# SINet: A Scale-insensitive Convolutional Neural Network for Fast Vehicle Detection

Xiaowei Hu, *Student Member, IEEE,* Xuemiao Xu, *Member, IEEE,*
Yongjie Xiao, Hao Chen, *Student Member, IEEE,* Shengfeng He, *Member, IEEE,*
Jing Qin, *Member, IEEE,* and Pheng-Ann Heng, *Senior Member, IEEE*

*Abstract*—Vision-based vehicle detection approaches achieve incredible success in recent years with the development of deep convolutional neural network (CNN). However, existing CNN-based algorithms suffer from the problem that the convolutional features are scale-sensitive in object detection task but it is common that traffic images and videos contain vehicles with a large variance of scales. In this paper, we delve into the source of scale sensitivity, and reveal two key issues: 1) existing RoI pooling destroys the structure of small scale objects; 2) the large intra-class distance for a large variance of scales exceeds the representation capability of a single network. Based on these findings, we present a scale-insensitive convolutional neural network (SINet) for fast detecting vehicles with a large variance of scales. First, we present a context-aware RoI pooling to maintain the contextual information and original structure of small scale objects. Second, we present a multi-branch decision network to minimize the intra-class distance of features. These lightweight techniques bring zero extra time complexity but prominent detection accuracy improvement. The proposed techniques can be equipped with any deep network architectures and keep them trained end-to-end. Our SINet achieves state-of-the-art performance in terms of accuracy and speed (up to 37 FPS) on the KITTI benchmark and a new highway dataset, which contains a large variance of scales and extremely small objects.

*Index Terms*—Vehicle detection, scale sensitivity, fast object detection, intelligent transportation system.

## I. INTRODUCTION

**A**UTOMATIC vehicle detection from images or videos is an essential prerequisite for many intelligent transportation systems. For example, vehicle detection from in-car videos (Fig. 1) is critical for the development of autonomous driving systems while vehicle detection from surveillance videos (Fig. 2) is fundamental for the implementation of
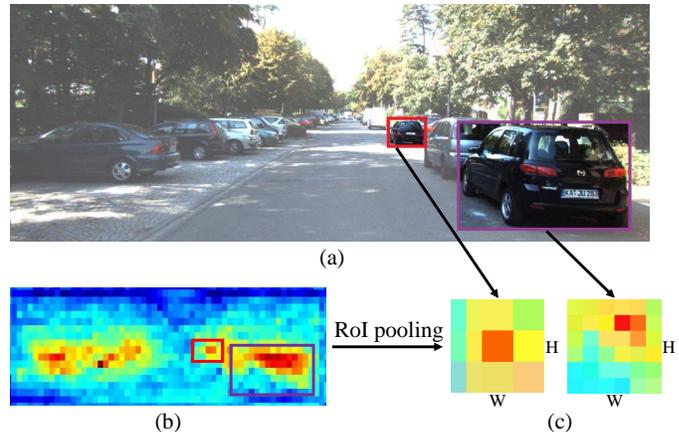
Fig. 1. The scale-sensitive problem. (a) An image includes both large and small vehicles. (b) The feature representations of small and large vehicles in deep layer are largely different. (c) Traditional RoI pooling introduces noise as it simply replicates the values on the feature map for the small vehicle.

intelligent traffic management systems. In this regard, over the past decade, a lot of effort has been dedicated to this field [1–20]. Some challenging benchmarks have also been proposed for evaluation and comparison of various detection algorithms [21]. On the other hand, in recent years, deep convolutional neural networks (CNNs) have achieved incredible success on vehicle detections as well as various other object detection tasks [22–30]. However, when applying CNNs to vehicle detection, one of the main challenges is that traditional CNNs are sensitive to scales while it is quite common that in-car videos or transportation surveillance videos contain vehicles with a large variance of scales (see the vehicles in Fig. 1 (a) and the input of Fig. 2). The underlying reason of this scale-sensitive problem is that it is challenging for a CNN to response to all scales with optimal confidences [31].

Existing CNN-based object detection algorithms attempt to make the network fit different scales by utilizing input images with multiple resolutions [23, 24, 26, 29, 31, 33, 34] or fusing multi-scale feature maps of CNN [22, 25, 28, 30, 35–40]. These methods, however, introduce expensive computational overhead and thus are still incapable of fast vehicle detection, which is essential for autonomous driving systems, real-time surveillance and prediction systems.

Instead of simply adding extra operations, we look into the detection network itself and scrutinize the underlying reasons of this scale-sensitive problem. We observe two main barriers. First, inadequate and/or imprecise features of small regions lead to the loss of detecting small objects (e.g., the red box in
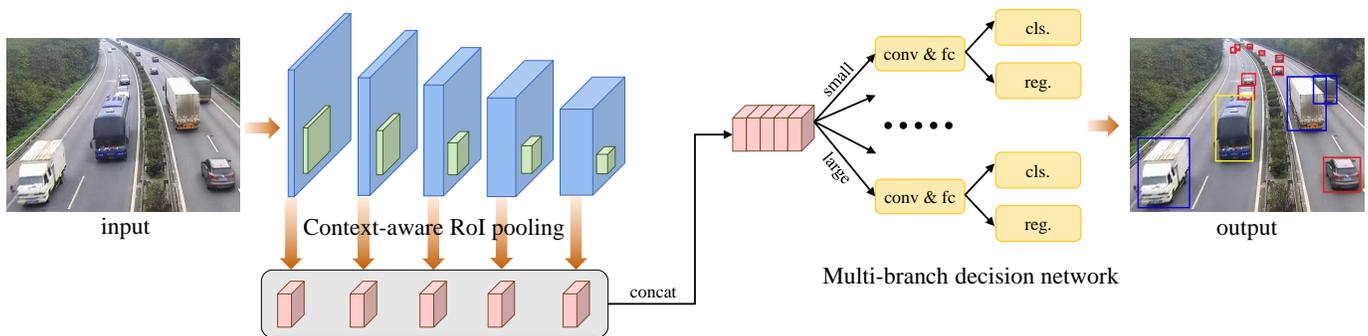
Fig. 2.    The schematic illustration of the pipeline of the proposed SINet: (i) we extract feature maps with multiple scales over the CNN [32] from the input image and get the proposals based on the CNN features [28]; (ii) each proposal on different layers is pooled into a fixed-size feature vector using the context-aware RoI pooling, in which the small proposals are enlarged by the deconvolution with bilinear kernels to achieve better representation (see section IV-B for details); (iii) we concatenate the features of proposals at each layer and feed them to the multi-branch decision network; and (iv) lastly, we fuse the predicted bounding boxes from all branches to produce the final detection results (car in red, bus in yellow, van in blue). Best viewed in color.

