

# Real-time Fusion Network for RGB-D Semantic Segmentation Incorporating Unexpected Obstacle Detection for Road-driving Images

Lei Sun<sup>1</sup>, Kailun Yang<sup>2</sup>, Xinxin Hu<sup>1</sup>, Weijian Hu<sup>1</sup> and Kaiwei Wang<sup>3</sup>

**Abstract**—Semantic segmentation has made striking progress due to the success of deep convolutional neural networks. Considering the demands of autonomous driving, real-time semantic segmentation has become a research hotspot these years. However, few real-time RGB-D fusion semantic segmentation studies are carried out despite readily accessible depth information nowadays. In this paper, we propose a real-time fusion semantic segmentation network termed RFNet that effectively exploits complementary cross-modal information. Building on an efficient network architecture, RFNet is capable of running swiftly, which satisfies autonomous vehicles applications. Multi-dataset training is leveraged to incorporate unexpected small obstacle detection, enriching the recognizable classes required to face unforeseen hazards in the real world. A comprehensive set of experiments demonstrates the effectiveness of our framework. On *Cityscapes*, Our method outperforms previous state-of-the-art semantic segmenters, with excellent accuracy and 22Hz inference speed at the full 2048×1024 resolution, outperforming most existing RGB-D networks.

**Index Terms**—Semantic scene understanding, RGB-D fusion, obstacle detection, autonomous driving.

## I. INTRODUCTION

ENVIRONMENT perception is a significant task for intelligent robots and systems in object classification, autonomous driving, and localization. In recent years, this field has witnessed remarkable progress thanks to deep Convolutional Neural Networks (CNNs) based semantic segmentation methods [1] [2] [3]. As an environment perception method to be applied in autonomous driving, safety, accuracy, and efficiency are the vital factors in semantic segmentation for upper-level navigational tasks. However, unexpected road hazards like debris, bricks, stones, and cargos become the most dangerous and difficult elements to detect in autonomous driving imagery. According to the AAA Foundation for Traffic Safety, debris on the road led to more than 200,000 crashes on U.S. roadways between 2011 and 2014, resulting in approximately 39,000 injuries and more than 500 deaths [4]. These obstacles are generally small in size but not fixed in shape and type, making detecting them a challenging subject that

\*This work has been partially funded through the AccessibleMaps project. This work has been supported by Hangzhou Surlmage Technology Co., Ltd.

<sup>1</sup>L. Sun, X. Hu and W. Hu are with State Key Laboratory of Modern Optical Instrumentation, Zhejiang University, China {leo\_sun, hxx\_zju, huweijian}@zju.edu.cn

<sup>2</sup>K. Yang is with Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, Germany kailun.yang@kit.edu

<sup>3</sup>K. Wang is with National Optical Instrumentation Engineering Technology Research Center, Zhejiang University, China wangkaiwei@zju.edu.cn

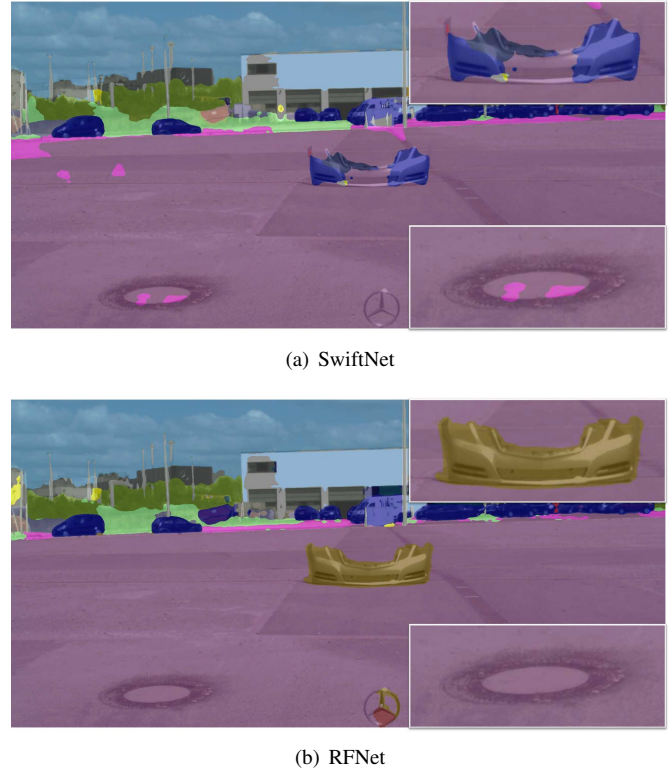


Fig. 1. Examples from the *Lost and Found* dataset and corresponding results of the methods: (a) SwiftNet (unexpected obstacle wrongly classified as car), (b) The proposed RFNet (clear and consistent segmentation).

