

# Vehicle Detection in Satellite Images by Parallel Deep convolutional Neural Networks

Xueyun Chen, Shiming Xiang, Cheng-Lin Liu, and Chun-Hong Pan.

National Laboratory of Pattern Recognition

Institute of Automation, Chinese academy of Sciences

Beijing 100190, China

Email: {xueyun.chen, smxiang, liucl, chpang}@nlpr.ia.ac.cn

**Abstract**—Deep convolutional Neural Networks (DNN) is the state-of-the-art machine learning method. It has been used in many recognition tasks including handwritten digits, Chinese words and traffic signs, etc. However, training and test DNN are time-consuming tasks. In practical vehicle detection application, both speed and accuracy are required. So increasing the speeds of DNN while keeping its high accuracy has significant meaning for many recognition and detection applications. We introduce parallel branches into the DNN. The maps of the layers of DNN are divided into several parallel branches, each branch has the same number of maps. There are not direct connections between different branches. Our parallel DNN (PNN) keeps the same structure and dimensions of the DNN, reducing the total number of connections between maps. The more number of branches we divide, the more swift the speed of the PNN is, the conventional DNN becomes a special form of PNN which has only one branch. Experiments on large vehicle database showed that the detection accuracy of PNN dropped slightly with the speed increasing. Even the fastest PNN (10 times faster than DNN), whose branch has only two maps, fully outperformed the traditional methods based on features (such as HOG, LBP). In fact, PNN provides a good solution way for compromising the speed and accuracy requirements in many applications.

**Keywords**—Remote Sensing; Object detection; Deep convolutional Neural Networks;

## I. INTRODUCTION

Detecting vehicle in high-resolution satellite images is a highly challenging task. Hidden in a tree, sheltered by a building or jammed together in the streets or parks in very close distances, vehicle detection always arouse human interests. Many works have been done, various features and methods have been used [1-8].

Hinz [3] built a hierarchical 3D-model to describe the prominent geometric features of the cars, Chen et al. [7] segmented road by its straight line contours and used SVM to detect vehicle in road. Liang et al. [8] combined the HOG descriptors and the selected Haar features and used Multiple Kernel Learning method to detect vehicle in wide area motion imagery. Zhao et al. [1] showed the boundary of the car body, the boundary of the front windshield, and the shadow were the key features for vehicle detection. Ali et al. [2] detected vehicle by adaboost method based on pose-indexed features and pose estimators. Grabner et al. [6] used boosting method based on Haar wavelets, HOG and LBP. They segmented the image into streets, buildings, trees, etc, discarding vehicle detections that

are not present on the streets. Kembhavi et al.[4] detected vehicles in San Francisco using the multi-scales HOG features computed on the color probability maps. They showed HOG outperformed SIFT. SIFT, HOG and LBP are the most popular features in object detection or image classification. SIFT is very similar to HOG, LBP is more suitable for texture feature such as face recognition [20].

Recently, machine learning intelligence has been reported enable to match human performance on recognition task of handwritten digits and traffic signs [9]. The machine is based on deep convolutional neural networks (DNN). Convolutional neural networks originates from the study on cats striate cortex by Hubel and Wiesel [12]. They first proposed the concept of receptive field. Fukushi [13] proposed Neocognitron, a hierarchy model of the multi-layer neural networks. He realized the concept of receptive field. LeCun [14-15] gave the normal form of convolutional neural networks (CNN)(LeNet-5, LeNet7), which has been used in many recognition tasks, include digit recognition, face detection. Rowley [17] detected face using a simple CNN which only had three layers but with three types receptive fields in the first layer. Garcia [16] realized a LeNet-5 structure CNN for face detection. They showed CNN outperformed Adaboost method [21] in CMU and MIT test sets obviously. Compared with CNN, DNN is more deep (6-10 layers) and wide (40-250 maps per layer). Training and test DNN are time-consuming works, the former often needs 1-2 days even with the help of a GPU card. For practical application like object detection in satellite images, speed and accuracy are the same important, we seek to simply DNN, increasing its speed while keep its accuracy in a high level.

Reminded by the LeNet-5 structure, where the 6 maps of the first max-pooling layer can supply  $2^6 - 1 = 127$  different map-connection manners to the higher layer theoretically. LeCun only used 15 manners of them. This implies that selection of the the map-connection manners are variable. In DNN, all the maps in the above convolutional layer has the full map-connections to the lower layer. This produced the maximal number of map-connections, and this is why that DNN is much slow than the traditional CNN. We simplify DNN by introducing the concept of branches, we divided all maps into several parallel branches, each branch has the same number of maps and the map-connections only occur in the internal maps of the branch.

