

## Object Tracking Based on Online Semi-Supervised SVM and Adaptive-Fused Feature

Ruxi Xiang<sup>1,2,3</sup>, Xifang Zhu<sup>1,2,3</sup>, Feng Wu<sup>1,2,3</sup>, Qinquan Xu<sup>1,2,3</sup>, Jianwei Li<sup>4</sup>

<sup>1</sup>College of Optoelectronic Engineering, Changzhou Institute of Technology, 213002 China

<sup>2</sup>Changzhou Institute of Modern Optical Technology, 213002 China

<sup>3</sup>Changzhou Key Laboratory of Optoelectronic Materials and Devices, 213002 China

<sup>4</sup>Key Laboratory of Optoelectronic Technology and Systems of Ministry of Education, Chongqing University, Chongqing, 400044 China

Email: zhuxfcz@yeah.net

**Abstract:** In order to improve the performance of tracking, we propose a new online tracking method based on classification and adaptive fused feature. We first label a few positive and negative samples, train the classifier by the online SSSM (Semi-Supervised Support Vector Machine) learning and these labelled samples, and then locate the position of the object from the next frame according to the trained classifier. In order to adapt more of the new samples, we need to update the classifier by finding new samples with high confident value obtained by the trained classifier and add them into the online SSSM. Finally we also update the object model by the online incremental PCA (Principal Component Analysis) because of background clutter, heavy occlusion and complicated object appearance changes. Compared with the basic mean shift tracking and the ensemble tracking method, experimental results show that our tracking method is able to effectively handle heavy occlusion and background clutter in some challenge videos including some thermal videos.

**Keywords:** Visual tracking, SSSM, feature fusion, incremental PCA, online.

### 1. Introduction

Recently, tracking problem formulated as a binary classification has been a promising direction used in some visual applications such as video surveillance and military field [1, 2]. The key idea of this kind of tracking method is to obtain an effective classifier from training samples which are labelled by online or offline method. In spite of speeding up the matching process by offline method, the classifier trained by offline method can't deal efficiently with changes in appearance or complex background, while the classifier, trained by the online method, on the other hand, can solve these problems effectively [3].

Many methods based on classifier have been proposed for visual tracking [1, 3-6]. S. Avidan [4], introduced the support vector tracker, combined the trained classifier by offline method with optical-flow feature to track the object. S. Avidan proposed the ensemble tracker which trained some weak classifiers by Adaboost method to obtain stronger classifier [5]. However, the ensemble tracker used the classification results with high confidence to update the classifier itself. Unfortunately, the classification error is gets more prominent over time and the tracking error increases as well. To further reduce the tracking error, co-tracking based on semi-supervised support vector machines [3] was proposed, with the underlying assumption that each feature is different and independent but completely independent features are difficult to obtain in reality. Co-tracking method makes use of two independent features to train the corresponding classifiers, produces the confidence maps for a new unlabelled sample using the trained classifiers, and then combines a final confidence map with the produced confidence maps. However, co-tracking method neglects the contribution of the each feature taking into consideration the fusion of two maps which were produced by the trained classifiers. To take full advantage of each feature, we combine all features to form an efficient feature which is helpful for distinguishing the object from the background by the distance between the candidate object and the model which is adaptively updated by online increment PCA [7].

Motivated by self-training and the fusion features idea [8], we present a new tracking method based on semi-supervised support vector machine and adaptive fused feature. We first label a few positive and negative samples, train the classifier by the online SSSM [9] learning and these labelled samples, and then locate the position of the object from the next frame according to the trained classifier. In order to adapt more of the new samples, we need to update the classifier by finding new samples with high confident value obtained by the trained classifier and add them into the online SSSM. Finally we also update the object model by the online incremental PCA because of background clutter, heavy occlusion and complicated object appearance changes. Compared with the basic mean shift tracking and the ensemble tracking method, experimental results show that our tracking method is able to effectively handle heavy occlusion and background clutter in some challenge videos including some thermal video.

The rest of this paper is organized as follows. Section 2 briefly reviews the online support vector machine and the self-training method. Section 3 describes the features fused by the new similarity MSSBRS and model object updated by incremental PCA. Our method is described in Section 4. Experimental results and comparisons are shown in Section 5. Finally, the conclusions are made in Section 6.

## 2. Background

### 2.1. Online support vector machine

Support Vector Machine (SVM) which is also a powerful learning tool for classification and regression samples easily solves small and high dimensional samples. When the number of the samples is very large the computation of SVM

will take considerable time. To reduce the computational time of training samples, Gert improves the basic SVM by incremental and decreasing leaning method [9].

