

# Spatial Context for Moving Vehicle Detection in Wide Area Motion Imagery with Multiple Kernel Learning

Pengpeng Liang<sup>a</sup> Dan Shen<sup>b</sup> Erik Blasch<sup>c</sup> Khanh Pham<sup>c</sup> Zhonghai Wang<sup>b</sup>  
Genshe Chen<sup>b</sup> Haibin Ling<sup>a</sup>

<sup>a</sup>Computer & Information Science Department, Temple University, Philadelphia, PA, USA

<sup>b</sup>Intelligent Fusion Technology, Inc, Germantown, MD, USA

<sup>c</sup>Air Force Research Lab, USA

{pliang,hbling}@temple.edu, {dshen,zwang,gchen@intfusiontech.com}, erik.blasch@rl.af.mil,  
Khanh.Pham@kirtland.af.mil

## ABSTRACT

Moving vehicle detection in wide area motion imagery is a challenging task due to the large motion of the camera and the small number of pixels on the target. At the same time, this task is very important for surveillance applications, and the result can be used for urban traffic management, accident and emergency responder routing. Also, the effectiveness of the context in object detection task can be further explored to increase target tracking accuracy. In this paper, we propose to use Spatial Context(SC) to improve the performance of the vehicle detection task. We first model the background of 8 consecutive frames with median filter, and get candidates by using background subtraction. The SC is built based on the candidates that have been classified as positive by Histograms of Oriented Gradient(HOG) with Multiple Kernel Learning(MKL). The region around each positive candidate is divided into  $m$  subregions with a fixed length  $l$ , then, the SC, a histogram, is built based on the number of positive candidates in each region. We use the publicly available CLIF 2006 dataset to evaluate the effect of SC. The experiments demonstrate that SC is useful to remove false positives, around which there are few positive candidates, and the combination of SC and HOG with multiple kernel learning outperforms the use of SC or HOG only.

**Keywords:** Vehicle detection, Spatial context, Multiple kernel learning

## 1. INTRODUCTION

Moving vehicle detection in wide area motion imagery is much more challenging than in traditional imagery, because of large amount of camera motion, small number of pixels on the targets, and the low frame rate of the video.<sup>1</sup> At the same time, more and more research focus on the understanding of such kind of imagery. A frame work to detect and track large number of cars in wide area image is introduced in,<sup>1</sup> and the performance of several state-of-the-art visual trackers are evaluated on Columbus Large Image Format (CLIF) dataset.<sup>2</sup> Zhao et al.<sup>3</sup> used the boundary of the car, the boundary of the front windshield and the shadow as features and integrated these features in the structure of the Bayesian network to detect vehicles. Prokaj et al.<sup>4</sup> used a similar method as<sup>1</sup> to perform vehicle detection, where the detected vehicles are refined by using the tracking result. Xiao et al.<sup>5</sup> used the road information from an additional co-registered geospatial information system (GIS) database as a constraint to refine the detected vehicles from the tracking results. Shi et al.<sup>6</sup> first used the vehicle detection result to construct trajectories, then they use these trajectories to estimate the road information which is used to remove false positives. Liang et al.<sup>7</sup> proposed to combine the HOG and Haar feature to improve the vehicle detection performance.

In order to use the framework introduced in,<sup>1</sup> we need to do background subtraction first. Due to the camera motion, we do registration for a set of consecutive frames. The more the consecutive frames we use, the smaller the common area is. So, in this task, the Gaussian mixture model<sup>8,9</sup> is not suitable, since it requires a relatively large number of frames to get a satisfying background estimate. In our approach, we choose to use the median background model. However, due to the large camera motion, 3D parallax, the registration and the background model is not robust enough to generate accurate foreground estimates, and there are a lot of false positives, so we need to classify these candidates as true positives and false positives.