Ciresan et al. [9-11] trained different DNNs by different preprocessed images, using the average of the outputs of all

DNNs. In most situations, such average only increases the accuracy by 0.5 – 1.0% [9], while adding the expense of time-consuming of a lot of times. We used four preprocessed images: gray, gradient, gradient in the image thresholding at 120, gradient in the negative image thresholding at 60. Instead of training four DNNs, we divide the maps of the first convolutional layer into four parts, allocating the same number of maps to the four images respectively, training the same PNN on the four images synchronously. Experiments show that our method increased the accuracy obviously, without more expense of time-consuming. PNN also outperformed the vehicle detection rate (about 70%) of [4] in the same challenging environments of San Francisco city.

## II. ARCHITECTURE OF PNN

The layers of PNN can be denoted as: (*Input*,  $C^1$ ,  $M^1$ , ...,  $C^{n_l}$ ,  $M^{n_l}$ ,  $H^1$ , ...,  $H^{n_h}$ , *Output*). Where  $C^i$  and  $M^i$  ( $i = 1, \dots, n_l$ ) are the convolutional and max-pooling layers respectively, they act as the feature extractor.  $H^l$  ( $l = 1, \dots, n_h$ ) are the hidden layers, they and the output layer compose the final Multi-Layer Perceptron (MLP) classifier. For convenience, suppose all convolutional and max-pooling layers have the same number of maps, denote  $C^i = (C_{j_1}^i, \dots, C_{j_{n_m}}^i)$ ,  $M^i = (M_{j_1}^i, \dots, M_{j_{n_m}}^i)$ , where  $n_m$  is the number of the mapping layer. Now divide all the maps of the convolutional and max-pooling layers into  $n_b$  branches as the following:

$$\left\{ \begin{array}{l} \text{Branch} = (B_1, \dots, B_{n_b}) \\ B_j = \left\{ (C_{j_k}^1, M_{j_k}^1, \dots, C_{j_k}^{n_l}, M_{j_k}^{n_l}) : \right. \\ \left. j = 1, \dots, n_b \right\} \end{array} \right\} \quad (1)$$

$b_w = \frac{n_m}{n_b}$  is the branch width.  $C_{j_k}^l = C_{(j-1)b_w+k}^l$ ,  $M_{j_k}^l = M_{(j-1)b_w+k}^l$ , suppose  $n_m$  can be divided by  $n_b$  without remainder.

Denote  $I, D^l, O^l$  as the definition domains of input, convolutional and max-pooling layers respectively.  $R^1 = [0, 255]$  is the gray range,  $R = [-1, 1]$  is the kernel function range. Then we have:

$$\left\{ \begin{array}{l} \text{Input} : I \rightarrow R^1 \\ C_{j_k}^l : D^l \rightarrow R \\ M_{j_k}^l : O^l \rightarrow R \end{array} \right\} \quad (2)$$

Denote  $flt_{j_k}^1$  as the filter which connects  $C_{j_k}^1$  with *input*,  $flt_{j_k}^l$  as the filter which connects  $C_{j_k}^l$  with  $M_{j_k}^{(l-1)}$ . Use  $\tanh$  as the kernel function,  $b_{j_k}^1, b_{j_k}^l$  are their biases:

$$\left\{ \begin{array}{l} C_{j_k}^1(x, y) = \tanh(b_{j_k}^1 + \sum_{(u,v) \in A^1} (flt_{j_k}^1(u, v) \\ \quad \times \text{Input}(x + u, y + v)) \\ C_{j_k}^l(x, y) = \tanh(b_{j_k}^l + \sum_{k=1}^{b_w} \sum_{(u,v) \in A^{l-1}} \\ \quad (flt_{j_k}^l(u, v) \times M_{j_k}^{(l-1)}(x + u, y + v)) \end{array} \right\} \quad (3)$$

where  $(x, y) \in D^l$ ,  $j = 1, \dots, n_b$ ,  $k = 1, \dots, b_w$ ,  $A^1$  and  $A^l$  are the definition domains of the filters.  $l = 2, \dots, n_l$ .

