

Progressive LiDAR Adaptation for Road Detection

Zhe Chen, Jing Zhang, and Dacheng Tao, *Fellow, IEEE*

Abstract—Despite rapid developments in visual image-based road detection, robustly identifying road areas in visual images remains challenging due to issues like illumination changes and blurry images. To this end, LiDAR sensor data can be incorporated to improve the visual image-based road detection, because LiDAR data is less susceptible to visual noises. However, the main difficulty in introducing LiDAR information into visual image-based road detection is that LiDAR data and its extracted features do not share the same space with the visual data and visual features. Such gaps in spaces may limit the benefits of LiDAR information for road detection. To overcome this issue, we introduce a novel Progressive LiDAR Adaptation-aided Road Detection (PLARD) approach to adapt LiDAR information into visual image-based road detection and improve detection performance. In PLARD, progressive LiDAR adaptation consists of two subsequent modules: 1) data space adaptation, which transforms the LiDAR data to the visual data space to align with the perspective view by applying altitude difference-based transformation; and 2) feature space adaptation, which adapts LiDAR features to visual features through a cascaded fusion structure. Comprehensive empirical studies on the well-known KITTI road detection benchmark demonstrate that PLARD takes advantage of both the visual and LiDAR information, achieving much more robust road detection even in challenging urban scenes. In particular, PLARD outperforms other state-of-the-art road detection models and is currently top of the publicly accessible benchmark leader-board.

Index Terms—Road Detection, LiDAR Processing, Computer Vision, Deep Learning, Autonomous Driving

I. INTRODUCTION

Robust urban road detection is critical for autonomous driving systems. Without adequate recognition of road areas, a self-driving vehicle could not make safe decisions to achieve reliable navigation. Over the years, segmentation techniques have been used to identify road areas in monocular images, and, more recently, the introduction of deep convolutional neural network(DCNN)-based image segmentation methods such as FCN [1] and DeepLab [2] has significantly improved the performance of visual image-based road detection.

Corresponding author: Jing Zhang.

This work was supported by Australian Research Council Projects FL-170100117, DP-180103424, IH-180100002, and National Natural Science Foundation of China (NSFC) under Grant 61806062.

Z. Chen, and D. Tao are with the UBTECH Sydney Artificial Intelligence Centre and the School of Computer Science, in the Faculty of Engineering and Information Technologies, at the University of Sydney, 6 Cleveland St, Darlington, NSW 2008, Australia (email: zche4307@uni.sydney.edu.au; dacheng.tao@sydney.edu.au).

J. Zhang is a visiting scholar with the School of Software and Advanced Analytics Institute, at the University of Technology Sydney, 15 Broadway, Ultimo NSW 2007, Australia. He is a lecturer with School of Automation at the Hangzhou Dianzi University (email: jing.zhang@uts.edu.au).

©2018 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

Despite progress (e.g. [3–10]), DCNNs may still underperform when there are visual noises such as variable illumination, over exposure, ambiguous appearances, and blurry images. To overcome these issues and improve road detection performance, many studies [11–13] have introduced LiDAR information to improve road detection. “LiDAR”¹ refers to the data acquired by measuring the distance to a target by illuminating the target with pulsed laser light[14]. Many studies have proved that LiDAR is robust to various visual noises and can complement monocular image data. For example, Caltagirone *et al.* [12] reported that 3D LiDAR point clouds provide sufficient information to detect roads and are robust to visual noises, making it possible to robustly detect roads using only LiDAR data. Moreover, the authors of [15] attempted to fuse LiDAR and visual information for road detection. However, existing approaches that utilize LiDAR data for road detection are still far from effective and provide only limited improvements over visual image-based road detection methods. Here, we investigate the difficulties encountered when utilizing LiDAR data for road detection and propose a novel and more effective approach to incorporate LiDAR information into a visual image-based road detection system.

By interrogating LiDAR information and visual information for road detection, we conclude that two major factors would pose difficulties for effective cooperation between the two types of information. First, since raw LiDAR data and raw visual image data are in different spaces, it is difficult to define a proper space to integrate both data types. For instance, in the KITTI road detection dataset [16, 17], the provided LiDAR data is defined in the 3D real-world space, while the visual images are defined on the 2D image plane. Although researchers can project the LiDAR data onto the 2D image plane using the calibration parameters, this may at the same time alter the road appearance in the LiDAR data, making road areas less distinguishable in the LiDAR data space. As a result, it will be difficult for a DCNN-based road detection model to learn a reliable road detection function based on the LiDAR data, let alone improve the visual image-based road detection models. Moreover, it is also difficult to appropriately integrate the features extracted from the LiDAR data and the visual features extracted from visual images. More specifically, since the road appearance in the LiDAR data is described by scattered points and road appearance in visual data is described by RGB values of pixels on the 2D image plane, it is highly likely that the features extracted from both data sources are also in different spaces. This gap in feature spaces could adversely impact the feature fusion performance and the final detection accuracy, thus existing feature fusion methods for road detection can hardly outperform state-of-the-art visual

¹an acronym of *light detection and ranging*

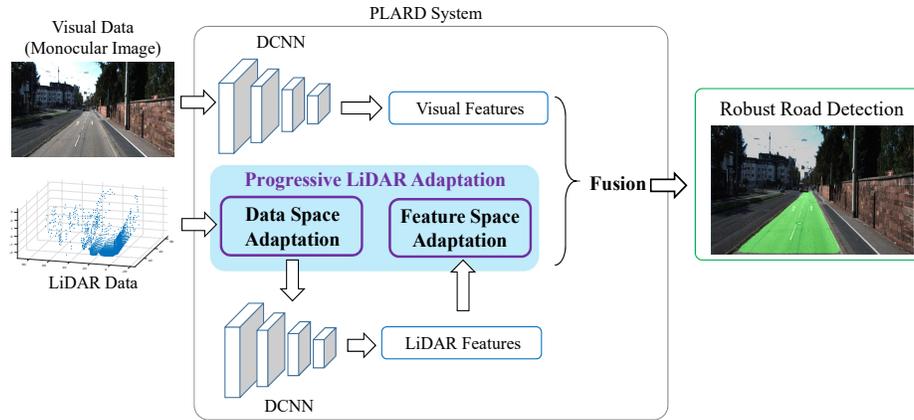


Fig. 1. Overview of the Progressive LiDAR Adaptation-aided Road Detection (PLARD) approach. We overcome the issue that LiDAR information and visual information are in different spaces when detecting road areas in urban scenes. In particular, the proposed progressive LiDAR adaptation consists of a data space adaptation step that adapts the view of raw LiDAR data to align the view of visual images and a feature space adaptation step that adapts the LiDAR features to visual features. By fusing the adapted LiDAR information and the visual information, PLARD achieves robust road detection.

image-based road detection algorithms.