has aroused interest among the robotics and computer vision community. For these reasons, it is desirable to develop a semantic segmentation based method incorporating pixel-wise unexpected obstacle detection.

Compared to expensive 3D sensors like LIDAR, RGB camera is a much lower cost solution with higher resolution. Based on RGB stereo camera, there have been some attempts to detect small obstacles with the help of geometry cues and CNNs [5], but only relying on apparent information in the RGB image alone is not sufficient for obstacle detection [6]. For example, manhole covers, and small obstacles can both cause gradient changes in the image. The traversable areas and obstacles in the depth map vary vastly in depth maps. Depth maps contain more location and contour information that can be used as a critical indicator of objects in real-world driving scenarios. In this sense, appropriately combining of appearance and depth is promising to improving the

performance [6] [7] [8]. But most accuracy-oriented RGB-D semantic segmentation works focus on indoor scenes [7] [9] [10], without assuring a fast inference speed that is necessary for autonomous vehicles.

On the other hand, the outstanding capacity of CNNs is based on a large amount of annotated data, especially for semantic segmentation tasks [11]. Current mainstream autonomous driving datasets generally assume only some fixed categories of objects in the scene, ignoring unforeseen hazards like unexpected small obstacles in the real world. For instance, *Cityscapes* [12] only divides objects to 19 classes, without defining any unexpected class. Multi-source training has been proven to effectively increase recognizable semantics without having to relabel the dataset [13]. However, previous multi-source training frameworks have only considered the heterogeneity in the label hierarchies of RGB data, missing the opportunity to leverage complementary depth information from different sources.

In this paper, we propose a framework that combines RGB-D semantic segmentation and obstacle detection. RFNet, a real-time fusion network for RGB-D semantic segmentation is elaborately designed. With our multi-dataset training strategy, our framework is able to classify 19 categories in *Cityscapes* incorporating pixel-wise unexpected small obstacle detection (see Figure 1). An extensive set of experiments shows the effectiveness and efficiency of the proposed framework for the semantic segmentation task. The main contributions of our work are threefold:

- We propose RFNet, a real-time fusion network for RGB-D semantic segmentation incorporating detection of unexpected obstacle, which achieves higher accuracy with fast inference compared to other state-of-the-art methods on the *Cityscapes* dataset.
- Depth complementary features are efficiently extracted in the proposed network, which improves the accuracy compared to the single RGB-stream architecture.
- Multi-dataset training and the depth stream in the architecture enable the network to work remarkably effective in detecting unexpected small objects.

## II. RELATED WORKS

### A. RGB-D Semantic Segmentation

High-quality dense depth maps from depth sensors like Kinect and RealSense boost the development of indoor semantic segmentation. Early attempt like [14] simply concatenated RGB and depth channels as a four-channel input, and fed it into a conventional RGB modal network. However, such method can not exploit complementary information from depth maps in most times [15]. Wang et al. [16] introduced depth-aware CNN which augmented conventional CNN with a depth similarity term, but it only works well with dense depth maps. Schneider et al. [17] designed a lightweight depth branch with GoogLeNet [18] and explored different points for merging the depth and RGB networks. In FuseNet [9] and RedNet [10], RGB images and depth maps are fed into two separate neural network branches respectively, which are fused before the upsampling. In [19], depth maps are pre-processed as HHA

features that encode horizontal disparity, height above ground and angle. Park et al. [20] proposed a multilevel feature fusion scheme by introducing multi-modal feature fusion to the RefineNet blocks. ACNet [7] achieved a breakthrough by proposing an attention complementary module to exploit complementary depth information efficiently. These studies prove that RGB-D semantic segmentation can achieve better segmentation results than single RGB-based methods. The major reason for this is that compared to the single RGB images, depth maps contain more location and contour information that benefit the context-critical semantic segmentation.

