



# FUSION OF MOTION SEGMENTATION WITH ONLINE ADAPTIVE NEURAL CLASSIFIER FOR ROBUST TRACKING

Slawomir Bak, Sundaram Suresh, François Bremond, Monique Thonnat

## ► To cite this version:

Slawomir Bak, Sundaram Suresh, François Bremond, Monique Thonnat. FUSION OF MOTION SEGMENTATION WITH ONLINE ADAPTIVE NEURAL CLASSIFIER FOR ROBUST TRACKING. Int. Joint Conf. on Computer Vision, Imaging and Computer Graphics Theory and Applications (VISAPP), Feb 2009, Lisbon, Portugal. inria-00496120v2

**HAL Id: inria-00496120**

**<https://inria.hal.science/inria-00496120v2>**

Submitted on 16 Jul 2010

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# FUSION OF MOTION SEGMENTATION WITH ONLINE ADAPTIVE NEURAL CLASSIFIER FOR ROBUST TRACKING

*Preparation of Camera-Ready Contributions to INSTICC Proceedings*

Sławomir Bąk

*Institute of Computing Science, Poznan University of Technology, ul. Piotrowo 2, Poznan, Poland  
sbak@skno.cs.put.poznan.pl*

Sundaram Suresh, Francois Brémond, Monique Thonnat

*INRIA Sophia Antipolis, PULSAR group  
2004, route des Lucioles, BP93  
06902 Sophia Antipolis Cedex - France  
firstname.surname@sophia.inria.fr*

**Keywords:** object tracking, neural network, Gaussian activation function, feature extraction, on-line learning, motion segmentation, reliability classification.

**Abstract:** This paper presents a method to fuse the information from motion segmentation with online adaptive neural classifier for robust object tracking. The motion segmentation with object classification identify new objects present in the video sequence. This information is used to initialize the online adaptive neural classifier which is learned to differentiate the object from its local background. The neural classifier can adapt to illumination variations and changes in appearance. Initialized objects are tracked in following frames using the fusion of their neural classifiers with the feedback from the motion segmentation. Fusion is used to avoid drifting problems due to similar appearance in the local background region. We demonstrate the approach in several experiments using benchmark video sequences with different level of complexity.

## 1 INTRODUCTION

Visual tracking of objects in complex environments is one of the most challenging problem in the machine vision field. Tracking algorithms are developed to faithfully determine the movement of image region in each video frame that matches with the given object. In general, developing a robust tracker is a challenging problem due to change in appearance, dynamic change in background, significant illumination variation, occlusion and scale change.

For the past two decades, many algorithms with different frameworks have been developed for object tracking. Among various available algorithms, detect-then-track, appearance-based and learning-based algorithms are widely used in the literature. In detect-then-track approaches (Stauffer and Brady, 2000), objects are detected and tracked effectively in real-time using the frame differencing or subtracting adaptively estimated background from the current frame. Appearance-based approaches (Cootes et al., 2001) create an object model from the first frame and incrementally follow the model in the subsequent frames. Learning-based algorithms use pattern recognition algorithms to learn the target objects in order to search

them in an image sequence (Avidan, 2004). A complete review on different algorithms in object tracking can be found in (Yilmaz et al., 2006).

In a practical scenario, the object of interest and its local surroundings in the video sequence change significantly, which affect the success of any tracking algorithm. The object/background dynamics change due to appearance, rapid illumination variation, occlusion and scale change. Hence, developing robust tracker to handle the aforementioned issues has been the focus of the most recent object tracking research and is also the main objective of this paper.

In this work the problem of object tracking in video sequence is converted into a binary classification problem and a discriminative model is developed to differentiate the object pixel from the background. Similarly, in (Collins et al., 2005) an adaptive online feature selection mechanism is presented. The algorithm selects the best features from a given feature set that can discriminate the object and the background effectively. A number of trackers are developed based on different sets of features. Then, the tracking method adaptively selects the trackers that are robust in the current situation. In (Nummiaro et al., 2003) an adaptive target model for efficient

mean-shift tracking is presented. The target model uses a histogram to handle the change in appearance of the object. This histogram is updated linearly only when the confidence level exceeds threshold value. Nevertheless, this adaptive tracker fails when there is an abrupt change in the appearance or illumination. In (Jepson et al., 2003), short-term and long term image descriptors are constantly updated and re-weighted using online-EM (expectation maximization) to handle abrupt change in object appearance. The above mentioned approaches require proper initialization of the object bounding box in the first frame and does not handle change in scale. In (Ho et al., 2004), the linear subspace of image space is constantly updated to achieve robust tracker under challenging imaging conditions. The subspace methods aim to maintain a robust foreground object and ignore the background completely. Since the method maintains the spatial integrity, it is suitable for rigid object tracking. Similarly, Support Vector Machine (SVM) is fused with an optic-flow-based tracker to improve vehicle tracking, (Avidan, 2004). Instead of minimizing an intensity difference function between successive frames, the SVM classification score is maximized to determine the position of the object in the next frame. Another approach, (Williams et al., 2005) extends the use of statistical learning algorithm for object localization. Here, a displacement expert is build to estimate displacement from the target region. A fully probabilistic relevance vector machine (RVM) is used to generate observations with Gaussian distributions which can be fused over time. Recently, in (Avidan, 2007), the appearance of object and its background are modeled using an ensemble of classifiers. Each one of the classifiers is trained to identify object and background classes. Then a strong classifier obtained by Ada-boost is used to locate the object in the next frame.