From (3) we can see that connections only occur between maps of the same branch. There is no connection between

branches. The max-pooling layer can be expressed as the following:

$$\left\{ \begin{array}{l} T_{j_k}^l : O^l \rightarrow D^l \\ T_{j_k}^l(x, y) = \underset{(x', y')}{\operatorname{argmax}} \left\{ \begin{array}{l} C_{j_k}^l(x', y') : \\ mx \leq x' < mx + m, \\ my \leq y' < my + m \end{array} \right\} \\ M_{j_k}^l(x, y) = C_{j_k}^l(T_{j_k}^l(x, y)) \end{array} \right\} \quad (4)$$

$T$  is a map from  $O^l$  to  $D^l$ ,  $m$  is a constant positive integer which determine the size of  $O^l$ , in this paper,  $m=2$ .

We denote  $\text{Dim}(\ast)$  is the number of all nodes in the layer  $\ast$ ,  $\text{Node}(\ast, i)$  as the value of the  $i$ -th node of the layer  $\ast$ ,  $n_c$  is the number of classes. Then we define  $H^l = (h_1^l, \dots, h_{\text{Dim}(H^l)}^l)$ ,  $1 \leq l \leq n_h$ , *Output* = ( $out_1, \dots, out_{n_c}$ ),  $w_{j_k}^l, w_{j_k}^o$  denote the weights of the hidden layers and output layer respectively,  $bia_j^l, bia_j^o$  denote the biases of the hidden layers and output layer respectively. We have:

$$\left\{ \begin{array}{l} h_j^1 = \tanh(bia_j^1 + \sum_{k=1}^{\text{Dim}(M^{n_l})} w_{j_k}^1 \text{Node}(M^{n_l}, k)) \\ h_j^l = \tanh(bia_j^l + \sum_{k=1}^{\text{Dim}(H^{l-1})} w_{j_k}^l h_k^{l-1}) \\ out_j = \tanh(bia_j^o + \sum_{k=1}^{\text{Dim}(H^{n_h})} w_{j_k}^o h_k^{n_h}) \end{array} \right\} \quad (5)$$

## III. TRAINING PNN

We define the training set as:  $\{(Input^{(q)}, Label^{(q)}) : q = 1, \dots, n_{sp}\}$ , where  $q$  is the sample-index,  $n_{sp}$  is the number of all samples in training set. In addition,  $Label^{(q)} = (lab_1^{(q)}, \dots, lab_{n_c}^{(q)})$ ,  $lab_i^{(q)} = 1$ , if  $Input^{(q)}$  belong to  $i$ -th class, otherwise equals to  $-1$ . We denote  $Output^{(q)} = PNN(W, Input^{(q)})$ , where  $W$  is the set of all filters, biases and weights in PNN. Then we have:

$$\left\{ \begin{array}{l} E = \frac{1}{2} \sum_{q=1}^{n_{sp}} \sum_{k=1}^{n_c} (out_k^{(q)} - lab_k^{(q)})^2 \\ W = \underset{W}{\operatorname{argmin}}(E) \end{array} \right\} \quad (6)$$

By the steepest descent Method, we have:

$$\Delta W(t+1) = -\varepsilon \frac{\partial E}{\partial W} + \gamma \Delta W(t) - \beta W(t) \quad (7)$$

Where  $\varepsilon$ = LearnRate,  $\gamma$ = Momentum,  $\beta$ =WeightDecay, LearnRate is a very small value, such as 0.001. Sometimes, Momentum is used to speed convergence, WeightDecay is used to limit the norms of the Weights.

### A. Back Propagation

We compute the error of every layer from output layer to the first convolutional layer by the back propagation algorithm, because  $\frac{d}{dx}(\tanh(x)) = 1 - (\tanh(x))^2$ , the error of the output and hidden layers are:

$$\left\{ \begin{array}{l} \delta out_j^{(q)} = (out_j^{(q)} - lab_j^{(q)})(1 - (out_j^{(q)})^2) \\ \delta h_j^{n_h(q)} = (\sum_{l=1}^{n_c} w_{lj}^{out} \delta out_l^{(q)})(1 - (h_j^{n_h(q)})^2) \\ \delta h_j^{l(q)} = (\sum_{k=1}^{n_c} w_{kj}^l \delta h_k^{(l+1)(q)})(1 - (h_j^{l(q)})^2) \end{array} \right\} \quad (8)$$