Compared to indoor scene depth maps from Kinect or RealSense, outdoor traffic scene depth maps are much more sparse. Li et al. [21] simply stacked smoothed depth maps with RGB images as a 4-channel input. Based on VGG [22], Kreso et al. [23] introduced a scale selection layer and used the depth maps as a guidance to produce a scale-invariant representation to free appearance from the scale. In [24], luminance information is used for depth map enhancement. Most recently, Deng et al. [25] proposed a Residual Fusion Block (RFB) to formulate the interdependencies of the encoders to extract cross-modal features based on ERFNet [26]. Low latency is crucial in autonomous driving applications, but most of these methods cannot meet the real-time constraint. In this paper, we propose a real-time fusion network to achieve swift inference while retaining a highly competitive performance among the state of the art for RGB-D segmentation.

### B. Unexpected Obstacle Detection for Self-driving Cars

Detecting unexpected small but potentially hazardous obstacles on the road is a vital task for autonomous driving, and this subject has always been a research hotspot. Generally these methods for detecting and localizing generic obstacles are based on stereo cameras integrated on self-driving cars. Among these methods, most are based on the generic geometric criteria. The Stixel algorithm [27] represents obstacles with a set of rectangular vertical obstacle segments, providing a robust representation of the 3D scene. Geometric point cluster methods like [28] and [29] exploit geometric relation between 3D points to detect and cluster obstacle points.

Because of the superiority in making use of visual appearance and context of images, CNNs are adopted in contemporary researches. Ramos et al. [8] presented a principled Bayesian framework to fuse the semantic segmentation predicted from a convolutional neural network and stereo-based detection results from the Fast Direct Planar Hypothesis Testing (FPHT) method. MergeNet [6] was proposed with a multi-stage training procedure involving weight sharing, separating learning of low and high level features from the RGB-D input and a refining stage which learns to fuse the obtained complementary features. But all these methods can only predict three main classes: free-space, obstacle, and background. To meet the demands of autonomous driving, we need a more universal approach that can enrich the detectable semantics beyond simple roads/obstacles separation. In this work, we address unexpected obstacle detection by incorporating it in a multi-source semantic segmentation framework to provide a unified pixel-wise scene understanding.

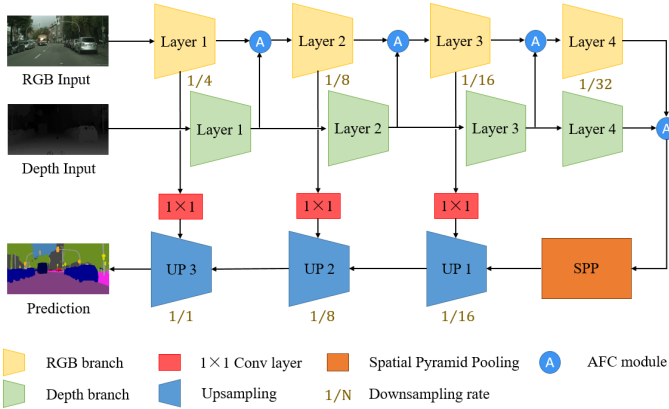


Fig. 2. Overview of RFNet: the proposed network architecture for real-time fusion-based RGB-D semantic segmentation.

### III. METHODOLOGY

#### A. Network Architecture

The entire network architecture of RFNet is shown in Figure 2. In the encoder part of the architecture, we design two independent branches to extract features for RGB and depth images separately—RGB branch as the main branch and Depth branch as the subordinate branch. In both branches, we choose ResNet-18 [30] as the backbone to extract features from inputs because ResNet-18 has moderate depth and residual structure, and its small operation footprint is compatible with real-time operation. After each layer of ResNet-18, the output features from Depth branch are fused to RGB branch after the Attention Feature Complementary (AFC) module. The spatial pyramid pooling (SPP) block gathers the fused RGB-D features from two branches and produces feature maps with multi-scale information. Finally, referred to SwiftNet [31], we design the efficient upsampling modules to restore the resolution of these feature maps with skip connections from the RGB branch.