To overcome these issues, here we propose a novel progressive LiDAR adaptation technique to make LiDAR information more compatible with visual information and to improve the visual image-based road detection more effectively. To achieve this, in the progressive LiDAR adaptation, we introduce proper transformation functions to successively adapt the LiDAR data space into the visual data space and adapt the LiDAR feature space into the visual feature space. Accordingly, the progressive LiDAR adaptation procedure consists of a data space adaptation step and a feature space adaptation step. The data space adaptation step transforms and aligns the LiDAR data space to the visual data space whilst still making the road areas easy-to-distinguish in the LiDAR data. Afterwards, through a cascaded fusion structure, the feature space adaptation step transforms the LiDAR feature space into the space that better complements and improves the visual features. By integrating the visual information with the adapted LiDAR information, we obtain a more robust road detection model: Progressive LiDAR Adaptation-aided Road Detection (PLARD) model. Fig. 1 shows an overview of our proposed system.

Using the well-known KITTI road detection benchmark [16], we perform comprehensive experiments for the proposed PLARD system to evaluate the effectiveness of different parts of the proposed technique and the overall performance gains over a visual image-based road detection system. Empirical results demonstrate that LiDAR information delivers more benefits for road detection via our proposed progressive LiDAR adaptation technique. Furthermore, on the test set of KITTI road detection benchmark, PLARD promisingly improves the road detection accuracy, outperforming other visual image-based road detection algorithms, LiDAR-based road detection algorithms, and the algorithms that fuse both information. In particular, PLARD achieves state-of-the-art performance on the publicly accessible leader-board. Indeed, an ensemble of 3 our PLARD models ranks the top on the leader-board at the time of writing this paper.

II. RELATED WORK

Road detection is beneficial to various other autonomous tasks [18–23]. Over the years, various algorithms have been developed to tackle the road detection problem [24–27]. For instance, model-based methods build shape [28, 29] or appearance models [30] to describe the road structure and then identify road areas in the input images. Learning-based methods then attempt to employ classifiers (such as SVM [31] and random forest [32]) to distinguish road from non-road areas. In practice, learning-based methods usually perform better than model-based methods.

By considering the road detection task as a semantic segmentation task, DCNNs have been demonstrated in recent years to be particularly useful for road detection. In particular, several typical algorithms have proven to be effective in semantic segmentation. For example, Long *et al.*[1] proposed the fully convolutional and upsampling layers to tackle the pixel-level semantic segmentation problems. Moreover, authors of [33, 34] achieved compelling semantic segmentation performance by introducing dilated convolution operations which can greatly enlarge receptive fields of convolutional kernels without reducing the resolution of the feature maps. By taking advantage of both the fully convolutional layer and dilated convolution operations, several studies [35–37] have achieved impressive performance on semantic segmentation benchmarks. These techniques have been widely used for detecting roads in urban scenes [7–9, 38].

In order to improve the effectiveness of DCNN-based road detection, several promising algorithms have been proposed. For instance, Mendes *et al.* [8] introduced a large contextual window and a network-in-network architecture to improve accuracy, while the study [39] introduced an efficient deep network following the “encoder-decoder” principle as discussed in “U-net” [40]. However, DCNNs are still susceptible to visual noises and they usually require excessively long processing times to guarantee better performance. As an example, the algorithm in [7] cost around 1s to process an image and

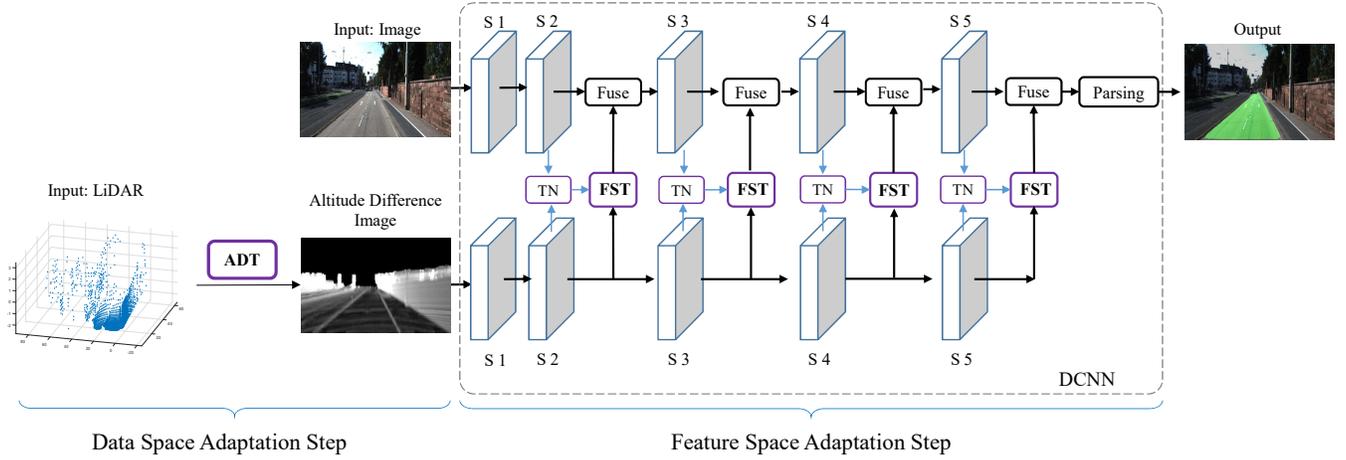


Fig. 2. Overall pipeline of the proposed PLARD system. In the data space adaptation step, we introduce the Altitude Difference-based Transformation (denoted “ADT”) to adapt the raw LiDAR data, obtaining a better aligned LiDAR data space where the roads are easier to be distinguished from other objects. In the feature space adaptation step, DCNNs are first employed to detect roads on visual and LiDAR data respectively. Then, Feature Space Transformation (denoted “FST”) modules are introduced to transform the LiDAR features and make them better complement and improve visual features. In each “FST” module, a Transformation Network (denoted “TN”) is employed to learn the transformation parameters. After feature transformation, the visual features and the adapted LiDAR features are fused via a cascaded structure. The cascaded fusion integrates features at all the convolution stages (denoted “S1-S5”) but the first stage. Lastly, in the parsing stage, PLARD performs classification on the integrated features, delivering robust road detection results.

the algorithm in [9] needed 2s, and both algorithms fail to achieve state-of-the-art performance, making them impractical for moving platforms like autonomous vehicles.