All these approaches require proper initialization and are limited to single object. Furthermore, these methods do not consider occlusions.

In this paper, the information from motion segmentation is fused with online adaptive neural classifier to handle aforementioned issues. The online learning sequential classifier is used to differentiate the object region from non-object region. The objects present in the video sequence are detected by using the motion segmentation information combined with the reliability classification. Then, the neural classifier is initialized to differentiate the object from its local background. In following frames the motion segmentation during adaptation phase helps in avoiding drifting problems due to similar appearance in the local background region.

The outline of the rest of the paper is the following. The used platform (SUP) is described in Section 2. The neural tracking algorithm with distinguished components are presented in Section 3. Section 4 illustrates experimental results and Section 5 contains some concluding remarks.

## 2 SUP PLATFORM

The platform for image sequence understanding called SUP (Scene Understanding Platform) is developed at the research group PULSAR at INRIA (Institut National de Recherche en Informatique et en Automatique), Sophia Antipolis. SUP is an environment for combining algorithms for video analysis which allows to flexibly combine and exchange various techniques at the different stages of video understanding process. Furthermore, SUP is oriented to help developers describing their own scenarios and building systems capable of monitoring behaviors, dedicated to specific applications. SUP takes as an input video stream (or even video streams if we consider several cameras), a geometric description of the unoccupied scene, a geometric description of models appearing on the scene and a set of behaviors of interest specified by experts of the application domain. Afterwards, it is possible to process the mentioned data by combining different modules including algorithms. In our approach the object classification module performs the initialization of tracked object.

### 2.1 Object Classification

The object classification module (Zúñiga et al., 2006) performs the initialization step in our approach. This classification method is intended for monocular video sequences, allowing to classify objects modeled independently from the camera position and objects orientation. A simple 3D object model representing by a parallelepiped is applied. This model is estimated using a set of 2D moving regions (obtained in a segmentation phase), the perspective matrix transform (obtained from camera calibration using the pin-hole camera model) and predefined 3D models of expected objects in the scene. The segmentation detects moving regions by subtracting the current image from the background image. The background image can be obtained by specifying image without objects. These moving regions (also called blobs) are merged to improve the classification performance by assembling 2D moving regions with better 3D model probability. These blobs with their associated class label are called mobiles. Each detected mobile is enclosed by a

2D bounding box and by the corresponding 3D parallelepiped. In our approach, these 2D bounding boxes are used as the motion information about visible objects in the scene.

### 3 NEURAL CLASSIFIER FOR TRACKING

The neural tracker aims at distinguish foreground objects from the background. In Fig. 1 the general model of the neural tracker is presented. First, mobiles generated by the object classification are used as motion information. If the 2D bounding box of a mobile relates to an area where none target exists, a new target is initialized. In this case the target goes through three steps: feature extraction, object/background separation and neural network training. Otherwise, if a target exists in the corresponding area, the tracking algorithm is applied. For each target the new position is predicted in the target surroundings during the localization step. Then, in the adaptation phase, the dimension of the target is corrected and the neural network is adapted using features computed in the target region. The following subsections describe the different components of the neural tracker.

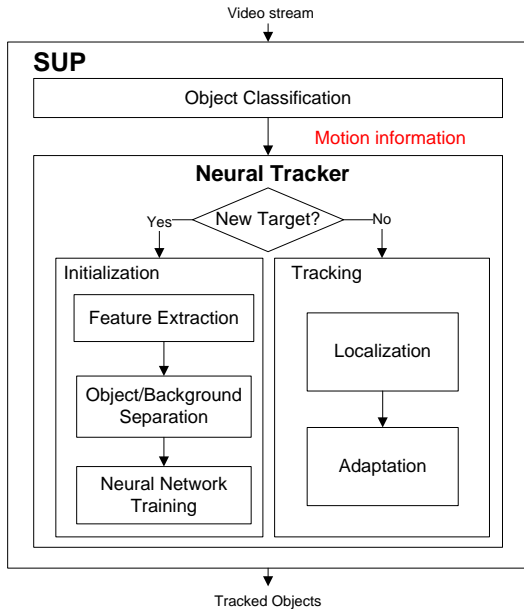


Figure 1: Model of the neural tracking system.

#### 3.1 Target Initialization

The initialization of new targets is based on the motion information obtained from the reliability classification. This information is given in a form of coordinates of 2D bounding boxes of mobiles with labels assigned in the classification phase. Features are extracted for each 2D bounding box together with its local background. The label of the target is used to initialize some constants in the neural network algorithm (e.g. maximal number of neurons in the neural network or size of extracted regions in feature extraction step).

#### 3.2 Feature Extraction

This step extracts features from target region and returns them as vectors. The neural network can use different types of features such as simple color features, region based features, local gradients or texture. In our work RGB color and texture features are applied. To obtain texture features the range filter is used. This method is based on morphological operations on intensity value. During this step the feature vector is obtained. It contains values representing features (e.g RGB values).

