# End-to-end trainable network for degraded license plate detection via vehicle-plate relation mining

Song-Lu Chen[a,b], Shu Tian[a], Jia-Wei Ma[a,b], Qi Liu[a,b], Chun Yang[a,b], Feng Chen[b,c] and Xu-Cheng Yin[a,b]

[a]*School of Computer and Communication Engineering, University of Science and Technology Beijing, Beijing 100083, China*
[b]*USTB-EEasyTech Joint Lab of Artificial Intelligence, University of Science and Technology Beijing, Beijing 100083, China*
[c]*EEasy Technology Company Ltd., Zhuhai 519000, China*

## ABSTRACT

License plate detection is the first and essential step of the license plate recognition system and is still challenging in real applications, such as on-road scenarios. In particular, small-sized and oblique license plates, mainly caused by the distant and mobile camera, are difficult to detect. In this work, we propose a novel and applicable method for degraded license plate detection via vehicle-plate relation mining, which localizes the license plate in a coarse-to-fine scheme. First, we propose to estimate the local region around the license plate by using the relationships between the vehicle and the license plate, which can greatly reduce the search area and precisely detect very small-sized license plates. Second, we propose to predict the quadrilateral bounding box in the local region by regressing the four corners of the license plate to robustly detect oblique license plates. Moreover, the whole network can be trained in an end-to-end manner. Extensive experiments verify the effectiveness of our proposed method for small-sized and oblique license plates. Codes are available at: https://github.com/chensonglu/LPD-end-to-end.

## 1. Introduction

License plate detection (LPD) has attracted great interest from academia and industry for many years owing to its importance in many practical applications, such as toll control, parking lot access, and traffic law enforcement. Accurate license plate detection is crucial for subsequent license plate recognition [37]. However, it remains a challenging task due to illumination variations, background changes, size variations, and viewpoint changes.

Before the deep learning era, most methods [36, 1, 7, 2] need to design handcrafted features for license plate detection. Recently, deep learning methods [4, 32, 31] have contributed to improving the license plate detection task. Many methods [4, 43, 41, 18] propose to localize the license plate directly from the input image, but these methods can not detect small-sized license plates properly, because the license plate is only a small part of the input image. There have been many previous works aiming at small-sized license plate detection by reducing the search area of the license plate using the vehicle proposal [32, 17, 9]. However, these methods can not handle large vehicles properly, such as buses and trucks, because their license plates are also a small part of them. Furthermore, [31] presents using the vehicle front region to further reduce the search area of the license plate. The vehicle front region is manually defined as the smallest region comprising the headlights and tires. However, [31] need to manually annotate the location of the vehicle front region, which is ambiguous and a waste of manpower.

Moreover, most previous methods [4, 43, 18, 41, 17, 31] simply consider the license plate in horizontal direction, which

is only applicable to limited scenes, such as highway bayonet charge and parking lot access. When it comes to more challenging scenes, such as on-road scenarios, they don't work for highly oblique license plates. Although in the literature [6, 32] there are some methods proposed to detect multi-oriented license plates, they are very complex due to adopting several separate networks.

In this work, we propose an end-to-end trainable network for degraded license plate detection via vehicle-plate relation mining, which can effectively detect the small-sized license plate and accurately localize the quadrilateral bounding box of the oblique license plate in real applications (e.g., on-road scenes). The method can detect the license plate in a coarse-to-fine manner. At the detection stage, we first estimate the location of the license plate based on the offset between the center of the license plate and the vehicle. Considering that the location obtained in this way is not always accurate, we refine the quadrilateral bounding box of the license plate in the local region around the license plate. The local region is simply obtained by expanding the background region around the license plate. In this way, many license plate regions of different vehicles in the input image can be obtained simultaneously, and they have various sizes and aspect ratios. To reduce the running time, all the estimated regions are scaled to the same size and aggregated together into LP patches, so all the license plates can be detected simultaneously in the LP patches. The aforementioned detection stages are combined to build an end-to-end network for license plate detection. Our method can greatly reduce the search area of the license plate, which can minimize false positives and improve the detection performance of small-sized and oblique license plates. Furthermore, estimating the local region can make LPD independent of the size of the vehicles, which is

advantageous to large vehicles.

Our main contributions can be summarized as:

- We propose a novel and applicable method for small-sized and oblique license plate detection by utilizing vehicle-plate relationships, where the license plate is precisely located in a coarse-to-fine scheme. Furthermore, the whole detection network is constructed as an end-to-end trainable network.

- We propose a novel method to estimate the local region around the license plate via vehicle-plate relation mining, which can greatly reduce the search area of the license plate.

- We propose a new method to localize the quadrilateral bounding box of the oblique license plate by regressing the four corners of the license plate.

The rest of this paper is organized as follows. Related work is described in Section 2. In Section 3, we describe our method in details. Section 4 presents comparative experiments and analyses. Final remarks are presented in Section 5.

## 2. Related Work

**Direct License Plate Detection** The following methods propose to localize the license plate directly from the input image. [4] detects the vehicle and the license plate with two independent branches to remove the effect that the vehicle suppresses the detection of the license plate. [43] presents a robust and efficient approach for license plate detection, which firstly accelerates the license plate localization using an effective image down-scaling method, then utilizes dense filters to extract candidate regions, and finally identifies the true license plates using a cascaded classifier. [41] presents to use multi-scale features to predict and regress the bounding box of the license plate. [18] proposes a method of detecting and recognizing the license plate, where the license plate is localized by Faster R-CNN [27] directly from the input image. However, these methods can not always detect small-sized license plates properly, because the license plate is only a small part of the input image[1].

**License Plate Detection with Vehicle Proposal** The following methods propose to reduce the search area of the license plate by the previous detection of vehicle, vehicle front region, or region around the license plate. In this way, it can improve the detection performance of small-sized license plates and reduce false positives of the license plate. [30] utilizes R-CNN[11] to generate vehicle proposals and then localizes the license plate in each vehicle region. [9] applies the Region Proposal Network (RPN) [27] to generate candidate vehicle proposals and then detects the license plate based on each proposal. [32] introduces a novel CNN framework capable of detecting the license plate in each predicted vehicle region. [17] utilizes YOLOv2 [25] to detect all the vehicles, then localizes the license plates in the vehicle patches simultaneously. However, these methods are still not favorable to large vehicles, such as trucks and buses, because their license plates are still only a small part of them. [31] proposes to detect the vehicle first, then detects the vehicle front region, and finally localizes the license plate in each vehicle front region. However, the vehicle front region needs to be annotated manually, which is a waster of manpower. [40] employs an attention-like method to estimate the local region around the license plate, then detects the license plate in the local region. However, in the literature [40], it utilizes low-resolution feature map and ROI pooling [10] for LPD, which causes loss of the spatial and semantic information, so the end-to-end model in [10] suffers significant performance degradation, especially for the large IOU threshold. Our method can maintain the semantic information by adopting the high-resolution feature map, and retain the spatial information by using space-invariant ROI warping [5].