Fig. 1 (b)). In particular, the commonly used RoI pooling [23] distorts the original structure of small objects, as it simply replicates the feature values to fit the preset feature length (as shown in the left example of Fig. 1 (c)). Second, the intra-class distance between different scales of vehicles is usually quite large. As illustrated in Fig. 1 (b), the red and purple boxes have different feature responds. This makes it difficult for the network to represent objects with different sizes using the same set of weights.

To cope with the above problems, we present a scale-insensitive convolutional neural network, named SINet, to detect vehicles with a large variance of scales accurately and efficiently. The network architecture is shown in Fig. 2. Object proposals are used on the feature map level to examine all the possible object regions, and the corresponding feature maps are fed to a decision network. Two new methods are proposed to overcome above-mentioned barriers. We first present a context-aware RoI pooling scheme to preserve the original structures of small scale objects. This new pooling layer involves a deconvolution with bilinear kernels which can maintain the context information and hence help produce features that are faithful to the original structure. These pooled features are then fed to a new, multi-branched decision network. Each branch is designed to minimize the intra-class distance of features, and therefore the network can more effectively capture the discriminative features of objects with various scales than traditional networks.

The proposed network achieves state-of-the-art performance on both detection accuracy and speed on the KITTI benchmark [21]. This method also shows a promising performance on detecting vehicles with low resolution input images, and brings detecting vehicles in low-resolution video surveillance into practice. Due to the lightweight architecture, real-time detection (up to 37 FPS) can be achieved on a 256×846 image. In order to demonstrate the proposed method in more practical scenes, we construct a new highway dataset, which contains vehicles with a vast variance of scales. To the best of our knowledge, it is the first dataset focuses on the highway scene. It contains 14388 well labelled images under different roads, time, weathers and traffic states. This dataset, as well as the source code of the SINet, are publicly available at https:

//xw-hu.github.io/. In summary, our contributions include:

- We present a context-aware RoI pooling layer, which can produce accurate feature maps for vehicles with small scales without extra space and time burdens. The proposed new pooling layer can be widely applied to existing architectures.
- We present a multi-branch decision network for vehicle detection. It can accurately classify vehicles with a large variance of scales without introducing extra computational cost.
- We construct the first large scale variance highway dataset, which provides a platform with practical scenes to evaluate the performance of various vehicle detection algorithms in handling target object with a large variance of scales.

## II. RELATED WORKS ON VEHICLE DETECTION

In this section, we give a brief introduction of the monocular vision vehicle detection methods, as our approach also belongs to the monocular vision detection. More comprehensive analysis of vehicle detection on monocular, stereo, and other vision-sensors can be found in [41].

Early works use the relative motion cues between the objects and background to detect the vehicles. Adaptive background models such as Gaussian Mixture Model (GMM) [3–5], Sigma-Delta Model [9] are widely used in vehicle detection by modeling the distribution of the background as it appears more frequently than moving objects. Optical flow is a common technique to aggregate the temporal information for vehicle detection [10] by simulating the pattern of object motion over time. Optical flow is also combined with symmetry tracking [8] and hand-crafted appearance features [7] for better performance. However, this kind of approach is unable to distinguish the fine-grained categories of the moving objects such as car, bus, van or person. In addition, these methods need lots of complex post-processing algorithms like shadow detection and occluded vehicle recognition to refine the detection results.

Then, the statistical learning methods based on the hand-crafted features are applied to detect the vehicles from the images directly. They first describe the regions of the image by some feature descriptors and then classify the image regions
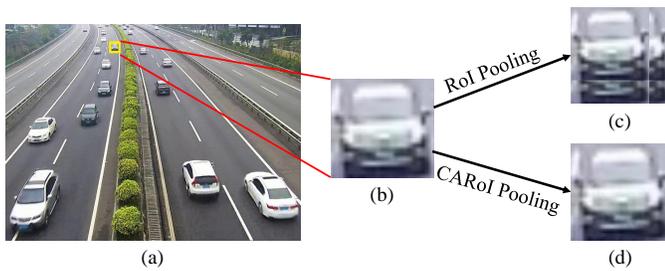
Fig. 3. The difference between RoI pooling and CARoI pooling. For the sake of clarity, we apply these two pooling layers on natural images instead of feature maps.

into different classes such as vehicle and non-vehicle. Features like HOG [11, 13], SURF [14], Gabor [13] and Haar-like [15, 16] are commonly used for vehicle detection followed by classifiers like SVM [11, 14], artificial neural network [13] and Adaboost [15, 16]. More advanced algorithms like DPM [12, 17] and And-Or Graph [1, 2] explore the underlying structures of vehicles and use hand-crafted features to describe each part of vehicles. These features, however, have limited ability of feature representation, which is difficult to handle complex scenarios.

Recently, features learned by the deep convolutional neural network show strong representation of the semantic meanings of objects, which make a great contribution to the state-of-the-art object detectors [23, 24, 27, 29]. Although these methods overperform lots of hand-crafted vehicle detection methods on the vehicle detection benchmark [21], the vehicles with a large variance of scales (Fig. 1 and Fig. 2) are still difficult to be detected accurately in a real-time manner due to the scale sensitive convolutional features. We will elaborate the scale-sensitive problem of current CNNs in the following section.

## III. WHY CURRENT CNNS ARE SCALE-SENSITIVE

It is well-known that CNNs are sensitive to scale variations in detection tasks [28]. In this section, we first carefully analyze its underlying reasons and then discuss how existing solutions solve this problem.

### A. Structure Distortion Caused by RoI Pooling

CNNs-based object detection algorithms fall into two categories. The first category is built upon a two-stage pipeline [23, 24, 26, 27, 33–36, 42, 43], where the first stage extracts proposals and the second stage predicts their classes. The second category aims to train an end-to-end object detector [37, 44, 45], which skips the object proposal detection and hence has relatively faster computational speed. Such a detector first implicitly divides the image into a grid, then simultaneously makes prediction for each square or rectangle in the grid, and finally figures out the bounding boxes of targeting objects based on the predictions of the squares or rectangles [44]. However, this grid-based paradigm cannot obtain comparable accuracy to two-stage detection pipeline, as the grids have too strong spatial constraints to predict small objects appeared as groups [43]. In this regard, most of the existing methods employ the two-stage detection pipeline.

In order to satisfy the input requirement of the classification networks, most two-stage object detection algorithms, e.g.,

SPP [25], Fast RCNN [23] and Faster RCNN [27], represent each proposal as a fixed-size feature vector by RoI pooling [23]. As shown in Fig. 1 (c), the RoI pooling divides every proposal into $H \times W$ sub-windows and uses the max pooling to extract one value for each sub-window so that the output can have a fixed-size of $H \times W$. If a proposal is smaller than $H \times W$, it is enlarged to $H \times W$ by simply replicating some parts of the proposal to fill the extra space. Unfortunately, such a scheme is not appropriate as it may destroy the original structures of the small objects (see Fig. 3 (c)). During the network training process, filling with replicated values not only leads to inaccurate representations in the forward propagation, but also accumulates errors in the backward propagation. The inaccurate representations and accumulated errors mislead the training and prevent the network from correctly detecting small scale vehicles. In our experiments, we find that this problem is critical for the low detection accuracy of small vehicles.