**RGB-D fusion module.** As discussed in the last part, the depth maps contain more contour and location information that benefit RGB semantic segmentation. In order to fuse RGB and depth information effectively, we design an RGB-D fusion module termed Attention Feature Complementary (AFC) module (shown in Figure 3) to make the network focus on learning more complementary informative features from RGB and Depth branches. As shown in Figure 3, in the AFC module, we leverage a SE block [32] as the channel attention method. SE block can learn to use global information to emphasize informative channels and suppress less useful channels, which helps the AFC module exploit informative features from both branches effectively.

With the multi-branch architecture, we have the RGB input feature maps  $X = [x_1, \dots, x_C] \in \mathbb{R}^{C \times H \times W}$  and depth input feature maps  $Y = [y_1, \dots, y_C] \in \mathbb{R}^{C \times H \times W}$ . First we use global average pooling as a channel descriptor based on channel attention mechanism, then we add a  $1 \times 1$  convolution layer with the same channels as input. This  $1 \times 1$  convolution layer is able to excavate correlations between channels. The followed sigmoid function is applied to activate the convolution result and constrain the value of the weight vector between 0 and

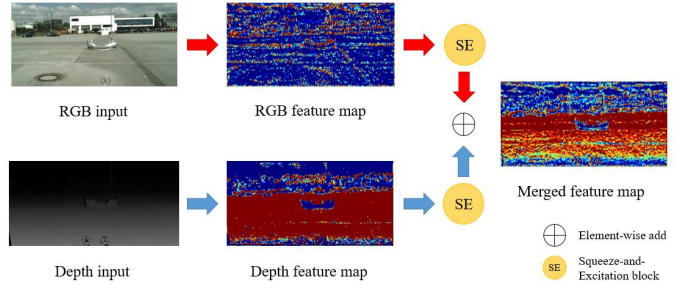


Fig. 3. AFC: Attention Feature Complementary module to exploit cross-model information from RGB and Depth inputs.

1. Next, we do outer product for the weight vector and input feature maps in both branches. Finally, by adding results from RGB branch and Depth branch, we have the resulted feature map  $Z \in \mathbb{R}^{C \times H \times W}$ , expressed as:

$$Z = X \otimes \sigma_1 [\phi_1 (X)] + Y \otimes \sigma_2 [\phi_2 (Y)] \quad (1)$$

Here,  $\phi$  denotes global pooling and  $1 \times 1$  convolution.  $\otimes$  and  $\sigma$  denote outer product and sigmoid function respectively. By applying such attention mechanism in RGB-D fusion, more informative features obtain higher values of weights, which helps us exploit complementary information from depth maps more effectively.

After four ResNet blocks and AFC module, the fused feature maps contain rich high-level semantic information. In order to increase the receptive field to cover pixels of large objects while maintaining a real-time speed, referred to [31] [33] [34], we adopt Spatial Pyramid Pooling (SPP) to average features over aligned grids with different granularities before the upsampling.

**Efficient upsampling module.** The purpose of the decoder is to upsample semantically rich visual features in coarse spatial resolution to the input resolution. We adopt a simple decoder that contains three simple upsampling modules with skip connections from the encoder. In the first two upsampling modules, low-resolution feature maps from the former block are upsampled with bilinear interpolation to the same resolution as feature maps from skip connection, then these two streams of feature maps are element-wisely added and finally mixed with a  $3 \times 3$  convolution. The third upsampling module is slightly different because we add a convolution layer and a 4-times bilinear interpolation at last to restore to the same resolution as the input. More precisely, the skip connection is routed before the second ReLU of the residual block because the current study shows that skip connection from any other stage impairs the accuracy [31].

#### B. Multi-Dataset Learning

As a data-driven technology, annotated labels are essential for semantic segmentation, but we can not annotate all classes in the real world. In order to utilize as much and diverse training data as possible and increase the number of recognizable classes from a few dozens to virtually anything that a scene can contain, multi-source learning is an effective

method. However, simply mixing two or more datasets for training may cause some problems. As shown in Figure 4, the heterogeneity in the annotation type and sample amount may cause overfitting to one of the data sources, leading to incomplete segmentation when simply mixing the datasets.