**Multi-Oriented License Plate Detection** [40] proposes a CNN-based MD-YOLO framework for multi-directional license plate detection via rotation angle prediction. [12] proposes to detect the license plate with a tightly bounding parallelogram via predicting the top-left, top-right and bottom-right corners of the license plate. [35] utilizes semantic segmentation to get the rotation angle of the license plate. [40, 12, 35] regard the oblique license plate as a parallelogram. However, it is not always accurate due to the perspective transformation of the license plate, because a highly oblique license plate is more like a arbitrary quadrilateral. [6] proposes to generate license plate candidates with RPN [27], then uses R-CNN [11] to regress the four corners of the license plate. [32] proposes to obtain the affine transformation parameters explicitly based on Spatial Transformer Networks (STN) [14], which can transform the tilted license plate into a horizontal direction. However, [6, 32] are complicated due to adopting several separate networks. Our method can localize the four corners of the license plate in an end-to-end manner.

**Vehicle Detection** Before the deep learning era, most methods [34] usually utilized information about symmetry, color, shadow, geometrical features (e.g., corners, horizontal/vertical edges), texture features and vehicle lights for vehicle detection. Recently, deep learning methods [3, 24] have contributed to improving vehicle detection. [38, 29, 39] design better anchor priors for vehicle detection, which can facilitate the matching between the anchor box and the ground truth box. [13, 15] utilize multi-scale features to be robust to scale change of the vehicles by adopting YOLOv3 [26]. [22, 21] introduce a backward feature enhancement network to generate high-recall proposals, then adopt a spatial layout preserving network to enhance tiny vehicle detection. In this work, we simply adopt vanilla SSD [20] for vehicle detection, which can detect various-sized vehicles by utilizing multi-scale features. We will employ more powerful vehicle detectors in future work.

---

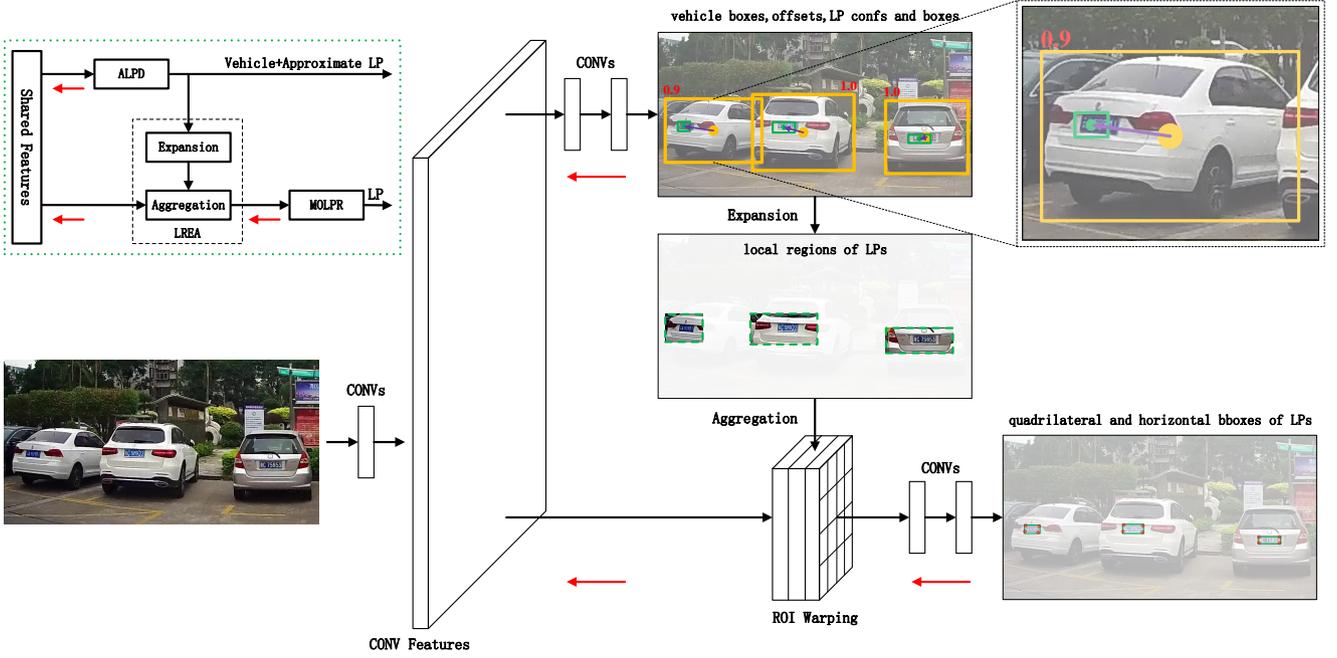[1]The average size of LPs is 0.26% of the full input image [31].

**Figure 1:** A thumbnail of the overall architecture is shown in the top-left corner (ALPD: approximate license plate detection; LREA: local region estimation and aggregation; MOLPR: multi-oriented license plate refinement). At the **ALPD** stage, first, the vehicle (orange rectangle) is detected, so the center of the vehicle (orange circle) is determined; second, the center of the license plate (green circle) is obtained based on the offset (purple arrow) between the center of the license plate and the vehicle; third, the size of the license plate is directly predicted from the input image. According to the center and size of the license plate, we can approximately estimate the license plate (green rectangle). Moreover, the probability of the vehicle containing a license plate (red number) is predicted simultaneously. An enlarged example is shown in the top-right corner. At the **LREA** stage, the local region of LP is obtained by expanding the background region around the license plate with a preset ratio, then all the expanded LP regions (green dashed rectangle) are aggregated into feature patches via ROI warping [5], in which all the LP regions have the same size and aspect ratio. Therefore, all the license plates can be detected simultaneously in the feature patches. At the **MOLPR** stage, the quadrilateral (red circle) and horizontal (green rectangle) bounding boxes of the license plate are detected simultaneously in the local region of LP. The network can be trained in an end-to-end manner, where the red arrows denote the backpropagation gradients.

## 3. Methodology

We propose an end-to-end trainable network for degraded license plate detection via vehicle-plate relation mining, which detects the license plate in a coarse-to-fine manner. The overall architecture is illustrated in Figure 1, where it firstly predicts the approximate location of the license plate utilizing spatial relationships between the license plate and the vehicle (Section 3.1), then estimates the local region by expanding the background region around the license plate followed by an aggregation operation (Section 3.2), and finally refines the quadrilateral bounding box of the license plate (Section 3.3).