### B. Intra-class Distance Caused by Scale Variations

The other important issue that causes scale sensitivity is the large intra-class distance between large and small scale objects. Once the features of each proposal are extracted, they are fed into a decision network for classification. Existing methods treat objects within the same class equally regardless of their scales. We argue that this may lead to inaccurate detection, as the intra-class distance between large and small scale objects may be as significant as the intra-class distance on their feature representations.

### C. Existing Solutions and Their Shortcomings

A lot of effort has been dedicated to solving this scale sensitivity issue. As mentioned, most existing solutions are designed based on two types of pyramid representations. The first one applies the concept of image pyramid (Fig. 4 (a)), which exploits the multi-scale input images to make the network fit all the scales [23, 26, 29, 31, 33, 34]. However, the main disadvantage of this representation is its large computational cost [36], prohibiting its application to real-time detection tasks.

The other representation is the feature pyramid, which exploits the information extracted from multi-layer feature maps. The first and straightforward attempt is to use the high resolution shallow layers to detect small objects, while using low resolution deep layers to detect large objects (as shown in Fig. 4 (b)). This strategy has been adopted by SSD [37], MSCNN [28], FCN [38] and SDP [30]. However, as the feature maps in the shallow layers lack of the semantic information, they usually fail to distinguish the small objects accurately.

In order to take full advantage of deep layer information to tackle scale variations, some researchers presented to combine multi-layer feature maps together to train a network (e.g., MultiPath [40] and HyperNet [35], see Fig. 4 (c)). However, due to the down-sampling operations used in the network, small objects cannot maintain sufficient spatial information in the deep layers, and thus they are still difficult to be detected. To better maintain the deep feature maps of small
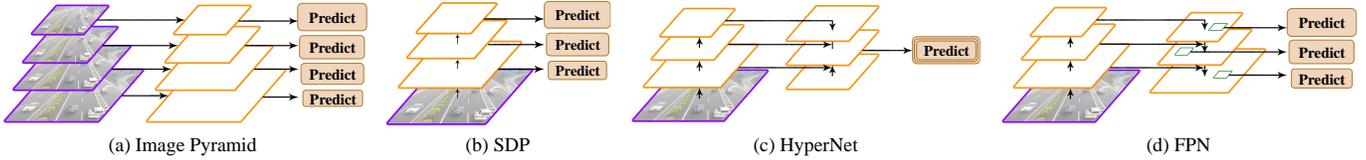
Fig. 4. (a) Multiple predictions based on multi-scale images. (b) Multiple predictions on multi-layer features. (c) Single prediction on concatenation features. (d) Multiple predictions on multi-layer features concatenating with low-layer features.

objects, another solution is proposed to use the high-resolution feature maps and the up-sampled deep feature maps together to predict small objects, such as [36], [39] (Fig. 4 (d)). The main problem of this solution is that the upsampling operation is performed on the entire feature map, which requires more memory resources and additional computational costs. While extra information leads to better accuracy, high computational cost is inevitable [46], which is unacceptable for our real-time vehicle detection task.

As a consequence, instead of introducing additional steps to solve the scale-sensitive problem, we aim to address this problem internally by introducing two simple solutions: a novel context-aware pooling and a multi-branch decision network, which lead to zero extra computational cost while effectively dealing with the scale sensitivity issue for real-time and accurate vehicle detection.

## IV. SCALE-INSENSITIVE NETWORK

### A. Overview

The architecture of the proposed scale-insensitive network (SINet) is illustrated in Fig. 2. Our SINet takes the whole image as input and outputs the detection result in an end-to-end manner. It first generates a set of convolutional feature maps [32] and obtains a set of proposals based on these feature maps using region proposal networks (RPN) [27], [28]. The RPN predicts the bounding boxes that have a large probability of containing objects and these predicted bounding boxes are called as proposals. Then, the proposed context-aware RoI pooling (CARoI pooling) is used to extract the features for each proposal. The CARoI pooling applies the deconvolution with bilinear kernels to enlarge the feature regions of the small proposals to avoid representing the small objects with the replicated values. The CARoI pooling is applied to multiple layers of the CNN and these pooled features at different layers are concatenated together to fuse the low-level detail information and the high-level semantic information to detect the objects [35]. After that, we split the SINet into multiple branches according to the sizes of the proposals, alleviating the training burden for the large intra-class variation of objects with different scales. In this case, we can improve the detection precision for both large objects and small objects. Lastly, we fuse all the predicted results from multiple branches into the final detection result. The deconvolution with bilinear kernels and the multi-branch decision network do not increase processing time because the former just deals with small proposals without enlarging the whole feature maps, and the latter processes the same number of proposals as traditional detection methods.

### B. Context-aware RoI Pooling

The context-aware RoI pooling (CARoI pooling) can adjust the proposals to the specified size without sacrificing important contextual information (as illustrated in Fig. 3 (d)).

In CARoI pooling, we have three cases to deal with. Firstly, if the size of a proposal is larger than the specified size, we shall extract the maximum value in each sub-window as original RoI pooling strategy (as described in Section III-A). Secondly, if the size of a proposal is smaller than the specified size, a deconvolution operation with bilinear kernel is applied to enlarge the proposal while keeping the circumstances from being impaired so that we can still extract discriminative features from the small proposals. The size of deconvolution kernel is dynamically determined by the proposal size and the predefined pooled size. Specifically, the kernel size is equal to the ratio between the specified size of pooled feature map and the size of each proposal. Thirdly, when the width of a proposal is larger than the pooled length and the height of this proposal is smaller than the pooled length, our CARoI pooling applies the deconvolution operation to enlarge the height of this proposal, splits the width of this proposal into several sub-windows (the number of the sub-windows is equal to the pooled length) and uses the maximum value of each sub-window as the most discriminative feature value.

Mathematically, we formulate the three cases mentioned above in the following equations. Let $y_k^j$ be the $j$-th output of CARoI pooling layer from the $k$-th proposal. The CARoI pooling computes $y_k^j = x_{i^*}$, where:

$$i^* = argmax_{i \in R(k,j)} x_i \tag{1}$$

$$x_i \in (X_k \otimes \sigma_k) \tag{2}$$

In above equations, $R(k, j)$ represents the index set of the sub-window where the output unit $y_k^j$ selects the maximum feature value. $x_i \in \mathbb{R}$ is the $i$-th feature value on the feature map. And we use the $X_k$ to represent a set of input features of $k$-th proposal. $\otimes$ denotes the deconvolution operation and $\sigma_k$ is the kernel of the deconvolution operation, which is determined by the scales of proposals. If the size of proposal is less than the pooled feature map size, this deconvolution kernel is equal to the ratio between the specified size of pooled feature map and the size of each proposal; otherwise, this deconvolution kernel is equal to one, which suggests this deconvolution operation doesn't take effect on the large proposals. After obtaining the discriminate features, the maximum values of these features in each sub-window are used to represent this proposal.