For classification, since the target in wide area motion imagery is very small as shown in Fig. 1, there is a limit to use shape, appearance, color, or texture models as used in.<sup>10</sup> Beyond using the appearance information, the effectiveness

of context information has been demonstrated in many recent works,<sup>11–13</sup> and is useful to improve the performance of the object detection task. Based on the idea that if a candidate is indeed a vehicle, there are usually other candidates which are also vehicles around it as shown in Fig. 1, we propose the Spatial Context (SC). Given the candidates obtained from background subtraction, we build SC based on the candidates classified as positive by Histograms of Oriented Gradient (HOG).<sup>14</sup> HOG is based on Scale-invariant feature transform (SIFT)<sup>15</sup> and uses oriented gradient histogram to delineate objects, which has been successfully used with support vector machines (SVM) in human detection.



Figure 1. Part of the original image, green and red boxes indicate the positive and negative candidates classified by HOG with GMKL respectively.

In order to benefit from both the appearance information and the context information, we use multiple kernel learning (MKL)<sup>16–19</sup> to combine these two kinds of features. The MKL can learn the trade-off between the appearance and context automatically. In this paper, we use Spatial Projected Gradient -Generalized MKL (SPG-GMKL)<sup>16</sup> which is very efficient and can handle millions of kernels.

The rest of the paper is organized as follows: In Section 2, we give a brief introduction to HOG descriptor. In Section 3, we describe the proposed Spatial Context (SC). In Section 4, the SPG-GMKL is introduced. In Section 5, we describe the dataset and give the experiment results. Finally, we conclude this paper in Section 6.

## 2. HISTOGRAMS OF ORIENTED GRADIENTS (HOG)

Histograms of Oriented Gradients (HOG)<sup>14</sup> use the distribution of the oriented gradients to describe a local part of an image. The gradient of a pixel in an image is calculated by masks, and for color images, HOG calculates gradients for each color channel, and chooses the one with the largest normal as the pixel’s gradient vector. Several masks were tested in,<sup>14</sup> including a 1-D point derivative (un-centered [-1,1], centered [-1,0,1] and cubic-corrected [1,-8,0,8,-1]), 2-D derivatives ( $\begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$ ,  $\begin{pmatrix} -1 & 0 \\ 0 & 1 \end{pmatrix}$ ) and a  $3 \times 3$  Sobel mask. Test results show that the simplest mask [-1,0,1] works best.

For each pixel,  $[-1, 0, 1]$  and  $[-1, 0, 1]^T$  are used to calculate a gradient horizontally and vertically respectively. Based on the horizontal and vertical gradients, the orientation and weight can be calculated by voting. The votes are accumulated into orientation bins over local spatial regions called cells composed by pixels. The weight of a vote is a function of the gradient magnitude and the function can be the magnitude itself, it’s square, it’s square root, etc..

In order to alleviate the effect of the variation of illumination and foreground-background contrast, effective local normalization is required. Normalization is done by grouping cells into larger spatial blocks and normalizing each block separately. In practice, the blocks are overlapped so that each cell response contributes to several components of the final descriptor. In,<sup>14</sup> four different normalization scheme are evaluated. They are (a)  $L1$ -norm,  $v \rightarrow v / (\|v\|_1 + \epsilon)$ ; (b)  $L1$ -sqrt,  $v \rightarrow \sqrt{v / (\|v\|_1 + \epsilon)}$ ; (c)  $L2$ -norm,  $v \rightarrow v / \sqrt{\|v\|_2^2 + \epsilon^2}$ ; (d)  $L2$ -Hys,  $L2$ -norm followed by limiting the maximum value of  $v$  to 0.2. It turns out that the performance of  $L2$ -Hys,  $L2$ -norm, and  $L1$ -sqrt are almost the same, whereas the  $L1$ -norm cannot compete with the other three schemes.

