Generalized Haar Filter based Deep Networks for Real-Time Object Detection in Traffic Scene

Keyu Lu¹, Jian Li¹, Xiangjing An¹, Hangen He¹

¹College of Mechatronic Engineering and Automation, National University of Defense Technology, Changsha 410073, Hunan, China

Email: keyu.lu@nudt.edu.cn; keyu.lu@alumni.ubc.ca

Abstract. Vision-based object detection is one of the fundamental functions in numerous traffic scene applications such as self-driving vehicle systems and advance driver assistance systems (ADAS). However, it is also a challenging task due to the diversity of traffic scene and the storage, power and computing source limitations of the platforms for traffic scene applications. This paper presents a generalized Haar filter based deep network which is suitable for the object detection tasks in traffic scene. In this approach, we first decompose a object detection task into several easier local regression tasks. Then, we handle the local regression tasks by using several tiny deep networks which simultaneously output the bounding boxes, categories and confidence scores of detected objects. To reduce the consumption of storage and computing resources, the weights of the deep networks are constrained to the form of generalized Haar filter in training phase. Additionally, we introduce the strategy of sparse windows generation to improve the efficiency of the algorithm. Finally, we perform several experiments to validate the performance of our proposed approach. Experimental results demonstrate that the proposed approach is both efficient and effective in traffic scene compared with the state-of-the-art.

Keywords: Generalized Haar filter; Deep networks; Object detection; Traffic scene

1 Introduction

Recent advances in self-driving vehicles and advance driver assistance system (ADAS) have attracted keen attention and interest from researchers and automobile manufacturers. Vision sensor plays an important role in these areas due to its faster response, lower price and power consumption compared with other popular sensors such as LiDAR and millimeter-wave radar. Moreover, vision sensor has the ability to capture rich information from traffic scene (such as luminance, color and texture) [1], which is beneficial for object detection.

Vision-based object detection is one of the fundamental functions for self-driving vehicle systems and advance driver assistance systems (ADAS), which need to detect the objects around and check whether they are dangerous to the host vehicle. However, object detection in traffic scene from image are still challenging job. On one hand, traffic scenes are diverse and the presence of objects (e.g. vehicles and pedestrians) in

traffic scenes are extremely flexible and unpredictable. on the one hand, most traffic scene applications have several special requirements such as real-time, portable (e.g. for ADAS device), low price and power consumption, which are quite different from object detection tasks in the ILSVRC [2] and COCO [3] competitions.

Recently, a great number of researchers have been interested in learning based approaches. Especially in recent years, the success of deep learning boosts the development of vision-based object detection. Convolutional Neural Network (CNN) [6] is one of the most popular forms of deep networks. By using the strategies of local receptive fields, weight sharing and spatial pooling [7], CNN has made a breakthrough in computer vision and image processing. For instance, AlexNet [6], which is a type of CNN, has made a startling achievement in the competition of ILSVRC-2012 [2] and demonstrated its superiority in image classification. However, with the increasing of network depth, CNN has met the bottleneck in training [8]. Recently, He *et al.* proposed the deep residual network [8] which employs shortcut connections to overcome the limitations in training a deeper CNNs.

Despite the fact that CNN has achieve a tremendous success in computer vision and image processing, there still exist two main problems that have to be solved when applied to objects detection tasks. Firstly, CNN has a multitude of convolutions that have to be calculated, it would be rather inefficient in object detection if using traditional dense sliding windows paradigm [4,5]. Secondly, CNN has the property of shift invariance, that is, it is less sensitive to the shift of objects in input image patch. Consequently, it can not achieve precise localization if it is directly applied to objects detection [9].

To overcome these limitations of CNN, Girshick *et al.* introduced the framework of region proposal based CNN (called R-CNN) [10] and successfully applied it to object detection. The main idea of region proposal based CNN is performing CNNs on candidate bounding-boxes (called "proposals") which have potential to contain objects [10]. Later, the updated version "Fast R-CNN" [9] is proposed to improve the efficiency of object detection. This method combines R-CNN with SPPnet [11]. In this way, computation of proposals can be sharing and runtime of object detection can be reduced dramatically. However, in both R-CNN and Fast R-CNN, the proposals are generated by Selective Search [12], which is quite inefficient and thus limits the detection speed of these methods. To overcome this limitation, Faster R-CNN [13] proposed to used Region Proposal Networks (RPNs) to generate proposals instead of Selective Search [12]. By this means, the stage proposal generation and CNN-based classification/regression can be performed under an unified framework and thus the detection speed can be boosted with the help of GPU.

However, region proposal based CNNs are complex and not easy to optimize [14]. For this reason, several works attempt to achieve real-time objects detection by regarding it as a regression problem. YOLO [14] is one of the pioneering works on deep regression networks based object detection. The approach is able to simultaneously outputs the location of bounding box, category and its confidence score for each object in the image at a extremely high frame rate. Nevertheless, it does not work well in small objects detection. In the same vein, Liu *et al.* proposed SSD (Single Shot MultiBox Detector) [15] for real-time objects detection. It equipped the deep regression networks with several new techniques such as multi-scale feature maps and default boxes [15].

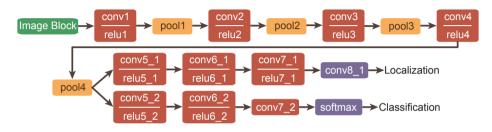


Fig. 1: Architecture of the Deep Networks.