Despite this progress in the visual image-based road detection, others proposed that LiDAR is robust to visual noises and they have attempted to detect the roads using LiDAR information. In [4], visual images were transformed from the perspective view into the birds-eye-view for road detection. Another study [12] took LiDAR point clouds rather than visual images as input, which can perform promising road detection in 3D real-world space. However, these studies did not effectively take advantage of both types of information, limiting their final detection performance. There are also studies [5, 15] that attempted to fuse both visual and LiDAR information for road detection, but existing fusion-based algorithms are either time-consuming or less effective than other state-of-the-art algorithms. In this study, we hypothesized that the difficulty of using LiDAR to improve road detection is due to the gaps between data spaces and feature spaces for LiDAR information and visual information. To overcome this issue, we introduce a progressive LiDAR adaptation technique that effectively adapts and integrates LiDAR information into the visual image-based road detection pipeline to boost the robustness and accuracy for road detection.

III. PLARD SYSTEM

A. Problem Definition

In this study, we formulate the road detection task as assigning pixels on the 2D image plane with a binary label indicating whether the pixel belongs to road areas. LiDAR information will be adapted for the road detection on the 2D image plane. Formally, suppose \hat{y} represents a ground-truth label, f is a road detection function, and W is the model

parameters of f . Taking LiDAR data L and visual data I as input, we tackle the road detection problem by optimizing the following objective:

$$\min_W \sum_i \sum_{x,y} \mathcal{L}(f(I_i, L_i; W), \hat{y})|_{x,y}, \quad (1)$$

where i indexes over training examples, x, y represent the horizontal and vertical offsets, respectively, on the image plane, and \mathcal{L} is a loss function.

B. Overview

Considering that gaps exist between the data and feature spaces, robust road detection is difficult to achieve solely by using a simple combination of LiDAR and visual information. To overcome this problem and improve road detection, we propose a progressive LiDAR adaptation technique to make the LiDAR information more compatible with visual information, thus more effectively incorporating both information types. Mathematically, we formulate that the road detection function f of PLARD system takes the following form:

$$f(I, L; W) = f_{parsing}(f_{fuse}(f_{vis}(I; W_{vis}), g(L; W_{lidar}))), \quad (2)$$

where g is the progressive LiDAR adaptation function, f_{vis} is the visual image-based road detection function, W_{vis} and W_{lidar} are parameters of the corresponding functions, f_{fuse} is a fusion operation, and $f_{parsing}$ is the final binary classification function that identifies road areas from fused features. Specifically, we implement f_{vis} by introducing a ResNet101[41] backbone and implement $f_{parsing}$ by a 2-class softmax function after a pyramid scene parsing module [42].

We implement the progressive adaptation function g by introducing two subsequent adaptation steps: data space adaptation step and feature space adaptation step. In the data space

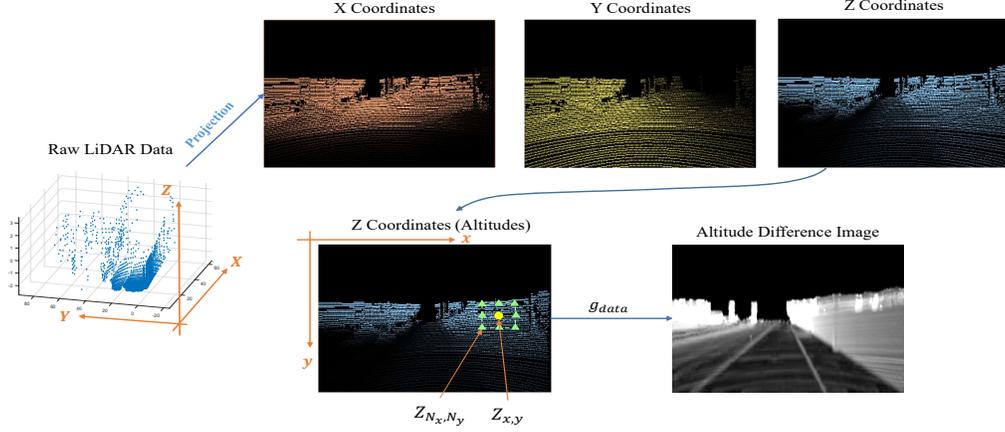


Fig. 3. An example of directly projected LiDAR data and altitude difference image. Using provided calibration parameters, LiDAR points described by 3-dimensional coordinate vectors in the real-world space can be projected onto the 2D image plane. Upper-right figures show the projected LiDAR points whose intensities represent the normalized X , Y , and Z coordinates. According to the definition of axes, Z can be considered as altitudes and then we can compute the absolute values of changes in altitudes with respect to spatial offsets between two locations (such as $Z_{x,y}$ and Z_{N_x, N_y} in the figure). With the computed differences in altitude, an altitude difference image (as illustrated in the bottom-right figure) is obtained. The altitude difference image can make roads easier to be distinguished from non-road areas by preserving the characteristics of road in the LiDAR data. Best view this figure in color.

adaptation, we transform raw LiDAR data from the 3D space to the 2D image plane while preserving the distinguishable characteristics of road areas. Then, in the feature space adaptation step, we introduce a learning-based module to transform the LiDAR features such that the transformed features better complement the visual features in road detection. Accordingly, we formulate the progressive LiDAR adaptation function g as follows:

$$g(L; W_{lidar}) = g_{feat}(f_{lidar}(g_{data}(L); W_{lidar})), \quad (3)$$

where g_{data} and g_{feat} represent the data space adaptation function and feature space adaptation function, respectively, and f_{lidar} is the LiDAR-based road detection function. The overview of PLARD is shown in Fig. 2.

In Section III-C1, we implement g_{data} by studying the altitude changes with respect to the 2D image plane. In Section III-C2, we implement g_{feat} by introducing a learning-based module into the DCNN architecture to transform the features extracted from LiDAR data into a space that better complement the visual features. Lastly, the implementation details of the PLARD approach is described in Section IV.

C. LiDAR Adaptation

In the data space adaptation step of the progressive LiDAR adaptation, we introduce a novel altitude difference-based transformation method to transform the LiDAR data space. In the feature space adaptation step, we introduce a transformation network to learn and transform the LiDAR feature space.

Both raw LiDAR data and raw visual data are in different spaces. Raw LiDAR data is composed of tens of thousands of points in the 3D real-world space and each LiDAR point is described by a 3-dimensional coordinate vector, while the visual data is composed of pixels on a 2D image plane and each pixel is described by a RGB value. Therefore, it is extremely

challenging to directly and smoothly integrate visual data with raw LiDAR data. Fortunately, with the help of calibration parameters, the 3D LiDAR points can be projected onto the 2D image plane, thereby obtaining an image with the projected LiDAR points. Although this obtained image can be used for road detection, we argue that road and non-road appearances in the projected LiDAR data will be less distinguishable from each other, thus diminishing the capacity of a road detection model to identify road areas, as illustrated in the upper-right corner of Fig. 3. Instead of using direct projection results, we propose an altitude difference-based transformation operation that better preserves the road characteristics to instantiate g_{data} and help adapt the 3D LiDAR data to the visual data space.