**Back-propagation.** Derivatives are diverted through CARoI pooling by back propagation to train the network. The partial derivative of loss $L$ respective to input variable $x_i$ is:

$$\frac{\partial L}{\partial x_i} = \sum_k \sum_j [i = i^*] \nabla_{\sigma_k} \left( \frac{\partial L}{\partial y_k^j} \right) \qquad (3)$$

where $i^*$ is the index described in Equation 1, which indicates the position of maximum values in each sub-window after the deconvolution. $\nabla_{\sigma_k}(\frac{\partial L}{\partial y_k^j})$ indicates the derivative of the deconvolution with respect to the loss $\frac{\partial L}{\partial y_k^j}$. This loss is propagated from following layers that are connected to CARoI pooling. This derivative will be accumulated by all RoIs and all positions ($\sum_k \sum_j$).

### C. Multi-branch Decision Network

As analyzed in Section III-B, another critical issue for CNN-based object detection is the large scale variations of targeting object, which is common in vehicle detection. To reduce the scale variance of objects, we present to split the proposals with different sizes into different branches and each branch is used to detect a set of objects with similar sizes. Each branch consists of one convolutional layer and one fully connected layer followed by two classifiers: one is used for classification; the other is used for bounding box regression (see [23] for details). Although we split the proposals into multiple branches, these proposals share the features extracted by some convolutional layers (as the blue boxes shown in Fig. 2).

The number of branches is empirically determined by considering the scale distribution of the dataset and the computational resources, which is discussed in section V-D. Here we take the two-branch decision network as an example but this technique can be easily extended to multi-branch decision network. In two-branch decision network, we mainly use the median value of all objects' scales in the training set as the reference threshold to split the proposals into the large branch or the small branch. During the training process, in order to make two branches share a portion of samples in the median scales and augment the size of training samples for each branch, the threshold for splitting proposals is dynamically changed in each training iteration. We simulate the threshold change by a Gaussian model, and the median value of all objects' scales is the mean value of the Gaussian model. In such a way, those proposals with the scales that are near to the median value of all objects' scales have opportunity to be categorized into the large and the small branches in the whole training procedure. In testing, we simply use median value to split the proposals.

### D. Implementation Details

**Network architecture.** In principle, our context aware RoI pooling and multi-branch decision network are general and can be built on any CNN architectures. In this paper, we test our algorithms based on the PVA network [42] and VGG network [47]. The kernel sizes of CARoI pooling are set to $6 \times 6$ in the PVA network and $7 \times 7$ in the VGG network.

We use the proposal extraction network (RPN) proposed by MS-CNN [28] to extract high-quality proposals from different layers of the CNN. Then we connect the multi-branch decision network at the end of the RPN to build the SINet as shown in Fig. 2. The whole network is trained in an end-to-end manner.

**Training strategies.** Stochastic gradient descent (SGD) is used to optimize our SINet. In order to make the training process stable, we first harness a small learning rate to train the RPN and then leverage a large learning rate to train the whole network end to end. We first set the learning rate as 0.0001 for 10k iterations with weight decay 0.0005 to train the RPN. Then, to train the whole network, we set the learning rate as 0.0005, reduce it by a factor of 0.1 at 40k and 70k iterations and stop learning after 75k iterations. If employing VGG net, we adjust the initial learning rate to 0.00005 and 0.0001 for the first and second stages respectively. To accelerate training and reduce overfitting [48], the weights of convolutional layers in VGG trained on ImageNet [49] or PVA trained on Pascal VOC [50] are used to initialize the RPN. Then, we utilize the well-trained weights of the fully connected layers in VGG and PVA to initialize the fully connected layers of the newly added multi-branch decision network. Other layers are initialized by random noise. There are four images in each batch. In addition, data argumentation methods and hard example mining strategies [28] are also used as in the MS-CNN.

**Inference.** In testing, for each input image, the network produces outputs of small and large objects in multiple branches. Then we combine them together and use non-maximum suppression (NMS) to refine the results. Instead of selecting only the bounding box with the maximum confidence from highly overlapping detection boxes, we choose several bounding boxes with the relatively high confidences among these boxes and average the coordinates of them. We call this strategy as *soft-NMS*, which is useful to improve the localization accuracy for occluded vehicles.

## V. EXPERIMENTS

In order to evaluate the effectiveness of the proposed SINet, we conduct experiments on two representative vehicle datasets: the KITTI dataset and a newly constructed large scale variance highway dataset (LSVH). The experiments are implemented on Ubuntu 14.04 with a single GPU (NVIDIA TITAN X) and 8 CPUs (Inter(R) Xeon(R) E5-1620 v3 @ 3.50GHz).

### A. Datasets and Evaluation Metrics

**KITTI dataset.** KITTI [21] is a widely used benchmark for vehicle detection algorithms. It contains various scales of vehicles in different scenes. It consists of 7481 images for training and 7518 images for testing. According to size, occlusion and truncation, the organizers classify these targeting vehicles into three difficulty levels: easy, moderate and hard; check [21] for detailed definition of these difficulty levels.

**LSVH dataset.** Highway is a typical road scene that contains vehicles with large scale variations, as the surveillance cameras usually cover a large and long view of the road. We construct a new large-scale variance highway dataset,
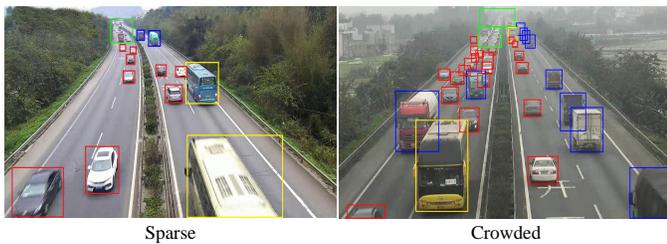
| Sparse | Crowded |

Fig. 5. Examples of our large scale variance highway (LSVH) dataset under the different scenes. Red, blue, orange and green boxes are labelled to indicate the "car", "van", "bus", and "don't care" regions respectively.