Despite the fact that it is able to be faster and more accurate compared with YOLO [14], its performance on small objects detection is still unsatisfactory.

In this paper, we introduce a practical and robust approach for real-time object detection in traffic scene. The approach requires less storage and computing resources, and thus it is feasible for traffic scene applications. The novelties of this work lie in four fold:

- (1) We present a local regression strategy for accurate objects detection. As we know, current regression based solutions (such as YOLO [14] and SSD [15]) employ global regression strategy to regress the bounding-box of each object from the whole image, which is actually a much more difficult task compared with local regression. Thus, global regression strategy often needs to be supported by a complex or large-scale network and yet its performance on small objects detection is still unsatisfactory. To solve these problem, we introduce the local regression strategy that performs tiny regression networks on small image patches to detect and locate objects. As the regression networks are less sensitive to the scaling and shift of the objects in image patches, a series of sparse sliding-windows can be obtained from multi-scale image pyramid. Finally, objects can be detected and located efficiently based on these sparse sliding-windows;
- (2) We introduce the generalized Haar filter based deep networks where each weight larger than 3×3 are constrained to the form of generalized Haar filter in training phase. Owing to the strong representation of Haar filters, the networks are able to achieve a high performance. Besides, the networks consume much less storage and computing resources compared with traditional deep networks, which make them possible to be utilized in traffic scene applications. In addition, the constraint of generalized Haar filter provides a form of regularization which is able to improve the generalization ability of the deep networks;
- (3) For object detection problems in traffic scene, a sparse windows generation method is proposed. The method first generates a series of sparse sliding-windows in multi-scale image pyramid by setting a specific stride according to the scale and shift tolerance of our deep network. In this way, it can be ensure that each object is completely contained in at least one window. Besides, in most object detection systems for traffic scene, camera is mounted on a fixed position (e.g. mounted on the top of a vehicle windshield). Consequently, objects (e.g. vehicles and pedestrians)

- produce the predetermined location-specific patterns in images. According to this assumption, perspective geometry is also utilized to further reduce the candidate windows:
- (3) We construct a tiny deep network that simultaneously outputs the bounding box, category and confidence score of detected object through two output channels: localization channel and classification channel. The network is efficient and consumes less resources. As our proposed method decomposes the global regression task into several easier local regression tasks, which can be handled effectively without the support of the complex or large-scale networks. Thus, our proposed approach can work efficiently and effectively by combining this tiny a deep network with local regression tasks.

The rest of this paper is organized as follows: In section 2, we introduce the architecture of our generalized Haar filter based deep networks and describe how to design their weights. In section 3, we propose the algorithm of sparse windows generation. Experimental results and corresponding discussion are presented in section 4, while conclusions and future works are given in section 5.

2 Generalized Haar Filter based Deep Networks

2.1 Architecture of the Deep Networks

To solve the problem of object detection, we construct a deep network that simultaneously outputs the bounding box, category and confidence score of detected object via two output channels: localization channel and classification channel (see Fig. 1). The location channel focuses on bounding box regression and outputs a 4-dimensional vector (dx_1, dx_2, dy_1, dy_2) which is used to describe a bounding box (see Fig. 2). The classification channel outputs the category and confidence score which are denoted by a 2-dimensional vector (l, s).

As shown in Fig. 1, our deep network consists of 11 convolution layers, 4 max-pooling layers and 1 softmax layer. To reduce the memory consumption, several low-level features are shared by localization and classification channels via conv1~pool4. As these two channels aim at different tasks, each of them has 4 independent layers to focus on different problems.

Our deep network is less sensitive to the scaling and shift of the input object owing to regression based localization channel. Consequently, instead of generating region proposals, objects can be detected and located efficiently based on sparse sliding-window paradigm and perspective geometry. More details on this issue are introduced in section 3.1 and related experiments are presented in section 4.

2.2 Generalized Haar Filter based Weights Design

Despite the fact that deep neural networks have achieved state-of-the-art performance in object detection, they consume considerable storage, computing resources and power [16]. Consequently, they are often unsuitable for power and memory constrained devices such

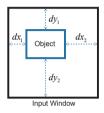


Fig. 2: Description of a bounding box using a 4-dimensional vector.

as vehicle-mounted embedded systems and mobile devices. To overcome this limitation, we introduce a novel convolution weights design method which is based on generalized Haar filter.

Haar-like filters have been successfully applied in object detection owing to their strong representation and high efficiency [4,17]. Instead of using a fixed number of rectangles and configuration types, generalized Haar filters are based on arbitrary configurations and number of rectangles [18]. In this work, the weights of our deep network are constrained based on generalized Haar filters.

Unlike original generalized Haar filters with arbitrary configuration and number of rectangles [18], in our work, the configuration and number of rectangles are obtained in a data-driven way. For a weight w_i of size $m \times m$ ($m \ge 3$), we constrain it to the following form:

$$w_i = \hat{w}_p \cdot k_i, \tag{1}$$

where $k_i \in R$ is a multiplication factor and \hat{w}_p is the p-th generalized Haar filter in the Haar filters space. As we know, there are 2^{m^2} types of configurations for the generalized Haar filters of size $m \times m$. In our case, \hat{w}_p and its negative form $-\hat{w}_p$ can be regarded as a same filter. Thus, the Haar filter space of size $m \times m$ contains 2^{m^2-1} filters. That is, $p \in [1, 2^{m^2-1}]$.