1) **Data Space Adaptation:** In the first stage, we introduce *altitude difference-based transformation* to implement data space adaptation. Altitude difference-based transformation is introduced based on the observation that road surfaces in the 3D space are flat and relatively smooth in the altitudes of LiDAR point clouds compared to other objects such as vehicles and buildings. After projecting LiDAR points onto the image plane, such smoothness can be maintained by recording the altitudes of the original 3D LiDAR points. As a result, road areas can be better distinguished in the projected LiDAR data according to the changes in altitudes on the image plane.

More specifically, on the LiDAR-projected 2D image plane, altitude difference-based transformation computes a pixel value $V_{x,y}$ located at (x,y) according to:

$$g_{data}(L)|_{x,y} = V_{x,y} = \frac{1}{M} \sum_{N_x, N_y} \frac{|Z_{x,y} - Z_{N_x, N_y}|}{\sqrt{(N_x - x)^2 + (N_y - y)^2}}, \quad (4)$$

where $Z(x,y)$ is the altitude of the LiDAR point projected on (x,y) , (N_x, N_y) denote positions in the neighbourhood of (x,y) , and M is the total number of considered neighborhood positions. If a neighboring pixel is not correlated to a 3D point, we simply ignore it. It is worth mentioning that Eq. 4 can be

viewed as a mean absolute value for gradients of altitudes of projected points with respect to the 2D image plane. Therefore, if an object is upright and sharp, its projected areas will have large altitude differences on the image plane. For example, in Fig. 3, road areas commonly have a small intensity on the altitude difference image while other objects generally have large intensities, distinguishing the roads from other objects. Suppose H and W are the height and width of the input image, respectively. According to Eq. 4, the complexity of altitude difference-based transformation is $O(MHW)$.

2) **Feature Space Adaptation:** In addition to the gap in data spaces, the features extracted from LiDAR data could also be inconsistent with the visual features extracted from images, since road areas may have different appearances in different data sources. This feature space inconsistency could further hamper the performance of integrating LiDAR features and visual features, thus limiting the overall benefit of introducing LiDAR information. We therefore attempt to transform the LiDAR feature space to make the LiDAR features better complement and improve visual features and visual image-based road detection performance. However, the main challenge of feature space adaptation is that we do not have a complete knowledge about how to properly transform the feature space. To tackle this issue, we introduce a learning-based module to find an appropriate space adaptation operation.

In general, we assume that the linear transformation can properly define a feature space adaptation operation, and we have:

$$g_{feat}(\mathbf{f}_{lidar}) = \alpha \mathbf{f}_{lidar} + \beta, \quad (5)$$

where α is a scalar vector, β is an offset vector, and \mathbf{f}_{lidar} is the LiDAR feature to be adapted:

$$\mathbf{f}_{lidar} = f_{lidar}(g_{data}(L); W_{lidar}). \quad (6)$$

To estimate α and β properly and achieve better feature space adaptation, we introduce a neural network, called the transformation network, to learn and adapt LiDAR features. Accordingly, the transformation network estimates α and β following:

$$\alpha = f_{\alpha}(\mathbf{f}_{lidar}, \mathbf{f}_{vis}; W_{\alpha}), \quad (7)$$

$$\beta = f_{\beta}(\mathbf{f}_{lidar}, \mathbf{f}_{vis}; W_{\beta}), \quad (8)$$

where f_{α} and f_{β} represent the neural network function for computing α and β , respectively, W_{α} and W_{β} are their corresponding weight parameters, and \mathbf{f}_{vis} is the visual feature:

$$\mathbf{f}_{vis} = f_{vis}(I; W_{vis}). \quad (9)$$

In a DCNN, we implement both f_{α} and f_{β} by using fully convolutional operations whose weight parameters are defined by W_{α} and W_{β} , respectively. \mathbf{f}_{lidar} and \mathbf{f}_{vis} are concatenated as input for both f_{α} and f_{β} . By optimizing the parameters W_{α} and W_{β} together with the overall road detection system, we can learn a proper transformation of feature spaces for the extracted LiDAR features, thus improving road detection more promisingly. Fig. 4 shows a detailed schematic of the feature space adaptation stage. Comparing to the overall DCNN architecture which may have over dozens of convolutional layers, the complexity of feature space adaptation is small

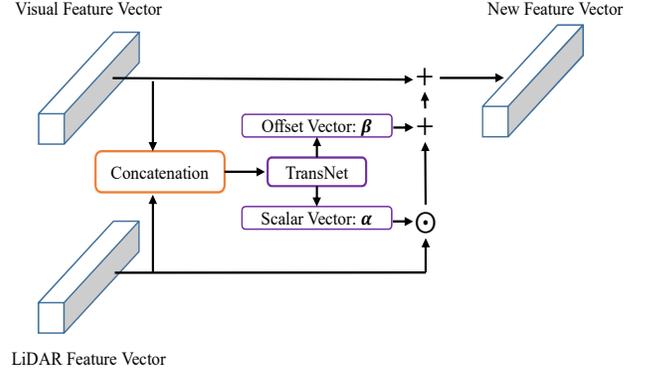


Fig. 4. Structure of a feature adaptation module in PLARD. Accepting as input the convolutional features of visual information and LiDAR information, the feature space adaptation introduces a transformation network (denoted “TransNet”) to learn and transform LiDAR features. More specifically, the transformation network outputs a scalar vector α and an offset vector β based on the concatenated visual and LiDAR features. By fusing the adapted LiDAR features with visual features via a residual structure, we obtain an improved feature vector. In this figure, \odot means element-wise multiplication and $+$ means addition in this figure.

because it only involves three 1×1 convolutional operations and three element-wise multiplication or addition operations. Suppose C is the channel number, H_k and W_k are height and width of the k -th convolutional stage, respectively, then the complexity of feature space adaptation for that stage is the complexity of 3 convolutional operations: $O(4C^2H_kW_k)$.

D. Cascaded Fusion for Adapted LiDAR Information

After LiDAR is adapted according to g_{data} and g_{feat} , we fuse the LiDAR and visual information to achieve more robust road detection. Specifically, we implement the fusion function f_{fuse} using an residual-based cascaded fusion structure. In this architecture, we make the adapted LiDAR features improve the visual features under a residual structure. By applying this residual-based fusion to every subsequent convolutional stage in a DCNN pipeline, we achieve robust road detection.