TABLE I
DATA DISTRIBUTION ON LSVH DATASET.

|  | Sparse | Crowded | Total |
|---|---|---|---|
| Image | 12979 | 1409 | 14388 |
| Car | 40025 | 35116 | 75141 |
| Bus | 4419 | 2480 | 6899 |
| Van | 10474 | 4812 | 15286 |
| Vehicle/Image | 4.23 | 30.10 | 6.76 |

which contains 16 videos captured under different scenes, time, weathers and resolutions, as shown in Fig. 5 and 6. As illustrated in Fig. 5 and Table I, the vehicle is classified into three categories (car, bus and van) under two scenes (sparse and crowded). We consider a video scene as a crowded scene in case that it contains more than 15 vehicles per frame on average; otherwise, it is considered as a sparse scene. For the vehicles that are too small to be recognized by human are labelled as "don't care", and these regions are ignored during training and evaluation. Specifically, an object whose height is less than 15 pixels will be ignored. There are 75141 cars, 6899 buses and 15286 vans in total in our LSVH. We use two strategies to split the videos into training/testing sets: (1): we separate each video into two parts, and select the first seventy percentages of each video as the training data and the remaining thirty percentages as the testing data; (2) we use the eight videos in "Sparse" as the training data, and the left four videos in "Sparse" as the testing data. To avoid retrieving similar images, we extract one frame in every seven frames of these videos as the training/testing images.

**Evaluation metrics.** We employ the well established average precision (AP) and intersection over union (IoU) metrics [50] to evaluate the performance; it has been widely used to assess various vehicle detection algorithms [21, 50]. For KITTI, we evaluate our method in all three difficulty levels. For LSVH, we evaluate the performance for car, bus and van under the scenes of sparse and crowded, respectively. The IoU is set to 0.7 in these two datasets, which means only the overlap between the detection bounding box and the ground truth bounding box greater than or equal to 70% is considered as a correct detection.

### B. Comparison with the State-of-the-arts

We compare the proposed SINet to the state-of-the-arts on both the KITTI dataset and our LSVH dataset. Table II shows the performance published on the KITTI website. In this experiment, the entire training set is used for training our models, and the tested results are uploaded to the KITTI website. We

TABLE II
RESULTS ON THE KITTI BENCHMARK. ALL METHODS ARE RANKED BASED ON THE "MODERATE".

| Model | Time/Image | Average Precision (%) | | |
|---|---|---|---|---|
|  |  | Moderate | Easy | Hard |
| **SINet_VGG (ours)** | 0.2s | **89.60** | 90.60 | 77.75 |
| **SINet_PVA (ours)** | **0.11s** | 89.21 | 91.91 | 76.33 |
| Deep3DBox [51] | 1.5s | 89.04 | 92.98 | 77.17 |
| SubCNN [29] | 2s | 89.04 | 90.81 | 79.27 |
| MS-CNN [28] | 0.4s | 89.02 | 90.03 | 76.11 |
| SDP+RPN [27, 30] | 0.4s | 88.85 | 90.14 | 78.38 |
| Mono3D [52] | 4.2s | 88.66 | 92.33 | 78.96 |
| 3DOP [53] | 3s | 88.64 | **93.04** | 79.10 |
| MV3D [54] | 0.45s | 87.67 | 89.11 | **79.54** |
| SDP+CRF (ft) [30] | 0.6s | 83.53 | 90.33 | 71.13 |
| Faster R-CNN [27] | 2s | 81.84 | 86.71 | 71.12 |
| MV3D (LIDAR) [54] | 0.3s | 79.24 | 87.00 | 78.16 |
| spLBP [55] | 1.5s | 77.39 | 87.18 | 60.59 |
| Reinspect [56] | 2s | 76.65 | 88.13 | 66.23 |
| Regionlets [57–59] | 1s | 76.45 | 84.75 | 59.70 |
| AOG [1, 2] | 3s | 75.94 | 84.80 | 60.70 |
| 3DVP [60] | 40s | 75.77 | 87.46 | 65.38 |
| SubCat [61] | 0.7s | 75.46 | 84.14 | 59.71 |
| YOLOv2 [45] | **0.03s** | 61.31 | 76.79 | 50.25 |
| YOLO [44] | **0.03s** | 35.74 | 47.69 | 29.65 |

compare our models with other 18 published methods. It is clear that our SINet achieves the highest accuracy on moderate case and fastest speed (except the one-stage deep learning based detectors, e.g. YOLO and YOLOv2, which are fast but with very low accuracy). Our method can achieve the same speed as YOLO and YOLOv2 by reducing the size of input images and the accuracy is still much better than these two methods (see Section V-C). For the computational efficiency among the two-stage deep learning based detectors, our SINet takes only $1/14$ of Deep3DBox [51].

Table III shows the performance on our LSVH dataset. It is obvious that both two variants of our SINet outperforms the MS-CNN baseline and Faster RCNN in terms of detection accuracy and efficiency. Our SINet also surpasses the one-stage detectors (YOLO and YOLOv2) by a significant margin for the accuracy. In particular, the SINet shows a good performance to detect the vehicles under the "Crowded" scene.

In Fig. 6, we visualize the vehicles detected by SINet on the images from KITTI dataset and our LSVH dataset. It is clear that our algorithm is effective to detect the vehicles with different orientations, scales, truncation levels under the different situations such as blurry, rainy, and occluded. Moreover, our SINet shows a strong ability on detecting vehicles with a large variance of scales, especially for small vehicles. This corroborates that the presented SINet has potential to be used as a powerful tool for intelligent transportation systems.

### C. Image Resolution Sensitivity

Since our SINet has a strong capability on feature representation for low resolution vehicles, it can also perform well on the low resolution images. This resolution insensitive property is actually very important for practical usage, and enabling fast

TABLE III
COMPARISON ON OUR LSVH DATASET. WE USE TWO STRATEGIES TO SPLIT THE DATASET (SEE SECTION V-A FOR DETAILS).

| Model | Time/Image | Strategy 1 | | | | | | | Strategy 2 | | | |
| | | Mean | Sparse | | | Crowded | | | Sparse | | | |
| | | | Car | Bus | Van | Car | Bus | Van | Mean | Car | Bus | Van |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **SINet_VGG (ours)** | 0.20s | **70.17** | **81.82** | **85.60** | **78.65** | **56.80** | **55.78** | 62.38 | **78.71** | 74.51 | **84.34** | **77.27** |
| **SINet_PVA (ours)** | **0.08s** | 70.04 | 81.40 | 84.39 | 77.39 | 53.76 | 54.06 | **69.23** | 77.69 | **74.79** | 81.35 | 76.92 |
| MS-CNN [28] | 0.23s | 63.23 | 79.94 | 83.71 | 76.79 | 51.74 | 32.95 | 54.26 | 72.66 | 71.13 | 74.34 | 72.50 |
| Faster RCNN [27] | 0.31s | 46.44 | 60.93 | 66.68 | 60.14 | 26.08 | 24.55 | 40.24 | 40.22 | 36.19 | 42.79 | 41.69 |
| YOLOv2 [45] | **0.03s** | 43.82 | 59.71 | 65.51 | 58.35 | 17.39 | 21.55 | 40.42 | 54.00 | 53.16 | 53.88 | 54.96 |
| YOLO [44] | **0.03s** | 16.53 | 23.06 | 31.13 | 22.44 | 3.87 | 8.35 | 10.32 | 23.78 | 22.97 | 24.52 | 23.85 |



Fig. 6. Examples of detection results by our SINet on the KITTI dataset (the first two rows) and our LSVH dataset (the last two rows). (Best viewed in color and at full size on a high-resolution display.)