### 3.1. Approximate License Plate Detection

At this stage, the vehicle is firstly detected, so the center of the vehicle is determined. After that, the network predicts the approximate location and size of the license plate, where the location is obtained based on the offset between the center of the license plate and the vehicle, and the size is directly predicted from the input image. Besides, the probability of the vehicle containing a license plate is predicted simultaneously. As shown in Figure 1, the location, and the size of the

license plate are not accurate in general cases, because they are directly predicted from the large input image, of which the license plate is only a small portion.

The ALPD network is based on SSD [20] for multi-task learning. The backbone network is shown in Table 1 without showing the ReLU activation function, and it is transformed from VGG-16 [33] followed by several extra layers. As for parameters, "k, s, d" mean kernel size, stride size, and dilation parameters [42] respectively. Moreover, we apply L2Norm [23] before combining the shallow features and the deep features, which avoids that the larger parameters "dominate" the smaller ones. As mentioned in Section 3.3, the first convolutional layer marked with "Δ" is the input layer of the MOLPR stage. Besides, layers marked with "*" are candidates for multi-scale detection.

The training objective of the ALPD network is defined as (1), which consists of the classification loss of the vehicle $L_{conf}(c)$, the regression loss of the vehicle $L_{loc}(l, g)$, the loss of whether the vehicle contains a license plate $L_{has\_lp}(v, lpc)$, the loss of the offset between the center of the licence plate and the vehicle $L_{off}(l, g, v)$, and the size loss of the license

**Table 1**
Backbone network of the ALPD network.

| Type | Filters | Parameters | Output |
|------|---------|------------|--------|
| Convolution△ | 64 | k:3,s:1 | $512 \times 512$ |
| Convolution | 64 | k:3,s:1 | $512 \times 512$ |
| Maxpool | - | k:2,s:2 | $256 \times 256$ |
| Convolution | 128 | k:3,s:1 | $256 \times 256$ |
| Convolution | 128 | k:3,s:1 | $256 \times 256$ |
| Maxpool | - | k:2,s:2 | $128 \times 128$ |
| Convolution | 256 | k:3,s:1 | $128 \times 128$ |
| Convolution | 256 | k:3,s:1 | $128 \times 128$ |
| Convolution | 256 | k:3,s:1 | $128 \times 128$ |
| Maxpool | - | k:2,s:2 | $64 \times 64$ |
| Convolution | 512 | k:3,s:1 | $64 \times 64$ |
| Convolution | 512 | k:3,s:1 | $64 \times 64$ |
| Convolution* | 512 | k:3,s:1 | $64 \times 64$ |
| Maxpool | - | k:2,s:2 | $32 \times 32$ |
| L2Norm | - | - | $32 \times 32$ |
| Convolution | 512 | k:3,s:1 | $32 \times 32$ |
| Convolution | 512 | k:3,s:1 | $32 \times 32$ |
| Convolution | 512 | k:3,s:1 | $32 \times 32$ |
| Maxpool | - | k:3,s:1 | $32 \times 32$ |
| Convolution | 1024 | k:3,s:1,d:6 | $32 \times 32$ |
| Convolution* | 1024 | k:1,s:1 | $32 \times 32$ |
| Convolution | 256 | k:1,s:1 | $32 \times 32$ |
| Convolution* | 512 | k:3,s:2 | $16 \times 16$ |
| Convolution | 128 | k:1,s:1 | $16 \times 16$ |
| Convolution* | 256 | k:3,s:2 | $8 \times 8$ |
| Convolution | 128 | k:1,s:1 | $8 \times 8$ |
| Convolution* | 256 | k:3,s:2 | $4 \times 4$ |
| Convolution | 128 | k:1,s:1 | $4 \times 4$ |
| Convolution* | 256 | k:3,s:2 | $2 \times 2$ |
| Convolution | 128 | k:1,s:1 | $2 \times 2$ |
| Convolution* | 256 | k:4,s:1 | $1 \times 1$ |

plate $L_{lp_{wh}}(l, g, v)$:

$$L_1(c, l, g, v, lpc) = \frac{1}{N}[L_{conf}(c) + L_{loc}(l, g) \\ + L_{has\_lp}(v, lpc) + L_{off}(l, g, v) + L_{lp_{wh}}(l, g, v)] \quad (1)$$

where $N$ is the number of the matched default boxes with the ground truth boxes of the vehicle, $c$ is the confidence of the vehicle, $l$ is the predicted box of the vehicle, $g$ is the ground truth box of the vehicle, $v$ is the ground truth of whether the vehicle contains a license plate, and $lpc$ is the predicted probability of the vehicle containing a license plate.

The learning target of vehicle detection is completely consistent with SSD [20], which predicts the vehicle presence confidence with cross-entropy loss (2) and regresses the bounding box of the vehicle with smooth L1 loss (3):

$$L_{conf}(c) = -\sum_{i=1}^{N}\sum_{p} \log(c_i^p) \qquad c_i^p = \frac{exp(c_i^p)}{\sum exp(c_i^p)} \quad (2)$$

$$L_{loc}(l, g) = \sum_{i=1}^{N}\sum_{m \in \{cx, cy, w, h\}} \mathbb{1}_{ij}^{p^+} smooth_{L1}\left(l_i^m - \overline{g}_j^m\right) \quad (3)$$

$$smooth_{L1} = \begin{cases} 0.5x^2 & |x| < 1 \\ |x| - 0.5 & otherwise \end{cases} \quad (4)$$

where the category $p$ is $\{vehicle, background\}$, the positive category $p^+$ is $vehicle$, and $\mathbb{1}_{ij}^{p^+} \in \{0, 1\}$ is the indicator of whether the $i$-th default box matches the $j$-th ground truth box. The smooth L1 loss [10] is defined as (4). Similar to SSD [20], the vehicle is regressed based on the center $(cx, cy)$ of the matched default box $(d)$ and its width $(w)$ and height $(h)$, as shown in (5).

$$\overline{g}_j^{cx} = \left(g_j^{cx} - d_i^{cx}\right)/d_i^w \qquad \overline{g}_j^{cy} = \left(g_j^{cy} - d_i^{cy}\right)/d_i^h \\ \overline{g}_j^w = \log\left(g_j^w/d_i^w\right) \qquad \overline{g}_j^h = \log\left(g_j^h/d_i^h\right) \quad (5)$$