Assume that training data  $(\mathbf{x}, \mathbf{y})$  with  $\mathbf{x} \in R^n$  and  $\mathbf{y} \in \{-1, 1\}$  is obtained from several frames, the SVM classifier is then represented as  $f(\mathbf{x}) = \mathbf{a}^T \phi(\mathbf{x}) + b$ , in which  $\mathbf{a}$  is a weight vector,  $b$  is a constant, and  $\phi(\mathbf{x})$  is a map function from the input space to feature space,  $N$  is the number of the training samples, and the classifier can be learned as

$$(1) \quad \min_a \left( \frac{1}{2} \|\mathbf{a}\|^2 + \omega \sum_{i=1}^N \lambda_i \right),$$

$$\text{s.t. } \mathbf{y}_i (\mathbf{a} \phi(\mathbf{x}_i) + b) \geq 1 - \lambda_i.$$

Here  $\lambda > 0$  is a slack parameter, and  $\omega$  is a penalty factor which represents the degree of punishment to the error. To further simplify the formula (1) can be rewritten as

$$(2) \quad \min_{0 \leq \alpha_i \leq \omega} W = \frac{1}{2} \sum_{i,j} \alpha_i \mathcal{Q}_{i,j} \alpha_j - \sum_i \alpha_i + b \sum_i y_i \alpha_i,$$

where  $\mathcal{Q}_{i,j} = \mathbf{y}_i \mathbf{y}_j k(\mathbf{x}_i, \mathbf{x}_j)$ , and  $k(\cdot)$  is a kernel function including Gaussian RBF

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2h^2}\right) \text{ in which } h \text{ is a variance of Gaussian RBF or}$$

polynomial kernel function  $k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + 1)^d$  in which  $d$  is a value of degree.

In this paper, we adopt the Gaussian RBF as our kernel function. The solution of the dual parameters  $\{a, b\}$  can be defined by minimizing (2). These solutions need to satisfy the KKT (Karush-Kuhn-Tucker) conditions expressed by (1) and (2) because it is a necessary condition for a solution in nonlinear to be optimal:

$$(3) \quad f_i = \frac{\partial W}{\partial \alpha_i} = \sum_j \mathcal{Q}_{i,j} \alpha_j + y_i b - 1 = y_i f(\mathbf{x}_i) - 1,$$

$$(4) \quad \frac{\partial W}{\partial b} = \sum_i y_i \alpha_i = 0.$$

Based on the partial derivatives  $f_i$ , the training samples can be classified into three categories, the samples on the margin are usually called support vectors ( $f_i = 0$ ), while samples exceeding the margin are error vector ( $f_i < 0$ ), the remaining samples, within the margin, are called reserve vectors ( $f_i > 0$ ).

In [9], incremental or decreasing SVM learning is a procedure that adds or removes samples from one vector of SVM according to the value  $f_c$  that is a partial derivative value adding the training yields margin each time. In the incremental procedure, new samples with  $f_c > 0$  are directly added to the reserve vector set due to the fact that they don't affect the solution. On the contrary, all other new samples are added to the set of either margin vector or error vector according to the value of  $f_c$  and  $a_c$ . that is a coefficient being incremental. In the decreasing procedure, the samples are removed according to the value  $f_c$  and  $a_c$ , No matter how the samples are added or removed, they should satisfy the KKT conditions.

## 2.2. Self-training method

Self-training, also called self-teaching or bootstrapping, is a classical training algorithm in semi supervised leaning [10]. In self-training frame, a classifier is first trained with only a few labelled samples, and then unlabelled samples are classified by the classifier. Typically unlabelled samples with the highest confidence are added to the training set, and the classifier is retrained.

## 3. Fusing feature and updating model

### 3.1. Feature fusion

#### 3.1.1. HOG

HOG (Histogram of Orientation Gradient Edge) is an effective representation of object appearance because it is insensitive to the change of illumination. We first pre-process the object by the guided filtering method [11], compute the edge (Fig.1b) of [12], and then reserve some pixels with greater than double the mean value. The orientation of the gradient edge of the object is quantified using  $N_h$  bins. In this paper,  $N_h$  is set 16.

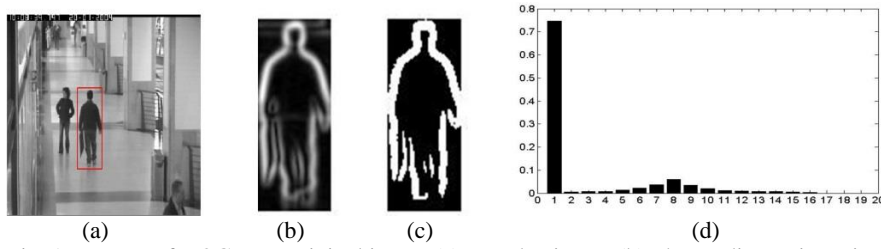


Fig. 1. Process of HOG: an original image (a); an edge image (b); the gradient orientation corresponding to the edge image (c); and a HOG with 16 bins (d)

#### 3.1.2. Fusing feature

During the tracking, each feature is different in contribution of the different frames, the more the contribution is, the bigger the weight is. On the contrary, the less the contribution is, the smaller the weight is, so it is extremely important and crucial to compute the weight of each feature.

Let  $fe_i$  and  $fe_{fu}$  represent the  $i$ -th feature and fusion feature of the object, respectively, the multiple features are fused by

$$(5) \quad fe_{fu} = \sum_{i=1}^{N_s} \eta_i fe_i, \quad \text{s.t.}, \quad \sum_{i=1}^{N_s} \eta_i = 1,$$

where  $\eta_i$  denotes the coefficient of the  $i$ -th feature, which is the normalized value of the similarity measure between the  $i$ -th feature and the object model. In this paper, the adopted similarity measure is the MSSBRS introduced in Section 3.2 that is more discriminative than the common measure, such as the Bhattacharyya similarity and the BRS. The selected features are different in visible video and thermal video, for example, the selected features are grayscale and HOG in thermal video and two features are selected from color and HOG in visible vide .