Mathematically, taking visual features \mathbf{f}_{vis} and the adapted LiDAR-based features \mathbf{f}_{lidar} as input, we implement the fusion function as follows:

$$\begin{aligned} & f_{fuse}^k(f_{vis}^k(f_{fuse}^{k-1}(I, L); W_{vis}^k), g^k(L; W_{lidar})) \\ &= f_{vis}^k(f_{fuse}^{k-1}(I, L); W_{vis}^k) + \lambda g^k(L; W_{lidar}), \end{aligned} \quad (10)$$

where k indicates features from the k -th convolutional stage of the DCNN in the road detection system and λ is a scalar parameter. It is worth noting that there are 5 convolutional stages when using ResNet-101 [41] as the backbone network. The right part of Fig. 2 shows the detailed structure of this fusion.

E. Overall Objective

During optimization, we train PLARD in an end-to-end manner to obtain parameters for both visual image-based and

LiDAR-based DCNNs. In addition to the objective for road detection with the fused LiDAR and image information, we also introduce a loss for the LiDAR-based DCNN such that the LiDAR-based road detection system directly detects roads from altitude difference images. Moreover, we follow the design of [42] and introduce an auxiliary loss in the visual image-based DCNN to facilitate convergence.

According to Eq. 1, training PLARD is to minimize a loss function w.r.t. all of its parameters. Suppose \mathcal{L}_{PLARD} is the overall loss for road detection with the fused LiDAR and image information, \mathcal{L}_{LiDAR} represents the loss only for the LiDAR-based DCNN, and \mathcal{L}_{aux} is the auxiliary loss in the visual image-based DCNN. Then, the loss function of PLARD can be written as:

$$\mathcal{L} = w_{parsing}\mathcal{L}_{parsing} + w_{lidar}\mathcal{L}_{lidar} + w_{aux}\mathcal{L}_{aux}, \quad (11)$$

where $w_{parsing}$, w_{lidar} and w_{aux} are the corresponding loss weights. The settings for $w_{parsing}$, w_{lidar} , and w_{aux} can be found in Sec. IV-C.

In this study, we exploit multinomial cross-entropy loss to define all these losses:

$$\mathcal{L}_{\{parsing,lidar,aux\}} = - \sum_{c=1}^2 (\hat{y}^c \log_{10}(y_{\{parsing,lidar,aux\}}^c)), \quad (12)$$

where \hat{y}^c is the ground-truth for category c and y^c is the prediction outputs. Specifically, $y_{parsing}^c$, y_{lidar}^c , and y_{aux}^c are computed according to Eq. 2, Eq. 6 (only for the final convolutional stage), and Eq. 10 (only for the convolutional stage $k = 4$), respectively.

IV. EXPERIMENT

In this section, we evaluate the effectiveness of PLARD on the KITTI road benchmark. We first adopt 5-fold cross-validation to illustrate the effectiveness of the individual parts in PLARD. Then, we evaluate PLARD on test set of the benchmark and compare it to other state-of-the-art road detection algorithms.²

A. Dataset

The KITTI road benchmark [16] is popular with road detection researchers due to its comprehensiveness. KITTI uses a wide variety of evaluation metrics to assess algorithm performance and also provides information captured by various sensors including visual cameras, LiDAR sensor, and GPS. KITTI contains 289 images for training and 290 images for testing, both containing three different road scene categories including Urban Marked roads (UM), Urban Multiple Marked lanes (UMM), and Urban Unmarked roads (UU). For fair evaluation, KITTI does not provide ground-truths for test images, and the number of submissions for online evaluation is limited. All the results are publicly accessible on its official website.

²Results on the test set are available on: http://www.cvlibs.net/datasets/kitti/eval_road.php.

B. Evaluation Metrics.

We follow the standard evaluation metrics used by KITTI, which are detailed in [43]. The metrics include the maximum F1-measure (MaxF), average precision (AP), precision rate (PRE), recall rate (REC), false positive rate (FPR), and false negative rate (FNR). The four latter measures are obtained at the working point of MaxF. According to KITTI's evaluation system, all the results are transformed into birds-eye-view space for evaluation and MaxF is selected as the metric for ranking the evaluated algorithms.

C. Implementation and Model Training

The implementation of PLARD is simple and straightforward. In this section, we discuss some implementation details that should be taken into account.

When dealing with altitude difference images, we set N_x and N_y as positions within a 7×7 window centered at (x, y) , thus the maximum value of M is 48 (excluding the centre). All the values in an altitude difference image are re-scaled within the range $[0, 255]$.

In PLARD, we employ PSPNet [42] as the visual image-based DCNN and use ResNet-101 [41] as the backbone. We also employ a 101-level DCNN to implement f_{lidar} and extract LiDAR features. To avoid excessive computational loads introduced by using two DCNNs, we make the channel numbers of the LiDAR-based DCNN 8 times smaller than the channel numbers of the visual image-based DCNN at the same level. In addition, we use hybrid convolutions [45] in the LiDAR-based DCNN to augment its expressive capacity with fewer channel numbers. During fusion, we use a uniform channel number, i.e. 256. The parameter λ in Eq. 10 is set to 0.1. Regarding loss weights in Eq. 11, we empirically set $w_{parsing}$, w_{aux} , and w_{lidar} as 1.0, 0.16, 0.4, respectively.

We resize all the images into 384 by 1280 during training and testing. We adopt the SGD algorithm to optimize all parameters with the learning rate gradually decaying from 1×10^{-4} to 1×10^{-6} . We train models using 80 epochs and use NVIDIA's GTX Titan GPUs for computation. The implementation of PLARD and the experimental settings are similar in both the ablation studies and the evaluation on the test set. The differences are as follows. First, to better illustrate the effectiveness of different parts of our model, we train models from scratch in the ablation studies. In contrast to ablation studies, we pre-train visual image-based DCNN in the PLARD system using external data [46] to improve robustness for the evaluation on test set. In addition, for the test set, we improve PLARD by adopting several data augmentation techniques, including multi-scale training and testing, random cropping, and disturbing the image brightness. Lastly, we extend the training period for the test set evaluation to three times longer than in the ablation study.