where  $1 \leq l < n_h$ . We denote the set  $Set^l(x, y) = \{(x', y', u, v) : x' + u = x, y' + v = y, (x', y') \in C^{l+1}, (u, v) \in A^{l+1}\}$ , then the error of the Max-Pooling layers are:

$$\begin{cases} \delta Node^{(q)}(M^{n_l}, i) = \sum_{k=1}^{Dim(H^1)} w_{kji}^1 \times \delta h_k^{1(q)} \\ \delta M_{jk}^{l(q)}(x, y) = \sum_{p=1}^{b_w} \sum_{(x', y', u, v) \in Set^l(x, y)} \delta c_{jp}^{(l+1)(q)}(x', y') f_{lpk}^{l+1}(u, v) (1 - (M_{jk}^{l(q)}(x, y))^2) \\ 1 \leq l < n_l \end{cases} \quad (9)$$

We define the map  $F_{jk}^{l(q)} : D^l \rightarrow O^l$  as:

$$F_{jk}^{l(q)}(x', y') = \begin{cases} (x, y), & \text{if } T_{jk}^{l(q)}(x, y) = (x', y') \\ (-1, -1), & \text{otherwise} \end{cases} \quad (10)$$

The error of the convolutional layer is:

$$\delta c_{jk}^{l(q)}(x, y) = \begin{cases} 0, & \text{if } F_{jk}^{l(q)}(x, y) = (-1, -1) \\ \delta M_{jk}^{l(q)}(F_{jk}^{l(q)}(x, y)), & \text{otherwise} \end{cases} \quad (11)$$

where  $1 \leq l \leq n_l$ ,  $1 \leq j \leq n_b$ ,  $1 \leq k \leq b_w$ .

### B. Weights Updating

Suppose momentum and WeightDecay are zero, the weights and biases of the output and hidden layers are updated as:

$$\begin{cases} w_{lj}^o(t+1) = w_{lj}^o(t) - \varepsilon \delta out_l^{(q)} h_j^{n_h(q)} \\ bia_l^o(t+1) = bia_l^o(t) - \varepsilon \delta out_l^{(q)} \\ w_{kj}^l(t+1) = w_{kj}^l(t) - \varepsilon \delta h_k^{l(q)} h_j^{(l-1)(q)} \\ bia_k^l(t+1) = bia_k^l(t) - \varepsilon \delta h_k^{l(q)} \\ w_{kj}^1(t+1) = w_{kj}^1(t) - \varepsilon \delta h_k^{1(q)} Node^{(q)}(M^{n_l}, j) \\ bia_k^1(t+1) = bia_k^1(t) - \varepsilon \delta h_k^{1(q)} \end{cases} \quad (12)$$

where  $2 \leq l \leq n_h$ , the filters and biases of the convolutional layers are updated as:

$$\begin{cases} f_{lpk}^l(u, v)(t+1) = f_{lpk}^l(u, v)(t) - \varepsilon \sum_{(x, y) \in D^l} M_{jk}^{(l-1)(q)}(x+u, y+v) \delta c_{jp}^{l(q)}(x, y) \\ b_{jp}^l(t+1) = b_{jp}^l(t) - \varepsilon \sum_{(x, y) \in D^l} \delta c_{jp}^{l(q)}(x, y) \\ f_{lpk}^1(u, v)(t+1) = f_{lpk}^1(u, v)(t) - \varepsilon \sum_{(x, y) \in D^1} Input^{(q)}(x+u, y+v) \delta c_{jp}^{1(q)}(x, y) \\ b_{jp}^1(t+1) = b_{jp}^1(t) - \varepsilon \sum_{(x, y) \in D^1} \delta c_{jp}^{1(q)}(x, y) \end{cases} \quad (13)$$

where  $2 \leq l \leq n_l$ ,  $1 \leq q \leq n_{sp}$ . In practical training, the samples are often input in batch, the batch size is in the range [10,100], and the weights update once after a batch input.

## IV. IMPLEMENTATION DETAIL

Vehicles can parked in any spatial place, they have colorful appearances, method based on gray is not suitable for vehicle locating. We rely on the gradients to locate vehicle. The normal gradient is the maximal norm of the gradients in RGB channel as computed in [18], In order to enhance the borders of the black and white vehicles, we computing gradients on the two thresholding images as shown in Figure 1.