Fig. 7. Image resolution sensitivity evaluation. These experiments were done on the KITTI training set.

computation by resizing an image to a small resolution. Fig. 7 illustrates our SINet is insensitive to the image resolution. The detection performance is robust with different sizes of input images. On the contrary, MS-CNN [28] is sensitive to the resolution of input images. A small input image decreases the accuracy dramatically, while increasing the input resolution leads to much more computational overhead.

### D. Ablation Analysis

We perform ablation analysis of SINet on the KITTI dataset to evaluate how different components affect the detection performance. As there is no ground truth provided for the testing set of KITTI, we follow [28] to split the training set into training and validation sets, and all of them are resized to $576 \times 1920$.

Table IV shows the experimental results. First, comparing with the baselines which are constructed by the MS-CNN framework with PVA (the $1^{st}$ row) or VGG (the $6^{th}$ row) network, the CARoI pooling dramatically improves the accuracy while no extra time is introduced, as shown in the $2^{nd}$ and $7^{th}$ row. Particularly, the improvements on "Moderate" and "Hard" categories are significant, which implies the recovered high resolution semantic features are very useful for detecting
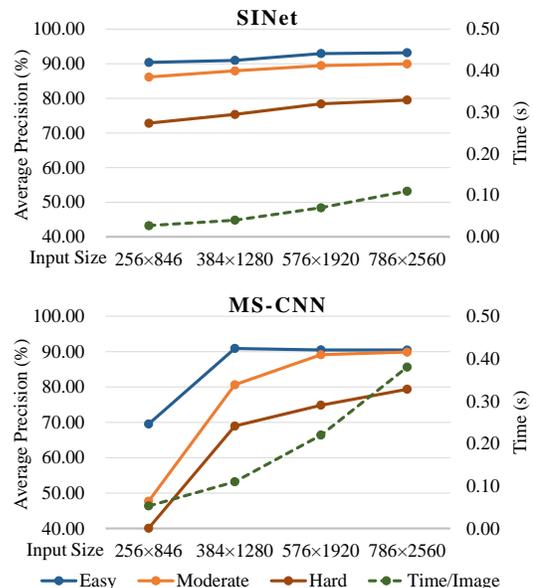
small objects. Moreover, our multi-branch decision network with two branches further enhances the accuracy while keeping the efficiency, as shown in the $3^{rd}$ and $8^{th}$ row. The soft-NMS post-processing contributes to the "Hard" category (the $4^{th}$ and $9^{th}$ row), which includes many occluded and truncated vehicles, demonstrating its effectiveness for occlusion and truncation cases. When we continue to increase the number of branches, the performance gain is limited but with more network parameters which occupy more memory. In this case, two-branch decision network is used in other experiments.

### E. Vehicle Scale Analysis.

We explore the detection performance of SINet on different scales of vehicles. This experiment is performed on our LSVH dataset which contain vehicles with a large variance of scales (as illustrated in Fig. 5). All vehicles are divided into three categories: "Small", "Medium" and "Large" based on the their scales. Specifically, the vehicles whose heights are greater than 15 pixels and smaller than 39 pixels belong to the "Small" category; the vehicles with the height between 39 pixels and 66 pixels are in "Medium" category; other vehicles with height greater than 66 pixels belong to "Large" category.

As shown in Table V, our SINet shows improvements on all scales of vehicles under the different scenes based on

TABLE IV
ABLATION ANALYSIS OF THE PRESENTED SINET ON THE KITTI VALIDATION SET. THE TIME IS EVALUATED ON A SINGLE NVIDIA TITAN X GPU (MAXWELL VERSION) WITH 12$GB$ MEMORY AND ONLY ONE IMAGE IS PROCESSED AT THE SAME TIME.

| Model | Network | Number of Branches | Post Processing | Average Precision (%) | | | Time/Image |
|---|---|---|---|---|---|---|---|
| | | | | Moderate | Easy | Hard | |
| MS-CNN_PVA | PVA | 1 | NMS | 82.69 | 91.85 | 69.74 | **0.07s** |
| +CARoI pooling | PVA | 1 | NMS | 88.39 | 92.52 | 75.70 | **0.07s** |
| +Multi-branch decision network | PVA | 2 | NMS | 89.36 | 92.96 | 78.11 | **0.07s** |
| **SINet_PVA** | PVA | 2 | Soft-NMS | 89.49 | 92.95 | 78.45 | **0.07s** |
| **SINet_PVA** | PVA | 3 | Soft-NMS | **89.53** | **93.31** | **78.53** | **0.07s** |
| MS-CNN_VGG [28] | VGG | 1 | NMS | 89.13 | 90.49 | 74.85 | 0.22s |
| +CARoI pooling | VGG | 1 | NMS | 90.07 | 95.30 | 79.31 | **0.20s** |
| +Multi-branch decision network | VGG | 2 | NMS | 90.22 | 95.82 | 80.02 | **0.20s** |
| **SINet_VGG** | VGG | 2 | Soft-NMS | **90.33** | **95.84** | **80.14** | **0.20s** |
| **SINet_VGG** | VGG | 3 | Soft-NMS | - | - | - | Out of memory |

TABLE V
VEHICLE SCALE ANALYSIS ON LSVH TESTING SET (STRATEGY 1, SUNNY). THE SIZE OF THE INPUT IMAGE IS $768 \times 1344$.

| Scene | | Sparse | | | | Crowded | | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | | SINet_PVA | MSCNN_PVA | SINet_VGG | MSCNN_VGG | SINet_PVA | MSCNN_PVA | SINet_VGG | MSCNN_VGG |
| Small | Car | 70.25 | 47.39 | **72.46** | 68.91 | 13.14 | 2.38 | **14.50** | 9.00 |
| | Bus | 55.06 | 30.22 | **59.54** | 57.37 | 20.10 | 3.48 | **22.36** | 5.37 |
| | Van | 48.34 | 25.84 | **51.40** | 48.78 | - | - | - | - |
| Medium | Car | 87.68 | 74.45 | **88.44** | 85.54 | 55.60 | 22.46 | **61.06** | 55.98 |
| | Bus | 86.02 | 70.46 | **90.31** | 86.13 | 25.45 | 10.48 | **29.04** | 12.18 |
| | Van | 73.63 | 57.67 | **76.08** | 75.20 | **17.16** | 3.87 | 10.16 | 3.58 |
| Large | Car | **84.72** | 78.94 | 83.13 | 77.27 | 82.26 | 59.18 | **83.25** | 80.51 |
| | Bus | 90.13 | 82.01 | **90.93** | 88.75 | **65.34** | 44.94 | 60.71 | 40.43 |
| | Van | 81.96 | 74.22 | **85.16** | 79.22 | **77.20** | 66.86 | 70.04 | 63.51 |

both PVA network and VGG network. The improvement on small vehicles is more significant compared with other sizes of vehicles, since the baseline methods introduce more artifacts and distortions (caused by the traditional RoI pooling) to small vehicles, which can be avoided by the CARoI pooling. Moreover, our SINet also achieves a dramatic improvement on the crowded scenes, especially for the vehicle with a small or medium scale. It shows that our approach is effective even under the complex situation, as shown in Fig. 6. However, the detection accuracy for the small scale crowded vehicles is still not satisfactory. This is because these objects are highly occluded, blurry and extremely small (Fig. 5 and Fig. 6).