D. Ablation Study

We first evaluate the effectiveness of different PLARD components and compare their performances to baseline methods using 5-fold cross-validation on the training set. Table I shows the average results for all the "UM", "UMM", and

TABLE I

ABLATION STUDIES FOR PLARD. RESULTS ARE BASED ON 5-FOLD CROSS VALIDATION USING THE TRAINING SET. FOUR COMPONENTS ARE COMPARED, INCLUDING “IMG” (VISUAL IMAGE-BASED ROAD DETECTION), “L-PROJ” (ROAD DETECTION USING DIRECTLY PROJECTED LIDAR DATA), “L-ADT” (ROAD DETECTION USING ALTITUDE DIFFERENCE TRANSFORMED LIDAR DATA), AND “FSA” (ROAD DETECTION BY APPLYING FEATURE SPACE ADAPTATION ON LIDAR FEATURES THROUGH CASCADED FUSION). BEST SCORES ARE HIGHLIGHTED IN **bold**.

Img	L-Proj	L-ADT	FSA	Speed (s/im)	Max F	AP	PRE	REC	FPR	FNR
✓				0.120	88.29	91.17	86.98	89.73	7.62	10.28
	✓			0.045	83.25	87.54	79.68	87.55	13.54	12.45
		✓		0.042	86.39	90.73	84.73	88.21	9.62	11.79
✓	✓			0.161	89.34	91.40	87.68	91.22	6.75	8.78
✓		✓		0.159	89.79	91.78	89.34	90.36	6.46	9.64
✓		✓	✓	0.160	92.85	93.14	93.16	92.55	3.10	7.45

TABLE II

OVERALL PERFORMANCE ON THE TEST SET OF KITTI ROAD DETECTION BENCHMARK. BEST SCORES ARE HIGHLIGHTED IN **bold**. “+”: A 3-MODEL ENSEMBLE WITH MULTI-SCALE TESTING.

Methods	Speed (s/im)	Input	MaxF	AP	PRE	REC	FPR	FNR
SPRAY[4]	0.045	Image	87.09 %	91.12 %	87.10 %	87.08 %	7.10 %	12.92 %
FCN_LC[8]	0.03	Image	90.79 %	85.83 %	90.87 %	90.72 %	5.02 %	9.28 %
HybridCRF[15]	1.5	Image + LiDAR	90.81 %	86.01 %	91.05 %	90.57 %	4.90 %	9.43 %
FTP [29]	0.28	Image	91.61 %	90.96 %	91.04 %	92.20 %	5.00 %	7.80 %
Up_Conv[39]	0.08	Image	93.83 %	90.47 %	94.00 %	93.67 %	3.29 %	6.33 %
LoDNN[12]	0.018	LiDAR	94.07 %	92.03 %	92.81 %	95.37 %	4.07 %	4.63 %
MultiNet[38]	0.17	Image	94.88 %	93.71 %	94.84 %	94.91 %	2.85 %	5.09 %
StixelNet II[44]	1.2	Image	94.88 %	87.75 %	92.97 %	96.87 %	4.04 %	3.13 %
RBNet[10]	0.18	Image	94.97 %	91.49 %	94.94 %	95.01 %	2.79 %	4.99 %
LidCamNet [11]	0.15	Image + LiDAR	96.03 %	93.93 %	96.23 %	95.83 %	2.07 %	4.17 %
NF2CNN	0.006	Image + LiDAR	96.70 %	89.93 %	95.37 %	98.07 %	2.62 %	1.93 %
PSPNet[42]	0.12	Image	96.29 %	93.71 %	96.22 %	96.35 %	2.09 %	3.65 %
PLARD	0.16	Image + LiDAR	96.83 %	93.98 %	96.79 %	96.86 %	1.77 %	3.14 %
PLARD+	1.5	Image + LiDAR	97.03 %	94.03 %	97.19 %	96.88 %	1.54 %	3.12 %

“UU” tasks. In particular, we investigate (1) the benefits of altitude difference-based transformation operations compared to directly projected LiDAR points; and (2) the improvements afforded by the proposed learning-based feature space adaptation module implemented via a cascaded fusion structure compared to the direct fusion implemented via a simple concatenation.

It can be seen that both the altitude difference-based transformation technique and the feature space adaptation technique achieve promising improvements compared to their counterparts. Specifically, training with the altitude difference image (L-ADT) surpasses the road detection model using directly projected LiDAR points (L-Proj) by around 3 points for “Max F”. In addition, fusing both visual and LiDAR information improves the visual image-based road detection baseline, indicating that LiDAR is helpful for boosting robustness. By further introducing the feature space adaptation procedure through a cascaded fusion structure, the final model (Img + L-ADT + FSA) achieves the highest performance among compared methods, demonstrating the effectiveness of the feature transformation and cascaded fusion structure. Note that the speed for LiDAR-based road detection does not include processing time for raw data projection and transformation.

E. Evaluation on the Test Set

1) *Quantitative Results*: Using all the training images in the benchmark, we evaluate the visual image-based PSPNet as the road detection baseline and we evaluate the performance of a single PLARD model as well as a 3-model PLARD ensemble using the multi-scale testing for augmentation. We compare the visual image-based PSPNet and our models with other state-of-the-art road detection algorithms, including SPRAY [4], FCN-LC [8], HybridCRF[15], FTP [29], Up-Conv [39], LoDNN[12], MultiNet [38], Stixel-Net II[44], RBNet[10], LidCamNet [11] and NF2CNN. All results were computed on the KITTI evaluation server, and the results of other studies are based on the scores reported on KITTI’s website. Overall algorithm performance is shown in Table II, and the detailed performance for different tasks, such as UM, UMM, and UU, is shown in Table III.

The overall results of PLARD and other state-of-the-art road detection systems are illustrated in details in Table II. The single PLARD model alone promisingly improves the visual image-based PSPNet and achieves superior MaxF score compared to other road detection algorithms. By further augmenting the PLARD with multi-scale testing and model ensemble, we achieve the best scores for most of the metrics, demonstrating the effectiveness of the proposed system for

TABLE III
PERFORMANCE ON DIFFERENT TASKS IN THE TEST SET OF KITTI ROAD DETECTION BENCHMARK. BEST SCORES ARE HIGHLIGHTED IN **bold**. “+”: A 3-MODEL ENSEMBLE WITH MULTI-SCALE TESTING.