For a trained deep network for vehicle detection, Fig. (4) illustrates the filter usage in the 3×3 Haar filter space. As shown in Fig. (4), filter usage is quite "sparse". In other words, a few generalized Haar filters tend to be used much more frequent than the rest of filters in 3×3 Haar filter space. Thus, these filters are more important and representative than the other filters. According to this fact, we try to reserve these "important" generalized Haar filters and remove the rest of filters from Haar filter space. In this way, the configuration and number of rectangles of generalized Haar filters are obtained in a data-driven way. Let Nr denote the number of filter we try to reserve. We sort the filters according to their usage count and select the top Nr filters (Nr=32 in our work). These selected filters are shown in Fig. (3). Then, the deep network is retrained by constraining the corresponding weights to these Nr filters with multiplication factors.

Generalized Haar filters are able to compress the deep networks by locking the relationship between each element in the weights larger than 3×3 . In this way, considerable storage resources can be saved. For each weight of our deep network, we only need to store a multiplication factor and a filter index. In this work, each multiplication factor (single-precision float-point) takes 4 bytes and each filter index (ranging from 1 to 32) takes less than 1 byte. Thus, only 5 bytes are needed for each weight which is larger



Fig. 3: The illustration of the selected generalized Haar filters.

than 3×3 . By contrast, in traditional deep networks, $4m^2$ bytes are consumed for each weight of size $m\times m$. Besides, despite the fact that weights are constrained according to Nr generalized Haar filters, each element of the weights is still able to keep a relatively high precision owing to the multiplication factors.

Unlike approaches in [17,4,18] where Haar filters are operated based on integral images. We calculate Haar filter based convolution by basic addition and multiplication operation. One reason is that most weights in deep networks are relatively small (3×3 and 5×5), accordingly, there is no obvious benefit in using integral images. Moreover, integral image of each channel is needed to be re-computed for each layer, which is inefficient. Besides, the configuration types and number of rectangles of the generalized Haar filter are diverse due to our data-driven filter selection strategy, consequently, lots of different computation rules are needed to be designed for these selected filters if calculated based on integral images. In this way, it difficult to achieve efficient batch calculation.

In our work, generalized Haar filters only contains two types of elements: -1 and +1. Thus, each generalized Haar filter can be regarded as a sign pattern matrix. Therefore, the weight w_i can be written as:

$$w_i = \hat{w}_p \cdot k_i = \operatorname{sign}(\hat{w}_p) \cdot k. \tag{2}$$

As we know, each convolution step can be regarded as a dot product operation. Let P_i denote the dot product between the Haar filter \hat{w}_p and the input patch x_i , that is:

$$P_i = \hat{w}_p \cdot x_i = \text{sign}(\hat{w}_p) \cdot x_i. \tag{3}$$

As a Haar filter can be regarded as a sign pattern matrix, equation (3) can be calculated by lookup table without multiplication. Then each convolution step can be transformed to the following form:

$$w_i \cdot x_i = (\hat{w}_p \cdot x_i) \cdot k_i = P_i \cdot k_i = k_i \cdot \sum P_i, \tag{4}$$

where " \sum " denotes the sum of the matrix elements. In this way, only one multiplication is needed for each convolution step. By contrast, in traditional deep networks, m^2 multiplication are needed. In a work, our approach consume much less computing resources,

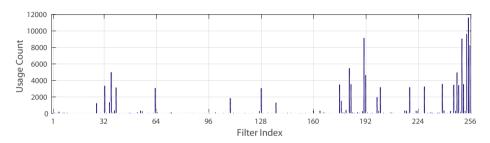


Fig. 4: Filter usage in 3×3 Haar filter space.

in this way, power consumption can be also reduced. Accordingly, our deep network is suitable for embedded devices and FPGA where power and available multipliers are limited.

In addition, the constraint of generalized Haar filter provides a form of regularization which is able to improve the generalization ability of the deep network. The loss function and the corresponding regularization term are introduced in the next subsection.

2.3 Multi-Task Training

Our deep network has two output channels: localization channel and classification channel. Each localization channel outputs a 4-dimensional vector $d=(dx_1,dx_2,dy_1,dy_2)$. Given the ground truth location vector $\hat{d}=(d\hat{x}_1,d\hat{x}_2,d\hat{y}_1,d\hat{y}_2)$, the localization loss can be defined as a squared loss form:

$$L_{loc}(d,\hat{d}) = \left\| d - \hat{d} \right\|^2. \tag{5}$$

As classification channel focuses on a typical classification problem, we define the classification loss L_{cla} as a traditional softmax-loss form.