### 3.2. MSSBRS

The BRD [13] (Bin-Ratio Dissimilarity) is a new similarity measure which considers the ratios between bin values of histograms, and The BRD is successful in the scene of classification. Nevertheless, it is prone to neglecting spatial information of the object. To improve the performance of distinguishing objects, we have introduced a modified BRD named the SBRS [14] which is better than the similarity Bhattacharyya and the BRS. We introduce a Modified Similarity SBRS (MSSBRS) which considers not only the spatial structure of the object but also the relation of the bins of the histogram of between the candidate object and the reference object.

Assume  $q_{\text{ref}}$  with second-order spatiogram [15] to represent the reference object, which is defined as  $q_{\text{ref}}(j)=[\mathbf{h}_j, \mathbf{m}_j, \mathbf{c}_j]$ , where  $\mathbf{h}_j$  is the  $j$ -th bin of the histogram, and  $\mathbf{m}_j, \mathbf{c}_j$  denote the mean vector and covariance of the coordinates of the pixels corresponding to the  $j$ -th bin respectively. Meanwhile, the histogram of the edge needs to satisfy  $\sum_{j=1}^M \mathbf{h}_j^2 = 1$ , where  $M$  is the number of the bins. Assume  $p_{\text{cal}}$  with second-order spatiogram, which defines as  $q_{\text{cal}}(j)=[\mathbf{h}'_j, \mathbf{m}'_j, \mathbf{c}'_j]$  to represent the candidate object. So the similarity  $\rho(q_{\text{ref}}, p_{\text{cal}})$  between reference spatiogram  $q_{\text{ref}}$  and candidate spatiogram  $p_{\text{cal}}$  is defined as

$$(6) \quad \rho(q_{\text{ref}}, p_{\text{cal}}) = \sum_{j=1}^M \beta_j d_j(\mathbf{h}_j, \mathbf{h}'_j),$$

where  $\beta_j$  is the weight of the  $j$ -th bin similarity, if the spatiogram is a 0-th order,  $\beta_j$  is set to 1, and if it is a second-order spatiogram,  $\beta_j$  can be defined as  $\beta_j = a_0 N(\mathbf{m}_j; \mathbf{m}'_j, k(\mathbf{c}_j + \mathbf{c}'_j))$ , in which the factor  $a_0$  ensures the similarity measure that satisfies the condition  $0 \leq \rho \leq 1$  and  $\rho(\mathbf{h}, \mathbf{h}) = 1$ ;  $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is a Gaussian function with a mean vector  $\boldsymbol{\mu}$  and a covariance matrix  $\boldsymbol{\Sigma}$  of the coordinates of the pixels. The bigger the control factor  $k$  is, the smoother the weight  $\beta_j$  is;  $d_j(\mathbf{h}_j, \mathbf{h}'_j)$  is a bin ratio similarity measure that improved the bin ratio dissimilarity [14], which is defined as

$$(7) \quad d_j(\mathbf{h}_j, \mathbf{h}'_j) = \frac{\|\mathbf{h} + \mathbf{h}'\|_2^2 \frac{\mathbf{h}_j \mathbf{h}'_j}{(\mathbf{h}_j + \mathbf{h}'_j)^2}}{N_e},$$

where  $N_e$  is the number of bins, and each pair of bins has at least one non-zero value in the two histograms  $\mathbf{h}_j, \mathbf{h}'_j$ , and  $\|\cdot\|_2$  is a  $L_2$  norm.

Given the reference object is represented as a set of  $B$  histograms  $\mathbf{q}=[\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_B]$ . Meanwhile, and the candidate object selected as the same size as the reference object is represented as a set of  $B$  histograms  $\mathbf{p}=[\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_B]$ , where  $\mathbf{q}_b$  and  $\mathbf{p}_b$  are the *spatiogram* of the  $b$ -th *fragment* respectively. In this paper, the fragments of the object are consisted of the fragments with three levels under different scales, and three levels of each scales are made up of the top level, the

middle level and the bottom, in which the top level is the image with equally divided into 16 rectangle fragments, and the second level is the image with equally divided into 4 rectangle fragments, and the bottom level is the original image. So all the fragments of each scale are concatenated as  $F = [F_1, F_b, \dots, F_{2^l}]$ ;  $\mathbf{q}$  is a reference histogram set corresponding fragments, and  $\mathbf{p}$  is a candidate histogram set. The similarity measure between  $\mathbf{q}$  and  $\mathbf{p}$  is expressed as

$$(8) \quad \rho(\mathbf{p}, \mathbf{q}) = \sum_{i=1}^B \lambda_b \rho_b \quad \text{s.t.}, \quad \sum_{b=1}^B \lambda_b = 1,$$

where  $\lambda_b$  represents the normalized weight corresponding to the  $b$ -th spatiogram,

which is defined as  $\lambda_b = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{d(\mathbf{p}_b, \mathbf{q}_b)}{2\sigma^2}}$  in which  $\sigma$  controls the normalized

weight of the  $\lambda_b$ . The similarity measure of the  $b$ -th spatiogram between  $\mathbf{q}$  and  $\mathbf{p}$ ,  $\rho_b$  is expressed as

$$(9) \quad \rho_b = \sum_{j=1}^M \beta_j^b d_j^b(\mathbf{h}_j^b, \mathbf{h}_j'^b).$$

### 3.3. Updating the model by incremental PCA

During the tracking, the tracking results with fixed object model are prone to producing a drift due to the fact that object or background appearance often changes. It is necessary to update online the object model in order to alleviate drift. PCA has been successful in face recognition and object tracking [7]. However, when the size of the object is very large, the complexity of these algorithms is very high. The main reason is that PCA computes the feature values and the feature vectors of training images. In this section, we introduce to update the object model by PCA which is similar to the method in [7].