Methods	UM		UMM		UU	
	Max F	AP	Max F	AP	Max F	AP
SPRAY[4]	88.14 %	91.24 %	89.69 %	93.84 %	82.71 %	87.19 %
FCN_LC[8]	89.36 %	78.80 %	94.09 %	90.26 %	86.27 %	75.37 %
HybridCRF[15]	90.99 %	85.26 %	91.95 %	86.44 %	88.53 %	80.79 %
FTP [29]	91.20 %	90.60 %	92.98 %	92.89 %	89.62 %	88.93 %
Up_Conv[39]	90.48 %	88.20 %	93.89 %	92.62 %	91.89 %	89.44 %
LoDNN[12]	92.75 %	89.98 %	96.05 %	95.03 %	92.29 %	90.35 %
MultiNet[38]	93.99 %	93.24 %	96.15 %	95.36 %	93.69 %	92.55 %
StixelNet II[44]	94.05 %	85.85 %	96.22 %	91.24 %	93.40 %	85.01 %
RBNet[10]	94.77 %	91.42 %	96.06 %	93.49 %	93.21 %	89.18 %
LidCamNet [11]	95.62 %	93.54 %	97.08 %	95.51 %	94.54 %	92.74 %
NF2CNN	96.09 %	88.40 %	97.77 %	93.31 %	95.47 %	86.98 %
PSPNet[42]	95.62 %	92.95 %	96.95 %	95.38 %	95.86 %	92.73 %
PLARD	96.34 %	93.43 %	97.53 %	95.61 %	96.13 %	93.00 %
PLARD ⁺	97.05 %	93.53 %	97.77 %	95.64 %	95.95 %	95.25 %



Fig. 5. Qualitative results of PLARD in different road scenes. Presented images show that PLARD is robust to severe illumination conditions. Best view in color.

robust road detection. In particular, our method achieves the highest MaxF and AP scores which are commonly used as the performance indicators of an approach.

In addition to the overall results, we also compare performance with respect to separate tasks (“UM”, “UMM”, and “UU”) in Table III. It can be observed that PLARD also outperforms other compared algorithms in most of the Max F and AP metrics. Especially, on “UM”, the augmented PLARD system outperforms NF2CNN, the best-performing algorithm among other methods, by around 1 point for MaxF. Moreover, for AP on the “UU” task, our method surpasses the second-ranking method LidCamNet[11] by around 2.5 points, demonstrating the generalizability of the proposed PLARD

system.

2) *Qualitative Results:* Fig. 5 shows some qualitative results of PLARD on test set of the benchmark. Columns from left to right show the road detection results of PLARD on UM, UMM, and UU tasks, respectively. It can be seen that our method is robust to severe illumination conditions like heavy shadows and over-exposed areas.

V. CONCLUSIONS

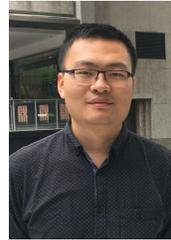
In this study, we introduce a novel road detection method called PLARD by progressively adapting LiDAR information to visual information. The PLARD first performs data space adaptation, which adapts LiDAR data to the 2D image space

to align with the perspective view by applying an altitude difference-based transformation. Then, the PLARD performs feature space adaptation to adapt the learned LiDAR features to visual features through cascaded fusion layers. By leveraging these two adaptation modules successively, PLARD takes advantage of both the visual and LiDAR information and is robust to various challenging aspects of urban scenes. We validate the PLARD model on the well-known KITTI road detection benchmark, where it outperforms other state-of-the-art road detection models and currently ranks the top of the leader-board, demonstrating its effectiveness and superiority over existing methods.

REFERENCES

- [1] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [2] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected crfs," in *ICLR*, 2015. [Online]. Available: <http://arxiv.org/abs/1412.7062>
- [3] P. Y. Shinzato, D. F. Wolf, and C. Stiller, "Road terrain detection: Avoiding common obstacle detection assumptions using sensor fusion," in *Intelligent Vehicles Symposium Proceedings, 2014 IEEE*. IEEE, 2014, pp. 687–692.
- [4] T. Kühnl, F. Kummert, and J. Fritsch, "Spatial ray features for real-time ego-lane extraction," in *VEHIT*. IEEE, 2012, pp. 288–293.
- [5] L. Xiao, B. Dai, D. Liu, T. Hu, and T. Wu, "Crf based road detection with multi-sensor fusion," in *Intelligent Vehicles Symposium (IV), 2015 IEEE*. IEEE, 2015, pp. 192–198.
- [6] C. C. T. Mendes, V. Frémont, and D. F. Wolf, "Vision-based road detection using contextual blocks," *arXiv preprint arXiv:1509.01122*, 2015.
- [7] D. Levi, N. Garnett, E. Fetaya, and I. Herzyliya, "Stixel-net: A deep convolutional network for obstacle detection and road segmentation." *BMVC*, 2015.
- [8] C. C. T. Mendes, V. Frémont, and D. F. Wolf, "Exploiting fully convolutional neural networks for fast road detection," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2016.
- [9] R. Mohan, "Deep deconvolutional networks for scene parsing," *arXiv/1411.4101*, 2014.
- [10] Z. Chen and Z. Chen, "Rbnet: A deep neural network for unified road and road boundary detection," in *International Conference on Neural Information Processing*. Springer, 2017, pp. 677–687.
- [11] L. Caltagirone, M. Bellone, L. Svensson, and M. Wahde, "Lidar-camera fusion for road detection using fully convolutional neural networks," *Robotics and Autonomous Systems*, 2018.
- [12] L. Caltagirone, S. Scheidegger, L. Svensson, and M. Wahde, "Fast lidar-based road detection using fully convolutional neural networks," in *Intelligent Vehicles Symposium (IV), 2017 IEEE*. IEEE, 2017, pp. 1019–1024.
- [13] L. Chen, J. Yang, and H. Kong, "Lidar-histogram for fast road and obstacle detection," in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 1343–1348.
- [14] Lidar. [Online]. Available: <https://en.wikipedia.org/wiki/Lidar>
- [15] L. Xiao, R. Wang, B. Dai, Y. Fang, D. Liu, and T. Wu, "Hybrid conditional random field based camera-lidar fusion for road detection," *Information Sciences*, 2017.
- [16] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *CVPR*, 2012.
- [17] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *International Journal of Robotics Research*, 2013.
- [18] L. Qi, M. Zhou, and W. Luan, "A dynamic road incident information delivery strategy to reduce urban traffic congestion," *IEEE/CAA Journal of Automatica Sinica*, vol. 5, no. 5, pp. 934–945, 2018.
- [19] L. Chen, X. Hu, W. Tian, H. Wang, D. Cao, and F. Wang, "Parallel planning: a new motion planning framework for autonomous driving," *IEEE/CAA Journal of Automatica Sinica*, pp. 1–12, 2018.
- [20] H. Kong, J.-Y. Audibert, and J. Ponce, "Vanishing point detection for road detection," in *Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2009, pp. 96–103.
- [21] Z. Chen, X. You, B. Zhong, J. Li, and D. Tao, "Dynamically modulated mask sparse tracking," *IEEE transactions on cybernetics*, vol. 47, no. 11, pp. 3706–3718, 2017.
- [22] Z. Chen, J. Li, Z. Chen, and X. You, "Generic pixel level object tracker using bi-channel fully convolutional network," in *International Conference on Neural Information Processing*. Springer, 2017, pp. 666–676.
- [23] Z. Chen, S. Huang, and D. Tao, "Context refinement for object detection," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 71–86.
- [24] Y. Xing, C. Lv, L. Chen, H. Wang, H. Wang, D. Cao, E. Velenis, and F.-Y. Wang, "Advances in vision-based lane detection: Algorithms, integration, assessment, and perspectives on acp-based parallel vision," *IEEE/CAA Journal of Automatica Sinica*, vol. 5, no. 3, pp. 645–661, 2018.
- [25] X. Han, J. Lu, C. Zhao, S. You, and H. Li, "Semi-supervised and weakly-supervised road detection based on generative adversarial networks," *IEEE Signal Processing Letters*, 2018.
- [26] J. Munoz-Bulnes, C. Fernandez, I. Parra, D. Fernández-Llorca, and M. A. Sotelo, "Deep fully convolutional networks with random data augmentation for enhanced generalization in road detection," in *International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2017, pp. 366–371.
- [27] D. Munoz, J. A. Bagnell, and M. Hebert, "Stacked hier-