### A. Object Locating

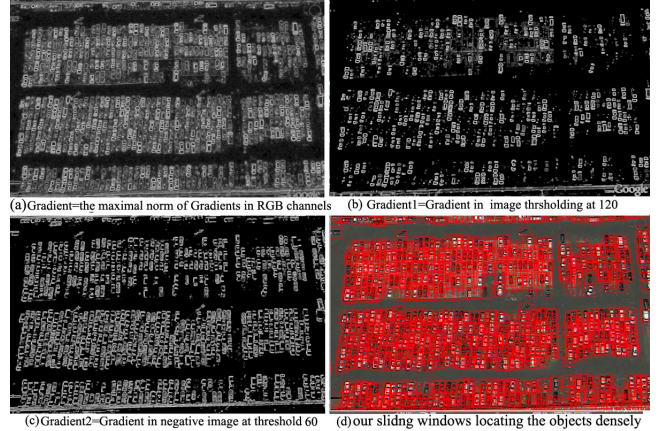


Fig. 1. Examples of locating vehicles in a large park. We compute three type gradients as (a), (b) and (c) respectively, locating the objects in the three images respectively, (d) the set of the location windows, there are 2918 location windows covering 544 vehicles, 99.2% vehicles are located correctly.

In Figure 1 (a), the border of the black vehicles are too dim to be locate correctly. In (c), all the borders of the dim black vehicles are enhanced. In (b), the borders of white vehicles are also enhanced, the backgrounds are flted, this helps to detect white vehicles under the trees. The final location windows is shown in (d). We locate the objects by Algorithm 1:

#### Algorithm 1 Object Locating

**Input:** The three gradient images, sliding window size, sliding step.  
**Output:** All location windows.

- 1: On each gradient image, generate the sliding window grid to cover the whole image.
- 2: For each sliding window  $W_p$  at position  $p = (x_0, y_0)$ , compute the geometric center  $p1 = (x_1, y_1)$  on  $W_p$ , center the  $W_p$  on  $(x_1, y_1)$ , denote it as  $W_{p1}$ .
- 3: Enlarge the size of  $W_{p1}$  twice, compute the new geometric center  $p2$  on the enlarged window. Center  $W_{p1}$  on the new center  $p2$ .
- 4: Output all location windows on the three gradient images.

The sliding window size is  $32 \times 32$ , the sliding step is 16. Some repetitive windows are filtered by a small distance limit (5 pixel). Our database has 63 images,  $1368 \times 972$  size, 6887 vehicles at all. Our method generates 197513 location windows, 3135 windows per image, 99.7% vehicles are located correctly. To achieve the same locating precision, the normal sliding window method needs 10400 sliding windows per image, our methods is more efficient in searching.

In order to get rotation-invariant and scale-invariant neural networks, we rotate every location window 11 times by:  $0^0$ ,  $4.5^0$ ,  $9^0$ ,  $\dots$ ,  $45^0$ , then shrink or enlarge the non-rotating images into multi-scalings: 0.8, 0.9, 1.0, 1.1, 1.2, 1.3. For each location window, we get four preprocessed images: Gray, Gradient, Gradient1 and Gradient2 (see the (a), (b), (c) of Fig. 1), these preprocessed images are normalized into  $48 \times 48$  size and [0,255] gray range. We store all these rotated, shrunk or enlarged preprocessed images in our database.

Fig. 2 shows some samples in our training database.

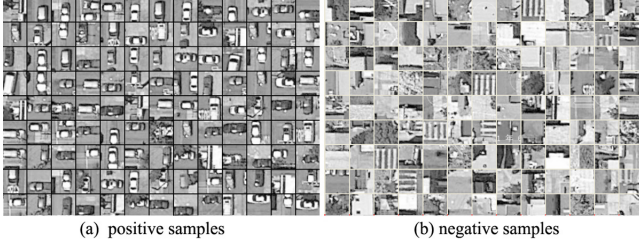


Fig. 2. Partial samples (rotation angle=0) from one image in training set

### B. Implementation of PNN

Fig. 3 shows the structure of PNN when using only Gray input. It has 9-layer: 48x48-80@C42x42,7-80@M21x21,2-80@C18x18,4-80@M9x9,2-80@C6x6,4-80@M3x3,2-300N-2N. It means: 48x48 input, a convolutional layer with 80 maps, 42x42 size and 7x7 filters, a max-pooling layer with 80 maps, 21x21 size and 2x2 fields,  $\dots$ , a max-pooling layer with 80 maps, 3x3 size and 2x2 fields, a full connected hidden layer with 300 nodes, an output layer with 2 nodes.

