

# Learning Scene-specific Object Detectors Based on a Generative-Discriminative Model with Minimal Supervision

Dapeng Luo<sup>a</sup>, Zhipeng Zeng<sup>a</sup>, Nong Sang<sup>b</sup>, Xiang Wu<sup>d</sup>, Longsheng Wei<sup>c,\*</sup>,  
Quanzheng Mou<sup>a</sup>, Jun Cheng<sup>c</sup>, Chen Luo<sup>e</sup>

<sup>a</sup>*School of Mechanical Engineering and Electronic Information, China University of Geosciences, Wuhan 430074 China*

<sup>b</sup>*School of Automation, Huazhong University of Science and Technology, Wuhan 430074 China*

<sup>c</sup>*School of Automation, China University of Geosciences, Wuhan 430074 China*

<sup>d</sup>*School of Automation, Nanjing University of Science and Technology, Nanjing 210094 China*

<sup>e</sup>*Huizhou School Affiliated to Beijing Normal University, Huizhou 516002, China*

---

## Abstract

One object class may show large variations due to diverse illuminations, backgrounds and camera viewpoints. Traditional object detection methods often perform worse under unconstrained video environments. To address this problem, many modern approaches model deep hierarchical appearance representations for object detection. Most of these methods require a time-consuming training process on large manual labelling sample set. In this paper, the proposed framework takes a remarkably different direction to resolve the multi-scene detection problem in a bottom-up fashion. First, a scene-specific objector is obtained from a fully autonomous learning process triggered by marking several bounding boxes around the object in the first video frame via a mouse. Here the human labeled training data or a generic detector are not needed. Second, this learning process is conveniently replicated many times in different surveillance scenes and results in particular detectors under various camera viewpoints. Thus, the proposed framework can be employed in multi-scene object detection applications with minimal supervision. Obviously, the initial scene-specific detector, initialized by sev-

---

\*Corresponding author

eral bounding boxes, exhibits poor detection performance and is difficult to improve with traditional online learning algorithm. Consequently, we propose Generative-Discriminative model to partition detection response space and assign each partition an individual descriptor that progressively achieves high classification accuracy. A novel online gradual optimized process is proposed to optimize the Generative-Discriminative model and focus on the hard samples: the most informative samples lying around the decision boundary. The output is a hybrid classifier based scene-specific detector which achieves decent performance under different scenes. Experimental results on six video datasets show our approach achieves comparable performance to robust supervised methods, and outperforms the state of the art self-learning methods under varying imaging conditions.

*Keywords:* Object detection, Unsupervised learning, Generative-Discriminative model

---

## 1. Introduction

With the development of intelligent surveillance systems, pedestrian and vehicle detection approaches have garnered profound interest from engineers and scholars. Many impressive works [1, 2, 3, 4, 5, 6, 7] have been published in the last several years. However, object detection and recognition remain a considerably difficult issue in highly populated public places such as airports, train stations and urban arterial roads, which have distributed multiple surveillance cameras with different viewpoints, illuminations and backgrounds in cluttered environments. One object class may show large inter-class variations under these different imaging conditions. A substantial amount of training data need to be collected and labeled to model one object category detector based on statistical learning. Without this, the detector, trained in constrained video environments, will deliver poor performance under different environmental conditions. How to robustly and stably locate objects in arbitrary video environments with minimal supervision is still an open issue.

Transfer learning methods [8, 9] can be used to learn different view-point scene detectors from pre-training models, which reduce the efforts involved in collecting samples and retraining in response to such variations in appearance. However, negative transfer often occurs when transferring very different scenes [10], significantly influencing the performance of target scene detectors

and limiting the application of the transfer learning.

The Multi-view object detection method is an alternative strategy to locate the object in various view-points, which minimizes the influence of changing imaging conditions. However, these approaches increase the discriminability of detectors through relatively complex additional stages involving view-invariant features [11, 12, 13, 14], the pose estimator [15, 16, 17, 18] or the 3D object model [19, 20, 21, 22, 23, 24], which make them computationally expensive and require an intensive training process on large datasets. These top-down approaches will result in considerable runtime complexity.

Let's focus on one surveillance camera used in one specific scenario. A common strategy is to train specific object detector from human collected and labelled training samples the scenario since each individual has limited variant poses in one scene. However, it is impossible to train scene-specific detectors for every scenario considering tedious human efforts and time costs, unless the training process requires minimal or even no supervision.

A method must be found to learn a scene-specific object detector without human intervention, extending to other scenes, making sure that every scene has its own detector and achieves satisfactory detection performance under different imaging conditions and viewpoints. Although the idea sounds attractive, this task is challenging because constructing an object model without prior knowledge is difficult, and there is no effective algorithm to automatically collect and label samples for training the detector on the fly.

Some studies on online-learning object detectors have been proposed and most of them adopt a similar framework, including an online learning detector and a validation strategy. However, these methods are not completely unsupervised learning methods, but only minimize manual effort, and are always initialized by several hundred human labeled training samples for one specific scene. Moreover, the number of manually labeled initial instances will increase proportionally considering the multi-scene object detection application. In addition, these approaches employ co-training methods [25, 26], the background subtraction [27], the generative model [28, 29] and tracking-by-detection approaches [30, 31] as the validation strategies to collect and label the online learning samples automatically. The performance of online-learning approaches are far from competitive with supervised methods because of the high label error in hard samples distributed around the decision hyperplane, compared to the human label in supervised training process.

Recently, some scene-specific object detection approaches [32, 33, 34] have been proposed to discover and label hard proposals for training an initial

generic detector without any manual annotation. Among these various solutions, Ye [32] has shown impressive results under challenging situations by optimizing a progressive latent model (PLM) with the difference of convex (DC) objective functions. However, it is a computational challenge to optimize a DC function straightforwardly throughout the sample space with an unsupervised manner. Furthermore, a careful initial generic detector is necessary to avoid local minima.

To overcome the limitations of existing works, in this paper, a Generative-Discriminative model (GDM) is proposed to partition the sample space into three disjoint subspaces: positive, negative and hard sample space. A Generative model learns the joint probability of positive and negative samples, while hard samples are classified by a discriminative model with solving a mixed integer programming problem. Similar to the concave-convex programming procedure, the GDM can be iteratively updated by an online gradual optimized process which is insensitive to initialization.

Our goal is to design an self-learning object detection framework, which can train object class detector in each particular scenario without human intervention. Instead of manually labeling several hundred initial training samples or employing a generic detector in existing methods, our scene-specific detector is obtained by simply marking several bounding boxes around the object in the first video frame via a mouse. This approach reduces human annotation effort to an effortless mouse operation within the first frame, which "determine" the interested object category in current surveillance video. Other than this, human labeled samples or general object detectors are not needed.

There are two processes in our framework:

In learning process, first, an initial sample set generated by affine transformation of these marked objects in the first frame. Second, a Generative-Discriminative model, trained by the initial sample set, runs as an initial detector on subsequence frames. Obviously, the initial detector has poor detection performance due to the incomplete initial training sample set. Third, for this case, we propose an online gradual optimized procedure to iterative optimize the Generative-Discriminative model (GDM) in an unsupervised manner. When the convergence condition is satisfied, the optimized process will stop and result in a hybrid classifier composed of a generative model and a discriminative model.

During object detection process, the two models work together to determine the location of real objects. The generative model is first used to detect

objects based on the sliding window strategy, and the detection responses, located near the classification boundary, will be further recognized by the discriminative model.

Our method, triggered by several bounding boxes, is minima supervised without human effort on collecting and labelling samples, and can be easily extended to other scenes, forming scene-specific objectors in various surveillance cameras with different viewing distances and angles. Thus, this is a bottom-up method to resolve the multi-scene object detection problem.

Moreover, both the use of Generative-Discriminative model (GDM) and online gradual optimized process make our method robust in tackling the most informative samples lying at the decision boundary, and the method achieves state of the art detection performance under different viewing angles, as shown in Fig. 1. By self-learning within three hours, our method produces satisfying results in three different viewpoints sequences of CAVIAR [31] and PETS2009 [35] datasets. In addition, we obtain competitive detection performance in the S2 sequences of PETS2009 dataset, compared to the Aggregated Channel Features (ACF) [6], one of the most robust supervised object detection approaches, trained by manually labeling 300 positive samples and 900 negative samples in the same viewpoint. To the best of our knowledge, this is the first time to demonstrate a scene-specific detector without additional manual annotated samples or generic detectors, compared to other object detection methods.

The main contributions of this paper are:

- 1) We present the Generative-Discriminative model (GDM), which partitions the sample space for improving the discriminability of scene-specific object classifiers while being efficient at run-time.

- 2) We present the online gradual optimized process which allows the detector, initialized by marking several bounding boxes around the object with poor detection performance, to successively improve classification accuracy and becoming more dedicated to challenging samples lying near the decision boundary.

- 3) We present a novel self-learning framework to train scene-specific object detectors by marking several bounding boxes around the object in the first frame, which can easily extend to various surveillance videos, and automatically achieve successful detection performance.