The feature vector is computed for each pixel in case of small objects (max. 1000 pixels). Otherwise we extract region based features like Region based Color Moments (RCM) to speed up the approach. In this case, the target is separated into rectangular regions. For each region the feature vector is calculated using mean values of pixels belonging to the region. These vectors are used for object/background separation.

#### 3.3 Object/Background Separation

Tracking an object will be efficient if we can separate the object region from the background accurately. The problem of tracking is converted into a binary classification problem and solved using neural networks. The online learning neural tracker presented in this paper is a generic approach, which is trained to separate the object region from the background. For the training step, object/background separation scheme is used. We represent the target as a 2D bounding box. The features inside the target can be labeled as ‘object class’ and features outside target region can be labeled as ‘background class’. Some of the features inside the target region will be similar to the background. Such features have to be labeled as ‘background class’ for better classification performance. For this purpose, we use a feature-based



object-background separation technique, see Fig. 2. The image region inside the yellow 2D bounding box represents the target region. The local background region is represented by the area between the yellow and the blue bounding boxes. This area is dependent on the size of the target. The probability density function of the feature in the target region and its local background is obtained to find the log-likelihood ratio of the sample belonging to the ‘object class’. The log-likelihood ratio  $L_i$  is obtained as

$$L_i = \log \frac{\max\{h_o(i), \epsilon\}}{\max\{h_b(i), \epsilon\}} \quad (1)$$

where  $h_o(i)$  and  $h_b(i)$  are the probabilities of  $i$ th sample belonging to the object and the background, respectively;  $\epsilon$  is a small non-zero value to avoid numerical instability. The class label  $y$  for a given  $i$ th sample is coded as

$$y_i = \begin{cases} 1 & \text{if } L_i > \tau_o \\ -1 & \text{otherwise} \end{cases} \quad (2)$$

where  $\tau_o$  is the threshold to decide on whether the feature belongs to target. Typical value of  $\tau_o$  is set at 0.65. The class label map obtained using the log-likelihood ratio is shown in Fig. 2 as a result of the object/background separation scheme.

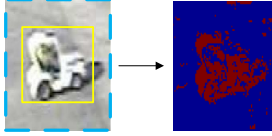


Figure 2: Object/background separation. Red pixels have features classified as target.

### 3.4 Neural Network Classifier

The neural network proposed in this paper adapts the temporal change in appearance of the object/background model. The basic building block of the neural classifier is the Radial Basic Function Network (RBFN). The RBFN is trained to recognize the object region and its local background by estimating the posterior probability of the feature vector belonging to the object or the background. The RBFN architecture consists of an input layer, a hidden layer and an output layer as shown in Fig 3.

The inter-connection weights only exist between the hidden and the output layer. Gaussians units are used in the hidden layer as activation functions because of their localization properties. Generally, the

output of the RBFN classifier with  $K$  hidden neurons has the following form:

$$\hat{y}_i = f(\mathbf{X}_i) = \sum_{j=1}^K \alpha_j \exp\left(\frac{-\|\mathbf{X}_i - \mu_j\|}{2\sigma_j^2}\right) \quad (3)$$

where  $\mu_j$  and  $\sigma_j$  is the center and width of the  $j$ th neuron, respectively;  $\alpha_j$  is the weight connection between the  $j$ th neuron and the output,  $\alpha_j \in \mathbb{R}$ ;  $\mathbf{X}_i$  is the feature vector of  $i$ th sample.

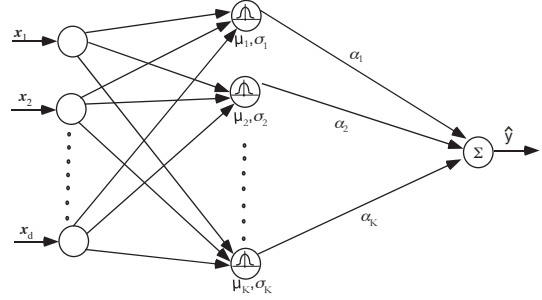


Figure 3: Architecture of RBFN used in this study.

RBFN classifier involves the allocation of new Gaussian hidden neurons, pruning the neurons and also adapting the neuron parameters. The RBFN begins with no hidden neurons (i.e., the network output is zero for the first sample). When the observation data are received sequentially, the network starts growing and shrinking by using some of them to make decision based on certain criterions. The input of the neural network are vectors obtained in the feature extraction phase. However, the output is compared with the class label map obtained by the object/background separation scheme to compute the classification error. We use a risk sensitive hinge loss function (emphasizing the error impact) to calculate the error in learning instead of simple mean square error. In (Suresh et al., 2008a) has been shown that the classifier developed using a risk sensitive hinge loss function can estimate the posterior probability more accurately than the mean square error loss function under sample imbalance. Growing criterion is based on the error value ( $e_i$ ) and a distance between the nearest neuron and the training sample. If the distance between the new training sample and the nearest neuron of the same class exceed a threshold a new hidden neuron will be added to the RBF classifier. Its parameters are set as follows:

$$\alpha_{K+1} = e_i \quad (4)$$

$$\mu_{K+1} = \mathbf{X}_i \quad (5)$$

$$\sigma_{K+1} = \kappa \|\mathbf{X}_i - \mu_{nr}\| \quad (6)$$

where  $\kappa$  is a positive constant which controls the overlap between the hidden neurons. The value of  $\kappa$  is 0.8 in our experiments.