- archical labeling,” in *European Conference on Computer Vision (ECCV)*. Springer, 2010, pp. 57–70.
- [28] M. Aly, “Real time detection of lane markers in urban streets,” in *Intelligent Vehicles Symposium*. IEEE, 2008, pp. 7–12.
- [29] A. Laddha, M. K. Kocamaz, L. E. Navarro-Serment, and M. Hebert, “Map-supervised road detection,” in *Intelligent Vehicles Symposium*. IEEE, 2016, pp. 118–123.
- [30] J. M. Alvarez, M. Salzmann, and N. Barnes, “Learning appearance models for road detection,” in *Intelligent Vehicles Symposium*. IEEE, 2013, pp. 423–429.
- [31] S. Zhou, J. Gong, G. Xiong, H. Chen, and K. Iagnemma, “Road detection using support vector machine based on online learning and evaluation,” in *Intelligent Vehicles Symposium*. IEEE, 2010, pp. 256–261.
- [32] L. Xiao, B. Dai, D. Liu, D. Zhao, and T. Wu, “Monocular road detection using structured random forest,” *International Journal of Advanced Robotic Systems*, vol. 13, no. 3, p. 101, 2016.
- [33] L. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. Yuille, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.” *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [34] C. Liang-Chieh, G. Papandreou, I. Kokkinos, K. Murphy, and A. Yuille, “Semantic image segmentation with deep convolutional nets and fully connected crfs,” in *ICLR*, 2015.
- [35] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” *arXiv preprint arXiv:1511.07122*, 2015.
- [36] G. Lin, A. Milan, C. Shen, and I. Reid, “Refinenet: Multi-path refinement networks for high-resolution semantic segmentation,” in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [37] F. Yu, V. Koltun, and T. Funkhouser, “Dilated residual networks,” in *Computer Vision and Pattern Recognition*, vol. 1, 2017.
- [38] M. Teichmann, M. Weber, M. Zoellner, R. Cipolla, and R. Urtasun, “Multinet: Real-time joint semantic reasoning for autonomous driving,” *arXiv preprint arXiv:1612.07695*, 2016.
- [39] G. L. Oliveira, W. Burgard, and T. Brox, “Efficient deep methods for monocular road segmentation.” in *IROS*, 2016.
- [40] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [42] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2881–2890.
- [43] J. Fritsch, T. Kuehnl, and A. Geiger, “A new performance measure and evaluation benchmark for road detection algorithms,” in *ITSC*, 2013.
- [44] N. Garnett, S. Silberstein, S. Oron, E. Fetaya, U. Verner, A. Ayash, V. Goldner, R. Cohen, K. Horn, and D. Levi, “Real-time category-based and general obstacle detection for autonomous driving,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 198–205.
- [45] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell, “Understanding convolution for semantic segmentation,” in *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2018, pp. 1451–1460.
- [46] G. Neuhold, T. Ollmann, S. Rota Bulò, and P. Kotschieder, “The mapillary vistas dataset for semantic understanding of street scenes,” in *International Conference on Computer Vision (ICCV)*, 2017. [Online]. Available: <https://www.mapillary.com/dataset/vistas>



Zhe Chen received the B.S. degree in Computer Science from University of Science and Technology of China, Hefei, China, in 2014. He is currently pursuing the Ph.D. degree at the UBTECH Sydney Artificial Intelligence Centre and the School of Computer Science, the Faculty of Engineering and Information Technologies, the University of Sydney. His current research interests include object recognition, computer vision, and deep learning. His studies was published in IEEE CVPR, ICONIP, and ECCV.



Jing Zhang received the B.S. degree from the Henan University and the Ph.D. from the University of Science and Technology of China (USTC). He used to work as a research fellow at the IFLYTEK Research. Then, he has been a lecturer at the School of Automation in the Hangzhou Dianzi University since 2017. Currently, he is a visiting scholar at the School of Software and Advanced Analytics Institute in the University of Technology Sydney. His research interests include computer vision and multimedia.

He published several papers at IEEE CVPR, ACM Multimedia, AAAI, IEEE TCSVT, Information Sciences, Neurocomputing, etc. He serves as a reviewer for a number of journals and conferences such as TIP, TCSVT, Information Sciences, Neurocomputing, ACM Multimedia.



Dacheng Tao (F'15) is Professor of Computer Science and ARC Laureate Fellow in the School of Computer Science and the Faculty of Engineering and Information Technologies, and the Inaugural Director of the UBTECH Sydney Artificial Intelligence Centre, at the University of Sydney. He mainly applies statistics and mathematics to Artificial Intelligence and Data Science. His research results have expounded in one monograph and 200+ publications at prestigious journals and prominent conferences, such as IEEE T-PAMI, T-IP, T-NNLS, T-CYB, IJCV,

JMLR, NIPS, ICML, CVPR, ICCV, ECCV, ICDM; and ACM SIGKDD, with several best paper awards, such as the best theory/algorithm paper runner up award in IEEE ICDM'07, the best student paper award in IEEE ICDM'13, the 2014 ICDM 10-year highest-impact paper award, the 2017 IEEE Signal Processing Society Best Paper Award, and the distinguished paper award in the 2018 IJCAI. He received the 2015 Austrian Scopus-Eureka Prize and the 2018 IEEE ICDM Research Contributions Award. He is a Fellow of the Australian Academy of Science, AAAS, IEEE, IAPR, OSA and SPIE.