The rest of this paper is arranged as follows: Section 2 briefly recalls related works. In Section 3, the analysis of our approach is provided. Generative model is described in section 4, discriminative model is explained in

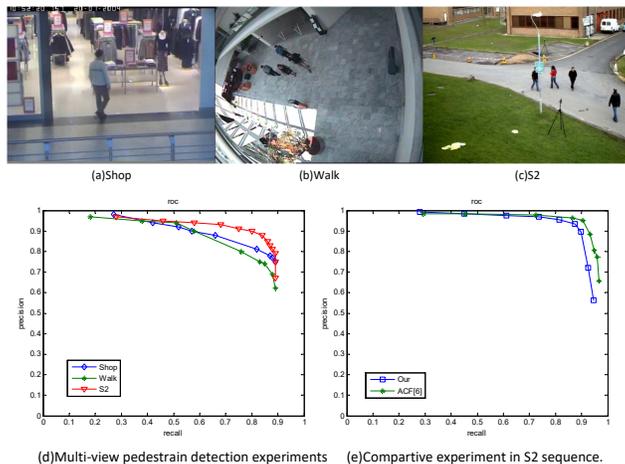


Figure 1: (a) Shop sequence from CAVIAR dataset; (b) Walk sequence from CAVIAR dataset; (c) S2 sequence from PETS2009 dataset (d) Experiments in multi-view video sequences; (e) Comparison with ACF[6].

section 5, and section 6 shows the online gradual optimized process. Experiments and results are presented in section 7, which is followed by the conclusions.

## 2. Literature Review

Object class detection is the core component in most computer vision tasks, and it has achieved prominent success for supervised learning based approaches, such as [1, 2, 3, 4, 5, 6, 7]. Here, we focus on the scene-specific object detection, online learning frameworks and multi-view object detection methods.

### 2.1. Scene-specific Object Detection

This paper proposes to address multi-scene object detection problem by bottom-up scene-specific object detection approaches. In essence, [32, 36, 37] have been widely used to train the scene-specific detector based on unsupervised object discovery techniques. These self-learning methods usually optimize the convex or semi-convex objective functions throughout the sample space. [32] even optimizes the difference of convex objective functions to improve the detection precise with a progressive latent model (PLM).

Despite the success of unsupervised object discovery techniques to improve the performance of the self-learning object detector, they are limited in two aspects. First, this is a mixed integer programming problem which poses a computational challenge when optimizing loss functions directly in the entire feature space with an unsupervised manner. Second, an effective initial generic detector is necessary to avoid falling into the local minima in optimized process.

## 2.2. Online learning Framework

Online learning frameworks have been published to adapt the detector to a particular scenario. However, it is difficult to use these approaches in the multi-scene object detection domain due to two major reasons:

First, traditional online learning object detection methods are inseparable from manual efforts. For instance, most of the online learning detectors are initialed by human labelled training samples. The detection performances are usually degraded when the number of manual labelled instances are minimized. Some works [29, 36] are proposed to specify a well-trained generic object detector to a specific scene. However, these works are particularly valuable in constrained applicant conditions. Recently, semi-supervised learning [38, 39], transfer learning [40, 9], and weak-supervised learning [41, 37, 42, 43], are employed to reduce the amount of labeled training data for object detector. How to minimize human effort in online object detection system is still a hot research topic.

Second, the new samples, which are used to train the detector online, need to be automatically collected and labeled. How to correctly label the new training samples is still a challenging topic. To date, various automatic annotation methods have been reported and can be broadly divided into four categories: 1) co-training based, 2) background subtraction based, 3) generative model based, and 4) tracking based.

In co-training based approaches [25, 26], two classifiers are trained simultaneously and labeled for each other. In background subtracted approaches [27], the foreground detector, based on a background model, is employed as an automatic labeler. Generative model based methods [28, 29] use the reconstructed model error to validate the detection responses and train the detector by a feedback process. Tracking based approaches [30, 31] collect and label online training samples by using tracking-by-detection methods which can interpolate missed object instances and false alarms as positive and negative samples, respectively.

However, the aforementioned methods have no special strategy to deal with the problematic samples located around decision boundary, the most informative and ambiguous part of the object categories. Our method employs online gradual optimized process and the Generative-Discriminative model to hierarchically process the detection responses, which makes the learning procedure focus on the hard samples and reduce the online labeling error.

### *2.3. Multi-view Object Methods*

Conventional object detection methods locating objects by considering single view cannot be employed in multi-view object detection since objects have wide variations in their poses, colors and shapes under multi-view imaging conditions. Thus, a common idea is to model object classes by collecting distinct views forming a bank of viewpoint-dependent detectors which can be used to predict the optimal viewing angle for detection.

In early stages, most multi-view detection approaches independently apply several single-view detectors and then combine their responses via arbitrary logic. Some impressive works [44, 45, 46, 47] have been reported in the domain of face detection, dealing with multiple viewpoints (frontal, semi-frontal and profile).

Following this progress, Thomas et.al. [22] no longer rely on single-view detectors working independently, but develop a single integrated multi-view detector that accumulates evidence from different training views. Several other approaches, such as [23, 24], build complex 3D part models containing connections between neighboring parts and overlapping viewpoints which achieve remarkable results in predicting a discrete set of object poses. The discrete views are usually treated independently, however, [48, 49, 50, 51] require evaluating a large number of view based detectors, resulting in considerable runtime complexity.

Recently, Pepik et.al. [52] proposed 3D deformable part models which extend the deformable part model to include viewpoint information and part-level 3D geometry information. This method establishes represented 3D object parts and synthesizes appearance models for viewpoints of arbitrary granularity on the fly, resulting in significant speed-up. Xu et.al. [53] proposed to accomplish multi-view learning with incomplete views by exploiting the connections between multiple views, enabling the incomplete views to be restored with the help of the complete views.

With respect to all the aforementioned approaches, our framework takes a markedly different direction: we establish scene-specific detectors in an unsupervised manner for each scenario, so as to prevent from integrating complex models and labeling substantial training samples in different viewpoints. To resolve the multi-scene object detection problem in bottom-up fashion is an improvement over other systems.

Other than the above mentioned approaches, there exist methods [54, 55] to detect unknown object classes from the motion segmentation. Although these methods learn foreground models in arbitrary scenarios without any a priori assumption, they are very different from our method, which address one object category detection problem: these methods can not recognize the detected responses because they use the cluster based global optimization procedure.

Our method initializes a scene-specific detector with several bounding boxes in the first frame, and eventually realizes a state of the art multiple object detection system without human label effort. Online gradual optimized strategy, which is proposed in this paper, is the key point which ensures our framework can be improved from a poor detection performance initial detector (see detail in section 6). Otherwise, the output of our framework is a Generative-Discriminative model which is very different from the other scene-specific object detectors. In this way, our method opens up the possibility for several different classifiers to work together to determine the locations of one object class in videos from self-learning.

### 3. Analysis of Our Method

Object detection can be viewed as a two-class classification problem [56]. However, in most applications, the sample space can be divided into three groups: positive, hard and negative samples, by two decision boundaries. Both hard positive and hard negative examples have a significant effect on enhancing the classification performance.

In this section, a Generative-Discriminative model has been proposed to describe the three sample spaces from one surveillance domain. A novel cost function will be employed to improve the model performance and speed up the convergence rate. After that, an example of the Generative-Discriminative model will be introduced to learn scene-specific objector with minimal supervision.

### 3.1. Generative-Discriminative Model (GDM)

A generative classifier  $G$ , with the inputs  $x$  and the label  $y$ , concentrates on capturing the generation process of  $x$  by modelling, which is robust to partial occlusion, viewpoint changes and significant intra-class variation of object appearances. Thus, the model is suitable to describe the positive sample space  $P_{c_+}$  and negative sample space  $P_{c_-}$ .

A discriminative model  $D$ , on the other hand, learns the difference between different categories. Thus,  $D$  can be employed to model the hard sample space  $P_{c_h}$ , so as to find the optimal classification of samples located between the positive and negative decision boundary:  $B_+$ ,  $B_-$ .

In other words, samples falling in  $P_{c_+}$  and  $P_{c_-}$  belong to labeled samples and hard samples are unlabeled samples. Thus, the cost function is as follows:

$$L(x) = \Upsilon \sum_{g(x) > B_+, g(x) < B_-} C_G(x, y) + \alpha \sum_{B_- \leq g(x) \leq B_+} C_D(x) \quad (1)$$

where  $C_G(x, y)$  and  $C_D(x)$  are the objective functions for labelled samples and unlabelled samples respectively.  $\Upsilon$  and  $\alpha$  are regularization factors.  $g(\cdot)$  is the classification function calculated by Generative model.