The probability of whether the vehicle contains a license plate can reduce false positives of the license plate. Very small-sized vehicles and vehicles with invisible license plate (occlusion, large vehicle pose, etc.) are recognized as without license plate; otherwise, the vehicles are labeled as containing a license plate. The probability is optimized by binary cross-entropy loss (6), where $\sigma$ is a sigmoid function to limit the confidence to $[0, 1]$.

$$L_{has\_lp}(v, lpc) = -\sum_{i=1}^{N}[v_i \cdot \log(\sigma\left(lpc_i\right)) \\ + \left(1 - v_i\right) \cdot \log\left(1 - \sigma\left(lpc_i\right)\right)] \quad (6)$$

The offset between the center of the license plate and the vehicle as well as the size of the license plate are estimated with smooth L1 loss [10], as shown in (7). The vehicle should contain a license plate; otherwise, the losses $L_{off}(l, g, v)$ and $L_{lp_{wh}}(l, g, v)$ are all set to 0 by setting $v_i = 0$:

$$L_{off}(l, g, v) = \sum_{i=1}^{N}\sum_{m \in \{off_{x,y}\}} \mathbb{1}_{ij}^{p^+} v_i smooth_{L1}\left(l_i^m - \overline{g}_j^m\right) \\ L_{lp_{wh}}(l, g, v) = \sum_{i=1}^{N}\sum_{m \in \{lp_{w,h}\}} \mathbb{1}_{ij}^{p^+} v_i smooth_{L1}\left(l_i^m - \overline{g}_j^m\right) \quad (7)$$

where $off_{x,y}$ is the offset between the center of the license plate and the vehicle in both x-direction and y-direction, and $lp_{w,h}$ is the width and height of the license plate. Based on the matched default box, it directly regresses the offset and limits the LP size to $(0, +\infty)$ by a logarithmic operation to prevent negative numbers, as shown in (8).

$$\overline{g}_j^{off_x} = g_j^{off_x}/d_i^w \qquad \overline{g}_j^{off_y} = g_j^{off_y}/d_i^h \\ \overline{g}_j^{lp_w} = \log\left(g_j^{lp_w}/d_i^w\right) \qquad \overline{g}_j^{lp_h} = \log\left(g_j^{lp_h}/d_i^h\right) \quad (8)$$

## 3.2. Local Region Estimation and Aggregation

Compared with detecting the license plate directly in the large input image, it is better to localize the license plate in a small local region. Based on the predicted center and size of the license plate obtained at the ALPD stage, the local region can be obtained by simply expanding the background region around the license plate with a preset ratio, which can not exceed the boundary of the corresponding vehicle to reduce redundant background noises.

After that, many license plate regions of different vehicles can be obtained simultaneously. However, these LP regions have different sizes and aspect ratios. To reduce the running time, all the estimated regions are aggregated via ROI warping [5] as feature patches, in which all components are scaled to the same size and aspect ratio. Therefore, all the license plates can be detected simultaneously in the feature patches. All the LP region features are extracted from the first convolutional layer of the ALPD network, as seen in Table 1, because the first convolutional layer has the same size as the input image, which preserves the spatial information and is favorable to the detection of small-sized license plates.

## 3.3. Multi-Oriented License Plate Refinement

At this stage, the quadrilateral and horizontal bounding boxes of the license plate are detected simultaneously in the local region around the license plate. The detection results are more accurate than those obtained at the ALPD stage, as illustrated in Figure 1.

**Table 2**
Backbone network of the MOLPR network.

| Type | Filters | Parameters | Output |
|------|---------|-----------|--------|
| Convolution | 512 | k:3,s:1 | $56 \times 56$ |
| Convolution* | 512 | k:3,s:1 | $56 \times 56$ |
| Maxpool | - | k:2,s:2 | $28 \times 28$ |
| Convolution | 512 | k:3,s:1 | $28 \times 28$ |
| Convolution* | 512 | k:3,s:1 | $28 \times 28$ |
| Maxpool | - | k:2,s:2 | $14 \times 14$ |
| Convolution | 512 | k:3,s:1 | $14 \times 14$ |
| Convolution* | 512 | k:3,s:1 | $14 \times 14$ |

The backbone network of this stage is shown in Table 2, where "k" means kernel size and "s" means stride size. Like SSD [20], layers marked by "*" are candidates for multi-scale detection.

The training objective of the MOLPR network is defined as (9), which consists of the classification loss of the horizontal license plate $L_{conf}(c')$, the regression loss of the horizontal license plate $L_{loc}(l', g')$, and the corner loss of the quadrilateral license plate $L_{corner}(l', g')$:

$$L_2(c', l', g') = \frac{1}{N'}[L_{conf}(c') + L_{loc}(l', g') + L_{corner}(l', g')] \quad (9)$$

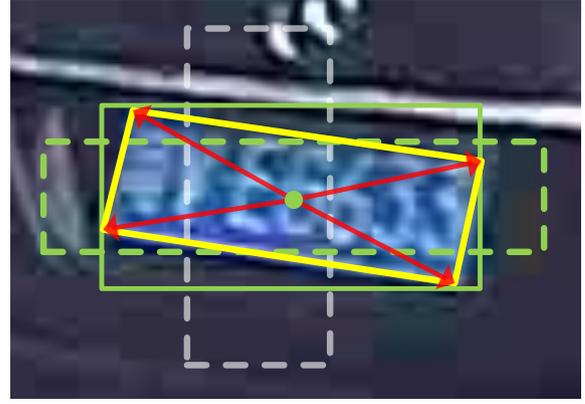where $N'$ is the number of the matched default boxes with



**Figure 2:** Four corners regression (red arrow) of the license plate. The matched default box (green dashed) is responsible for regressing the quadrilateral bounding box (yellow solid), where the default box is evaluated by IOU with the horizontal ground truth box (green solid). The four corners of the license plate are obtained based on the center of the matched default box. The irrelevant default box (grey dashed) is ignored because of low IOU.

the ground truth boxes of the license plate, $c'$ is the confidence of the license plate, $l'$ is the predicted box of the license plate, and $g'$ is the ground truth box of the license plate. The losses of the horizontal license plate $L_{conf}(c')$ and $L_{loc}(l', g')$ are similar to vehicle detection, as shown in (2) and (3).