Given  $\mathbf{ch}_1 = [\mathbf{hi}_1, \mathbf{hi}_2, \dots, \mathbf{hi}_{n_{in}}]$  and  $\mathbf{ch}_2 = [\mathbf{hi}_{n_{in}+1}, \mathbf{hi}_{n_{in}+2}, \dots, \mathbf{hi}_{n_{in}+m_o}]$ , where  $\mathbf{hi}_i$  is an observation histogram with  $d$  bins, and  $m_o, n_{in}$  is the number of the training histograms  $\mathbf{ch}_1$  and  $\mathbf{ch}_2$ , respectively;  $\overline{\mathbf{ch}_1}$ ,  $\Sigma_1$  is the mean vector and the scatter matrix  $S_{\mathbf{ch}_1}$  of  $\mathbf{ch}_1$  respectively, and  $\overline{\mathbf{ch}_2}$ ,  $\Sigma_2$  is the mean vector and the Scatter matrix  $S_{\mathbf{ch}_2}$  of  $\mathbf{ch}_2$  respectively;  $\mathbf{ch} = [\mathbf{ch}_1 \ \mathbf{ch}_2]$  is their concatenation,  $S_{\mathbf{CH}}$  is the scatter matrix of CH, which is expressed as

$$(10) \quad S_{\mathbf{CH}} = S_{\mathbf{ch}_1} + S_{\mathbf{ch}_2} + \frac{n_{in}m_o}{n_{in} + m_o} (\overline{\mathbf{ch}_2} - \overline{\mathbf{ch}_1})(\overline{\mathbf{ch}_2} - \overline{\mathbf{ch}_1})^T.$$

Its detailed proof can refer to [7]. After the decomposition of the incremental PCA, we select the corresponding vectors according to the first five big eigen values, and then we make them assemble the feature set. The reference histogram is approximated as

$$(11) \quad \mathbf{H}_b \approx \sum_{k=1}^{n_u} \mathbf{V}(:, k, b)^T \text{cih}(k, b) \mathbf{V}(:, k, b) + \overline{\mathbf{h}_b},$$

where  $\mathbf{V}(:, k, b)$  in (11) is the  $k$ -th eigen vector of the  $b$ -th fragment in the scatter matrix  $S_{\mathbf{ch}}$ ;  $\overline{\mathbf{h}_b}$  is the mean histogram of the  $b$ -th fragment;  $\text{cih}(k, b)$  is the  $k$ -th eigen-

values of the  $b$ -th fragment of the object;  $n_u$  is the number of the selected feature histogram (In this paper, we set  $n_u = 5$ . In that we consider the time and space complexity and the performance of the tracking method. If  $n_u > 5$ , the time and space complexity is high, hence it affects the speed of the tracking method, on the contrary, the performance of the tracking method is affected). To further adaptive the change of the object, we also consider the initial histogram, therefore, the new object model is expressed as

$$(12) \quad \mathbf{H}_{0,b} = \gamma \mathbf{H}_{\text{init},b} + (1 - \gamma) \mathbf{H}_{1,b},$$

where  $\gamma$  is a learning rate, which is computed by the similarity between the current and the initial object. (We set from 0.6 up to 0.8 due to the fact that the appearance changes are prone to neglecting if is bigger than 0.8, on the contrary, the appearance changes are excessively considered if is smaller than 0.6.)

#### 4. Object tracking based on online semi-SVM and adaptive fused feature

Self-training has been a widely used method in semi-supervised learning, in which a classifier is firstly trained with a small amount of labelled samples, and then the unlabelled samples are classified by the classifier. If the confident values of these unlabelled samples are high, we need to add these unlabelled samples into the training set as the label samples and update the classifier by retraining these label samples. Our method combines the incremental SVM with self-learning to track the object. To improve the drifting problem resulted from the self-learning, we need to update the object mode on line.

The larger the number the more it leads to degradation in tracking results, while the smaller also lead to unreliable results since the number of the labelled samples is important for training. In semi-supervised learning, only a few samples are labelled. In this paper, we find that our method can perform well with only 15 labelled positive samples and 60 negative labelled samples. The positive samples are a histogram set including some feature histograms of the objects and the negative sample is also a histogram set of including some feature histograms of the areas which do not overlap with the objects. After the initialization, we will locate the object and update the classifier.

##### 4.1. Locating the object and updating the classifier

Given a self-learning classifier and an unlabelled frame, we need to compute the location of the object within the frame. Specifically, the location of the unlabelled frame is classified by the classifier using the sliding window technology, and then the confidence map is generated. One common method of searching the maximum value of the confidence map is a gradient ascent algorithm which is known at the prior location of the object. However, this method suffers from the spatial constraint. To simplify the procedure, we only compute the global maximum value in the confidence map.