In addition, the distance  $Dis(B_+, B_-)$  between the positive and negative decision boundaries has a huge impact on the model. The smaller the distance, the more accurate the generative model is to describe the positive and negative samples. Furthermore, a smaller distance also means fewer hard samples, which provides a convergence condition in optimizing a Generative-Discriminative model. Thus, the cost function can be formulated as follows:

$$L(x) = \Upsilon \sum_{x \in P_{c_+}, P_{c_-}} C_G(x, y) + \alpha \sum_{x \in P_{c_h}} C_D(x) + \lambda \cdot Dis(B_+, B_-) \quad (2)$$

where  $\lambda$  is the regularization fact of the distance item  $Dis(\cdot)$ .

The cost function  $L(x)$  has three latent variables: the pseudo labels of unlabelled samples, the positive boundary  $B_+$  and the negative boundary  $B_-$ , which cannot be solved efficiently with the traditional tools of optimization. In this paper, we propose the online gradual optimized process to solve the objective function in a dynamically changing sample space (see detail in section 6).

Many generative and discriminative models can be employed in the GDM, such as Autoencoder networks, conditional random fields, hidden Markov model and Bayesian networks. In our case, we propose an Online Selector Fern(OSF) generative model and an iterative SVM(ISVM) discriminative model, considering real-time requirements in video object detection applications. This results in a hybrid model which is both highly effective and computationally efficient (running at over 60 fps for  $768 \times 576$  images).

### 3.2. Self-learning Framework

Based on the Generative-Discriminative model, a self-learning framework is proposed to train a scene-specific object detector. The initial training sample set is prepared by affine transformation of the several selected patches in the first surveillance video frame. Thus, the human labeled training data or a generic detector are not needed. An OSF generative model is trained by the initial training data as the initial detector. Obviously, the detector has poor detection performance due to the incomplete training sample set. As shown in Fig. 2, the detector will be improved by a propagation optimal process.

Thus, we propose the online gradual optimized process to train a Generative-Discriminative model in fully-autonomous fashion:

First, we define positive and negative decision boundaries in the OSF classifier. The detected responses are then collected as positive, negative and hard samples, respectively. To ensure the high correct label rate, the OSF classifier has large initial margin setting between the two boundaries.

Second, an ISVM discriminative model is proposed to solve a mixed integer programming problem and label the hard samples in the dynamically changing subspace.

Third, the labelled hard samples are used to online train the OSF classifier and gradually minimize the margin between positive and negative boundaries, improving the ability to express the positive and negative sample spaces.

This process is repeated till convergence. The output is a Generative-Discriminative model, composed of an OSF classifier and an ISVM model.

In detection process, the scanning window strategy is employed. Most image patches are classified by the OSF classifier and only a small fraction located between the dual boundaries, collected as hard samples, are classified by the ISVM model, which makes our approach robust while being efficient at run-time.

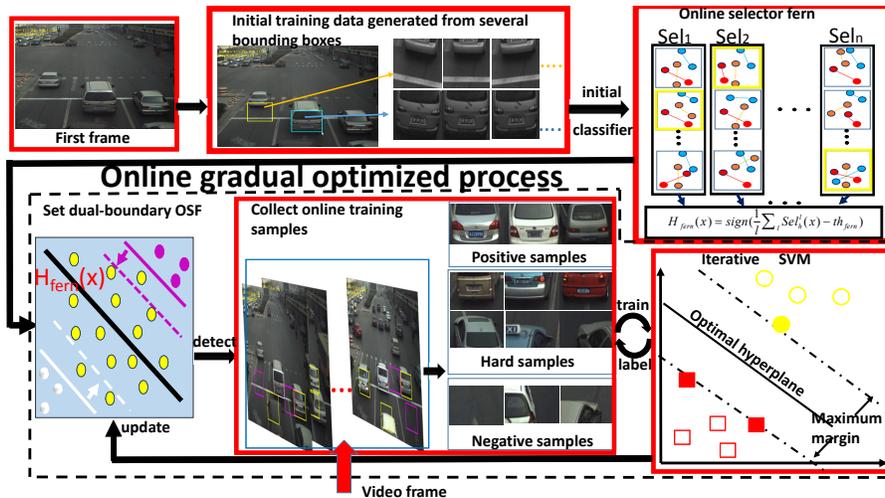


Figure 2: Approach overview

Next section, we will first introduce the Boosting fern classifier and its online learning algorithm realized by an online selection operator. The iterative SVM will be presented in section 5. Following, the online gradual optimal process is shown in section 6.

#### 4. Online Selector Fern(OSF) Generative Model

Traditional fern classifiers are widely used in the object tracking field [57], due to efficiency and high performance in tracking affine transformation planar objects. [58, 59, 60] extended fern classifier methods to detect objects appearing in the image under different orientations and view points. In this paper, we propose the online selector fern (OSF) integrated fern classifiers and an online feature selection strategy. Fern classifiers play the role of weak classifier and can be boosted into a strong model by online selection operators. Furthermore, the generative model will evolve into the dual-boundaries OSF, when the positive and negative boundaries are initialized around the decision hyperplane of the generative model.

First, we briefly describe the boosting fern model [59], second, the OSF generative model will be introduced to further improve the ability of the

model to express labelled samples, and third, the positive and negative boundaries will be initialized around the decision hyperplane of the OSF model, which evolves into the dual-boundary OSF classifier with the least model parameters.

#### 4.1. Boosting Fern

Let  $S^i = (f^i, C^i), i = 1, 2, \dots, I$  denotes training samples, which is an image patch set and their labels. Where  $I$  is the number of samples,  $C \subseteq \{c_+, c_-\}$  is the sample labels,  $f$  is an  $J$ -dimension feature vector which is used to describe the samples:

$$f = (f_1, f_2, \dots, f_J) \quad (3)$$

We employ the Local Binary Feature (LBF) [57] to map an image sample to a Boolean feature space by comparison between two random position intensities of a sample. Fern is denoted as a feature sub-space random sampled from the  $J$ -dimension feature space. The  $l$ th fern can be denoted as:

$$F_l = f_{l,1}, f_{l,2}, \dots, f_{l,s}, \quad l = 1, 2, \dots, L \quad (4)$$

Where  $s$  is the number of features in a fern. For a training sample, this gives an  $s$ -digit binary code to describe the sample appearances. In other words, each Fern maps 2D image samples to a  $2s$ -dimensional feature space, as shown in Fig. 3 (a). Accordingly, apply the fern  $F_l$  to each labelled training sample and learning the posterior distribution as histogram in each class  $P(F_l|c_+)$  and  $P(F_l|c_-)$ , as shown in Fig. 3(b).

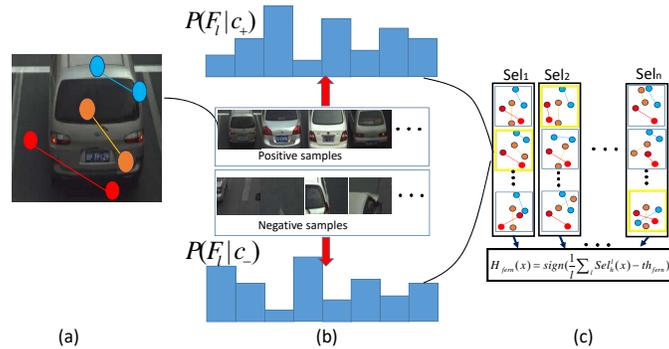


Figure 3: (a) LBF features;(b) posterior distribution of  $F_l$ ;(c) online selector of Boosting fern.

A single fern does not give the accurate estimation of the generative model in a specific surveillance scene, but we can build an ensemble of ferns by randomly choosing different subsets of LBF features. A weak classifier is defined by the co-occurrence of a random Fern  $F_l$ 's observation, which is further linearly combined to an ensemble classifier.

$$h_l(x) = \frac{P(F_l(x) = k|c_+) + \varepsilon}{P(F_l(x) = k|c_+) + P(F_l(x) = k|c_-) + \varepsilon} \quad (5)$$

where  $\varepsilon$  is a smoothing factor. Thus, the classifier is estimated approximately:

$$H_{fern}(x) = \text{sign}\left(\frac{1}{L} \sum_l h_l(x) - th_{fern}\right) \quad (6)$$

$th_{fern}$  is the threshold of each weak fern classifier, which is usually set to 0.5. Next section will introduce the weak classifier selection strategy and the online selector fern algorithm.

