as well as the result of segmentation are given in Figure 3. Grey pixels correspond to the foreground. No morphological operations, typically used to remove spurious one pixel effects and make the object solid, have been performed on the images shown in Figure 3 a and b. Figure 3 c, however, shows the same segmentation result when a morphological open and then morphological close operation are applied. The support region for the operations is a two pixel wide square.

A detail of the first frame of the original sequence (shown in figure 3(a)) containing several small objects as well as the segmentation result for that part of the frame with and without morphological operations applied is shown enlarged in Figure 4. Light grey pixels are classified as foreground due to the BNNs recognizing that these are new values not yet stored, while the dark grey ones are stored but classified as foreground based on the learned PDFs.

The neural networks used in the experiments are fairly simple. The simulation application implements BNNs containing 20 processing, two summation and one output neuron per pixel in the classification subnet. The activation and replacement subnet attribute for additional 20, i.e. 41 processing units respectively, bringing up the total of neurons used per pixel to 84. The input neurons of the classification shown in Figure 1 just map the input to the output and need not be implemented as such.

The learning rate (β) of the networks was set to 0.005 and the smoothing parameter (σ) for the classification subnet used was set to 10. The activation threshold (θ) of the activation subnet was set to 0.95.

The performance of the simulation application allows for efficient experimenting. It is capable of processing a single frame of size 720×480 in 2.25 seconds on average, which translates to 8 frames of 160×120 pixels per second or 2.2 frames per second (fps) for images sized 320×240 pixels, on a 3.0 GHz Pentium IV based system.

While the experiments conducted prove the capa-

bility of the system to segment the foreground objects, the main goal of the research presented is to achieve real-time segmentation, using a hardware-based solution. In a hardware implementation the delay of the network (segmentation time) corresponds to the time needed by the signal to propagate through the network and time required to update it. In a typical FPGA implementation this can be done in less than 20 clock cycles, which corresponds to a 2ms delay through the network, for a FPGA core running at 100ns clock rate. Thus, the networks themselves are capable of achieving a throughput of some 500 fps, which is more than sufficient for real-time segmentation of video sequences.

6 CONCLUSION AND FURTHER RESEARCH

Object segmentation is a fundamental task in several important domains of video processing. The complexity of captured and stored video material is on the rise and current motion based segmentation algorithms are not capable of handling high-resolution sequences in real time. The possibility of resolving this problem through a highly parallelized approach is the focus of the research presented.

The research resulted in a new motion based object segmentation and background modeling algorithm, proposed here. It is parallelized at sub-pixel level. The basis of the approach is employment of a novel neural network architecture designed specifically to serve as a model of background in video sequences and a Bayesian classifier to be used for object segmentation. The new Background Modeling Neural Network is an unsupervised classifier, differing from the approaches published before. The proposed model is independent of the features used and general since it does not impose restrictions in terms of the probability density functions estimated.

A PC based system has been developed to evaluate the algorithm using a complex maritime sequence.

The results obtained through these experiments are illustrated in the paper via representative frames.

Full exploitation of the algorithm's parallelism can be achieved only if the system is implemented in hardware, allowing for highly-parallelized execution.

The speed of PC based system and the projected speed of a hardware component implemented as an FPGA is discussed.

Future work will proceed in several directions: Use of features different than RGB values will be explored to evaluate the impact of the choice of features on the performance of the system. Methods to enhance the segmentation, other than morphological transformations, will be explored (e.g. single frame color-based segmentation, depth cues from stereo sequences). Finally, development of a FPGA based system which would achieve real time segmentation of HDTV and QuadHDTV sequences will be explored.

REFERENCES

- Boult, T., Micheals, R., X.Gao, Lewis, P., Power, C., Yin, W., and Erkan, A. (1999). Frame-rate omnidirectional surveillance and tracking of camouflaged and occluded targets. In *Proc. of IEEE Workshop on Visual Surveillance*, pp. 48-55.
- Ellis, T. and Xu, M. (2001). Object detection and tracking in an open and dynamic world. In *Proc. of the Second IEEE International Workshop on Performance Evaluation on Tracking and Surveillance (PETS'01)*.
- Karmann, K. P. and von Brandt, A. (1990). Moving object recognition using an adaptive background memory. In *Timevarying Image Processing and Moving Object Recognition*, 2, pp. 297-307. Elsevier Publishers B.V.
- Li, L., Huang, W., Gu, I., and Tian, Q. (2003). Foreground object detection from videos containing complex background. In *Proc. of the Eleventh ACM International Conference on Multimedia (MULTIMEDIA'03)*, pp. 2-10.
- Li, L., Huang, W., Gu, I., and Tian, Q. (2004). Statistical modeling of complex backgrounds for foreground object detection. In *IEEE Trans. Image Processing*, vol. 13, pp. 1459-1472.
- Parzen, E. (1962). On estimation of a probability density function and mode. In *Ann. Math. Stat.*, Vol. 33, pp. 1065-1076.
- Specht, D. F. (1991). A general regression neural network. In *IEEE Trans. Neural Networks*, pp. 568-576.
- Stauffer, C. and Grimson, W. (2000). Learning patterns of activity using real-time tracking. In *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, pp. 747-757.
- Ya, L., Haizhou, A., and Guangyou, X. (2001). Moving object detection and tracking based on background subtraction. In *Proc. of SPIE Object Detection, Classification, and Tracking Technologies*, pp. 62-66.