The quadrilateral bounding box of the license plate is obtained by regressing the four corners of the license plate, as illustrated in Figure 2. The corner loss of the license plate is optimized with smooth L1 loss [10], as shown in (10):

$$L_{corner}(l', g') = \sum_{i=1}^{N'} \sum_{m \in \{tl, tr, br, bl\}} \mathbb{1}_{ij}^{p'+} smooth_{L1} \left( l_i'^m - \overline{g}_j'^m \right) \quad (10)$$

where the positive category $p'^+$ is *license plate* and $m \in \{tl, tr, br, bl\}$ are the four corners (top-left, top-right, bottom-right, bottom-left) of the license plate. The regression target is shown in (11), where the four corners of the license plate are directly regressed based on the center $(cx, cy)$ of the matched default box $(d')$ and its width $(w)$ and height $(h)$.

$$\overline{g}_j'^X = \left( g_j'^X - d_i'^{cx} \right) / d_i'^w \quad X \in \{tlx, trx, brx, blx\}$$
$$\overline{g}_j'^Y = \left( g_j'^Y - d_i'^{cy} \right) / d_i'^h \quad Y \in \{tly, try, bry, bly\} \quad (11)$$

## 3.4. End-to-End Trainable Detection Network

By integrating the aforementioned detection stages, we develop an end-to-end trainable network for degraded license plate detection. Combining (1) and (9) together, the loss of the whole network is shown in (12), where $\alpha$ is simply set to 1 to balance these loss terms.

$$L = L_1(c, l, g, v, lpc) + \alpha L_2(c', l', g') \quad (12)$$

**Table 3**
Ablation study (AP) of different datasets and IOU thresholds.

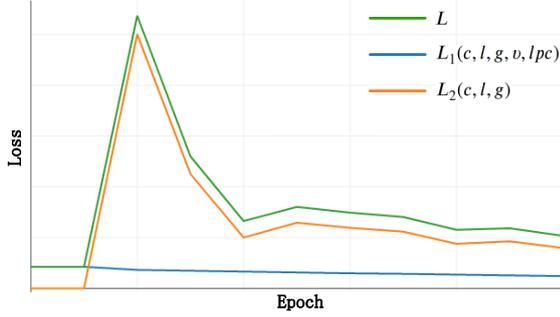| Method | LREA & MOLPR | has_lp confidence | vehicle boundary | IOU=0.5 | | IOU=0.75 | |
| | | | | TILT720 Test | TILT1080 Test | TILT720 Test | TILT1080 Test |
|---|---|---|---|---|---|---|---|
| Ours (ALPD) | | | | 76.71% | 77.71% | 26.27% | 35.27% |
| Ours (E2E) | √ | | | 88.11% | 86.20% | 53.73% | 53.37% |
| | √ | √ | | 88.95% | 87.29% | 54.46% | 55.94% |
| | √ | √ | √ | *89.19%* | *87.67%* | *54.51%* | *56.92%* |



**Figure 3**: Training loss. $L_1(c, l, g, v, lpc)$ is the loss of the ALPD network, $L_2(c', l', g')$ is the loss of the MOLPR network, and $L$ is the total loss of the end-to-end network.

During end-to-end training, the ALPD network can be firstly optimized to estimate the local region of the license plate, then the entire network will be optimized simultaneously. Concretely, during the first few epochs, $L_1$ goes down and $L_2$ remains zero because the untrained ALPD network can not estimate the location of the license plate; then $L_2$ goes up dramatically because the ALPD network can approximately estimate the license plate after training for some epochs, and the MOLPR network starts learning to regress the four corners of the license plate; finally the total loss $L$ goes down steadily because the ALPD and MOLPR network are optimized simultaneously.

# 4. Experiments

We mainly follow SSD [20], including the data augmentation strategies, such as random crop and color distortion, etc. We adopt SSD512 as the baseline network of the ALPD stage, which is initialized with the ILSVRC CLS-LOC dataset [28]. The backbone network of the MOLPR stage is trained from scratch, where the input size is $56 \times 56 \times 64$ (height $\times$ width $\times$ channel). Our model is trained for 60K iterations using the Adam optimizer [16]. The momentum parameters are set to $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Batch size and weight decay are set to 32 and $5 \times 10^{-4}$ respectively. The learning rate is first initialized to $10^{-4}$ and then decreased 10 times at 20K and 40K iterations. All the experiments are carried on a PC with 4 NVIDIA Titan Xp GPUs.

## 4.1. Datasets

**TILT720.** We employ an automobile data recorder to collect on-road videos with a size of $720 \times 1280$. After keyframe extraction and careful annotation, a total of 1033 images are obtained. All visible vehicles are labeled with a horizontal bounding box, and their corresponding license plates are annotated with a quadrilateral bounding box. The horizontal bounding box of the license plate is regarded as the tightest boundary of the quadrilateral bounding box. For simplicity, we name the dataset TILT720 (mulTi-oriented lIcense pLate deTection dataset 720p). All images are randomly divided into the training-validation set and test set by 9:1.

**TILT1080.** Similar to the TILT720 dataset, we obtain the TILT1080 dataset with another automobile data recorder. The TILT1080 dataset contains 4112 images, and all images have a size of $1080 \times 1920$. All images are randomly divided into the training-validation set and test set by 9:1.

## 4.2. Evaluation Protocols

**Horizontal Bounding Box.** We adopt the general AP (Average Precision) to evaluate the horizontal bounding box of the license plate. To be specific, we follow the 11-points computation of the VOC2007 [8] with different IOU thresholds (0.5 and 0.75). If it is not specified, the IOU threshold is set to 0.5.

Moreover, we hope that the expanded region at the LREA stage could completely contain the license plate for the next MOLPR stage. To evaluate it, we define a new evaluation criterion $C_{recall}$, as shown in (13):

$$C_{recall} = \frac{1}{M} \sum_{i=1}^{M} ER_i \cap LP_i = LP_i \quad (13)$$

where $ER_i$ means the $i$-th expanded region, $LP_i$ denotes the $i$-th license plate, and $M$ is the number of LP ground truths.

**Quadrilateral Bounding Box.** We adopt the classical precision (P), recall (R) and $F_1$-score (F) to evaluate the quadrilateral bounding box of the license plate:

$$P = \frac{TP}{TP + FP} \quad R = \frac{TP}{TP + FN} \quad F = \frac{2PR}{P + R} \quad (14)$$

where TP means true positive, FP means false positive, and FN means false negative. A quadrilateral bounding box is

considered as correct when the IOU with a quadrilateral ground truth box is greater than the threshold (0.5 or 0.75) under the confidence threshold 0.5.

### 4.3. Ablation Study

As demonstrated in Table 3, we adopt the ALPD network as the baseline model. The network only achieves 26.27% and 35.27% on the test set of TILT720 and TILT1080 with IOU threshold 0.75, because it fails to accurately localize the license plate from the large input image.