#### 4.2. Online Selector Fern(OSF)

From section 4.1, the fern based weak classifier is more discriminant than single feature based weak classifier because every fern is a group of features in fixed size of  $s$ . Moreover, the online learning process of each weak classifier is simplified by updating the post probability of every fern. However, in the online learning process, we must find a method to select the most discriminant fern classifier and minimize the first item in equation (2), the cost function of a Generative-Discriminative model. We denote a fern based weak classifier set  $h_l^{set} = h_{F_1}(x), h_{F_2}(x), h_{F_3}(x), \dots, h_{F_M}(x)$ , and a selector:

$$Sel_l = h_{F_m}(x), m \in 1, 2, \dots, M \quad (7)$$

where  $m$  is chosen according to an optimization criterion. In this paper, we use the criterion of picking the one that minimizes the following Bhattacharyya distance between the distributions of the class  $c_+$  and background  $c_-$ :

$$B_l = 2 \sum_{k=0}^{2^r} \sqrt{P(F_l(x = k)|c_+)P(F_l(x = k)|c_-)} \quad (8)$$

As shown in Fig. 3(c), a fixed set of  $N$  selectors are initialized randomly, each has a fixed set of  $M$  fern based weak classifiers. When the weak classifiers of each selector receive a new training sample, the classifiers are updated by changing the post probability distribution of each fern according to the location in a  $2^s$  dimension feature space partitioned by the fern. The weak classifier with the smallest Bhattacharyya distance is selected. This procedure is repeated for all selectors. A strong classifier is obtained by linear combination of selectors.

$$H_{fern}(x) = \text{sign}\left(\frac{1}{L} \sum_l Sel_h^l(x) - th_{fern}\right) \quad (9)$$

#### 4.3. Dual-boundaries OSF Classifier

Let  $H_{fern}^0(x)$  denote the initial OSF classifier based detector, which is firstly applied to the target video by sliding window search. All detected visual examples  $d_i (i \in 1, 2, \dots, N)$  are collected and divided into positive sample set,  $Set_{pos}$ , hard sample set,  $Set_{hard}$  and negative sample set,  $Set_{neg}$  based on the confidence level calculated by the OSF classifier.

$$\begin{cases} d_i \in Set_{pos} & \text{if } H_{fern}^0(d_i) > \beta + \frac{\theta}{2} \\ d_i \in Set_{hard} & \text{if } \beta + \frac{\theta}{2} > H_{fern}^0(d_i) > \beta - \frac{\theta}{2} \\ d_i \in Set_{neg} & \text{if } H_{fern}^0(d_i) < \beta - \frac{\theta}{2} \end{cases} \quad (10)$$

$\beta$  is the decision hyperplane.  $\beta + \frac{\theta}{2}$  and  $\beta - \frac{\theta}{2}$  are the positive and negative boundaries around  $\beta$  and have large margin in the initial stage. Thus, the most detection responses will located between the two boundaries and collected as hard samples with uncertain labels. The online selector fern classifier becomes an dual-boundary OSF.

Simplifying the parameters about the positive and negative boundaries to a scalar threshold  $\beta$  greatly reduces the complexity of the GDM. Experiments show that the model performs well even in clutter surveillance videos, as shown in Fig. 5.

## 5. Iterative SVM

A standard SVM classifier for two-class problem can be defined as:

$$\left\{ \begin{array}{l} \min_{W,b,\xi} \frac{1}{2} \|W\|^2 + c \sum_{i=1}^N \xi_i, \\ s.t. \forall i : y_i(W^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0, i = 1, 2, \dots, N \end{array} \right. \quad (11)$$

where  $x_i \in R^n$  is a training sample,  $y_i \in \{-1, 1\}$  is the label of  $x_i$ .  $c > 0$  is a regularization constant.

The traditional SVM belongs to the supervised discriminative model. However, semi-supervised SVM [61] or mi-SVM (Multiple Instance SVM) [62] can be trained by labelled and unlabelled samples simultaneously. In this paper, labelled samples, collected from several bounding boxes in the first frame, are much less than the traditional semi-supervised algorithms. Therefore, an iterative SVM algorithm is proposed to gradually label high-confidence hard samples in subspace.

As shown in Fig. 2, the initial training samples, generated by affine transformation of the several bounding boxes in the first frame, can be used as the initial labelled sample set  $L_0 = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$ , and the hard samples takes the place of unlabeled sample set  $U = x_n, x_{n+1}, \dots, x_g$  which is automatically collected from the detection responses of online selector fern based detector.

$$\left\{ \begin{array}{l} \min_{\hat{y}_j} \min_{W,b,\xi} \frac{1}{2} \|W\|^2 + c \sum_{i=1}^N \xi_i + d \sum_{j=1}^M \xi_j, \\ s.t. \forall i : y_i(W^T x_i + b) \geq 1 - \xi_i, \xi_i \geq 0, x_i \subseteq L, \\ s.t. \forall j : \hat{y}_j(W^T x_j + b) \geq 1 - \xi_j, \xi_j \geq 0, x_j \subseteq U \end{array} \right. \quad (12)$$

where  $\hat{y}_j$  is the pseudo labels of unlabelled samples, and  $d$  is the regularization constant of the unlabelled term.

Thus, the manual labelled samples are no longer needed in the optimized process. The SVM model can be trained in an unsupervised manner, as follows:

First, the initial SVM classifier, denoted as  $H_{SVM}^0(x)$ , is trained by the same initial sample set, generated by affine transformation of the bounding boxes in the first frame, which ensures the model is correctly initialized. The

HOG feature [2] is employed to train  $H_{SVM}^0(x)$ , which is different from the LBF feature used in the fern classifier training process. Thus, two different types of features can be integrated in our self-learning framework, which is crucial to improve the detection performance in cluttered environments [3].

Second, we perform the  $H_{SVM}^0(x)$  on the unlabeled sample set  $U$ . The predicted labels are denoted as  $lab_u = y_n(0), y_{n+1}(0), \dots, y_g(0)$ .

Third, denote  $th_{SVM}^P$  and  $th_{SVM}^N$  as the positive and negative thresholds according to SVM classification confidence. The labelled positive sample set can be updated by adding some hard samples with high classification confidence (more than  $th_{SVM}^P$ ), while the labelled negative sample set can be updated by adding some hard samples with low classification score (less than  $th_{SVM}^N$ ). This is a more conservative manner to update the labelled sample set when compared to the convenient semi-supervised SVM algorithm. Using the updated training set  $L_t$ , we train a new model, and perform classification again on  $U$ . The predicted labels are denoted as  $lab_u = y_n(1), y_{n+1}(1), \dots, y_g(1)$ .

If all the hard example labels are unchanged, the algorithm stops after the  $k$ th iteration. The SVM model and the predicted labels  $lab_u = y_n(k), y_{n+1}(k), \dots, y_g(k)$  of the hard sample set are the final output. If changed, perform the second and third step for the  $(k + 1)$ th iteration.

This is a heuristic optimization procedure similar to mi-SVM [62] without the constraint of at least one positive sample in the unlabelled sample set. In addition, the optimization process performs in a subspace, embedding in the online gradual optimal process, produces surprising experimental results without fine tuning, as shown in Fig. 8.

The iterative optimized process is triggered by an online hard sample collected module without manual annotation, as shown in Fig. 2. When the online selector Ferns are retrained by the convergent iterative SVM and labeled hard sample set according to a online gradual optimized process, we obtain a Generative-Discriminative model, consisting of a  $H_{fern}(x)$  and a  $H_{SVM}(x)$ , more dedicated to the problematic samples. As a consequence, the updated classifier puts more emphasis on the most distinctive parts of the object that reduces the global classification error.

## 6. Online Gradual Optimized Procedure

With the online gradual optimized procedure, a poor performance detector is acceptable in the beginning of the learning process, and will im-

prove from iteratively learning the hard samples located close to the decision boundary, which is the key idea behind our proposed framework.

The optimized procedure is divided into three steps. The first step optimize the generative model by solving  $C_G(x, y)$ , the second step labels hard samples by optimizing the discriminative model  $D$ , and the third step reduce the margin gradually between the positive and negative boundaries.

In essence, this procedure similar to an active learning process [56, 57] which reduce gradually the amount of human assistance of annotating hard samples by updating classifier models for the most distinctive parts of the object domain. However, to obtain labels for these samples, the active learner has to ask an oracle (e.g., a human expert) for labels. In our framework, the human expert can be instead by the discriminative model which has learned the hard samples in advance. Thus, the distance  $Dis(B_+, B_-)$  can be updated by evaluating the performance of the generative classifier  $G$  in real time.

From equation (10), the parameter  $\theta$ , determining the margin between the dual boundaries, can be minimized by the followed equation:

$$\theta = 1 - \nu \zeta_{fern} \quad (13)$$

where  $\nu$  is a sensitivity parameter which controls the learning speed of dual-boundary OSF classifier (set to 0.85 in our experiments).  $\zeta_{fern}$  measures the performance of the OSF classifier which makes the margin reduce process adaptive to the classifier learning process.  $\zeta_{fern}$  can be computed by:

$$\zeta_{fern} = \frac{\sum_{d_i \in Set_{hard}} [(H_{fern}(d_i) - \beta) \times sign(H_{SVM}(d_i))]}{\sum_{d_i \in Set_{hard}} |(H_{fern}(d_i) - \beta)|} \quad (14)$$

The overall online gradual optimized process is shown in Tab. 1. The output is a Generative-Discriminative model integrated by a ISVM model and a OSF classifier with dual decision boundaries. When the Generative-Discriminative model is used to detect individual class objects in a video, most candidate windows, generated by the sliding window strategy, are classified by the OSF classifier. Only a small fraction, located between the positive and negative boundaries, are dedicatedly classified by the ISVM model. This hybrid classifier system exploits the strengths of the individual classifier models by first performing a sample space partitioning and, second, assigning to each partition region an individual classifier that achieves high classification accuracy while being efficient at run-time. Moreover, the self-learning process is fully-autonomous and can be extended to other surveillance scenes

---

Table 1: ONLINE GRADUAL OPTIMIZED PROCESS

---

**Input:** The initial training sample set  $L_0 = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$  generated from the affine warping of several bounding boxes in the first video frame. Denote empty hard sample set  $Set_{hard} \in \{\}$ . Initialised online dual-boundary OSF classifier  $H_{fern}^0(x)$ , and some initialised parameters including  $\beta = 0.5$ ,  $\theta_0 = 1$  and  $\nu = 0.85$ .