## VI. CONCLUSION

In this paper, we present a scale-insensitive network, denoted as SINet, for fast detecting vehicles with a large variance of scales. Two new techniques, context-aware RoI pooling and multi-branch decision network, are presented to maintain the original structures of small objects and minimize the intra-class distances among objects with a large variance of scales. Both of the techniques require zero extra computational effort. Furthermore, we construct a new highway dataset which contains vehicles with large scale variance. To our knowledge, it is the first large scale dataset focuses on the highway scene. Our SINet achieves state-of-the-art performance on both accuracy and speed on KITTI benchmark and our LSVH dataset. Further investigations include evaluating the SINet

on more challenging datasets and integrating it into some intelligent transportation systems.

## REFERENCES

[1] B. Li, T. Wu, and S.-C. Zhu, "Integrating context and occlusion for car detection by hierarchical and-or model," in *ECCV*, 2014.

[2] T. Wu, B. Li, and S.-C. Zhu, "Learning and-or model to represent context and occlusion for car detection and viewpoint estimation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 9, pp. 1829–1843, 2016.

[3] C. Stauffer and W. E. L. Grimson, "Adaptive background mixture models for real-time tracking," in *CVPR*, 1999.

[4] Z. Chen, T. Ellis, and S. A. Velastin, "Vehicle detection, tracking and classification in urban traffic," in *ITSC*, 2012.

[5] C. Premebida and U. Nunes, "A multi-target tracking and gmm-classifier for intelligent vehicles," in *ITSC*, 2006.

[6] E. Martinez, M. Diaz, J. Melenchon, J. A. Montero, I. Iriondo, and J. C. Socoro, "Driving assistance system based on the detection of head-on collisions," in *Intelligent Vehicles Symposium*, 2008.

[7] J. Cui, F. Liu, Z. Li, and Z. Jia, "Vehicle localisation using a single camera," in *Intelligent Vehicles Symposium*, 2010.

[8] S. Kyo, T. Koga, K. Sakurai, and S. Okazaki, "A robust vehicle detecting and tracking system for wet weather conditions using the imap-vision image processing board," in *ITSC*, 1999.

[9] M. Vargas, J. M. Milla, S. L. Toral, and F. Barrero, "An enhanced background estimation algorithm for vehicle detection in urban traffic scenes," *Vehicular Technology*, vol. 59, no. 8, pp. 3694–3709, 2010.

[10] Z. Sun, G. Bebis, and R. Miller, "On-road vehicle detection using optical sensors: A review," in *ITSC*, 2004.

[11] Q. Yuan, A. Thangali, V. Ablavsky, and S. Sclaroff, "Learning a family of detectors via multiplicative kernels," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 3, pp. 514–530, 2011.

[12] H. T. Niknejad, A. Takeuchi, S. Mita, and D. McAllester, "On-road multivehicle tracking using deformable object model and particle filter with improved likelihood estimation," *IEEE Transactions on Intelligent Transportation Systems*, vol. 13, no. 2, pp. 748–758, 2012.

[13] Z. Sun, G. Bebis, and R. Miller, "Monocular precrash vehicle detection: Features and classifiers," *IEEE Transactions on Image Processing*, vol. 15, no. 7, pp. 2019–2034, 2006.

[14] J.-W. Hsieh, L.-C. Chen, and D.-Y. Chen, "Symmetrical surf and its applications to vehicle detection and vehicle make and model recognition," *IEEE Transactions on Intelligent Transportation Systems*, vol. 15, no. 1, pp. 6–20, 2014.

[15] W.-C. Chang and C.-W. Cho, "Online boosting for vehicle detection," *Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 40, no. 3, pp. 892–902, 2010.

[16] S. Sivaraman and M. M. Trivedi, "A general active-learning framework for on-road vehicle recognition and tracking," *IEEE Transactions on Intelligent Transportation Systems*, vol. 11, no. 2, pp. 267–276, 2010.

[17] J. J. Yebes, L. M. Bergasa, R. Arroyo, and A. Lázaro, "Supervised learning and evaluation of kitti's cars detector with dpm," in *Intelligent Vehicles Symposium Proceedings*, 2014.

[18] Q. Wang, J. Gao, and Y. Yuan, "A joint convolutional neural networks and context transfer for street scenes labeling," *IEEE Transactions on Intelligent Transportation Systems*, 2017.

[19] Y. Yuan, Z. Xiong, and Q. Wang, "An incremental framework for video-based traffic sign detection, tracking, and recognition," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 7, pp. 1918–1929, 2017.

[20] Q. Wang, J. Gao, and Y. Yuan, "Embedding structured contour and location prior in siameed fully convolutional networks for road detection," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 1, pp. 230–241, 2018.

[21] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *CVPR*, 2012.

[22] S. Bell, C. Lawrence Zitnick, K. Bala, and R. Girshick, "Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks," in *CVPR*, 2016.

[23] R. Girshick, "Fast r-cnn," in *ICCV*, 2015.

[24] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *CVPR*, 2014.

[25] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.

[26] K. He, X. Zhang, S. Ren, and J.Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.

[27] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NIPS*, 2015.

[28] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *ECCV*, 2016.

[29] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Subcategory-aware convolutional neural networks for object proposals and detection," in *WACV*, 2017.

[30] F. Yang, W. Choi, and Y. Lin, "Exploit all the layers: Fast and accurate cnn object detector with scale dependent pooling and cascaded rejection classifiers," in *CVPR*, 2016.

[31] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. LeCun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," *arXiv preprint arXiv:1312.6229*, 2013.

[32] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[33] A. Shrivastava, A. Gupta, and R. Girshick, "Training region-based object detectors with online hard example mining," in *CVPR*, 2016.

[34] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," *arXiv preprint arXiv:1611.05431*, 2016.

[35] T. Kong, A. Yao, Y. Chen, and F. Sun, "Hypernet: towards accurate region proposal generation and joint object detection," in *CVPR*, 2016.