When the new training sample does not satisfy the criterion for adding a new hidden neuron, the network parameters are adapted using a Decoupled Extended Kalman Filter (DEKF) (Puskorius and Feldkamp, 1994). More details on learning algorithm can be found in (Suresh et al., 2008b).

In order to maintain a compact network and remove a non-performing neuron, a pruning strategy is incorporated in the algorithm. Pruning of neurons ensures that the neurons that have been added in the past and have not contributed significantly to the network performance are removed from the network. This strategy consists in removing neurons which contributions is less than a threshold for  $M$  consecutive observations. The parameter  $M$  depends on the number of samples in each class.

After training phase, each sample is tested using the final neural network to calculate the probability map, see Fig 4.

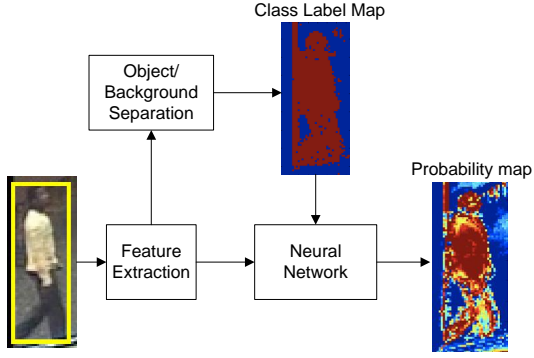


Figure 4: The neural network learning phase is composed of three components

### 3.5 Localization

Object localization phase is needed to calculate a displacement of the target. It is achieved by computing the distance between weighted centers of the target in current and previous frame. Let  $c^i$  be the weighted object center obtained from the corresponding probability map  $\hat{p}^i$  in  $i^{th}$  frame. Then, using the neural network and the location of the target from  $i^{th}$  frame we calculate the posterior probability map  $\hat{p}^{i+1}$  in  $(i+1)^{th}$  frame. Here we assume that there is an overlap between the target region in the subsequent frames. Otherwise, the tracker fails to detect the correct center of the object. It can happen considering fast moving objects. In this case the tracker should use the motion information to compare generated 2D bounding boxes with already existing targets. Assum-

ing an overlap, the new center is estimated as

$$c^{i+1} = \left[ \frac{\sum_j I_j^{i+1} \hat{p}_j^{i+1} \hat{p}_j^i}{\sum_j \hat{p}_j^{i+1} \hat{p}_j^i} \right] \quad (7)$$

where  $I_j^{i+1}$  is the location (coordinates) of the  $j^{th}$  feature. Next, the 2D bounding box is moved by the vector computed using both centers. The probability map is updated considering the new location of the target and the procedure is repeated until the distance between both centers ( $\Delta c$ ) can be neglected. The pseudo code of the procedure is presented in Listing 1.

---

#### Listing 1 Pseudo code of the localization algorithm.

---

```

Localization
begin
  do
    compute the probability map
    compute the new center
    compute the distance
    move the object
  while(  $\Delta c > 2$  pixels)
end.

```

---

### 3.6 Adaptation

This phase is responsible for the correction of the size of the object bounding box and the adaptation of the neural network to handle changes in appearance of the object/background model. This procedure is based on the probability map obtained from the neural classifier and the motion information generated by object classification module.

Firstly, we use the probability map to find the current object dimension. For this purpose, we define the class label map based on the estimated posterior probability. Then, the morphological operations (close + open) are applied to remove noises. After noise removal, the map is used for size determination based on a calculation of continuous regions in width and height dimension.

Secondly, the object classification generates mobiles which are compared with existing targets. If a mobile relates to an area where none target exists then a new one is initialized. Otherwise, the mobile is used to adapt the target located in this area. The history of dimensions of the targets and mobiles in previous frames is also used. Here, the information coming from motion segmentation is fused with the estimated posterior probability map generated by the neural classifier. The pseudo code of the information fusion is given in Listing 2.

The algorithm starts by checking if there is a mobile in the related target area. We do not adapt the

---

**Listing 2** Pseudo code for information fusion at adaptation level.

---

```
Adaptation(framei)
begin
  if no mobile /* the object classification fails */
    do not adapt the neural network and the target size
  else
    [ $c^t, s^t$ ]i - apply morphological operations
    [ $c^m, s^m$ ]i - mobile from the object classification
    compare [ $c^t, s^t$ ]i with [ $c^t, s^t$ ]i-1 and [ $c^m, s^m$ ]i with [ $c^m, s^m$ ]i-1
    if(  $\Delta s^t > 40\%$  and  $\Delta s^m > 40\%$  )
      do not adapt the neural network and the target size
    else
      /* fusion of neural classifier and motion segmentation */
      compare [ $c^t, s^t$ ]i with [ $c^m, s^m$ ]i
      if(  $\Delta s_i > 10\%$  )
        compare [ $c^t, s^t$ ]i with the history of mobiles (max 5 previous mobiles)
        if( min(  $\Delta s$  ) > 20% )
          [ $c^t, s^t$ ]i := [ $c^m, s^m$ ]i /* use the center and the size of current mobile */
          update the history of mobiles
        adapt the neural network using [ $c^t, s^t$ ]i
        [ $c^{new}, s^{new}$ ]i - apply morphological operations
        compare [ $c^{new}, s^{new}$ ]i with [ $c^t, s^t$ ]i
        if(  $\Delta s > 10\%$  ) /* the neural network fails */
          retrain the neural network
      end.
end.
```