---

**Output:**  $H_{hyb}(x)$  composed by a  $H_{fern}(x)$  and a  $H_{SVM}(x)$ .

---

Training initial SVM model  $H_{SVM}^0(x)$  from the initial training sample set  $L_0$ .

t=0

while ( $\theta > 0.3$ )

- using  $H_{fern}^t(x)$  to detect object in video
- collect  $d_i \in Set_{hard}$ , and calculated the number of collected hard samples  $N_{hard}$
- if ( $N_{hard} > 100$ )

$H_{SVM}^{t+1}(x) = \text{ISVM}(H_{SVM}^t(x), Set_{hard}, L_0)$   
; optimize the discriminative model in eq (2)

$Lab_u = H_{SVM}^{t+1}(Set_{hard})$   
; label the hard samples set by  $H_{SVM}^t(x)$

$H_{fern}^{t+1}(x) = \text{OSF}(H_{fern}^t(x), Set_{hard}, Lab_u)$   
; optimize the generative model in eq (2)

classify hard sample by  $H_{fern}^{t+1}(x)$

calculate  $\zeta_{fern}$  by equation (14)

calculate  $\theta$  by equation (13)  
; optimize the distance item in eq (2)

$Set_{hard} \in \{\}$ ; clear the  $Set_{hard}$

t = t+1

-end if

repeat

or object class detection tasks. Resolving the multi-scene object detection problem is very important, and can be settled by combining multiple self-learning scene-specific detectors.

## 7. Experiments And Comparisons

The proposed method is evaluated on vehicle and pedestrian detection problems, which play a key role in current intelligent surveillance systems. For the vehicle detection task, the GRAM-RTM dataset [63] and the Vehicle dataset are used to evaluate our approach. The Vehicle dataset has been captured and labeled by ourselves. The two datasets, composed of 6 video sequences with different view points and resolution levels, show a real urban road scene with multiple vehicles at the same time. As shown in Fig. 4, four sequences: Hx, Yk, Hi and Hn, are used from these datasets, which have 6415, 1663, 7520 and 1085 frames, respectively, of different resolutions:  $1224 \times 1024$ ,  $512 \times 288$ ,  $800 \times 480$  and  $2448 \times 2048$ . Hx has 912 GT instances of the vehicle, whereas Yk and Hi have 344 and 2089 GT instances of the vehicle, respectively. To evaluate the pedestrian detection performance, six sequences: TownCentreXVID.avi (Town), PNNLParkingLot2.avi (PNN), WalkByShop1front.avi (Shop), OneShopLeave2Enter (Enter), Meet-Crowd.avi (Walk), and S2.L1View-001.avi (S2), are used from the well-known public Towncenter [64], PNN-Parking-Lot2/Pizza [36], CAVIAR [31] and PETS2009 [35] datasets, which have different appearances due to different imaging view points, as shown in Fig. 4. The Ground-Truth is available at [31, 35, 64, 36]. Note that the Town sequences are much longer than the other sequences with an average of about sixteen people visible at any frame. This is a challenged dataset to most pedestrian detection methods.

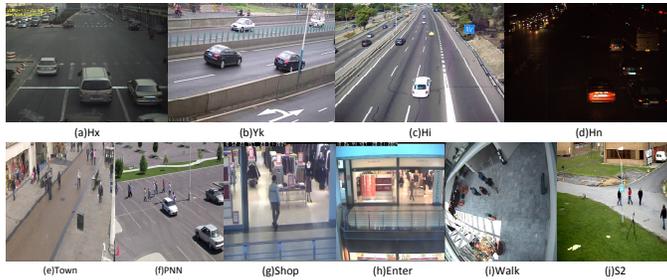


Figure 4: Video sequences.

In each experiment, we trigger the video-specific object learning algorithm by several bounding boxes in the first frame. Thus, this method can be conveniently extended into each video sequence. In terms of the necessary parameters that define our Generative-Discriminative model, we would like to point out in all the experiments, described in the following section, the dual-boundary OSF classifier have 10 selectors with initial margin parameter  $\theta_0$  setting to 1. Each selector uses 10 Random fern classifiers with 6 binary local features. When online optimizing the ISVM model, HOG feature is employed to describe the samples which are divided into cells of size  $16 \times 16$  pixels and each group of  $8 \times 8$  cells is integrated into a block. In other parameters, in online gradually optimized process,  $\nu$ , controlling the updating speed of the online dual-boundary OSF classifier, is set to 0.85.

In traditional object detection methods, the detection performance varies considerably due to resolution difference in different videos. It is difficult for these traditional methods to set a appropriate detection scale to satisfy all the possible resolutions in different videos. But in our framework, we can conveniently ensure the optimal detection scales for every testing video according to the bounding box in the first frame, which describes the accurate object size in surveillance videos. From this, we increase 11 different scales and achieve robust detection performance in each testing video.

In experiments, our approach has demonstrated state of the art detection performance in each testing video after no more than 5 hours self-training process, learning about 400-1500 samples without any human effort, as shown in Tab. 6.

### *7.1. Online Generative-Discriminative model Structure and Parameter Test*

We initially analyze the object detection performance of the proposed framework on the sequence Hx from the Vehicle dataset, which will reveal the influence of different hybrid classifier strategies and human annotation. As shown in Fig. 5, these classifiers were initialized by the same sample set generated by affine warping of several bounding boxes in the first frame of the sequence Hx, but have different learning processes and classifier systems. Here, we evaluate five different alternatives:

**Human fern:** the vehicles are detected only by the OSF classifier, which is online trained by 850 human labeled frames, including 246 positive samples and 500 negative samples collected from the sequence Hx.

**Human fern SVM:** the detector is a Generative-Discriminative model consisting of a dual-boundary OSF classifier and an ISVM model, which is

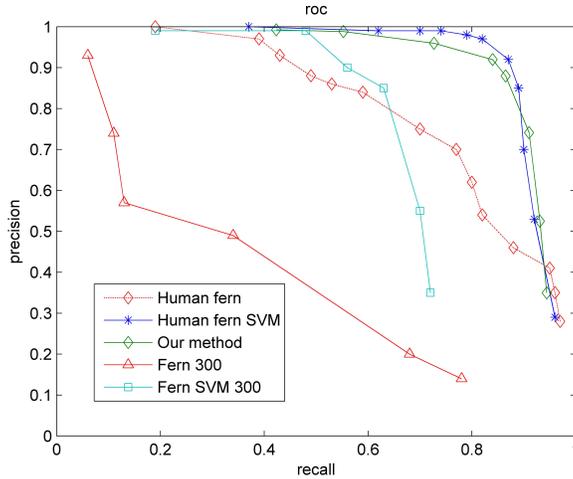


Figure 5: Experiments on analyzing our framework structure.

supervised by 850 frames, including manually annotated 246 positive samples and 500 negative samples.

**Our approach:** The detector is a Generative-Discriminative model composed of a dual-boundary OSF classifier and an ISVM classifier. The OSF classifier is used first to detect objects by the sliding window searching strategy and then the ISVM focuses on recognizing the hard samples distributed around the decision boundary. Note the Generative-Discriminative model is obtained by self-learning **1460** frames, about **336** positive samples and **687** negative samples, all collected and labeled automatically; some collected samples are shown in Fig. 6.

**Fern 300:** the detector is the only OSF classifier which is supervised by 300 frames under human guidance, learning about 160 positive samples and 160 negative samples.

**Fern SVM 300:** the detector is a Generative-Discriminative model which is obtained by human annotated 300 frames, about 160 positive samples and 160 negative samples.

The ROC curves are shown in Fig. 5. The Generative-Discriminative model based detector significantly outperforms single classifier. The results clearly demonstrate the ability of our Generative-Discriminative model to handle the most distinctive parts of the object category.