To adapt to the changes of the object or background, it is vital to update the classifier and the object model on line. The main problem is how to select online samples comprising of positive samples and negative samples, that is, we determine whether or not to add the current sample to training set and update the classifier. In this paper, we determine whether or not to add the current sample to the feature set according to its value of the global maximum in the confidence map, if the value is higher than the given threshold (threshold is selected manually and its value is 0.3 because of considering effective and accuracy of the classifier), the current sample is added and then the classifier is updated. The negative samples are selected according to whether or not adding positive sample. When the current object is added, we select the  $k$  highest peaks which do not overlap with the object in the confidence map in order to select the most important negative samples. However, the number of the samples can become very large in time, moreover, they make our method of studying the changes of the object or the background work much slower. To improve the tracking velocity, we maintain the number of training samples  $N_i$  (we set  $N_i = 15$  in our experiments because of considering effective and accuracy of the classifier). If the number of the trained samples is larger than  $N_i$ , the old sample will be removed from the SVM by decreasing method.

#### 4.2. Pseudo-code of our method

In this section, the detailed procedure of the presented algorithm will be given as the following:

*Initialization:*

For  $i=1:n_{in}$

**Step 1.** Acquire the  $i$ -th labelled object either manually or by other object tracking method

**Step 2.** Compute the histogram of the intensity and the HOG

**Step 3.** Fuse feature histogram through the Equation (5)

**Step 4.** Acquire a positive sample and desired negative samples

End

Train the classifier using incremental/decremental SVM

$L=1$ ;

#### **Online self-tracking with adaptive update model**

For  $i=n_{in}+1: n_{in}+N_{fr}$

**Step 1.** Build a confidence map making use of the trained classifier

**Step 2.** Locate the object employing the maximum value of the confidence map

**Step 3.** If maximum value is bigger than threshold1 and similarity between the current object

(a) the current object is added to the positive sample set

(b) the images corresponding to the first four minimum value of the confidence map are added to the negative sample set

$L=L+1$

(c) update the SSSVM classifier

(d) the current object is added to the feature set



```

    Else
        Track the object using the basic particle filter
    End
Step 4. If  $L > 20$ 
    Remove old samples from the SVM by the decreasing SVM
    Update the object model by (12) and retrain the classifier.
     $L = 1$ ;
    End
End

```

## 5. Experimental results and discussion

To evaluate the performance of the proposed method, we perform our method in some visible and thermal video sequences. The visible sequences are selected from two vehicle sequences and the CAVIAR database<sup>1</sup>, and the thermal sequences are selected from OTCBVS Benchmark Dataset<sup>2</sup>. In our experiments, two selected features are intensity histogram and HOGE, and then it fuses them by MSSBRS to represent the tracked object and local background. Of course, the other feature, such as optical flow can be easily selected. The selected features are intensity and HOGE for thermal sequences, and two features from  $r$ ,  $g$ ,  $b$ , HOGE are selected according to the method in [8] for visible sequences. The scale which may produce the maximum score was selected when scale changes are bound to  $\pm 10\%$  in each frame.

All the experiments are carried out with a core 4 Duo 2.9 GHz processor with 2 GB RAM under Matlab R2010b. To demonstrate the performance of our method, we compare our tracker with the ensemble tracking and the basic mean-shift tracking in some invisible and thermal sequences.

### 5.1. Results in visible sequences

The first experiment can be seen on a video sequence of a white car crossing the crossroads under heavy fog and heavy wind. The sequence is 150 frames long and the size of each image is  $576 \times 768$  pixels. The sequence is challenging because the car can be easily affected by the heavy fog. Obviously, it is difficult to track the car accurately. In spite of this, our method is still able to track the car throughout the entire sequence. The ensemble tracker [5] is able to track the car for the first 100 frames, however the following frame is clear drift. The main reason is that the ensemble tracking method suffers from the fog, where our method can solve this problem. Fig. 2 shows the results of the ensemble tracker and our method, which contain the tracking results and the corresponding confidence map, and it can see that our confidence maps have a clear peak at the object center (Fig. 3), so it leads to more stable tracking.

---

<sup>1</sup>[http://i21www.ira.uka.de/image\\_sequences/](http://i21www.ira.uka.de/image_sequences/)