---

neural network and the target size in the case of absence of mobiles. If the mobile exists we try to fuse mobile information with output of the neural classifier. In order to achieve this we compare both information to make decision about the true dimension of the object. First, morphological operations are applied to estimate the current target size. Let us assume that the center and the size at frame  $i$  are represented by  $[c, s]_i$ . The result of morphological operations is given as  $[c^t, s^t]_i$ . While the mobile is represented by  $[c^m, s^m]_i$ . Next, both target and mobile size are compared with target and mobile size at previous frame, respectively. If the size change ( $\Delta s$ ) in both cases is greater than 40% we do not allow to adapt the target size and also the parameters of the neural network are not modified. Here, we try to avoid abrupt changes obtained from the object classification which happens in a case of shadows, crowds or another disturbances coming from the motion information. Then, the fusion rules is used to decide whether it is necessary to correct the dimension and the localization of the target. We check the difference between the target and the mobile size ( $\Delta s_i$ ) in current frame. If it is greater than 10% we compare the target size with the history of mobiles. If the history suggest that the mobile size is correct (at previous frames, the mobile size was always greater than current target size) we accept the

new size and the new center. After establishing of the target dimension the neural network is trained using samples from the target and its local background region. The new probability map is computed and the new dimension of the target is determined. If the difference between new and previous size is greater than 10% we assume that the neural network fails. All neurons are removed and the network is retrained.

After the adaptation step, the tracker checks if there are any occlusion in the scene. Occlusion events are presented in Listing 3. We distinguish two types of occlusion. First, the static occlusion means that the moving target is occluded by static item. This decision is made only if the density value of the probability map decreases more than 10%. However, the dynamic occlusion is detected when the target crosses another target and the common area is greater than 60% of the area belonging to one of them. We assume that the first target is occluded if the density value is less than the density of second target.

## 4 EXPERIMENTAL RESULTS

To validate the tracking approach we have tested many video sequences. This section presents tracking results only for challenging video sequences that

---

**Listing 3** Occlusion events.

---

```
OcclusionEvent(Target A)
  Constraint{A.density < 10%}
  Alarm("Occlusion")
  State_Assignment{A.state = OCCLUDED}

DynamicOcclusionEvent(Target A, Target B)
  Constraint{Crossing.Area > 60% of A.Area}
  Alarm("Dynamic Occlusion")
  State_Assignment{
    if(A.density < B.density)
      A.state = DYNAMIC_OCCLUDED}
```

---

illustrate the advantage and limitations of the tracker. The online learning phase enhances the ability to track under changing background and illumination conditions, changing in appearance and scale and improper initialization. Furthermore we propose solutions for some difficulties such as occlusions or shadows occurrence. We have also tested our approach on weakly contrasted sequences where the detection of the true dimension of the object is challenging problem.

#### 4.1 Gerhome video sequence

The experiments were performed on the data obtained from the Gerhome laboratory which is a realistic site reproducing the environment of a typical apartment of an elderly person. The laboratory is located in the CSTB (Centre Scientifique de Techniques du Bâtiment) at Sophia Antipolis. The Laboratory promotes research in the domain of activity monitoring and assisted living (Zouba et al., 2008). In Fig. 5, an example sequence is presented. At the first frame a target is initialized as a bounding box with a label which is obtained from the object classification module. During initialization for each new target the neural network is created and an identity is assigned. Next, this neural network is trained using features computed from the target and the local background region. At frame 67 we can observe important issues. The motion information is not always correct which leads to track noise due to illumination change ('2-PERSON'). Nevertheless if in subsequent frames the target is not confirmed by the motion information coming from the object classification, the target is assumed to be a noise and afterwards is removed. However a more important issue is that the tracked person is split into two targets ('0-PERSON' and '3-PERSON'), caused also by noise coming from motion information. We do not apply any merging algorithm for tracked targets because it is very difficult to de-

cide whether few targets in fact represent one real object or several different objects. The neural network is not also helpful in that case because parts of a real object could have a completely different appearance model which prevents from merging such kind of targets. At frame 93 we can see '0-PERSON' marked as occluded target. This information is correct because the woman is occluded by a cupboard door. Here, the rules which are based on the density value of the probability map returned by the neural network are applied to make a decision about presence of occlusion. However, the density value is not always reliable. Two targets with similar size of bounding box but with different silhouette could have completely different density value. The neural network in our approach does not use information about finer object shapes which should be considered in future work. At frame 133 we show that the neural tracker is able to capture true dimension of the object in following frames.

We also tested our approach on a long-term sequence. An elderly woman in her apartment (real world scene with occlusions and illumination changes) was tracked during 1 hour and 23 minutes (50.000 frames). A woman left for short periods of time the observed room (and came back) 7 times. The tracker was confused only 22 times (id switched).