For a weight w_i $(m \times m, m \ge 3)$ in our deep network, the training goals is to obtain a Haar filter index p and a multiplication factor k_i . Each w_i is constrained according to least squares principle, that is:

$$\min \left\| (w_i - \hat{w}_r \cdot \lambda_r)^2 \right\|_1, \quad r = 1, 2, \dots, 2^{m^2 - 1}.$$
 (6)

For each Haar filter in the Haar filter space, the corresponding multiplication factor λ_r can be obtained by:

$$\frac{d}{d\lambda_r} \left\| \left(w_i - \hat{w}_r \cdot \lambda_r \right)^2 \right\|_1 = 0. \tag{7}$$

From equation (7), λ_r can be calculated:

$$\lambda_r = \frac{\sum w_i \cdot \hat{w}_r}{\sum \hat{w}_r^2}.$$
 (8)

Having obtained λ_r , we constrain each weight by adding a regularization term to its original loss function $Co_i(w_i, w_i^*)$:

$$C_i(w_i, w_i^*) = Co_i(w_i, w_i^*) + \varphi \cdot \min_r \left\| (w_i - \hat{w}_r \cdot \lambda_r)^2 \right\|_1.$$
 (9)

This regularization term can not only reduce the storage consumption of our deep network, but also improve the generalization ability of the deep network. Further experiments are presented in section 4.

As the deep network are trained via stochastic gradient descent (SGD), the loss function $C_i(w_i, w_i^*)$ is needed to be transformed to a differentiable form. According to [19], a maximum function for $Z = \{z_1, \ldots, z_n\}$ can be smoothly approximated as:

$$F_{\max}(Z) = \frac{\sum_{r=1}^{n} z_r^{q+1}}{\sum_{r=1}^{n} z_r^q},$$
(10)

where $p \ge 1$, $\sum z_r^q \ne 0$ and $z_i \ge 0$. The continuity and differentiability have been proofed by [19]. Inspired by this work, we first define the intermediate variable z_r as:

$$z_{r} = \exp(-\left\| (w_{i} - \hat{w}_{r} \cdot \lambda)^{2} \right\|_{1})$$

$$= \exp(-\left\| (w_{i} - \hat{w}_{r} \cdot \frac{\sum w_{i} \cdot \hat{w}_{r}}{\sum \hat{w}_{r}^{2}})^{2} \right\|_{1}).$$
(11)

Then the loss function of each weight can be transformed to a differentiable form:

$$C_i(w_i, w_i^*) = Co_i(w_i, w_i^*) - \varphi \cdot \ln \left(\frac{\sum_{r=1}^{2^{m^2 - 1}} z_r^{q+1}}{\sum_{r=1}^{2^{m^2 - 1}} z_r^q} \right).$$
(12)

Having obtained the differentiable form of the loss function of each weight, the updated weight w_i^{t+1} can be obtained using stochastic gradient descent (SGD). Then, the Haar filter index p and a multiplication factor k_i can be obtained by:

$$\begin{cases}
p = \arg\min_{r} \left\| \left(w_i^{t+1} - \hat{w}_r \cdot \frac{\sum w_i^{t+1} \cdot \hat{w}_r}{\sum \hat{w}_r^2} \right)^2 \right\|_{1} \\
k_i = \frac{\sum w_i^{t+1} \cdot \hat{w}_p}{\sum \hat{w}_p^2}
\end{cases}$$
(13)

The procedure for training our deep network is demonstrated in algorithm 1. For the i-th weight of our deep network, we first construct the current weight using the filter index and multiplication factor which are updated in previous iteration. Then, the target w_i^* are obtained via a standard forward propagation approach. After that, based on

the loss function $C_i(w_i^t, w_i^*)$ (see equation (12)), the gradients can be computed using a standard backward propagation approach. Having obtained the gradients, weight is updated using stochastic gradient descent (SGD). Lastly, using the updated weight, filter index and multiplication factor can be updated via equation (13).

Algorithm 1 Parameter update of the *i*-th weight.

Input: A minibatch of inputs x_i ;

Loss function $C_i(w_i^t, w_i^*)$;

Generalized Haar filter space \hat{w}_r ,

 $(r=1,2,\ldots,2^{m^2-1});$

Current filter index \tilde{p} ;

Current multiplication factor \tilde{k}_i ;

Current learning rate l_t .

Output: Updated filter index p;

Updated multiplication factor k_i ;

Updated learning rate l_{t+1} .

Procedure:

1: Constructing current weight:

$$w_i^t = w_{\tilde{p}} \cdot \tilde{k}_i;$$

2: Obtaining target using standard forward propagation:

$$w_i^* = \text{Forward}(w_i^t, x_i);$$

3: Gradients are computed via standard backward propagation: $\frac{\partial C_i}{\partial w^t} = \text{Backward}(w_i^t, w_i^*);$

$$\partial w_i^v = 0$$
 stic gradient descent:

4: Updating weight using stochastic gradient descent:

$$w_i^{t+1} = \text{UpdateWeight}(w_i^t, \frac{\partial C_i}{\partial w_i}, l_t);$$

5: Updating filter index:

$$p = rg \min_{r} \left\| \left(w_i^{t+1} - \hat{w}_r \cdot rac{\sum w_i^{t+1} \cdot \hat{w}_r}{\sum \hat{w}_r^2}
ight)^2
ight\|_1;$$
tion factor:

6: Updating multiplication factor:

$$k_i = \frac{\sum w_i^{t+1} \cdot \hat{w}_p}{\sum \hat{w}_p^2};$$

7: Updating learning rate:

$$l_{t+1} = \text{UpdateLearningRate}(l_t, t).$$

3 Sparse Windows Generation

Due to the regression based localization channel, our deep network is less sensitive to the scaling and shift of the input object. Consequently, instead of traditional dense sliding-window paradigm, we employ sparse sliding-window strategy to achieve real-time object detection in traffic scene. Besides, in most object detection systems for traffic scene, camera is mounted on a fixed position (e.g. mounted on the top of a vehicle windshield). Accordingly, objects (e.g. vehicles and pedestrians) produce the predetermined location-specific patterns in images. Thus, the potential appearance of objects in the images can be obtained using the perspective geometry of the given scene. Based on

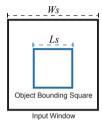


Fig. 5: Illustration of object size ratio.

the perspective geometry, a series of sparse windows can be generated according to the scale and shift tolerance of our deep network. Finally, the locations and categories of objects in the given image can be obtained efficiently by performing our deep network on these sparse windows.