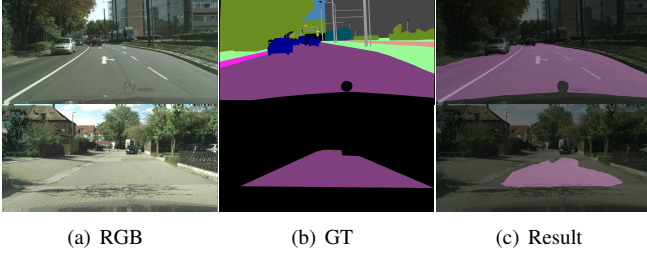


Fig. 4. RGB images and ground truth from *Cityscapes* (first row) and *Lost and Found* (second row) respectively. The last column shows the inference result if simply training on two datasets without consideration of the heterogeneity in the annotation style.

This is because classes in different datasets may conflict with each other. For example, the annotation type of these classes is different, or a certain class is a subclass of a class in another dataset. To facilitate multi-source learning with such heterogeneity, we design some training strategies. Formally, we have datasets  $D_1, \dots, D_c, \dots, D_n$ . Class set  $A$  contains classes that do not conflict with each other, and class set  $B$  contains the rest. For these conflicted classes, we refer to dataset  $D_c$  as a standard annotation. Let us denote an image by  $x$ , and the corresponding human annotation for  $x$  is provided and denoted by  $y$ , where  $y(m, n) \in 1, \dots, C$  is the label of pixel  $x(m, n)$ , and  $C$  is the total number of classes.  $l$  denotes the total number of images in all datasets and  $\phi$  denotes the segmentation model. We train on the joint multi-datasets with the loss function shown below:

$$loss = \frac{1}{l} \sum_{i=1}^l [L_A(\phi(x_i), y_i) + \lambda L_B(\phi(x_i), y_i)] \quad (2)$$

where  $L_A(\cdot, \cdot)$  and  $L_B(\cdot, \cdot)$  denotes cross entropy loss function for class set  $A$  and  $B$  respectively.  $\lambda$  is a hyper-parameter balancing the weights of different classes, and in our work we set  $\lambda$  as following:

$$\lambda = \begin{cases} 1, & x_i \in D_c \\ 0, & x_i \notin D_c \end{cases} \quad (3)$$

For instance, in this work we leverage *Cityscapes* [12] and *Lost and Found* [5]. Note that *Cityscapes* has annotations on 19 classes except for those unexpected small obstacles. We make some modifications to the loss function while training. The road class in *Cityscapes* and small obstacle class in *Lost and Found* do not conflict with classes in other datasets, which belong to class set  $A$ . The rest classes in *Cityscapes* are conflicted with the background class in *Lost and Found*, so they are divided into class  $B$ . In this situation, we assume *Cityscapes* as standard dataset  $D_c$ . In the training stage, background class and free-space in *Lost and Found* should not be counted in the loss function. In our situation, the

ignorance of background class makes coarse-annotated free-space class in *Lost and Found* helpful for improving the training data amount, so we also include free-space in the final loss. With the presented multi-dataset training strategy, our RFNet learns to predict 19 classes from *Cityscapes* and the critical unexpected small obstacle class from the *Lost and Found* dataset.

Although unexpected small obstacle is a generalized conception, which is not limited to obstacle types in *Lost and Found*, the definition of this particular set of classes allows us to meet the demand by exploiting the power of deep learning methods. For example, learning that all kinds of obstacles have some common contextual property, being of small dimensions and surrounded at least partly by free-space. Thereby, the network is able to generalize far beyond its training data with the multi-source learning strategy when facing innumerable possible corner cases.

## IV. EXPERIMENTS

### A. Datasets

In this work, two RGB-D semantic segmentation datasets: *Cityscapes* and *Lost and Found* are exploited.

*Cityscapes* [12] is a large-scale RGB-D dataset that focuses on semantic understanding of urban street scenes. It contains 2975/500/1525 images in the training/validation/testing subsets, both with finely annotated labels on 19 classes. The images cover 50 different cities with a full resolution of  $2048 \times 1024$ .

The *Lost and Found* [5] dataset consists of 2014 annotated frames from 112 stereo video sequences, along with coarse annotations of free-space areas and fine-grained annotations of the small obstacles on the road. Among them, training set and validation set contain 814 and 1200 images with a resolution of  $2048 \times 1024$ , covering different small obstacles present at long distance with non-uniform road textures/appearances and pathways with many non-obstacle class objects acting as distractors.