Our method has comparable performance to the supervised Generative-

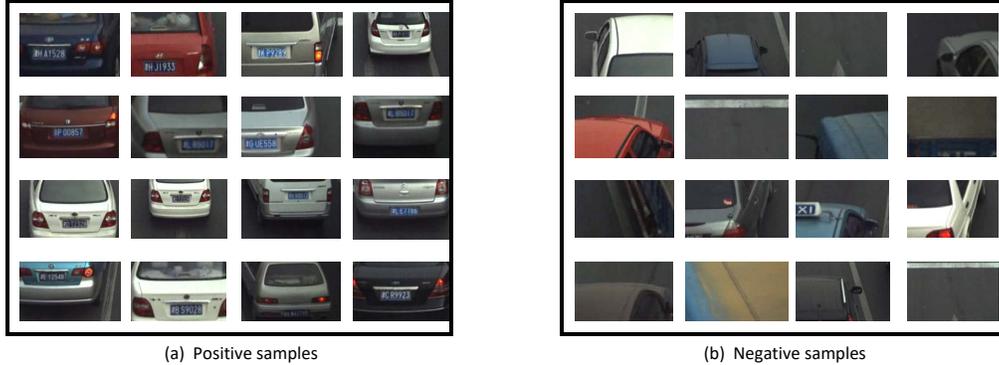


Figure 6: Online collected training data.

Discriminative model (Human fern SVM), demonstrating our framework, with a high label correct rate, can improve detection performance by focusing learning on the problematic samples located near the decision boundary.

In addition, the positive and negative thresholds  $th_{SVM}^P$  and  $th_{SVM}^N$ , denoted in the ISVM model, are critical for the "hard sample" learning. The two parameters are mutually opposite numbers. We change the value of  $th_{SVM}^P$  from 0.2 to 1.0. and applied to vehicle detection on sequence Hx. The ROC curve is shown in Fig. 7. The F-measure is calculated in Tab. 2. From the experimental point of view, when the  $th_{SVM}^P$  is about 0.8 (remains the same in subsequent experiments), the detection performance is satisfactory.

Table 2: COMPARISON WITH DIFFERENT  $th_{SVM}^P$

Sequence	0.2	0.4	0.6	0.8	1.0
Hx	0.7718	0.8010	0.7795	0.8779	0.7544

### 7.2. Comparison with Scene-specific object detection Methods

[32, 36, 37] are the widely used scene-specific object detection methods in recent years. [32] achieves the best performance in multiple databases due to the progressive latent model and the graph-based label propagation. Compared with [32], the GDM improves the precision by nearly 10 percent in the Towncentre dataset, and achieves the best detection performance in the PNN-Parking-Lot2/Pizza dataset, as shown in Fig. 8 and Tab. 3. The main

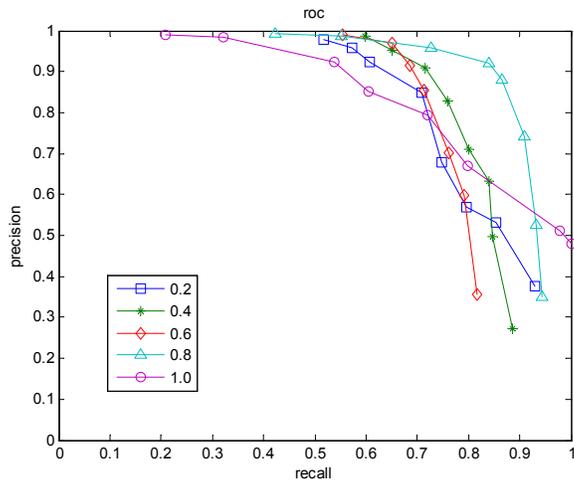


Figure 7: Parameter  $th_p$  test.

reason is that our method proposes to optimize a multiple latent variables model in subspace, and can autonomously learn a scene-specific detector through a poor initial detector. Detection results are shown in Fig. 11.

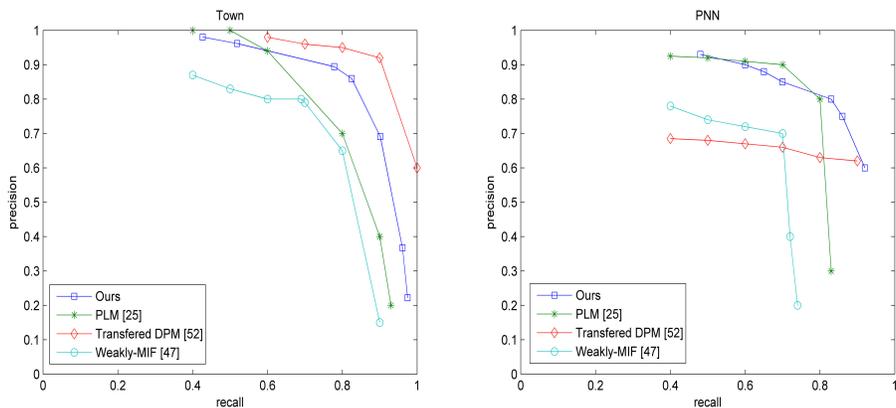


Figure 8: Comparison with scene-specific object detection methods

### 7.3. Comparison with Online and offline Learning Methods

In this section, the proposed method is compared with three online learning object detection methods [31, 26, 28], which have achieved satisfactory results among the online learning object detection frameworks. More

Table 3: COMPARISON WITH SCENE-SPECIFIC OBJECT DETECTION METHODS

Sequence	Shu[47]	Cinbis[52]	Ye[25]	Ours
Town	0.7422	0.9098	0.7466	0.8419
PNN	0.7000	0.7342	0.8000	0.8147

specifically, the proposed approach is compared with four supervised methods: Boosted fern [59], ISVM [61], ACF [6] and FernSVM, a supervised Generative-Discriminative model. 300 positive samples and 900 negative samples are collected and labelled manually from S2 sequence to train the offline learning classifiers. In the Hx sequence of the Vehicle Dataset, the number of supervised positive and negative training data become 200 and 500, respectively.

The ROC curves are shown in Fig. 9 and Fig. 10. The according F-measures are calculated in Tab. 4 and Tab. 5. The ACF method outperforms the proposed method on Hx and S2 sequences. Our approach outperforms the other online object detection methods and achieves detection performance competitive with the supervised methods. We demonstrate the primary reason for these results is our special strategy to address the hard samples, required to detect objects in cluttered environments. Detection results are shown in Fig. 11.

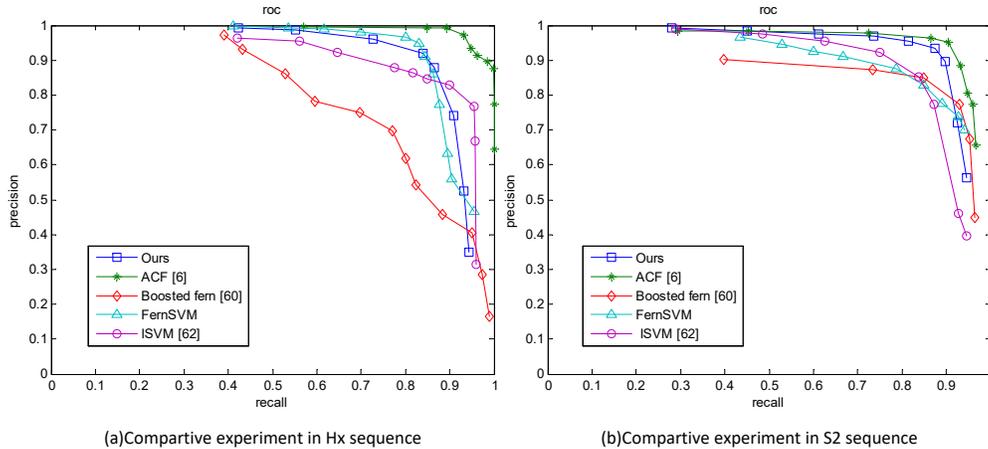


Figure 9: Comparison with supervised learning methods

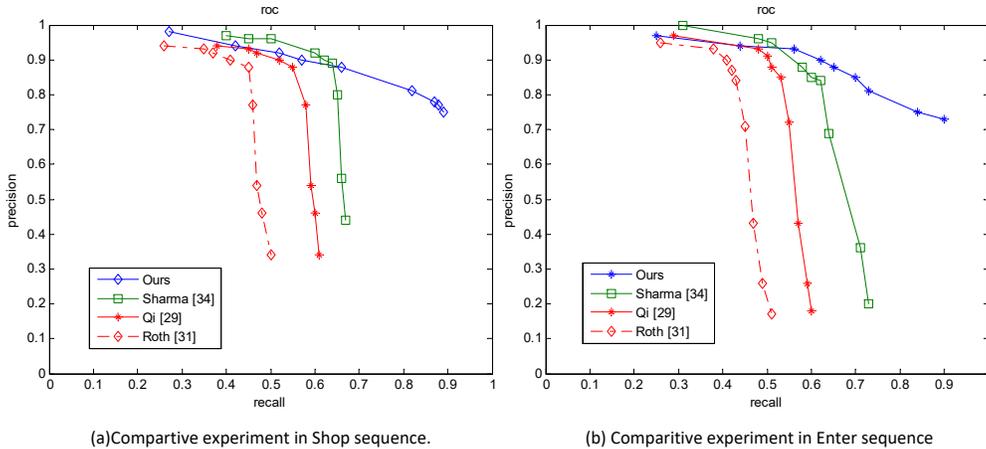


Figure 10: Comparison with online learning methods

Table 4: COMPARISON WITH ONLINE LEARNING METHODS

Sequence	Roth[31]	Qi[29]	Sharma[34]	Ours
Shop	0.5955	0.6769	0.7446	0.8225
Enter	0.5634	0.6454	0.6992	0.8061