3.1 Sparse Sliding-Window Strategy

We first generate a set of input windows U_s using sparse sliding-window strategy. In order to obtain U_s which is able to fit the objects in different size, we define a size ratio for each input window. As shown in Fig. 5, let Ws denote the size of each input window and Ls denote the size of each object bounding square. The size ratio between object bounding square and the corresponding input window is represented by Rs:

$$Rs = \frac{Ls}{Ws}. (14)$$

We then generate image pyramid by resizing the given image to different scales. For each resized image in the image pyramid, we assume that each deep network is responsible for objects with Rs ranging from $0.5 \sim 0.7$ (objects with Rs beyond this range would be detected in other image scale in the image pyramid). By setting the stride of sliding-windows to 0.3, it can be ensure that each object is completely contained in at least one window. Thus, our approach is able to combine the benefits of regression based method and sliding-window based method. Instead of applying global regression on whole image, our approach performs local regression for bounding box localization. Consequently, the approach has the potential to detected smaller objects compared with global regression based method [14,15]. Besides, the approach maintains a relatively high efficiency thanks to the sparse sliding-window paradigm. Further experiments are presented in section 4.

3.2 Perspective Geometry

Having obtained the set of input windows U_s by sparse sliding-window strategy. We then generate a set of input windows U_p according to the perspective geometry. Let (x_{3D},y_{3D},z_{3D}) denote a 3D point position in world coordinates, (x_{2D},y_{2D}) denotes

the corresponding 2D point position in pixel coordinates. According to the perspective geometry, they satisfy:

$$z \begin{bmatrix} x_{2D} \\ y_{2D} \\ 1 \end{bmatrix} = M \begin{bmatrix} x_{3D} \\ y_{3D} \\ z_{3D} \\ 1 \end{bmatrix}, \tag{15}$$

where M is a camera projection matrix which is the product of an intrinsic matrix and an extrinsic matrix:

 $M = \begin{bmatrix} f_x & 0 & u_0 & 0 \\ 0 & f_y & v_0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mathbf{R} & \mathbf{T} \\ \mathbf{0^T} & 1 \end{bmatrix}.$ (16)

In equation (16), R is a rotation matrix which is the result of three rotations around the world coordinate axes. As in most object detection systems for traffic scene, the angles of camera rotations are relatively small. Accordingly, the camera projection matrix M can be approximated as:

$$M = \begin{bmatrix} m_{11} & 0 & m_{13} & m_{14} \\ 0 & m_{22} & m_{23} & m_{24} \\ 0 & 0 & 1 & m_{34} \end{bmatrix}. \tag{17}$$

Then, the equation (15) can be transformed to:

$$\begin{cases}
 x_{2D} = \frac{m_{11} \cdot x_{3D} + m_{13} \cdot z_{3D} + m_{14}}{z_{3D} + m_{34}} \\
 y_{2D} = \frac{m_{22} \cdot y_{3D} + m_{23} \cdot z_{3D} + m_{24}}{z_{3D} + m_{34}}
\end{cases}$$
(18)

In this work, each input window of our deep network has the same height and width. We use d_{2D} to represent the height or width of the input window located at (x_{2D},y_{2D}) in pixel coordinates (center aligned), and let d_{3D} denote the corresponding height or width in world coordinates. Then d_{2D} can be formulated as:

$$d_{2D} = \frac{m_{11} \cdot d_{3D}}{z_{3D} + m_{3A}}. (19)$$

As we know, the locations of most objects in traffic scene are limited (e.g. vehicles and pedestrians will not appear on sky region). In this work, the location and size of each input window of our deep network is described by a triplet (x_{2D}, y_{2D}, d_{2D}) in pixel coordinates. Let $[x_{3D}^{\min}, x_{3D}^{\max}]$ and $[y_{3D}^{\min}, y_{3D}^{\max}]$ respectively represent the location ranges of the corresponding windows in x_{3D} and y_{3D} axes of the world coordinates. By solving the equation set that consists of equation (18) and (19), we can find 4 boundary planes in $x_{2D}y_{2D}d_{2D}$ space which jointly limit the locations and sizes of input windows:

$$\begin{cases}
d_{2D} = x_{2D} \cdot g_x(x_{3D}^{\min}) - g_x(x_{3D}^{\min}) \cdot m_{13} \\
d_{2D} = y_{2D} \cdot g_y(y_{3D}^{\min}) - g_y(y_{3D}^{\min}) \cdot m_{23} \\
d_{2D} = x_{2D} \cdot g_x(x_{3D}^{\max}) - g_x(x_{3D}^{\max}) \cdot m_{13}
\end{cases},$$

$$d_{2D} = y_{2D} \cdot g_y(y_{3D}^{\max}) - g_y(y_{3D}^{\max}) \cdot m_{23}$$
(20)