<http://groups.inf.ed.ac.uk/vision/CAVIAR/CAVIARDATA1/>

<sup>2</sup> <http://www.cse.ohio-state.edu/OTCBVS-BENCH/bench.html>

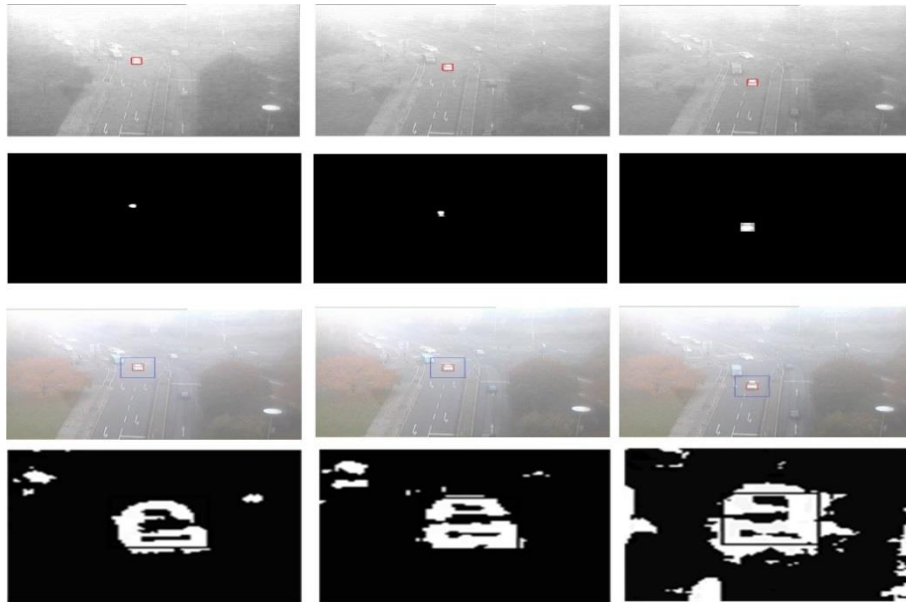


Fig. 2. Comparison of our method and the ensemble tracker method. The first two rows are the tracking results and the corresponding confidence maps of the proposed method, and later two rows are the tracking results and the corresponding confidence maps of the ensemble tracker

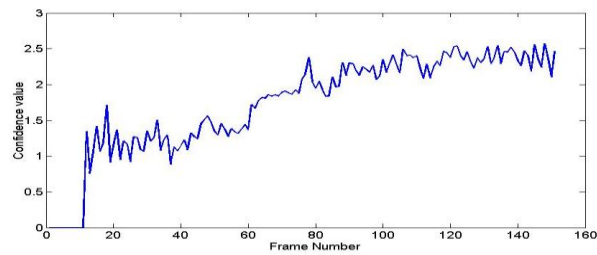


Fig. 3. Confidence value of the car under the heavy fog

## 5.2. Results in thermal sequences

Our method can track not only the object in visible sequences but also the object in thermal sequences. Here we show results of two experiments<sup>3</sup>; one is the boat in the sea and the other is the pedestrian across the big tree. The only modification is the feature space that is used to represent the object because the thermal object has few efficient cues, so the efficient features are selected as the intensity and the HOG.

In the first thermal sequence, a boat at sea, which is over 800 frames is tracked, because the background of the sequence suffers from the clutter and it is represented using the mixture histogram which is fused between the intensity histogram with 16 bins and the HOG with 16 bins. The tracking results of our method and the basic mean shift tracking method -are shown in Fig. 4. We observe that our method is able to track robustly the boat because we adopt the adaptive

<sup>3</sup> <http://www.cse.ohio-state.edu/OTCBVS-BENCH/bench.html> and <http://www.cs.technion.ac.il/~idol/>

object model, while the basic mean shift tracking method [16] is slight drift due to the fact that the thermal object does not provide enough information for tracking. Meanwhile, we also observe that the confidence map has a clear peak and helps the tracking.

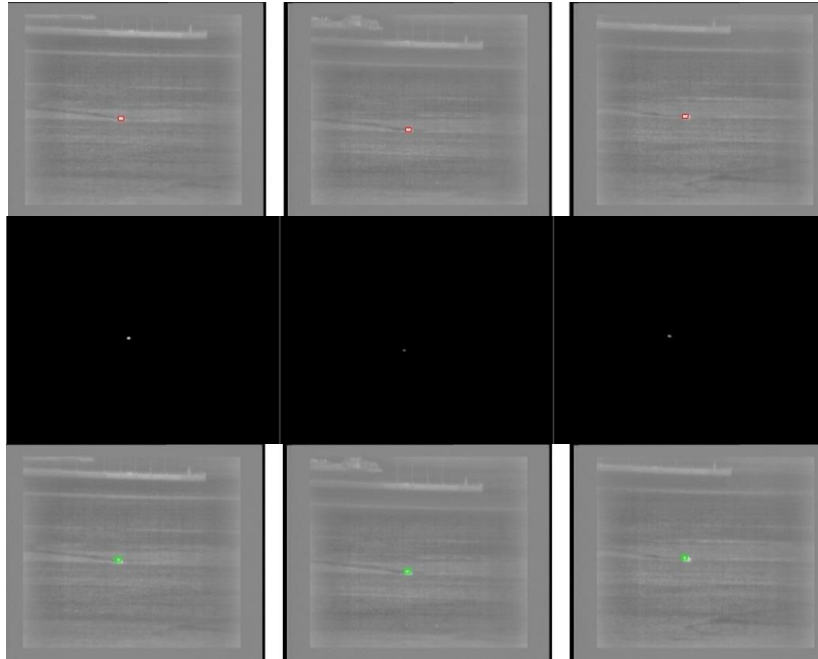


Fig. 4. Results of our algorithm and the basic mean shift algorithm on a thermal video of a boat. The top row is the results of our method, the middle row is the confidence map of our method, and the bottom row is the results of the basic mean shift algorithm. The frame number is 33,283, and 734 from the left to right respectively. Note that the confidence maps are smoother and have a clear peak at the object centre

Fig. 5a shows the results of our method compared with the basic mean shift method on the boat sequence. As can be seen from the Fig. 5a, our method can track more stable and more consistent than the basic mean shift method. Furthermore, it is close to the ground truth. Fig 5b shows the confidence value of the tracked object in the first 300 frames.