#### 7.4. Multi-view Object Detection in Video Sequences

Further tests check if the system can self-adjust to view point changes. Our system, initialized from several bounding boxes, was applied to the video sequences Yk and Hi from Vehicle and GRAM-RTM datasets. For the Yk sequence, the detector self-learns 156 positive samples and 323 negative samples. For the Hi sequence, 244 positive samples and 598 negative samples were automatically collected and labelled to online train the video-specific detector. As shown in Fig. 12, the new detector achieves state of the art performance without any human intervention. Following this action, our framework is used for multi-view pedestrian detection in CAVIAR and PETS2009 dataset. The number of online learning samples, self-learning duration and detection speed, different in each video, as shown in Tab. 6. Thus, the trained scene-specific detector can detect object in real time on a standard PC (Intel Core i5 3.2 GHz with 4 GB RAM). The detection performance is shown in Fig. 12. It is valuable to note all self-learning processes are fully unsupervised without any prior knowledge or constraint. Our method can be easily employed in other surveillance scenes and form a bottom-up

Table 5: COMPARISON WITH SUPERVISED LEARNING METHODS

Method	S2		Hx	
	F-Measure	FPS	F-Measure	FPS
Ours	0.9036	62	0.8779	10
ACF[6]	0.9271	14	0.9518	5
ISVM[62]	0.8438	2	0.8624	4
Boosted fern[60]	0.8496	92	0.7320	11
FernSVM	0.8391	62	0.8771	10

multi-view object detection method.

Table 6: DISCRIPTION OF ONLINE LEARNING PROCESS

Sequence	Positive samples	Negative samples	Duration times(s)	Detection speed(FPS)
Yk	156	323	815	36
Hi	244	598	2751	15
Hx	336	687	25196	10
Shop	281	479	1121	19
Walk	410	311	2728	34
S2	504	994	9787	62

## 8. Conclusions and Discussions

This paper presents a self-learning video object detection framework. In this framework, a Generative-Discriminative model can be trained by online gradual optimized process without any human labelled samples or generic detectors. This framework can be easily employed in multiple surveillance scenarios and will result in scene-specific detectors more dedicated by a hierarchical optimized process to the problematic samples. Consequently, the Generative-Discriminative model puts more emphasis on the most distinctive parts of the object category, which leads to the reduction of the global classification error. This process simulates the autonomous learning of human. Experimental results show our approach achieves high accuracy on multi-scene vehicle and pedestrian detection tasks.



Figure 11: Some detection results using our approach in Vehicle, GRAM-RTM, Town-center, PNN-Parking-Lot2/Pizza, CAVIAR and PETS2009 Datasets. Note that the Hx, Yk, Hi and Hn sequences have high resolutions. Thus, a ROI region of purple box was settled for improving the detection speed. The detection results were shown in the last three rows.

Future investigations will integrate the self-learning detector with an on-line learning tracker to form a scene-specific multiple object detection and tracking system which will simultaneously increase the performance of the detection and tracking.

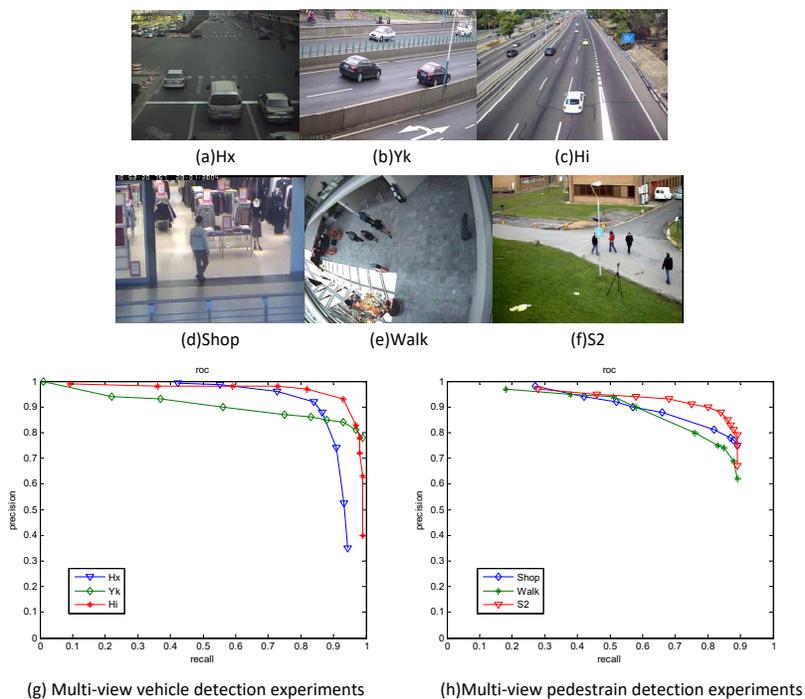


Figure 12: Multi-view vehicle and pedestrian detection experiments.

## Acknowledgment

This work was supported by the National Natural Science Foundation of China (61302137, 61603357, 61271328 and 61603354), Wuhan huanghe Elite Project, Fundamental Research Funds for the Central Universities Young Teacher Promotion Program-Outstanding Youth Foundation, China University of Geosciences (Wuhan)(CUGL170210), Fundamental Research Funds for National University, China University of Geosciences (Wuhan) (1610491B06).

## References

- [1] P. Viola, M. Jones, Rapid object detection using a boosted cascade of simple features, in: IEEE Conference on Computer Vision and Pattern Recognition, Vol. 1, 2001, pp. 511–518.
- [2] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: IEEE Conference on Computer Vision and Pattern Recognition, Vol. 1, 2005, pp. 886–893.

- [3] X. Wang, T. X. Han, S. Yan, An hog-lbp human detector with partial occlusion handling, in: IEEE International Conference on Computer Vision, 2009, pp. 32–39.
- [4] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, D. Ramanan, Object detection with discriminatively trained part-based models, *IEEE Trans. Pattern Anal. Mach. Intell.* 32 (9) (2010) 1627–1645.
- [5] B. Li, B. Tian, Y. Li, D. Wen, Component-based license plate detection using conditional random field model, *IEEE Trans. Intell. Transp. Syst.* 14 (4) (2013) 1690–1699.
- [6] P. Dollár, R. Appel, S. Belongie, P. Perona, Fast feature pyramids for object detection, *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (8) (2014) 1532–1545.
- [7] J. Redmon, A. Farhadi, Yolo9000: Better, faster, stronger, arXiv preprint arXiv:1312.6229, 2017.
- [8] X. Wang, M. Wang, W. Li, Scene-specific pedestrian detection for static video surveillance, *IEEE Trans. Pattern Anal. Mach. Intell.* 36 (2) (2014) 361–374.
- [9] J. Pang, Q. Huang, S. Yan, S. Jiang, L. Qin, Transferring boosted detectors towards viewpoint and scene adaptiveness, *IEEE Trans. Image Process* 20 (5) (2011) 1388–1400.
- [10] M. T. Rosenstein, Z. Marx, L. P. Kaelbling, T. G. Dietterich, To transfer or not to transfer, in: *Advances in Neural Information Processing Systems*, Vol. 898, 2005.
- [11] A. Torralba, K. P. Murphy, W. T. Freeman, Sharing visual features for multiclass and multiview object detection, *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (5).
- [12] B. Leibe, A. Leonardis, B. Schiele, Robust object detection with interleaved categorization and segmentation, *Int. J. Comput. Vision* 77 (1-3) (2008) 259–289.
- [13] J. R. R. Uijlings, K. E. A. V. D. Sande, T. Gevers, A. W. M. Smeulders, Selective search for object recognition, *Int. J. Comput. Vision* 104 (2) (2013) 154–171.