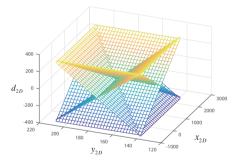


Fig. 6: Illustration of 4 boundary planes that jointly limit the locations and sizes of input windows in $x_{2D}y_{2D}d_{2D}$ space.

where:

$$\begin{cases}
g_x(x_{3D}) = \frac{m_{11} \cdot d_{3D}}{m_{11} \cdot x_{3D} - m_{13} \cdot m_{34} + m_{14}} \\
g_y(y_{3D}) = \frac{m_{11} \cdot d_{3D}}{m_{22} \cdot y_{3D} - m_{23} \cdot m_{34} + m_{24}}
\end{cases}$$
(21)

As shown in Fig. 6, in the $x_{2D}y_{2D}d_{2D}$ space, possible input windows are distributed in the region of inverted pyramid which is enclosed by 4 boundary planes. We use U_p to represent the set of input windows in this region, then the final sparse windows U_f can be obtained by:

$$U_f = U_p \cap U_s. \tag{22}$$

As shown in Fig. 7, we illustrate the final sparse windows in each image of the image pyramid using the TME Motorway dataset [20], which is a challenging and widely-used dataset in vehicle detection. Note that, all of the sparse windows have the same size $(48 \times 48 \text{ in this work})$, which means that they can be directly utilized as the input of the deep network. In this way, the sparse windows are avoided to be resized and the efficiency of the approach is ensured.

4 Experimental Results

The experiments in this section mainly focus on vehicle detections in traffic scene, which is an important issue in a wide range of applications such as autonomous driving, autonomous navigation and advance driver assistance systems (ADAS) [21]. In this work, the performance of our approach is evaluated on a broadly used datasets: TME motorway dataset [20], which is designed for vehicle detection and localization in challenging traffic scene with various lighting conditions and complex traffic situations [20].

In the following experiments, the output bounding-boxes of our approach are refined by mean shift and non-maximum suppression, and we use an intersection-over-union (IoU) threshold of 0.7 to determine the correctness of detection. All the experiments in this section are performed on a GTX1080 GPU.

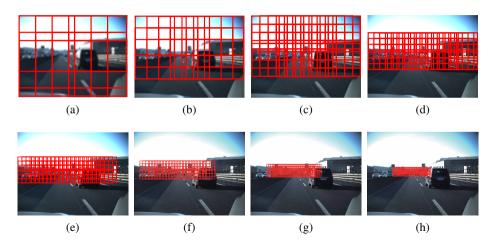


Fig. 7: Final sparse windows in each image of the image pyramid.

4.1 Generalized Haar filter based weights vs traditional weights

As our deep network has two output channels: classification channel and localization channel. We respectively use the classification and localization error of training and test as the performance metrics for the following experiments. The classification error of both training and test is defined as:

$$Er_{cla} = \frac{FP + FN}{N},\tag{23}$$

where N is the total number of training or test samples, FP and FN respectively denotes the number of positive and negative samples which are incorrectly classified.

For the localization channel, the error of both training and test is defined as:

$$Er_{loc} = \frac{\sum_{i=1}^{N} \left\| d_i - \hat{d}_i \right\|^2}{4N},$$
 (24)

where the vector $d=(dx_1,dx_2,dy_1,dy_2)$ represents the output of the localization channel, and $\hat{d}=(d\hat{x}_1,d\hat{x}_2,d\hat{y}_1,d\hat{y}_2)$ is the ground truth location vector.

Fig. 8 illustrates how the errors of classification and localization changed when increasing the number of training iterations. According to this figure, our deep network is able to achieve a high performance when training is convergent. Despite the fact that the deep network with generalized Haar filter based weights (hereafter called G-Haar weights) has a slightly larger training error and needs more training iterations to reach the state of convergence, it is able to produce a less test error in both classification and localization tasks when training is convergent. In other words, the deep network with G-Haar weights has a stronger generalization ability than that with traditional weights.

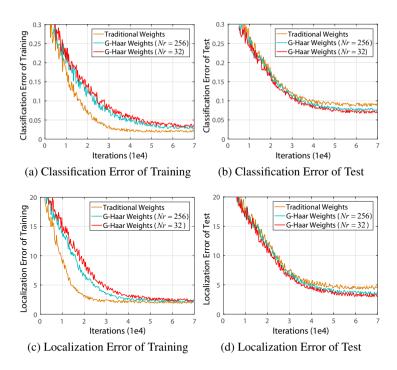


Fig. 8: Comparisons of generalized Haar filter based weights (G-Haar weights) and traditional weights on TME motorway dataset.

This is owing to the regularization effect of G-Haar weights. Besides, the network employing G-Haar weights with less Nr (the number of selected filters) needs more training iterations and tend to have stronger generalization ability when training iteration is convergent.

4.2 Storage and computing resources consumptions

We investigate the storage and computing resources consumptions in this subsection. The experiment is also performed on TME motorway dataset [20]. In the experiment, the number of selected filters Nr is 32, and all the weights are stored using single-precision floating-point format (32 bits). In the deep network (see Fig. 1), the size of $conv1\sim conv5_x$ is 3×3 , and the rest of convolution kernels have the size of 1×1 . As the weights of size 1×1 consume much less resources than that of size 3×3 , we only evaluate the dimensions of $conv1\sim conv5_x$ and their effects on resources consumptions in table 1.