#### 4.2 ETISEO video sequence

ETISEO was a Video Understanding Evaluation project which ended successfully in the December 2006. The ETISEO data is a public resource intended for research communities (Nghiem et al., 2007). We present a sequence from metro scene (ETI-VS1-MO-7-C1) where a man leaves his bag on a metro station, Fig 6. This is an example of video with weakly contrasted objects. As we can see at frame 1, two targets are initialized. Let us focus on the target '0-PERSON'. The contrast between the man's trouser and the background is so weak that legs are not detected using a background segmentation algorithms. In this case the neural network has also problems with capturing legs because of similarity in the background and leg appearance. When the person drops the bag the neural network recognizes the bag as a part of the person which is shown at frame 140. The person and the bag are recognized as a single target until the object classification module returns the information about new classified object '5-LUGGAGE'. However in following frames the bag is occluded by the person legs which leads the tracker to decide that it is a noise and should be removed. At frame 241 the bag is initialized again and a new identity is assigned ('7-LUGGAGE'). When the person goes away the bag is

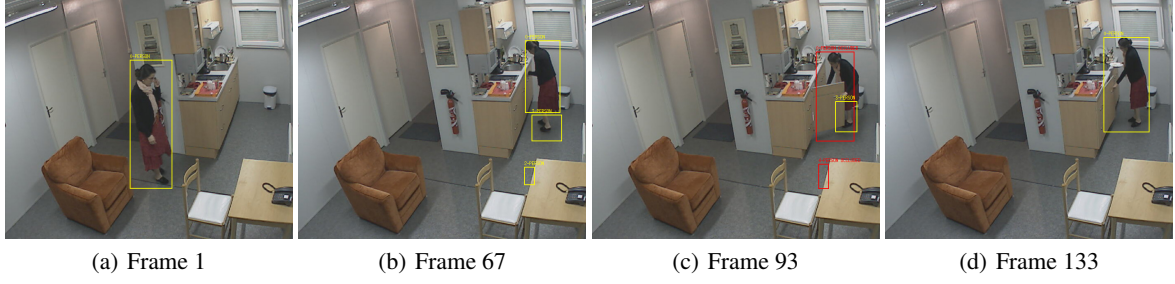


Figure 5: Tracking results from Gerhome laboratory. Person is split into two parts.

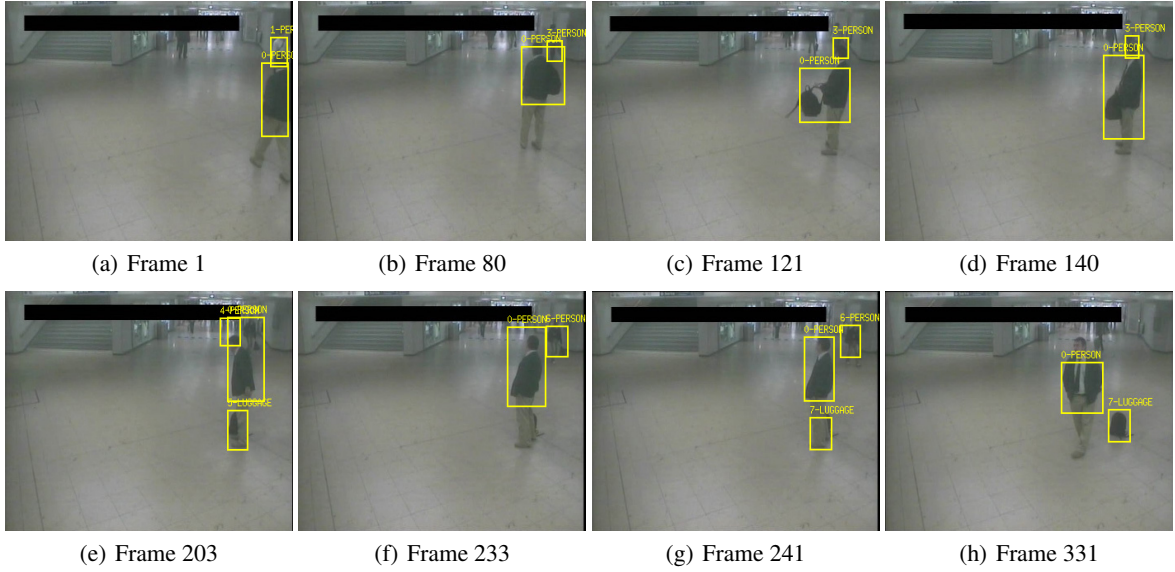


Figure 6: Tracking results for ETISEO data with weakly contrasted objects.

still tracked which enables to recognize the scenario ‘abandoned luggage’. At last frame we can see that the person legs are still not detected as a target but despite that the tracker is able to track the person using only the body trunk.

### 4.3 TRECVID video sequence

We tested our approach on TREC Video Retrieval Evaluation (organized by NIST, TRECVID 2008) data obtained from Gatwick Airport surveillance system. Firstly we present a sequence where two people are crossing each other. The neural tracker has the ability to manage with dynamic and static occlusions which are presented in Fig. 7. We can see that two targets are detected, ‘0-PERSON’ and ‘1-PERSON’, respectively. For the next two frames (32, 43) we indicate that the proposed algorithm can adapt the bounding box and capture the true dimension effectively. When two targets are crossing (frame 54) dynamic

occlusion is detected. To make decision which targets are occluded the neural network is used. For each target the probability map of the overlapping area is computed and the most probable one is chosen. The most probable target means the target which has the largest density value of the probability map. During dynamic occlusion localization stage does not use probability map to localize target but its history which is based on the target displacement. Also the adaptation process is suspended. The history contains the information about the velocity and the direction of the target motion represented by displacement vector. Using that information our tracker is able to approximate location of the occluded target. We resume using the neural network only if we confirm separated objects from the motion information. At frame 63 the shadows problem is presented. The person with identity 2 is initialized. Nevertheless in subsequent frames the target is not confirmed by the motion features which cause that the target is assumed to be a