[36] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," *arXiv preprint arXiv:1612.03144*, 2016.

[37] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *ECCV*, 2016.

[38] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 640–651, 2017.

[39] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta, "Beyond skip connections: Top-down modulation for object detection," *arXiv preprint arXiv:1612.06851*, 2016.

[40] S. Zagoruyko, A. Lerer, T.-Y. Lin, P. O. Pinheiro, S. Gross, S. Chintala, and P. Dollár, "A multipath network for object detection," *arXiv preprint arXiv:1604.02135*, 2016.

[41] S. Sivaraman and M. M. Trivedi, "Looking at vehicles on the road: A survey of vision-based vehicle detection, tracking, and behavior analysis," *IEEE Transactions on Intelligent Transportation Systems*, vol. 14, no. 4, pp. 1773–1795, 2013.

[42] K.-H. Kim, S. Hong, B. Roh, Y. Cheon, and M. Park, "Pvanet: Deep but lightweight neural networks for real-time object detection," *arXiv preprint arXiv:1608.08021*, 2016.

[43] Y. Li, K. He, J. Sun *et al.*, "R-fcn: Object detection via region-based fully convolutional networks," in *NIPS*, 2016.

[44] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *CVPR*, 2016.

[45] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *CVPR*, 2017.

[46] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama *et al.*, "Speed/accuracy trade-offs for modern convolutional object detectors," *arXiv preprint arXiv:1611.10012*, 2016.

[47] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[48] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *NIPS*, 2014.

[49] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International*

*Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.

[50] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.

[51] A. Mousavian, D. Anguelov, J. Flynn, and J. Kosecka, "3d bounding box estimation using deep learning and geometry," *arXiv preprint arXiv:1612.00496*, 2016.

[52] X. Chen, K. Kundu, Z. Zhang, H. Ma, S. Fidler, and R. Urtasun, "Monocular 3d object detection for autonomous driving," in *CVPR*, 2016.

[53] X. Chen, K. Kundu, Y. Zhu, A. G. Berneshawi, H. Ma, S. Fidler, and R. Urtasun, "3d object proposals for accurate object class detection," in *NIPS*, 2015.

[54] X. Chen, H. Ma, J. Wan, B. Li, and T. Xia, "Multi-view 3d object detection network for autonomous driving," *arXiv preprint arXiv:1611.07759*, 2016.

[55] Q. Hu, S. Paisitkriangkrai, C. Shen, A. van den Hengel, and F. Porikli, "Fast detection of multiple objects in traffic scenes with a common detection framework," *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 4, pp. 1002–1014, 2016.

[56] R. Stewart, M. Andriluka, and A. Y. Ng, "End-to-end people detection in crowded scenes," in *CVPR*, 2016.

[57] C. Long, X. Wang, G. Hua, M. Yang, and Y. Lin, "Accurate object detection with location relaxation and regionlets relocalization," in *ACCV*, 2014.

[58] X. Wang, M. Yang, S. Zhu, and Y. Lin, "Regionlets for generic object detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 10, pp. 2071–2084, 2015.

[59] W. Y. Zou, X. Wang, M. Sun, and Y. Lin, "Generic object detection with dense neural patterns and regionlets," *arXiv preprint arXiv:1404.4316*, 2014.

[60] Y. Xiang, W. Choi, Y. Lin, and S. Savarese, "Data-driven 3d voxel patterns for object category recognition," in *CVPR*, 2015.

[61] E. Ohn-Bar and M. M. Trivedi, "Learning to detect vehicles by clustering appearance patterns," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 5, pp. 2511–2521, 2015.

**Xiaowei Hu** received the B.Eng. degree in the Computer Science and Technology from South China University of Technology, China, in 2016. He is currently working toward the Ph.D. degree with the Department of Computer Science and Engineering, The Chinese University of Hong Kong. His research interests include computer vision and deep learning. Mr. Hu is a recipient of the Hong Kong Ph.D. Fellowship.

**Xuemiao Xu** received her B.S. and M.S. degrees in Computer Science and Engineering from South China University of Technology in 2002 and 2005 respectively, and Ph.D. degree in Computer Science and Engineering from The Chinese University of Hong Kong in 2009. She is currently a professor in the School of Computer Science and Engineering, South China University of Technology. Her research interests include object detection, tracking, recognition, and image, video understanding and synthesis, particularly their applications in the intelligent transportation system.
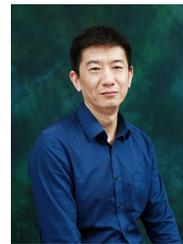
**Yongjie Xiao** received the B.Eng. degree in the Computer Science and Technology from South China University of Technology, China, in 2016. He is currently working toward the Master degree with the School of Computer Science and Engineering, South China University of Technology. His research interests include intelligent transportation, object detection and deep learning.

**Hao Chen** received the Ph.D. degree in Computer Science and Engineering from The Chinese University of Hong Kong, China, in 2017. He is currently a post-doctoral fellow in The Chinese University of Hong Kong. His research interests include medical image analysis, deep learning, and health informatics. Dr. Chen was a recipient of the Hong Kong Ph.D. Fellowship.

**Shengfeng He** obtained his B.Sc. degree and M.Sc. degree from Macau University of Science and Technology and his Ph.D. degree from City University of Hong Kong. He is an Associate Professor in the School of Computer Science and Engineering at South China University of Technology. He was a Research Fellow at City University of Hong Kong and a visiting Ph.D. student at Georgia Institute of Technology. His research interests include computer vision, image processing, computer graphics, and deep learning.

**Jing Qin** received his Ph.D. degree in Computer Science and Engineering from the Chinese University of Hong Kong in 2009. He is currently an assistant professor in School of Nursing, The Hong Kong Polytechnic University. He is also a key member in the Centre for Smart Health, SN, PolyU, HK. His research interests include innovations for healthcare and medicine applications, medical image processing, deep learning, visualization and human-computer interaction and health informatics.

**Pheng-Ann Heng** received his B.Sc. from the National University of Singapore in 1985. He received his MSc (Comp. Science), M. Art (Applied Math) and Ph. D (Comp. Science) all from the Indiana University of USA in 1987, 1988, 1992 respectively. He is a professor at the Department of Computer Science and Engineering at The Chinese University of Hong Kong (CUHK). He has served as the Director of Virtual Reality, Visualization and Imaging Research Center at CUHK since 1999 and as the Director of Center for Human-Computer Interaction at Shenzhen Institute of Advanced Integration Technology, Chinese Academy of Science/CUHK since 2006. He has been appointed as a visiting professor at the Institute of Computing Technology, Chinese Academy of Sciences as well as a Cheung Kong Scholar Chair Professor by Ministry of Education and University of Electronic Science and Technology of China since 2007. His research interests include AI and VR for medical applications, surgical simulation, visualization, graphics and human-computer interaction.