As shown in table 1 (column of Mem.), the networks using G-Haar weights are able to dramatically reduce memory resources (about $0.8n^2$ times, n=3 in this work). This is due to the fact that only a filter index and a multiplication factor are needed to stored for each G-Haar based convolution kernel, which totally consumes 5 bytes of

conv1	conv2	conv3	conv4	conv5_x	G-Haar	Mem.	Mul./St.	Cla. Err.	Loc. Err.
3×64	64×128	128×256	256×256	256×128	×	5.97 MB	9	0.147	5.19
3×64	64×256	256×512	$512{\times}1024$	1024×128	×	32.13 MB	9	0.135	4.78
3×256	256×512	512×512	$512{\times}1024$	1024×512	×	67.74 MB	9	0.092	4.21
3×128	128×256	256×512	512×1024	1024×1024	×	96.04 MB	9	0.054	3.22
3×64	64×128	128×256	256×256	256×128	✓	901.96 KB	1	0.095	4.33
3×64	64×256	$256{\times}512$	$512{\times}1024$	1024×128	\checkmark	4.51 MB	1	0.072	3.48
3×256	256×512	512×512	512×1024	1024×512	\checkmark	9.58 MB	1	0.067	2.97
3×128	128×256	256×512	512×1024	1024×1024	\checkmark	13.68 MB	1	0.042	2.32

Table 1: Storage and computing resources consumptions of our deep network.

memory space. By contrast, there are n^2 weights needed to stored for each traditional convolution kernel of $n \times n$, thus, $4n^2$ of memory space is required.

Besides, as multipliers are major computing resource consumed by deep networks, we use the number of multiply operations required for each convolution step to measure the consumption of computing resource for the deep networks. As we know, each traditional convolution kernel of $n \times n$ needs n^2 multiply operations for each convolution step. Nevertheless, for any convolution kernel using G-Haar weights, each convolution step can be transformed to the form of equation (4). In this way, only one multiply operation is required for each G-Haar convolution step. Accordingly, as shown in table 1 (column of Mul./St.), the computing resource consumption of the deep networks using G-Haar weights is only $1/n^2$ of that using traditional weights (n=3 in this experiment).

In addition, power consumption is directly influenced by storage and computing resources utilizations [22]. As our deep network with G-Haar weights can markedly reduce storage and computing resources (including memory accesses), it would have a great increase in power-efficiency. This merit is especially meaningful when the deep network implements on embedded systems or mobile devices (such as FPGA and ARM), which are quite sensitive to power consumption.

4.3 Comparing with state-of-the-art methods

For further analysis, our proposed approach is evaluated via comparing with some state-of-the-art methods. In the following experiments, "Ours" represents our complete method that employs G-Haar weights and the stage of sparse windows generation. "Ours (Traditional Weights)" and "Ours (Without Sparse Windows)" denote the variants of our approach that utilize traditional weights instead of G-Haar weights and without the stage of sparse windows generation, respectively. The experiment settings of our approaches is the same as that in the row 7 of table 1.

As shown in Fig. 9, thanks to the strong representation of deep convolution neural networks, regression based deep networks (such as YOLO [14], SSD500 [15] and our approach) tend to perform better than traditional hand-craft methods. At the same time, as exhaustive sliding windows can be avoided when using regression based deep

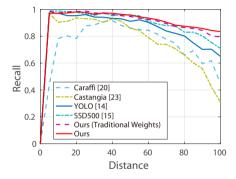


Fig. 9: Performance evaluation in different vehicle distances on TME motorway dataset.

networks and convolution can be efficiently computed in GPU, as shown in table 2, regression based deep networks are able to achieve real-time performance in object detection task in traffic scene. Besides, table 2 also demonstrates that the stage of sparse windows generation is able to reduce unnecessary computation and improve the efficiency dramatically.

Moreover, quantitative evaluation in Fig. 9 and qualitative results in Fig. 10 indicate that our proposed method is able to achieve better performance on small object detection compared with other regression based deep networks such as YOLO [14] and SSD500 [15] which utilize global regression strategy. SSD500 [15] employs multi-scale feature maps to detect objects in different scales and it perform better than YOLO [14]. However, its performance on small objects detection is still unsatisfactory due to the fact that global regression, that is regressing the bounding-box of each object from the whole image, is a much more difficult task compared with our local regression strategy. Beside, input images (e.g. 1024×768 in TME motorway dataset [20]) have to be resized to 500×500 in SSD500 [15], which leads to the lacking of resolution of small objects such as vehicles and pedestrians that are far away from camera. Our proposed method decomposes the global regression task into several easier local regression tasks, and detect multi-scale objects from image pyramid. In this way, the resolution of small objects is ensured and thus better performance can be achieve.

Table 2: Runtime analysis on TME motorway dataset.

Method	Average Runtime
Caraffi [20]	0.1s
Castangia [23]	0.05s
YOLO [14]	0.022
SSD500 [15]	0.043
Ours (Without Sparse Windows)	32s
Ours	0.025



Fig. 10: Qualitative results on TME motorway dataset in comparison with two state-of-the-art approaches which are based on deep regression networks. The rows of "L", "M" and "S" respectively mean that the average object sizes are large, middle and small.