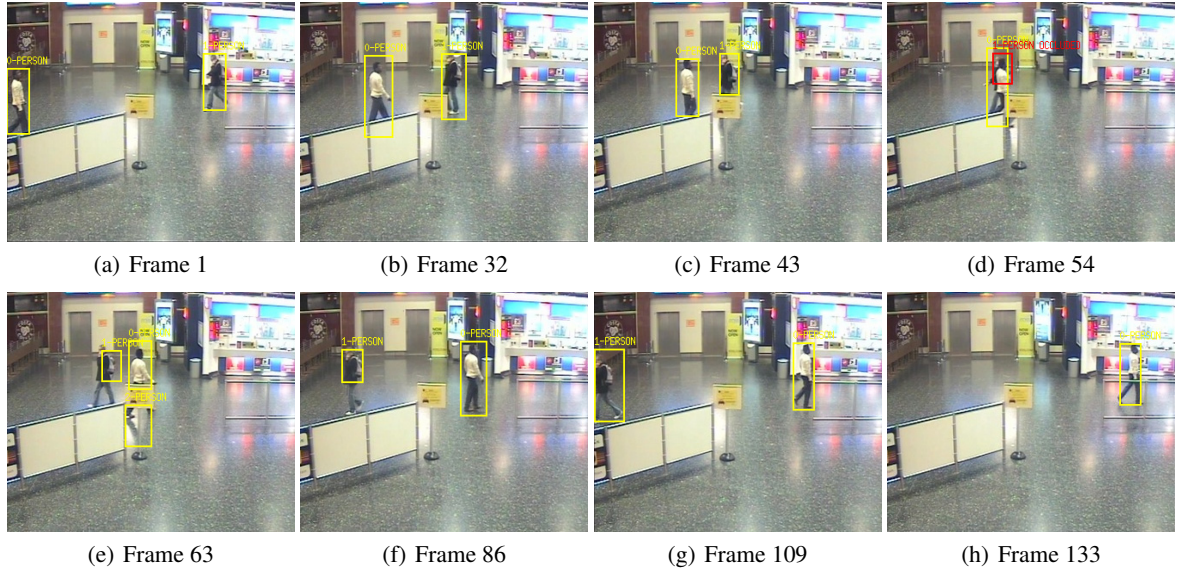


Figure 7: Tracking results for TREC Video Retrieval Evaluation data for crossing people.

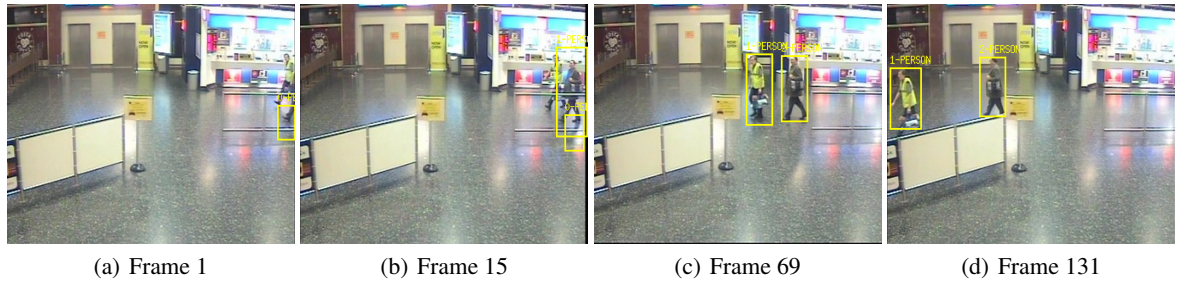


Figure 8: Tracking results for TREC Video Retrieval Evaluation data. Incorrect initialization.

noise and afterwards is removed. Moreover at frame 133, an example presents a situation where tracking is very difficult because of similarity in the background and in the object appearance. In that case the motion information is used to improve the tracking robustness. It is worthy noting that using only an appearance information coming from the neural network is not always possible to differentiate the object from the background accurately. Consequently we can not obtain robust tracking using only appearance information which leads us to fuse the neural classifier with the motion information.

However mobiles are not always correct which is showed in Fig 8. At the first frame the target is badly initialized because of the shadows problem. Incorrect initialization once again leads to track noise. Hence, using mobiles from following frames we correct the tracking by retraining the neural network or by removing the target which has been recognized as a noise and reinitialized as a new one. At frames 69 and

131 we show that despite bad initialization the tracker is able to capture correct dimension of the target.