**LREA & MOLPR.** By adding the LREA stage and MOLPR stage, we obtain the end-to-end detection network, as illustrated in Figure 1. With a large IOU threshold, it improves almost 20% on the test set because of localizing the license plate in the local region, which proves the effectiveness of our method. With a small IOU threshold, it can also improve about 10% on the test set.

**has_lp confidence.** The ALPD network will inevitably estimate the approximate location and size of the license plate, no matter there is a visible license plate or not. However, the license plate is not always visible, especially for the very small-sized vehicle, occluded vehicle, and large-posed vehicle. The confidence of whether the vehicle contains a license plate can reduce false positives of the license plate, and it is fixed to 0.5. In this way, the invisible license plate can be filtered out and the precision is improved.

**vehicle boundary.** The expanded region at the LREA stage can be limited by the predicted vehicle boundary to avoid redundant background noises. It can further improve the performance on both test sets with different IOU thresholds.

**Table 4**
The influence to vehicle detection (AP) with IOU threshold 0.5. "E2E" means end-to-end.

| Method | TILT720 Test | TILT1080 Test |
|---|---|---|
| SSD | 87.82% | 87.51% |
| Ours (ALPD) | 87.85% | 87.50% |
| Ours (E2E) | 87.83% | 87.52% |

Moreover, as demonstrated in Table 4, either the ALPD network or the end-to-end network (E2E), our methods have no influence on vehicle detection compared with vanilla SSD [20]. As can be seen, there is still a lot of room to improve the performance of vehicle detection, and we will employ more powerful vehicle detectors in future work.

### 4.4. Experiments with Expansion Ratio

We hope the expanded region, obtained at the LREA stage, can fully contain the license plate for the next MOLPR stage. The expanded region is obtained based on the center and size of the license plate predicted at the ALPD stage, which takes the center of the license plate as the center and expands the width and height with the same ratio. To verify the effect of different expansion ratios at the LREA stage, we conduct comparative experiments on the trainval set and

test set of TILT720. Apart from the expansion ratio at the LREA stage, all other settings of our end-to-end network are the same.

**Table 5**
The effect of different expansion ratios at the LREA stage.

| Expansion Ratio | trainval | | test | |
|---|---|---|---|---|
| | AP(%) | $C_{recall}$(%) | AP(%) | $C_{recall}$(%) |
| 1 | 46.83 | 5.60 | 40.35 | 4.40 |
| 2 | 90.75 | 96.18 | 81.14 | 94.80 |
| 3 | *90.76* | 96.43 | *89.19* | 96.00 |
| 4 | 90.60 | 97.39 | 87.23 | 96.40 |
| 5 | 90.41 | 97.44 | 80.62 | 97.60 |
| $+\infty$ | 88.76 | **97.73** | 75.57 | **98.00** |

As shown in Table 5, it achieves the best AP performance with expansion ratio 3 on the trainval set. When the expansion ratio is less than 3, the AP increases gradually; when the expansion ratio is greater than 3, the AP decreases gradually. Due to the restriction of the vehicle boundary, the expansion ratio $+\infty$ represents the vehicle region. As can be seen, a too small or too large expansion ratio leads to significant performance degradation, because a small region can not fully contain the license plate, and a large region is not favorable to small-sized license plate. Moreover, the $C_{recall}$ increases as the expansion ratio increases. However, when the expansion ratio is greater than 1, the $C_{recall}$ has little improvement, so we set it to 3 by default. The results on the test set further validate that an expanded region of appropriate size is favorable to license plate detection.

### 4.5. Experiments with Horizontal Bounding Box

We compare Faster R-CNN [27], YOLOv2 [25] and SSD [20] with our proposed method, and the input size of all these methods is 512. The backbone network of Faster R-CNN, SSD, and our method is VGG-16 [33], while the backbone network of YOLOv2[2] is unchanged. Besides, we conduct comparative experiments with LPD methods [32] and [4] as well as scene text detection method TextBoxes [19]. Except for [32][3], all other methods are trained by ourselves with the trainval set of TILT720 and TILT1080 respectively. As shown in Table 6, our method (E2E) achieves the best performance with different datasets and IOU thresholds. Moreover, as shown in Figure 4, with a larger IOU threshold, our method (E2E) can obtain a larger performance gap than other methods. For example, with the IOU threshold 0.5, our method (E2E) is 2.56% and 1.33% better than SSD [20] on the test set of TILT720 and TILT1080 respectively; with the IOU threshold 0.75, our method (E2E) is 7.45% and 3.04% better than SSD [20]. Moreover, we enlarge the expanded region at the LREA stage to the whole vehicle region and uti-

---

[2]The backbone network of YOLOv2 has 19 convolutional layers and is comparable to VGG-16.

[3]The authors have publicly released their trained models for license plate detection. Please refer to https://github.com/sergiomsilva/alpr-unconstrained for more details.

**Table 6**
Comparative experiments (AP) with horizontal bounding box. "E2E" means end-to-end, "FC" means four corners, and "VP" means vehicle proposal.

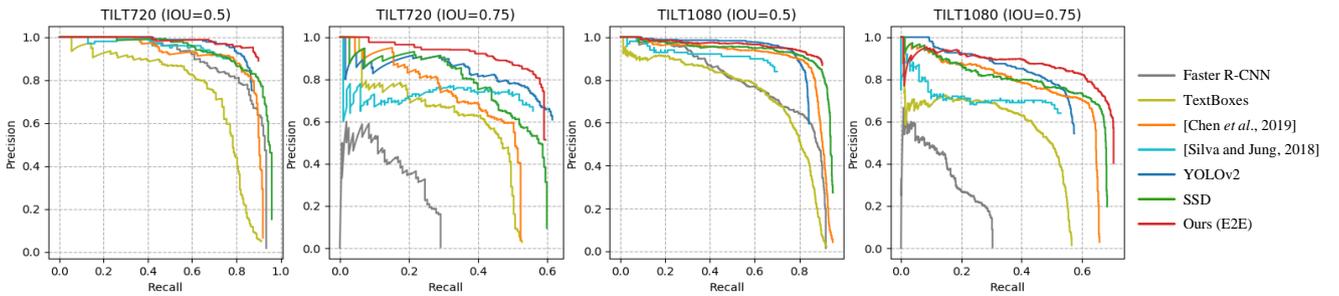| Method | IOU=0.5 | | IOU=0.75 | |
| --- | --- | --- | --- | --- |
| | TILT720 Test | TILT1080 Test | TILT720 Test | TILT1080 Test |
| Faster R-CNN [27] | 81.65% | 73.88% | 13.63% | 14.29% |
| YOLOv2 [25] | 80.80% | 79.58% | 51.66% | 49.32% |
| SSD [20] | 86.63% | 86.34% | 47.06% | 53.88% |
| TextBoxes [19] | 69.67% | 67.56% | 37.24% | 38.66% |
| Method [32] | 74.67% | 64.78% | 42.67% | 38.61% |
| Method [4] | 84.05% | 82.05% | 45.35% | 53.42% |
| Ours (E2E+VP) | 75.57% | 74.29% | 34.26% | 35.52% |
| Ours (E2E) | ***89.19%*** | ***87.67%*** | ***54.51%*** | ***56.92%*** |



**Figure 4:** The precision-recall curve of different methods, datasets, and IOU thresholds. Our method achieves the best performance, especially for the larger IOU threshold.

lize the vehicle proposal for the next MOLPR stage (E2E+VP). As shown in Table 6, the detection performance is significantly degraded compared with detecting LP in the local region (E2E), because the vehicle proposal is too large and not favorable to small-sized license plate detection.