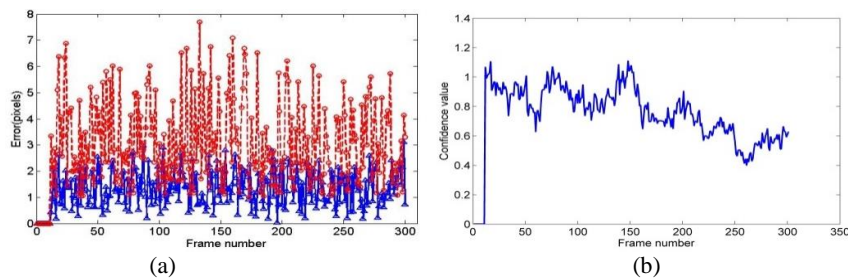


Fig. 5. The tracking error about our method and the basic mean shift method of the first 300 frames in the boat sequence, Note that solid blue is the tracking error of our method and dashed red is the tracking error of the basic mean shift algorithm (a); the confidence value of the first 300 frames (b)

So far, we only take into account the partial occlusions of the object. With respect to the complete occlusions, we can track the object using the basic particle filter tracker [17]. Specifically, as long as the confidence value is higher than the given threshold, the tracking position preserve unchanged until the next value is lower than the threshold, and then we begin with tracking the object via the basic particle filter tracker. In our experiments, we find that the object is occluded when the confidence value is less than 0.5, namely, we may set the threshold to 0.5. While tracking object, the confidence value is firstly computed by the trained classifier and then judge if it is above the threshold (that is to say, the object reappears if the threshold value is high, as is shown Fig. 6c. On the contrary ,the basic particle filter tracker will stop performing it if the confidence value of the current frame is above the threshold ,and then our method begin to track the object from the next frame . When the target is occluded, our method can accurately track the object than other tracking method, as is shown in Fig. 6. Fig. 7 shows the confidence value of all complete 151 frames, and finds that the confidence value between the frame 138 and the frame 143 is very low.

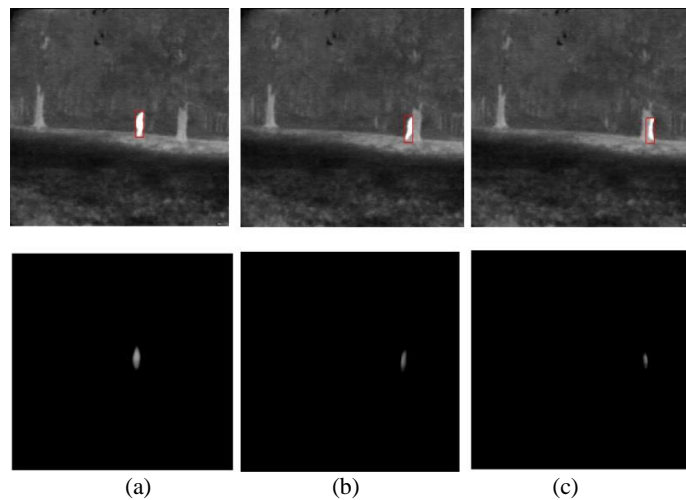


Fig. 6. Tracking results and the confidence map corresponding the frame number 93,136, and 144 of our tracker. Note that the confidence maps are smoother and have a clear peak at the object centre

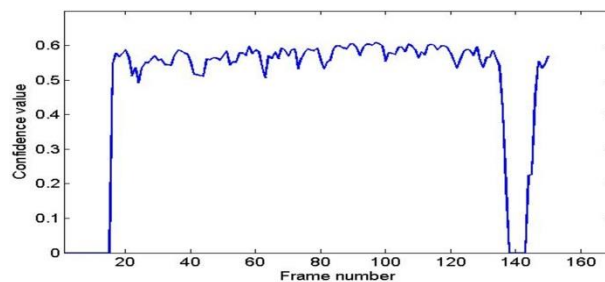


Fig. 7. The confidence value of all complete 151 frames. Note that the confidence value between the frame 138 and the frame 143 is very low

Fig. 8 shows the tracking results of additional sequences applying proposed our method.