### 3. SPATIAL CONTEXT(SC)

Given a candidate  $p$  in frame  $I$ , which has been classified as positive by HOG<sup>14</sup> with GMKL,<sup>16</sup> we first divide the region around  $p$  into  $m$  subregions with length  $l$ . Each subregion is a fan (e.g., center with rotating blades) and all the subregions have the same center angle  $\theta$ .  $m$  is determined by  $\theta$ , and  $m = 360/\theta$ . Next, we calculate the number of candidates which are also classified as positive by HOG and GMKL. Then, we build a histogram with  $m$  bins, where each bin corresponds to a specific subregion. The value of each bin is the number of positive candidates in the corresponding subregion. The value of the  $k$ th bin of candidate  $p$  is

$$h_p(k) = \#\{i \in C \wedge i \neq p \wedge i \in R(k)\} \quad (1)$$

where  $C$  is the set of the candidates classified as positive by HOG with GMKL, and  $R(k)$  is the subregion corresponding to the  $k$ th bin. The approach to build SC is described in Algorithm 1.

---

**Algorithm 1** Building SC for a Candidate

---

**Require:**

$p$ : a candidate  
 $C$ : the set of candidates of  $I$

**Ensure:**

$h$ : the SC, a histogram, for  $p$

- 1: **for** each bin  $k$  of  $h$  **do**
- 2:    $count \leftarrow 0$
- 3:   **for** each  $i \in C$  **do**
- 4:     **if**  $i \in C \wedge i \neq p \wedge i \in R(k)$  **then**
- 5:        $count \leftarrow count + 1$
- 6:     **end if**
- 7:   **end for**
- 8:    $h(k) \leftarrow count$
- 9: **end for**
- 10: **return**  $h$

---

### 4. MULTIPLE KERNEL LEARNING WITH SPECTRAL PROJECTED GRADIENT (SPG)

The objective of multiple kernel learning is to learn the optimal combination of the given base kernels. In,<sup>17</sup> the optimal kernel is approximated as  $\mathbf{K}_{opt} = \sum_k d_k \mathbf{K}_k$  where  $\mathbf{d}$  corresponds to the trade-off among the base kernels. However, the sum of base kernels is too constraining, and is just the concatenation of individual kernel's feature space. In order to relieve this constraint, in,<sup>18</sup> the generalized multiple kernel learning (GMKL) is proposed, so that the combination of base kernels can be more flexible. The primal formulation for optimizing the GMKL is the following:

$$\begin{aligned} \min_{\mathbf{d}} \quad & T(\mathbf{d}) \quad \text{subject to} \quad \mathbf{d} \geq 0 \\ \text{where} \quad & T(\mathbf{d}) = \min_{\mathbf{w}, \mathbf{d}} \frac{1}{2} \mathbf{w}^t \mathbf{w} + \sum_i l(y_i, f(\mathbf{x}_i)) + r(\mathbf{d}) \end{aligned}$$

where  $w$  are the weights,  $l$  is the loss function, and both the regularizer  $r$  and the kernel can be any general differentiable functions of  $\mathbf{d}$  with a continuous derivative. The optimization of  $\mathbf{d}$  is carried out in the outer loop, while the optimization of  $\mathbf{d}$  is performed in the inner loop.

Though the GMKL in<sup>18</sup> allows more flexible combination of kernels, the projected gradient descent optimizer is inefficient, since the computation of the step size and a reasonably accurate gradient direction are all expensive. In order to alleviate this problem, using Spectral Projected Gradient (SPG) descent to solve the GMKL is proposed in.<sup>16</sup> The SPG has the following advantages:<sup>18</sup> (a) second order information is explored to determine the step size; (b) employs a non-monotone step size selection criterion requiring fewer function evaluations; (c) it is robust to gradient noise, and (d) it can take steps when far away from the optimum. The pseudo code of SPG-GMKL is given in Algorithm 2.