### 4.6. Experiments with Quadrilateral Bounding Box

We conduct comparative experiments on the test set of TILT720 with different IOU thresholds. As shown in Table 7, our method (E2E) achieves the best $F_1$-score with different IOU thresholds. SSD [20] has poor performance because it can only detect the horizontal bounding box of the license plate. Furthermore, we upgrade vanilla SSD [20] and make it capable of detecting the four corners of the license plate (SSD+FC), which simulates the MOLPR stage and can directly localize the horizontal and quadrilateral bounding box of the license plate in the input image. SSD+FC achieves better performance than vanilla SSD [20], especially for the larger IOU threshold. However, due to directly detecting LP from the large input image, SSD+FC still lags behind our end-to-end method (E2E). Besides, our method with vehicle proposal (E2E+VP) still suffers great performance degradation, which proves the effectiveness of detecting LP in the local region of our method (E2E). All methods, except for our method (E2E), suffer low recall, because these methods localize the license plate in a relatively large region (input image or vehicle region), which leaves the confidence of LP

at a low level. Our method can localize the license plate in a small region around the license plate, which reduces background noises and can detect the license plate with high confidence. The aforementioned experiments prove that our method can precisely localize the quadrilateral bounding box of the oblique license plate. Some qualitative detection results are illustrated in Figure 5.

## 5. Conclusion

In this work, we propose an end-to-end trainable network for small-sized and oblique license plate detection via vehicle-plate relation mining, which detects the license plate in an end-to-end scheme. First, we propose a novel method to estimate the local region around the license plate using spatial relationships between the license plate and the vehicle, which can greatly reduce the search area and precisely detect very small-sized license plates. Second, we propose a new method to localize the quadrilateral bounding box by regressing the four corners of the license plate to robustly detect oblique license plates. Finally, based on the aforementioned methods, we develop an end-to-end trainable network for degraded license plate detection. Extensive experiments verify the effectiveness of our method, especially for a large IOU threshold. In future work, we would further promote the detection performance of vehicle detection to reduce false negatives of the license plate.

**Table 7**

Comparative experiments with quadrilateral bounding box on the test set of TILT720. "E2E" means end-to-end, "FC" means four corners, and "VP" means vehicle proposal.

| Method | IOU=0.5 | | | IOU=0.75 | | |
|---|---|---|---|---|---|---|
| | P | R | F | P | R | F |
| SSD [20] | 98.66% | 58.80% | 73.68% | 65.10% | 38.80% | 48.62% |
| Method [32] | 88.79% | 76.00% | 81.90% | 53.27% | 45.60% | 49.14% |
| Ours (E2E+VP) | 86.14% | 69.60% | 76.99% | 45.05% | 36.40% | 40.27% |
| Ours (SSD+FC) | 97.47% | 61.60% | 75.49% | 75.32% | 47.60% | 58.33% |
| Ours (E2E) | 90.61% | 88.80% | ***89.70%*** | 60.41% | 59.20% | ***59.80%*** |



**Figure 5:** Detection results of TILT720 (first row) and TILT1080 (second row). Our method can correctly localize small-sized and oblique license plates as well as license plates of large buses and trucks.

# References

[1] Anagnostopoulos, C., Anagnostopoulos, I., Psoroulas, I.D., Loumos, V., Kayafas, E., 2008. License plate recognition from still images and video sequences: A survey. IEEE Trans. Intell. Transp. Syst. 9, 377–391.

[2] Arafat, M.Y., Khairuddin, A.S.M., Khairuddin, U., Paramesran, R., 2019. Systematic review on vehicular licence plate recognition framework in intelligent transport systems. IET Intell. Transp. Syst. 13, 745–755.

[3] Cao, L., Jiang, Q., Cheng, M., Wang, C., 2016. Robust vehicle detection by combining deep features with exemplar classification. Neurocomputing 215, 225–231.

[4] Chen, S.L., Yang, C., Ma, J.W., Chen, F., Yin, X.C., 2019. Simultaneous end-to-end vehicle and license plate detection with multi-branch attention neural network. IEEE Trans. Intell. Transp. Syst. , 1–10, Online.

[5] Dai, J., He, K., Sun, J., 2016. Instance-aware semantic segmentation via multi-task network cascades, in: IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Las Vegas, NV, USA. pp. 3150–3158.

[6] Dong, M., He, D., Luo, C., Liu, D., Zeng, W., 2017. A cnn-based approach for automatic license plate recognition in the wild, in: British Machine Vision Conference, BMVA Press, London, UK.

[7] Du, S., Ibrahim, M., Shehata, M.S., Badawy, W.M., 2013. Automatic license plate recognition (ALPR): A state-of-the-art review. IEEE Trans. Circuits Syst. Video Techn. 23, 311–325.

[8] Everingham, M., Gool, L.J.V., Williams, C.K.I., Winn, J.M., Zisserman, A., 2010. The Pascal Visual Object Classes (VOC) Challenge. International Journal of Computer Vision 88, 303–338.

[9] Fu, Q., Shen, Y., Guo, Z., 2017. License plate detection using deep cascaded convolutional neural networks in complex scenes, in: Proceedings of the 24th International Conference on Neural Information Processing, Springer, Guangzhou, China. pp. 696–706.

[10] Girshick, R., 2015. Fast R-CNN, in: IEEE International Conference on Computer Vision, IEEE, Santiago, Chile. pp. 1440–1448.

[11] Girshick, R., Donahue, J., Darrell, T., Malik, J., 2014. Rich feature hierarchies for accurate object detection and semantic segmentation, in: IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Columbus, OH, USA. pp. 580–587.

[12] Han, J., Yao, J., Zhao, J., Tu, J., Liu, Y., 2019. Multi-oriented and scale-invariant license plate detection based on convolutional neural networks. Sensors 19, 1175.