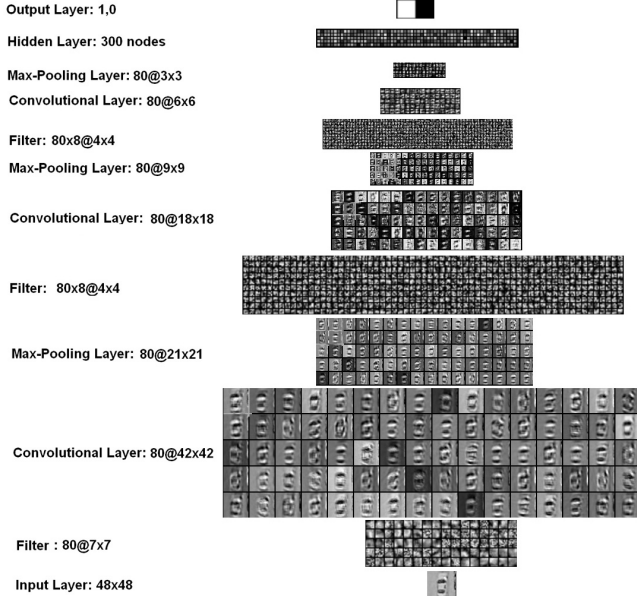


Fig. 3. The Structure of PNN ( $b_w=8$ ), input only Gray image

We wrote the PNN codes according to the formulas (3-13), we trained PNN on the GPU card, initial weights were set by a uniform random distribution in the range  $[-0.05, 0.05]$ , all initial biases were set to zero. LearnRate=0.001, Momentum=0, WeightDecay=0, batch size=50. Training ended when the validation error was near to zero. The samples with the same scaling and rotation angle were trained in one epoch, we changed scalings or rotation angles in the store sequence when a new epoch began.

## V. EXPERIMENT

Our database include 63 images, 6887 vehicles, 197513 samples from google earth at San Francisco city. 31 images, 3874 vehicles, 94858 samples are used as training set, other

32 images and 102655 samples are used as test set. We denote False Alarm Rate (FAR) and Detection Rate (DR) as:

$$\begin{cases} \text{FAR} = \frac{\text{number of false alarms}}{\text{number of vehicles}} \times 100\% \\ \text{DR} = \frac{\text{number of detected vehicles}}{\text{number of vehicles}} \times 100\% \end{cases} \quad (14)$$

To be fair and objective, some overlapped False Alarms are fused into one alarm.

TABLE I. FAR OF PNN (INPUT ONLY GRAY IMAGE)

$b_w$	test(s)	train(h)	Detection Rate				
			95%	90%	85%	80%	75%
2	35.21	23.10	55.7%	36.4%	23.5%	17.2%	12.9%
4	49.05	28.52	54.0%	34.6%	22.4%	16.5%	12.1%
8	56.43	34.87	53.4%	32.5%	21.3%	15.9%	11.4%
10	75.20	37.58	52.8%	30.4%	20.8%	14.9%	10.7%
16	110.7	46.80	50.9%	28.5%	18.6%	13.6%	9.35%
20	144.6	54.30	47.3%	26.7%	17.3%	12.3%	8.75%
40	237.3	75.23	44.6%	25.6%	16.0%	11.2%	8.25%
80	364.5	102.9	41.2%	23.0%	14.4%	10.5%	7.83%

Table 1 shows the influence of  $b_w$ . If  $b_w=1$ , the PNN training will not converge. In 85% detection rate, from  $b_w=2$  to  $b_w=80$ , the fastest speed is about 10 times of the slowest speed, with FAR increases no more than 10%. The test and training time unit is second and hour on GPU cards.

TABLE II. FAR OF PNN ( $b_w=80$ )

Input Data	Detection Rate				
	95%	90%	85%	80%	75%
Gray	41.2%	23.0%	14.4%	10.5%	7.83%
Gradient	43.5%	24.8%	15.7%	11.3%	8.20%
Gradient1	48.0%	26.5%	17.5%	12.8%	9.45%
Gradient2	44.6%	24.7%	15.9%	11.4%	8.26%

Table 2 shows the results of different inputs. In contrast to our expectation, the gradient input does not perform better than gray input, this is because the gradient image lost many details information of the object texture. Gradient2 performs better than Gradient1. Fig. 1 shows that Gradient2 contains more information than Gradient1.