---

**Algorithm 2** SPG-GMKL

---

```
1:  $n \leftarrow 0$ 
2: Initialize  $\mathbf{d}^0$  randomly
3: repeat
4:    $\alpha^* \leftarrow \text{SolveSVM}_c(\mathbf{K}(\mathbf{d}^n))$ 
5:    $\lambda \leftarrow \text{SpectralStepLength}$ 
6:    $\mathbf{p}^n \leftarrow \mathbf{d}^n - \mathbf{P}(\mathbf{d}^n - \lambda \nabla \mathbf{W}(\mathbf{d}^n))$ 
7:    $s^n \leftarrow \text{Non-Monotone}$ 
8:    $\epsilon \leftarrow \text{TuneSMPrecision}$ 
9:    $\mathbf{d}^{n+1} \leftarrow \mathbf{d}^n - s^n \mathbf{p}^n$ 
10: until converged
```

---

## 5. EXPERIMENT

### 5.1 Dataset

We use Columbus Large Image Format (CLIF) 2006<sup>2</sup> dataset to evaluate the proposed Spatial Context (SC). The scene of this dataset is a flyover of the Ohio State University (OSU) from a large format monochromatic electro-optical platform which is comprised of a matrix of six cameras and the size of each image is 2672(width) by 4008(height) pixels. Since our focus is vehicle detection, we use a  $2672 \times 1200$  subregion which contains an expressway as shown in Figure 2 and Figure 3. The subregion not only includes horizontal and vertical roads, but also an overpass. We labeled the vehicles in 102 frames of camera 3, and there are 9364 vehicles in total which are used as test data. We also labeled 1730 candidates obtained from background subtraction from 16 frames of camera 1 which are used as training data.

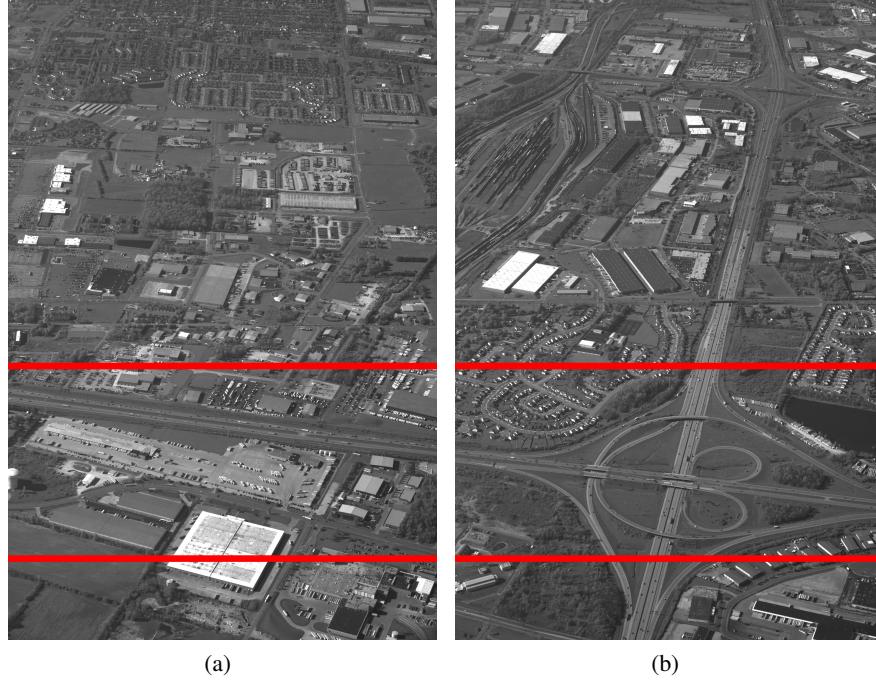


Figure 2. (a) and (b) are the original images, where the region between the two red lines is the subregion.

### 5.2 Quantitative Evaluation

For quantitative evaluation, we use the distance between a positive candidate and the groundtruth to decide whether it is a true positive. Given the groundtruth  $G = \{g_1, g_2, \dots, g_S\}$ , if the candidate  $c$  is true positive, there exists  $g \in G$  and the distance between the center of  $g$  and the center of  $c$  is not greater than ten pixels. Since we classify the candidates obtained



(a)



(b)

Figure 3. (a) and (b) are subregions for 2(a) and 2(b) respectively.

from background subtraction, we use two different ways to calculate the recall. One way is to use the number of the actual groundtruth,  $S$ , i.e., the actual number of vehicles in each image; another way is to use the number of candidates which are indeed vehicles as groundtruth,  $S'$ . We give precision-recall curves for both. Without classification, the performance of background subtraction is not satisfying, the precision is only 0.398 at the recall 0.854. For parameters of SC, we test several configurations of  $m$  and  $l$ , the difference is tiny, and we use 6 and 200 for  $m$  and  $l$  respectively in the following experiment.