5 Conclusions and Future Works

In this paper, we presents a novel network system for object detection tasks in traffic scene. In this system, we introduce a local regression strategy for accurate objects detection. Compared with traditional global regression based object detection, the local regression task is easier to handled without the support of the complex or large-scale networks. According to this fact, we handle the local regression tasks by using several tiny deep networks which simultaneously output the bounding boxes, categories and confidence scores of detected objects. In order to satisfy the storage, power and computing source requirements of the platforms for traffic scene applications, the weights of the deep networks are constrained to the form of generalized Haar filter in training phase. Furthermore, to achieve the real-time performance, we introduce the strategy of sparse windows generation to reduce the runtime of our system.

In the experiments, we first evaluate the performance and generalization ability of our generalized Haar filter based weights by comparison with traditional weights. Then the consumptions of storage and computing resources are evaluated. Finally, the effectiveness and efficiency of our approach are validated in comparison with some recently published state-of-the-art methods. Experimental results demonstrate that the proposed approach is proved to be efficient, robust and source saving in challenging traffic scene.

As the proposed approach is suitable for the platforms which have strict limitations on memory, power and computing sources, our future works will focus on implementing the proposed approach (path forward) to FPGA (Field Programmable Gate Array) for advance driver assistance system (ADAS). Owing to the generalized Haar Filter based weights, only one multiply operation is required for each convolution step. This make it possible to construct more parallel pipelines in FPGA where the number of multipliers is limited. Moreover, each local regression task is an independent computation. Thus, all of the local regression tasks can run in parallel in FPGA. In this way, the system can achieve real-time response without extra effort.

Acknowledgements

This work is supported by the National Natural Science Foundation of China (Grant No. 61473303).

References

- K. Lu, J. Li, X. An, H. He, Vision sensor-based road detection for field robot navigation, Sensors 15(11) (2015) 29594–29617.
- Large Scale Visual Recognition Challenge (ILSVRC), http://www.image-net.org/challenges/LSVRC/, [Online; accessed Oct. 23, 2016].
- MS COCO Visual Recognition Challenges, http://mscoco.org/, [Online; accessed Oct. 23, 2016].
- P. Viola, M. J. Jones, Robust real-time face detection, International Journal of Computer Vision (IJCV) 57(2) (2004) 137–154.
- P. F. Felzenszwalb, R. B. Girshick, D. McAllester, D. Ramanan, Object detection with discriminatively trained part-based models, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) 32(9) (2010) 1627–1645.
- 6. A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, Advances in Neural Information Processing Systems 25 (2012) 1–9.
- T. Wang, D. J. Wu, A. Coates, A. Y. Ng, End-to-end text recognition with convolutional neural networks, in: International Conference on Pattern Recognition (ICPR), 2012, pp. 1051

 4651.
- 8. K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, arXiv preprint arXiv:1512.03385.
- R. Girshick, Fast r-cnn, in: IEEE International Conference on Computer Vision (ICCV), 2015, pp. 1440–1448.
- R. Girshick, J. Donahue, T. Darrell, J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, in: IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2014, pp. 580–587.
- K. He, X. Zhang, S. Ren, Spatial pyramid pooling in deep convolutional networks for visual recognition, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) 37(9) (2014) 1904–1916.
- 12. J. Uijlings, K. van de Sande, T. Gevers, A. Smeulders, Selective search for object recognition, in: International Journal of Computer Vision (IJCV), 2013, pp. 154–171.
- S. Ren, K. He, R. Girshick, J. Sun, Faster r-cnn: Towards real-time object detection with region proposal networks, IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI).

- 14. J. Redmon, S. Divvala, R. Girshick, A. Farhadi, You only look once: Unified, real-time object detection, arXiv preprint 1506.02640.
- 15. W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, A. C. Berg, Ssd: Single shot multibox detector, arXiv preprint 1512.02325.
- 16. S. Han, H. Mao, W. J. Dally, Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, Fiber 56(4) (2016) 3–7.
- 17. K. Zhang, L. Zhang, M.-H. Yang, Real-time compressive tracking, in: European Conference on Computer Vision (ECCV), 2012, pp. 864–877.
- P. Dollár, Z. Tu, H. Tao, S. Belongie, Visual tracking with online multiple instance learning, in: IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2009, pp. 983–990.
- 19. Z. Yu, An efficient method for solving nonlinear unconstrained min-max problem, Master's thesis. Xi'An University of science and technology, Xi'An, China (2011).
- C. Caraffi, T. Vojíř, J. Trefný, J. Šochman, J. Matas, A system for real-time detection and tracking of vehicles from a single car-mounted camera, in: IEEE Intelligent Transportation Systems Conference, 2012, pp. 975–982.
- 21. Z. Sun, G. Bebis, R. Miller, Monocular precrash vehicle detection: features and classifiers, IEEE Transactions on Image Processing (TIP) 32(9) (2006) 2019–2034.
- 22. M. Horowitz, Computing's energy problem (and what we can do about it), in: IEEE International Solid State Circuits Conference, 2014, pp. 10–14.
- L. Castangia, P. Grisleri, P. Medici, A. Prioletti, A. Signifredi, A coarse-to-fine vehicle detector running in real-time, in: IEEE Intelligent Transportation Systems Conference, 2014, pp. 691–696.