[13] Hu, X., Xu, X., Xiao, Y., Chen, H., He, S., Qin, J., Heng, P., 2019. Sinet: A scale-insensitive convolutional neural network for fast vehicle detection. IEEE Trans. Intell. Transp. Syst. 20, 1010–1019.

[14] Jaderberg, M., Simonyan, K., Zisserman, A., Kavukcuoglu, K., 2015. Spatial transformer networks, in: Annual Conference on Neural Information Processing Systems, Montreal, Quebec, Canada. pp. 2017–2025.

[15] Kim, K., Kim, P., Chung, Y., Choi, D., 2018. Performance enhancement of yolov3 by adding prediction layers with spatial pyramid pooling for vehicle detection, in: IEEE International Conference on Advanced Video and Signal Based Surveillance, Auckland, New Zealand. pp. 1–6.

[16] Kingma, D.P., Ba, J., 2015. Adam: A method for stochastic optimization, in: Proceedings of the 3rd International Conference on Learning Representations, San Diego, CA, USA.

[17] Laroca, R., Severo, E., Zanlorensi, L.A., Oliveira, L.S., Gonçalves, G.R., Schwartz, W.R., Menotti, D., 2018. A robust real-time automatic license plate recognition based on the YOLO detector, in: International Joint Conference on Neural Network, IEEE, Rio de Janeiro, Brazil. pp. 1–10.

[18] Li, H., Wang, P., Shen, C., 2019. Towards end-to-end car license plates detection and recognition with deep neural networks. IEEE Trans. Intell. Transp. Syst. 20, 1126–1136.

[19] Liao, M., Shi, B., Bai, X., Wang, X., Liu, W., 2017. Textboxes: A fast text detector with a single deep neural network, in: Proceedings of the 31st AAAI Conference on Artificial Intelligence, San Francisco, California, USA. pp. 4161–4167.

[20] Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S.E., Fu, C.Y., Berg, A.C., 2016. SSD: single shot multibox detector, in: Proceedings of the 14th European Conference on Computer Vision, Springer, Amsterdam, The Netherlands. pp. 21–37.

[21] Liu, W., Liao, S., Hu, W., 2019. Towards accurate tiny vehicle detection in complex scenes. Neurocomputing 347, 24–33.

[22] Liu, W., Liao, S., Hu, W., Liang, X., Zhang, Y., 2018. Improving tiny vehicle detection in complex scenes, in: IEEE International Conference on Multimedia and Expo, San Diego, CA, USA. pp. 1–6.

[23] Liu, W., Rabinovich, A., Berg, A.C., 2015. Parsenet: Looking wider to see better. CoRR abs/1506.04579.

[24] Mo, Y., Han, G., Zhang, H., Xu, X., Qu, W., 2019. Highlight-assisted nighttime vehicle detection using a multi-level fusion network and label hierarchy. Neurocomputing 355, 13–23.

[25] Redmon, J., Farhadi, A., 2017. YOLO9000: Better, faster, stronger, in: IEEE Conference on Computer Vision and Pattern Recognition, IEEE, Honolulu, HI, USA. pp. 6517–6525.

[26] Redmon, J., Farhadi, A., 2018. Yolov3: An incremental improvement. arXiv preprint arXiv:1804.02767 .

[27] Ren, S., He, K., Girshick, R.B., Sun, J., 2015. Faster R-CNN: towards real-time object detection with region proposal networks, in: Annual Conference on Neural Information Processing Systems, Montreal, Quebec, Canada. pp. 91–99.

[28] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L., 2015. ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision 115, 211–252.

[29] Sang, J., Wu, Z., Guo, P., Hu, H., Xiang, H., Zhang, Q., Cai, B., 2018. An improved yolov2 for vehicle detection. Sensors 18, 4272.

[30] SG, K., HG, J., HI, K., 2017. Deep-learning-based license plate detection method using vehicle region extraction. Electronics Letters 53, 1034–1036.

[31] Silva, S.M., Jung, C.R., 2017. Real-time brazilian license plate detection and recognition using deep convolutional neural networks, in: Conference on Graphics, Patterns and Images, IEEE, Niterói, Brazil. pp. 55–62.

[32] Silva, S.M., Jung, C.R., 2018. License plate detection and recognition in unconstrained scenarios, in: Proceedings of the 15th European Conference on Computer Vision, Springer, Munich, Germany. pp. 593–609.

[33] Simonyan, K., Zisserman, A., 2015. Very deep convolutional networks for large-scale image recognition, in: Proceedings of the 3rd International Conference on Learning Representations, San Diego, CA, USA.

[34] Sun, Z., Bebis, G., Miller, R., 2006. On-road vehicle detection: A review. IEEE Trans. Pattern Anal. Mach. Intell. 28, 694–711.

[35] Tian, J., Wang, G., Liu, J., 2019. Semantic region proposals for adaptive license plate detection in open environment. J. Electronic Imaging 28, 023017.

[36] Tian, Y., Song, J., Zhang, X., Shen, P., Zhang, L., Gong, W., Wei, W., Zhu, G., 2016. An algorithm combined with color differential models for license-plate location. Neurocomputing 212, 22–35.

[37] Wang, J., Huang, H., Qian, X., Cao, J., Dai, Y., 2018. Sequence recognition of chinese license plates. Neurocomputing 317, 149–158.

[38] Wang, Y., Liu, Z., Deng, W., 2019. Anchor generation optimization and region of interest assignment for vehicle detection. Sensors 19, 1089.

[39] Wu, Z., Sang, J., Zhang, Q., Xiang, H., Cai, B., Xia, X., 2019. Multi-scale vehicle detection for foreground-background class imbalance with improved yolov2. Sensors 19, 3336.

[40] Xie, L., Ahmad, T., Jin, L., Liu, Y., Zhang, S., 2018. A new CNN-based method for multi-directional car license plate detection. IEEE Trans. Intell. Transp. Syst. 19, 507–517.

[41] Xu, Z., Yang, W., Meng, A., Lu, N., Huang, H., Ying, C., Huang, L., 2018. Towards end-to-end license plate detection and recognition: A large dataset and baseline, in: Proceedings of the 15th European Conference on Computer Vision, Springer, Munich, Germany. pp. 261–277.

[42] Yu, F., Koltun, V., 2016. Multi-scale context aggregation by dilated convolutions, in: International Conference on Learning Representations, San Juan, Puerto Rico.

[43] Yuan, Y., Zou, W., Zhao, Y., Wang, X., Hu, X., Komodakis, N., 2017. A robust and efficient approach to license plate detection. IEEE Trans. Image Processing 26, 1102–1114.