To demonstrate the usefulness of SC, the performance of SC, HOG and the combination of SC and HOG are evaluated. The precision-recall curves are given in Fig. 4. From the results, we can see that HOG outperforms SC, but the combination of SC and HOG performs best, especially at a high recall rate.

### 5.3 Qualitative Evaluation

Some qualitative results are given in Fig. 5 and Fig. 6. From the results, we can see that SC is useful to remove false positives when there are few candidates that are classified as positive by HOG with GMKL around them. Usually, this kind of candidates are away from the road. So, SC is useful to remove the false positives that are away from the road. At the same time, the number of misclassified vehicles on the road is very small.

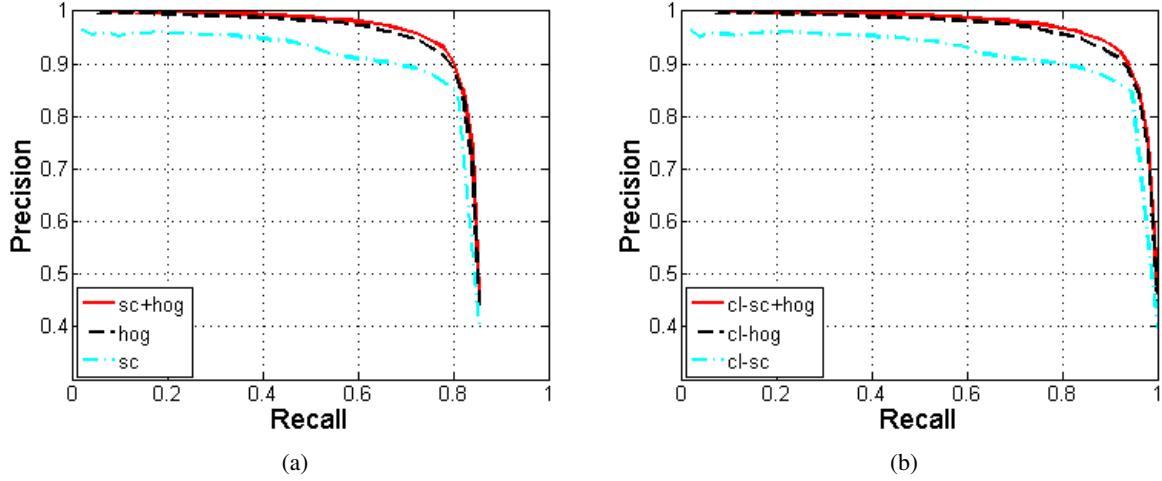


Figure 4. The precision-recall curve for SC, HOG and the combination of SC and HOG.(a) is plotted based on  $S$ , (b) is plotted based on  $S'$ .

## 6. CONCLUSION

In this paper, we propose Spatial Context (SC) to improve the vehicles detection task in wide area motion imagery. In order to make use of both the appearance and context information, we use SPG-GMKL which is very efficient to combine SC and HOG. Experiments on the CLIF 2006 dataset verified that SC is useful to remove false positives that are away from the road.