In Fig. 9 we present more complex sequence in which there are many objects crossing each other. At first frame targets are initialized. We can observe that both, '2-PERSON' and '3-PERSON' in fact should be represented by two separated targets. But at this level of information, where we do not use any face detection algorithm it is impossible to separate these objects. Unfortunately, at frame 15, a woman '4-PERSON' is lost because of noise coming from motion information. It is caused by a shadow effect and occlusions (by the barrier and '0-PERSON'). The woman is reinitialized as a new target later. However, at frame 52 is shown that the tracker is able to detect few objects crossing each other. Both '0-PERSON' and '1-PERSON' are recognized as occluded targets because their density is less than density of the '3-PERSON' but none of these targets is lost. At frame 110 is shown that the tracker is able to track objects

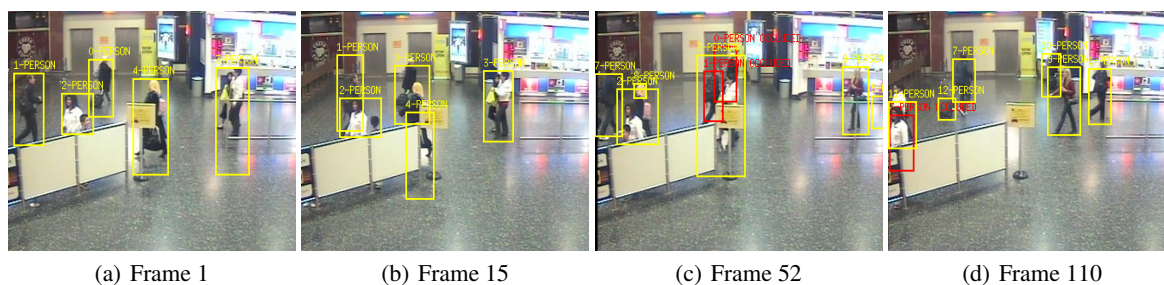


Figure 9: Tracking results for TREC Video Retrieval Evaluation data. Complex sequence.

effectively during very complex sequence. The target ‘1-PERSON’ is crossing many other targets and its identity remains unchanged.

## 5 Conclusions

We have presented an online learning neural tracker to handle dynamic changes in object/background appearance, illumination and scale. The problem of tracking is treated as a binary classification problem, and online learning radial basis function classifier is used to differentiate the object and the background region. The information from motion segmentation is fused with this neural classifier for robust object tracking. The results indicate that the tracker is robust under normal illumination variation, appearance and scale changes. Occlusions are also handled using the motion and the probability map information.

Furthermore, consideration of different features like shape or silhouette might be beneficial. Additional research can also be carried out in order to handle separated parts of the same object and improve the initialization step.

## REFERENCES

- Avidan, S. (2004). Support vector tracking. *IEEE Trans. Pattern Anal. Machine Intell.*, 26(8):1064–1064.
- Avidan, S. (2007). Ensemble tracking. *IEEE Trans. Pattern Anal. Machine Intell.*, 29(2):261–271.
- Collins, R. T., Liu, Y., and Leordeanu, M. (2005). Online selection of discriminative tracking features. *IEEE Trans. Pattern Anal. Machine Intell.*, 27(10):1631–1643.
- Cootes, T., Edwards, G., and Taylor, C. (2001). Active appearance models. *IEEE Trans. Pattern Anal. Machine Intell.*, 25(5):681–685.
- Ho, J., Lee, K., Yang, M., and Kriegman, D. (2004). Visual tracking using learned linear subspaces. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, volume 1, pages 782–789.
- Jepson, A. D., Fleet, D. J., and El-Maraghi, T. (2003). Robust online appearance model for visual tracking. *IEEE Trans. Pattern Anal. Machine Intell.*, 25(10):1296–1311.
- Nghiem, A. T., Brémond, F., Thonnat, M., and Valentin, V. (2007). Etiseo, an evaluation project for video surveillance systems. In *IEEE International Conference on Advanced Video and Signal based Surveillance*, London, U.K.
- Nummiaro, K., Koller-Meier, E., and Van Gool, L. (2003). An adaptive color-based particle filter. *Image vision computing*, 21(1):99–110.
- Puskorius, G. and Feldkamp, L. (1994). Neurocontrol of nonlinear dynamical systems with kalman filter trained recurrent networks. *IEEE Trans. on Neural Networks*, 5(2):279–297.
- Stauffer, C. and Brady, J. (2000). Learning patterns of activity using real-time tracking. *IEEE Trans. Pattern Anal. Machine Intell.*, 22(8):747–757.
- Suresh, S., Sundararajan, N., and Saratchandran, P. (2008a). Risk sensitive loss functions for sparse multi-category classification problems. *Information Sciences*, 178(12):2621–2638.
- Suresh, S., Sundararajan, N., and Saratchandran, P. (2008b). A sequential multi-category classifier using radial basis function networks. *Neurocomputing*, 71(7-9):1345–1358.
- Williams, O., Blake, A., and Cipolla, R. (2005). Sparse bayesian learning for efficient visual tracking. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 27(8):1292–1304.
- Yilmaz, A., Javed, O., and Shah, M. (2006). Object tracking: A survey. *ACM Computing Surveys*, 38(4):1–45.
- Zouba, N., Boulay, B., Brémond, F., and Thonnat, M. (2008). Monitoring activities of daily living (adls) of elderly based on 3d key human postures. In *Proc. of International Cognitive Vision Workshop*, pages xx–xx, Greece.
- Zúñiga, M., Brémond, F., and M.Thonnat (2006). Fast and reliable object classification in video based on a 3d generic model. In *Proceedings of the International Conference on Visual Information Engineering (VIE2006)*, pages 433–440, Bangalore, India.