TABLE III. FAR OF FOUR METHODS (INPUT MULTI-IMAGES)

Method	Detection Rate				
	95%	90%	85%	80%	75%
PNN( $b_w=80$ )	39.63%	21.54%	13.84%	10.08%	7.52%
PNN( $b_w=40$ )	43.01%	24.03%	14.82%	10.67%	7.94%
HOG+SVM	65.21%	40.21%	28.71%	21.81%	15.42%
LBP+SVM	74.35%	46.82%	32.20%	24.72%	17.37%

In Table 3, we input multi-images: Gray, Gradient, Gradient1 and Gradient2. We divided the 80 maps of the first convolutional layer of PNN into four equal parts, allocate 20 maps for each image. Comparing Table 2 and Table 3, we can see the result of multi-images is much better than any single image. It shows that the multi-images are complementary to each other, they have produced good resonance effect.

HOG feature is computed here as [18], where Gaussian smoother parameter  $\sigma = 2$ , derivative mask is  $[-1, 0, 1]$ , spacial orientation bins is 9, cell size is 8x8, each overlapped block include four cells. LBP feature is computed as [19], where  $P=8$ ,  $R=2$ , using 58 uniform patterns and 1 nonuniform pattern. The detection window is divided into  $1 \times 1 + 2 \times 2 + 3 \times 3 + 4 \times 4 + 5 \times 5 = 55$  blocks. We used rbf kernel in SVM, all other parameters are optimized.



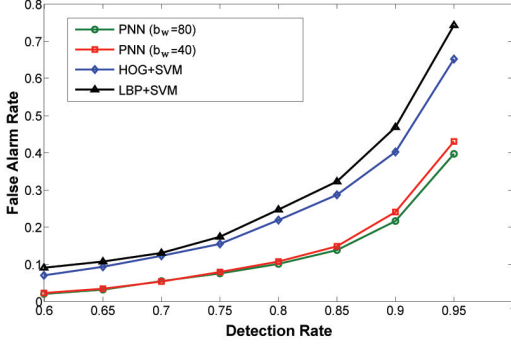


Fig. 4. False Alarm Rates of four methods on our vehicle test set.

Fig. 4 shows the difference between  $b_w=80$  and  $b_w=40$  are very subtle. Fig. 5 shows cars in different orientations

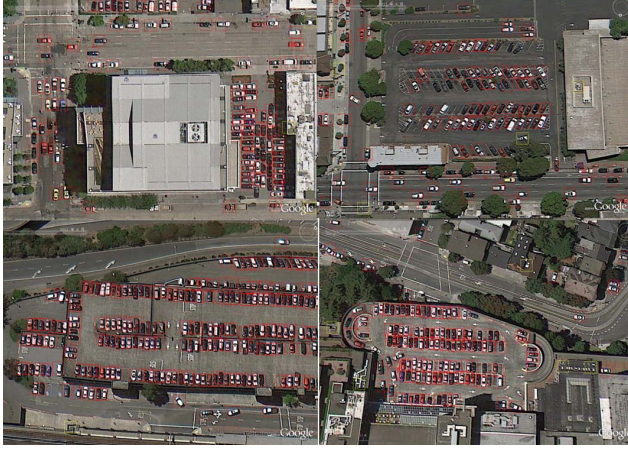


Fig. 5. Some detection results by PNN ( $b_w=40$ ) in San Francisco city. Red frames are the right detection, the yellow frames are the false alarm.

and places are detected correctly, it shows that our PNN has got scale-invariant and rotation-invariant power by training on samples with different scaling and rotation angles. However some black cars are missed, it implies that illumination and color play important roles in vehicle detection.

## VI. CONCLUSION

Practical applications like vehicle detection in satellite images often have high requirements in speed and accuracy. We proposed parallel deep convolutional neural network (PNN), which divides all maps of DNN into different branches, no direct map-connections between branches. The conventional DNN becomes a special form of PNN which has only one branch. Experiments on vehicle detection showed PNN can be 10 times faster than DNN, while keeping a good accuracy. All PNNs outperformed the traditional methods based on features. They also outperformed the vehicle detection result of [4] in the same challenging environments of San Francisco city. Via adjusting the branch width, PNN provides a good solution for compromising both speed and accuracy requirements in practical applications.