## REFERENCES

- [1] Reilly, V., Idrees, H., and Shah, M., “Detection and tracking of large number of targets in wide area surveillance.” in *European Conference on Computer Vision*, pp. 186–199 (2010).
- [2] O. Mendoza-Schrock, J. A. Patrick, and E. P. Blasch, “Video image registration evaluation for a layered sensing environment.” in *IEEE Nat. Aerospace Electronics Conf. (NAECON)* (2009).
- [3] T. Zhao and R. Nevatia, “Car detection in low resolution aerial image,” in *ICCV* (2001).
- [4] J. Prokaj, M. Duchaineau, and G. Medioni, “Inferring tracklets for multi-object tracking,” in *Workshop of Aerial Video Processing joint with IEEE CVPR* (2011).
- [5] J. Xiao, H. Cheng, H. S. Sawhney, and F. Han, “Vehicle detection and tracking in wide field-of-view aerial video,” in *CVPR* (2010).
- [6] X. Shi, H. Ling, E. Blasch, and W. Hu, “Context-driven moving vehicle detection in wide area motion imagery,” in *Int'l Conf. on Pattern Recognition (ICPR)* (2012).
- [7] P. Liang, G. Teodoro, H. Ling, E. Blasch, G. Chen, and L. Bai, “Multiple kernel learning for vehicle detection in wide area motion imagery .” in *International Conference on Information Fusion* (2012).
- [8] C. Stauffer and W. E. L. Grimson, “Adaptive background mixture models for real-time tracking.” in *Computer Vision and Pattern Recognition*, pp. 2246–2252 (1999).
- [9] P. KaewTraKulPong and R. Bowden, “An improved adaptive background mixture model for real-time tracking with shadow detection.” in *European Workshop on Advanced Video Based Surveillance Systems* (2001).
- [10] M.-E. Nilsback and A. Zisserman, “A visual vocabulary for flower classification,” in *Computer Vision and Pattern Recognition* (2), pp. 1447–1454 (2006).
- [11] G. Heitz and D. Koller, “Learning spatial context: Using stuff to find things,” in *ECCV* (1) (2008).
- [12] A. Jain, A. Gupta, and L. S. Davis, “Learning what and how of contextual models for scene labeling,” in *ECCV* (4) (2010).
- [13] S. K. Divvala, D. Hoiem, J. Hays, A. A. Efros, and M. Hebert, “An empirical study of context in object detection,” in *CVPR* (2009).

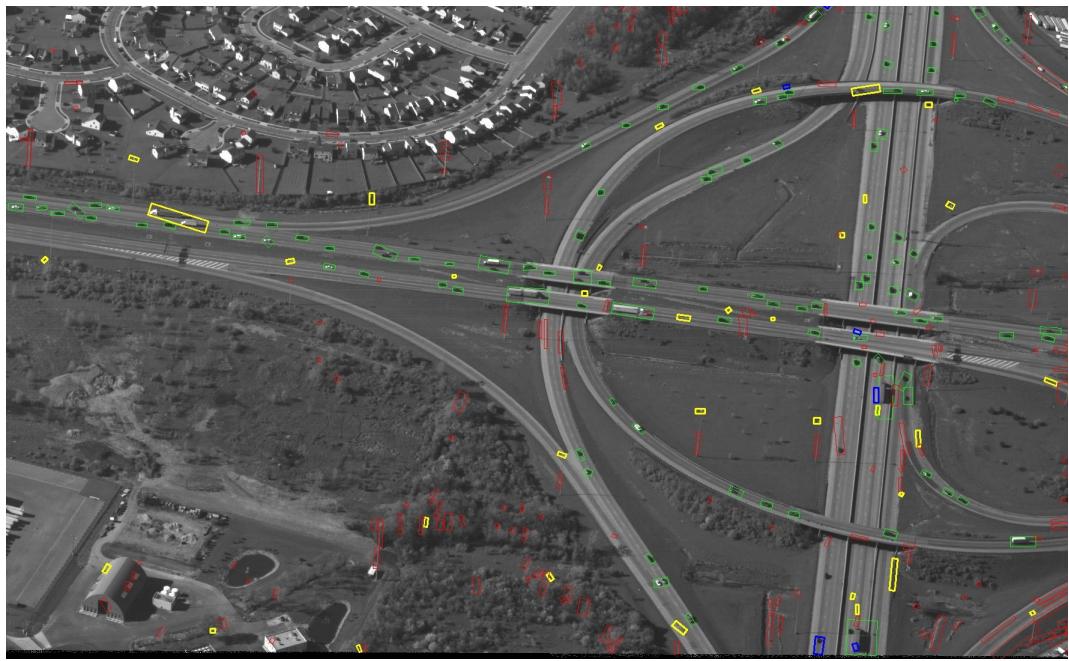


(a)