- [14] B. C. Ko, J. H. Jung, J. Y. Nam, View-independent object detection using shared local features, *Journal of Visual Languages & Computing* 28 (2015) 56–70.
- [15] M. Viola, M. J. Jones, P. Viola, Fast multi-view face detection, in: *IEEE Conference on Computer Vision and Pattern Recognition*, Vol. 11, 2003, p. 276286.
- [16] B. Wu, R. Nevatia, Cluster boosted tree classifier for multi-view, multi-pose object detection, in: *IEEE International Conference on Computer Vision*, 2007, pp. 1–8.
- [17] B. Wu, H. Ai, C. Huang, S. Lao, Fast rotation invariant multi-view face detection based on real adaboost, in: *IEEE International Conference on Automatic Face and Gesture Recognition*, 2004, pp. 79–84.
- [18] X. Zhu, D. Ramanan, Face detection, pose estimation, and landmark localization in the wild, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 2879–2886.
- [19] D. Hoiem, C. Rother, J. Winn, 3d layoutcrf for multi-view object class recognition and segmentation, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- [20] N. Razavi, J. Gall, L. Van Gool, Backprojection revisited: Scalable multi-view object detection and similarity metrics for detections, *Lecture Notes in Computer Science* 6311 (2010) 620–633.
- [21] E. Seemann, B. Leibe, B. Schiele, Multi-aspect detection of articulated objects, in: *IEEE Conference on Computer Vision and Pattern Recognition*, Vol. 2, 2006, pp. 1582–1588.
- [22] A. Thomas, V. Ferrar, B. Leibe, T. Tuytelaars, B. Schiel, L. Van Gool, Towards multi-view object class detection, in: *IEEE Conference on Computer Vision and Pattern Recognition*, Vol. 2, 2006, pp. 1589–1596.
- [23] D. Q. Zhang, S. F. Chang, A generative-discriminative hybrid method for multi-view object detection, in: *IEEE Conference on Computer Vision and Pattern Recognition*, Vol. 2, 2006, pp. 2017–2024.

- [24] S. Savarese, F. F. Li, 3d generic object categorization, localization and pose estimation, in: IEEE International Conference on Computer Vision, 2007, pp. 1–8.
- [25] O. Javed, S. Ali, M. Shah, Online detection and classification of moving objects using progressively improving detectors, in: IEEE Conference on Computer Vision and Pattern Recognition, Vol. 1, 2005, pp. 696–701.
- [26] Z. Qi, Y. Xu, L. Wang, Y. Song, Online multiple instance boosting for object detection, *Neurocomputing* 74 (10) (2011) 1769–1775.
- [27] V. Nair, J. J. Clark, An unsupervised, online learning framework for moving object detection, in: IEEE Conference on Computer Vision and Pattern Recognition, Vol. 2, 2004, pp. 317–325.
- [28] P. M. Roth, H. Grabner, H. Bischof, D. Skocaj, A. Leonardist, Online conservative learning for person detection, in: Proceedings of IEEE Workshop on VS-PETS, 2005, pp. 223–230.
- [29] X. Wang, G. Hua, T. X. Han, Detection by detections: Non-parametric detector adaptation for a video, in: IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 350–357.
- [30] P. Sharma, C. Huang, R. Nevatia, Unsupervised incremental learning for improved object detection in a video, in: IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 3298–3305.
- [31] P. Sharma, R. Nevatia, Efficient detector adaptation for object detection in a video, in: IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 3254–3261.
- [32] Q. Ye, T. Zhang, W. Ke, Q. Qiu, J. Chen, G. Sapiro, B. Zhang, Self-learning scene-specific pedestrian detectors using a progressive latent model, in: IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 509–518.
- [33] F. Xiao, J. L. Yong, Track and segment: An iterative unsupervised approach for video object proposals, in: IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 933–942.

- [34] S. Kwak, M. Cho, I. Laptev, J. Ponce, C. Schmid, Unsupervised object discovery and tracking in video collections, in: IEEE International Conference on Computer Vision, 2015, pp. 3173–3181.
- [35] J. Ferryman, A. Shahrokni, Pets2009: Dataset and challenge, in: Performance Evaluation of Tracking and Surveillance (PETS-Winter), 2009 Twelfth IEEE International Workshop on, 2009, pp. 1–6.
- [36] G. Shu, A. Dehghan, M. Shah, Improving an object detector and extracting regions using superpixels, in: IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 3721–3727.
- [37] R. G. Cinbis, J. Verbeek, C. Schmid, Weakly supervised object localization with multi-fold multiple instance learning, *IEEE Trans. Pattern Anal. Mach. Intell.* 39 (1) (2017) 189–203.
- [38] A. Levin, P. A. Viola, Y. Freund, Unsupervised improvement of visual detectors using co-training., in: IEEE International Conference on Computer Vision, 2003, pp. 626–633.
- [39] Y. Yang, G. Shu, M. Shah, Semi-supervised learning of feature hierarchies for object detection in a video, in: IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 1650–1657.
- [40] M. Wang, W. Li, X. Wang, Transferring a generic pedestrian detector towards specific scenes, in: IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 3274–3281.
- [41] A. Prest, C. Leistner, J. Civera, C. Schmid, V. Ferrari, Learning object class detectors from weakly annotated video, in: IEEE Conference on Computer Vision and Pattern Recognition, 2012, pp. 3282–3289.
- [42] Y. Liu, Y. Wang, A. Sowmya, F. Chen, Soft hough forest-erts: Generalized hough transform based object detection from soft-labelled training data, *Pattern Recognition* 60 (2016) 145–156.
- [43] K. Kumar Singh, F. Xiao, Y. Jae Lee, Track and transfer: Watching videos to simulate strong human supervision for weakly-supervised object detection, in: IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 3548–3556.

- [44] Z. G. Fan, B. L. Lu, Fast recognition of multi-view faces with feature selection, in: IEEE International Conference on Computer Vision, Vol. 1, 2005, pp. 76–81.
- [45] J. Ng, S. Gong, Multi-view face detection and pose estimation using a composite support vector machine across the view sphere, in: International Workshop on Recognition, Analysis, and Tracking of Faces and Gestures in Real-Time Systems, 1999, pp. 14–21.
- [46] M. Weber, W. Einhauser, M. Welling, P. Perona, Viewpoint-invariant learning and detection of human heads, in: IEEE International Conference on Automatic Face and Gesture Recognition, 2000, pp. 20–27.
- [47] S. Z. Li, Z. Zhang, Floatboost learning and statistical face detection, IEEE Trans. Pattern Anal. Mach. Intell. 26 (9) (2004) 1112–1123.
- [48] M. Stark, M. Goesele, B. Schiele, Back to the future: Learning shape models from 3d cad data, in: British Machine Vision Conference, 2010, pp. 1–11.
- [49] R. J. López-Sastre, T. Tuytelaars, S. Savarese, Deformable part models revisited: A performance evaluation for object category pose estimation, in: IEEE International Conference on Computer Vision Workshops, 2011, pp. 1052–1059.
- [50] C. Gu, X. Ren, Discriminative mixture-of-templates for viewpoint classification, in: European Conference on Computer Vision, 2010, pp. 408–421.
- [51] C. M. Christoudias, R. Urtasun, T. Darrell, Unsupervised feature selection via distributed coding for multi-view object recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, 2008, pp. 1–8.
- [52] B. Pepik, M. Stark, P. Gehler, B. Schiele, Multi-view and 3d deformable part models, IEEE Trans. Pattern Anal. Mach. Intell. 37 (11) (2015) 2232–2245.
- [53] C. Xu, D. Tao, C. Xu, Multi-view learning with incomplete views, IEEE Trans. Image Process 24 (12) (2015) 5812–5825.

- [54] H. Celik, A. Hanjalic, E. A. Hendriks, Unsupervised and simultaneous training of multiple object detectors from unlabeled surveillance video, *Computer Vision and Image Understanding* 113 (10) (2009) 1076–1094.
- [55] I. Huerta, M. Pedersoli, J. González, A. Sanfeliu, Combining where and what in change detection for unsupervised foreground learning in surveillance, *Pattern Recognition* 48 (3) (2015) 709–719.
- [56] Y. Amit, 2D object detection and recognition: Models, algorithms, and networks, Vol. 99, MIT Press, 2002.
- [57] M. Ozuysal, P. Fua, V. Lepetit, Fast keypoint recognition in ten lines of code, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2007, pp. 1–8.
- [58] M. Villamizar, F. Moreno-Noguer, J. Andrade-Cetto, A. Sanfeliu, Efficient rotation invariant object detection using boosted random ferns, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2010, pp. 1038–1045.
- [59] M. Villamizar, J. Andrade-Cetto, A. Sanfeliu, F. Moreno-Noguer, Bootstrapping boosted random ferns for discriminative and efficient object classification, *Pattern Recognition* 45 (9) (2012) 3141–3153.
- [60] D. Levi, S. Silberstein, A. Bar-Hillel, Fast multiple-part based object detection using kd-ferns, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 947–954.
- [61] A. R. Shah, C. S. Oehmen, B.-J. Webb-Robertson, Svm-hustlean iterative semi-supervised machine learning approach for pairwise protein remote homology detection, *Bioinformatics* 24 (6) (2008) 783–790.
- [62] S. Andrews, I. Tsochantaridis, T. Hofmann, Support vector machines for multiple-instance learning, *Advances in Neural Information Processing Systems* 15 (2) (2003) 561–568.
- [63] R. Guerrero-Gómez-Olmedo, R. J. López-Sastre, S. Maldonado-Bascón, A. Fernández-Caballero, Vehicle tracking by simultaneous detection and viewpoint estimation, in: *International Work-Conference on the Interplay Between Natural and Artificial Computation*, 2013, pp. 306–316.

- [64] B. Benfold, I. Reid, Stable multi-target tracking in real-time surveillance video, in: IEEE Conference on Computer Vision and Pattern Recognition, 2011, pp. 3457–3464.