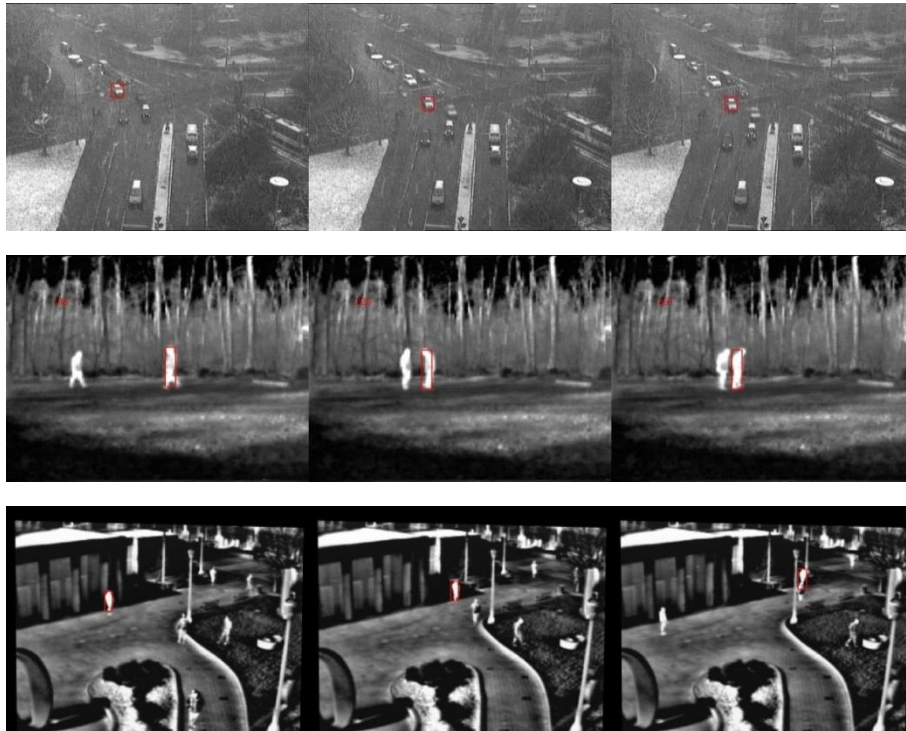


Fig. 8. Tracking results of our method under the visible sequence and the infrared sequences. The first row shows the results of our method on frames 114,211 and 226 in car under the heavy snow. The second row shows the results of our method on frames 143,225, and 391 from the OTCBVS sequence. The third row shows the results of our method on frames 22,124, and 259 from the OTCBVS sequence. Note that our method can robustly track the object

## 6. Conclusion

We treat object tracking as a binary classification and propose a robust tracking algorithm based on online semi-SVM with the self-training framework and adaptive fused features whose fused coefficients are computed by new similarity MSSBRS. The object model is updated using the incremental PCA. The proposed method only labels a small amount data to train the classifier during the initialization, and then the tracked object is classified according to the confidence value, the classifier is updated online by the self-training framework and new samples with high confident value. With regards to the complete occlusion, we use the basic particle filter algorithm to solve it. However, the object is represented by a fixed object scale, the proposed method is not able to deal with large changes in object scale. In future, we would like to make the proposed method adapt to these changes.

**Acknowledgments:** This work was supported by the Research Foundation of Changzhou Institute of Technology (Grant No YN1208), the Natural Science Foundation of China (Grant No 61475027), Collaborative innovation fund of Jiangsu province (XYN1408) and the Research Foundation of Changzhou Institute of Modern Optical Technology (Grant No CZGY005).

## References

1. Lu, Z., L. J. P. Van Der Maaten. Preserving Structure in Model-Free Tracking. – IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. **36**, 2014, No 4, pp. 756-769.
2. Wu, Y., M. Pei, M. Yang, J. Yuan, Y. Jia. Robust Discriminative Tracking via Landmark-Based Label Propagation. – IEEE Transactions on Image Processing, Vol. **24**, 2015, No 5, pp.1510-1523.
3. Tang, F., S. Brennan, Q. Zhao, H. Tao et al. Co-Tracking Using Semi-Supervised Support Sector Machines. – In: Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2007, pp. 1-8.
4. Avidan, S. Support Vector Tracking. – IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. **26**, 2004, No 8, pp. 1064-1072.
5. Avidan, S. Ensemble Tracking. – IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. **29**, 2007, No 2, pp. 261-271.
6. Grabner, C. L. A. H. H. Semi-Supervised Online Boosting for Robust Tracking. – In: Proc. of European Conference on Computer Vision, 2008, pp. 234-247.
7. Ross, D. A., J. Lim, R. S. Lin, M. H. Yang. Incremental Learning for Robust Visual Tracking. – International Journal of Computer Vision, Vol. **77**, 2008, pp.125-141.
8. Wang, J. Q., Y. Yagi. Integrating Color and Shape-Texture Features for Adaptive Real-Time Object Tracking. – IEEE Transactions on Image Processing, Vol. **17**, 2008, No 2, pp. 235-240.
9. Cauwenberghs, G., T. Poggio. Incremental and Decremental Support Vector Machine Learning. – In: Proc. of Neural Information Processing Systems, 2000, pp. 409-415.
10. Zhu, X. Semi-Supervised Learning Literature Survey. Computer Sciences Technical Report, 2006.
11. He, K., J. Sun, X. Tang. Guided Image Filter. – IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. **35**, 2013, No 6, pp. 397-1409.
12. Geusebroek, J. M., R. Van Den Boomgaard, A. W. M. Smeulders et.al. Color Invariance. – IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. **23**, 2001, No 12, pp. 1338-1350.
13. Xie, N., H. Ling, W. Hu et. al. Use Bin-Ratio Information for Category and Scene Classification. – In: Proc. of IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2010, pp. 2313-2319.
14. Xiang, R., J. Li, X. Wang, F. Youjia. Probabilistic Tracking Method Based on the Spatial Bin-Ratio Information. – International Journal of Optomechatronics, Vol. **5**, 2011, No 4, pp. 360-377.
15. Birchfield, T., S. Rangarajan. Spatial Histograms for Region-Based Tracking. – ETRI Journal, Vol. **29**, 2007, pp. 697-699.
16. Comaniciu, D., V. Ramesh, P. Meer. Kernel-Based Object Tracking. – IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. **25**, 2003, No 5, pp. 564-577.
17. Gordon, N. J., D. J. Salmond, A. Smith. Novel-Approach to Nonlinear Non-Gaussian Bayesian State Estimation. – In: IEE Proc. of Radar and Signal Processing. Vol. **140**, 1993, pp. 107-113.