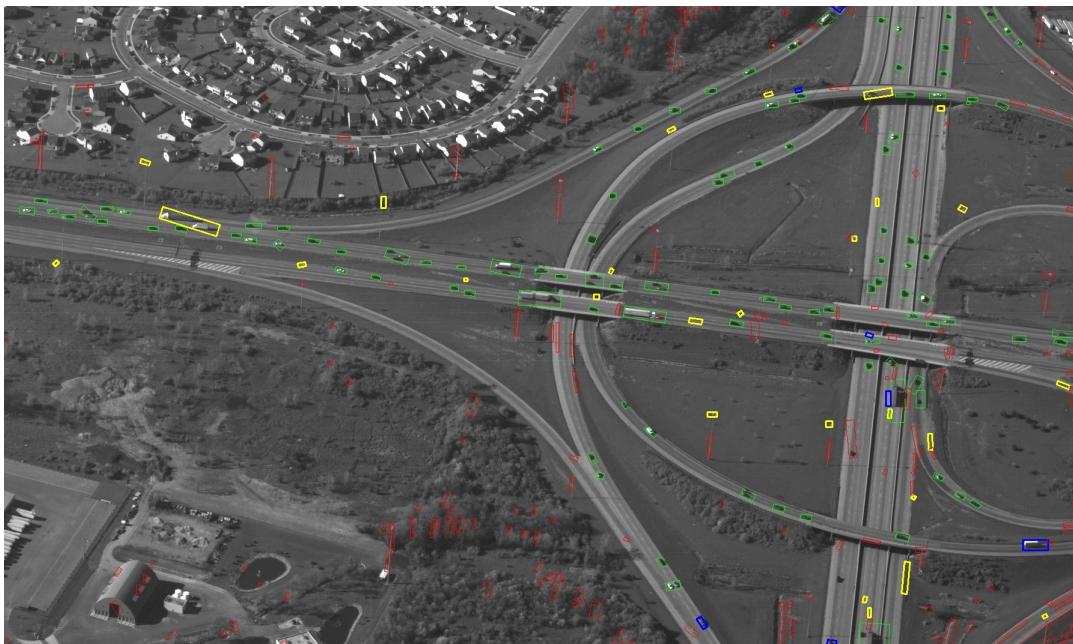


(b)

Figure 5. The classification result. The green, yellow, red, and blue bounding boxes indicate true positive, false positive, true negative, false negative respectively. (a), (b) are the results of SC, TC+HOG respectively.



(a)



(b)

Figure 6. The classification result. The green, yellow, red, and blue bounding boxes indicate true positive, false positive, true negative, false negative respectively. (a), (b) are the results of SC, TC+HOG respectively.

- [14] N. Dalal and B. Triggs, “Histograms of oriented gradients for human detection.” in *Computer Vision and Pattern Recognition*, pp. 886–893 (2005).
- [15] D. G. Lowe, “Distinctive image features from scale-invariant keypoints,” *International Journal of Computer Vision*, vol. 60, no. 2, pp. 91–110 (2004).
- [16] A. Jain, S. V. N. Vishwanathan, and M. Varma, “Spg-gmkl: generalized multiple kernel learning with a million kernels,” in *KDD* (2012).
- [17] M. Varma and B. R. Babu, “More generality in efficient multiple kernel learning,” in *ICML* (2009).
- [18] M. Varma and D. Ray, “Learning the discriminative power-invariance trade-off,” in *ICCV* (2007).
- [19] M. Gönen and E. Alpaydin, “Multiple kernel learning algorithms,” *Journal of Machine Learning Research*, vol. 12, pp. 2211–2268 (2011).