## ACKNOWLEDGMENT

This work was supported in part by the National Basic Research Program of China (973 Program) Grant 2012CB316300 and the Strategic Priority Research Program of the Chinese Academy of Sciences (Grant XDA06030300).

## REFERENCES

- [1] T. Zhao R. Nevatia, Car Detection in Low Resolution Aerial Images, Proc. ICCV , vol.1, pp. 710-717, 2001.
- [2] K. Ali, F. Fleuret, D. Hasler, and P. Fua, A Real-Time Deformable Detector, IEEE Trans. PAMI , 34(2):225-239, February 2012.
- [3] S. Hinz, Detection and Counting of Cars in Aerial Images, Proc. ICIP, 2003.
- [4] A. Kembhavi, D. Harwood, L. S. Davis, Vehicle Detection Using Partial Least Squares, IEEE Trans. PAMI , 33(6):1250-1265 June 2011.
- [5] L. Eikvil, L. Aurdal and H. Koren, Classification-based Vehicle Detection in High-resolution Satellite Images, Journal of Photogrammetry and Remote Sensing, 64(1):65-72, January 2009.
- [6] H. Grabner, T. Nguyen, B. Gruber, and H. Bischof, On-Line Boosting-Based Car Detection from Aerial Images, ISPRS J. Photogrammetry and Remote Sensing, 63(3):382-396, 2008.
- [7] L. Chen , Z. Jiang, J. Yang, Y. Ma, A Coarse-to-fine Approach for Vehicles Detection from Aerial Images, International Conference on Computer Vision in Remote Sensing (CVRS), pp. 221-225, 2012.
- [8] P. Liang, G. Teodoro, H. Ling, E. Blasch, G. Chen, L. Bai, Multiple Kernel Learning for vehicle detection in wide area motion imagery, 15th International Conference on Information Fusion (FUSION), pp. 1629-1636, 2012.
- [9] D. C. Cireşan, U. Meier and J. Schmidhuber, Multi-column Deep Neural Networks for Image Classification , Proc. CVPR 2012.
- [10] D. C. Cireşan, U. Meier, L. M. Gambardella, and J. Schmidhuber, Convolutional Neural Network Committees for Hand-written Character classification, In International Conference on Document Analysis and Recognition, pp. 1250-1254, 2011.
- [11] D. C. Cireşan, U. Meier, J. Masci, L. M. Gambardella, and J. Schmidhuber, Multi-Column Deep Neural Network for Traffic Sign Classification, Neural Networks, January 23, 2012.
- [12] D. H. Wiesel and T. N. Hubel, Receptive fields of single neurones in the cats striate cortex , J. Physiol., 148:574-591, 1959. 2
- [13] K. Fukushima, Neocognitron: A self-organizing neural network for a mechanism of pattern recognition unaffected by shift in position , Biological Cybernetics, 36(4):193-202, 1980.
- [14] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, Gradient-based learning applied to document recognition , Proceedings of the IEEE, 86(11):2278-2324, November 1998.
- [15] Y. LeCun, F.-J. Huang, and L. Bottou, Learning methods for generic object recognition with invariance to pose and lighting, Proc. CVPR, 2004.
- [16] C. Garcia and M. Delakis, Convolutional Face Finder: A Neural Architecture for Fast and Robust Face Detection, IEEE Trans. PAMI, 26(11):1408-1423, November 2004.
- [17] H. A. Rowley, S. Baluja, and T. Kanade, Neural Network-Based Face Detection, IEEE Trans. PAMI, 20(1):23-38, January 1998
- [18] N. Dalal and B. Triggs, Histograms of oriented gradients for human detection Proc. CVPR ,vol. 1, pp. 886-893, 2005.
- [19] T. Ojala, M. Pietikainen, T. Maenpää, Multiresolution Gray Scale and Rotation Invariant Texture Classification with Local Binary Patterns, IEEE Trans. PAMI, 24(7):971-987, July 2002.
- [20] C. Huang, S. Zhu, K. Yu, Large Scale Strongly Supervised Ensemble Metric Learning, with Applications to Face Verification and Retrieval <http://arxiv.org/abs/1212.6094>, 2011.
- [21] P. Viola and M. Jones, Rapid Object Detection Using a Boosted Cascade of Simple Features, Proc. Intl. Conf. Computer Vision and Pattern Recognition, vol. 1, pp. 511-518, 2001.