Both disparity images from *Cityscapes* and *Lost and Found* are obtained by using the semi-global matching algorithm [35], which is a sophisticated method for the estimation of a dense disparity map from a rectified stereo image pair.

### B. Implementation Details

The models were implemented on a single 2080Ti GPU with CUDA 10.0, CUDNN 7.6.0, and PyTorch 1.1. Adam [36] is used for optimization with the learning rate set to  $4 \times 10^{-4}$ , where cosine annealing learning rate scheduling policy [37] is adopted to adjust learning rate with a minimum value of  $1 \times 10^{-6}$  in the last epoch. The weight decay is set to  $1 \times 10^{-4}$ . We initialize the ResNet-18 in both RGB branch and Depth branch with pre-trained weights from ImageNet [38], and initialize the rest part of the model with kaiming initialization [39]. More precisely, we average the weights for RGB inputs to match the shape of one-channel depth image in the Depth branch, as research works [7] [17] show that RGB pre-trained weights also boost depth image feature extraction. For pre-trained parameters, we update them with a 4 times

smaller learning rate and apply 4 times smaller weight decay. Because the left and bottom part of the disparity images are not applicable due to the restrictions of semi-global matching algorithm, we crop these pixels and resize images back to the original resolution with bilinear upsampling. The rest of the data augmentation operations consist of scaling with random factors between 0.5 and 2, random horizontal flipping, and random cropping with an output resolution of  $768 \times 768$ . We train all the models for 200 epochs with a batch size of 8.

### C. Results and Analysis

**Ablation Study.** We perform the ablation study on our RFNet to explore the influence of different architecture variants and fusion schemes on the network accuracy where the results are shown in Table I. Results in this section are obtained by evaluating on the blended validation set of *Cityscapes* and *Lost and Found*, which includes all images from both validation datasets. All backbones in these models are initialized with ImageNet pre-trained weights.

In the table, the single RGB method only exploits the RGB branch of RFNet. Here, compared to SwiftNet [31], the only difference is the SE block after each block of the ResNet-18. It is a control group to determine if the depth information helps improve the accuracy, which achieves a mean Intersection over Union (mIoU) of 69.20%. In the RGB-D-Stack method, we stack depth maps with respective RGB images to form a 4-channel input to the single branch of RFNet. The low accuracy of the method (65.20% in mIoU) proves that depth information is not exploited effectively in this way. We also design RGB-D-Fusion (concatenation), where the only difference of this method to the RGB-D-Fusion (element-wise add) in our RFNet is that RGB feature maps and depth feature maps are concatenated to a higher dimension feature maps and restore to the original dimension after a  $1 \times 1$  convolution. Results show that this method (68.67%) performs clearly worse than RFNet (72.22%). This is because in a compact network like RFNet, concatenation is a more inefficient way to make use of the depth information.

To eliminate the cause that more parameters in two-branch RFNet make it perform better, we design and train the RGB-RGB-Fusion method. The difference of the RGB-RGB-Fusion method to RFNet is that inputs are duplicate RGB images instead of RGB-D images, and after each AFC module, the element-wise added feature maps are divided by 2. The accuracy (69.37%) is much lower than RFNet and approximately the same as the single RGB method, proving the benefit of fusion in RFNet is not simply owing to the increased parameters. We also perform an experiment to explore the influence of the proposed multi-dataset training strategy. It turns out that without multi-dataset training strategy, our RFNet gets nearly 20% lower IoU because of the class conflictions in two datasets. Finally the proposed RFNet with the proposed multi-dataset training strategy achieves a mIoU of 72.22%, which is significantly better than the baseline (single RGB architecture) and other fusion-based variants, demonstrating the effectiveness of our fusion scheme bridged by the designed attention complementary modules.

**Numerical Performance Comparison.** Based on our multi-dataset training, we create a benchmark to compare our RFNet with the other two real-time networks: ERF-PSPNet [33] (a light-weight network), SwiftNet [31] (whose network architecture is very similar to our RFNet). The first two networks only take RGB input. Table II shows IoU of all 20 classes in the new multi-source setting. Our RFNet achieves higher accuracy in most of the classes. Compared to SwiftNet, RFNet improves accuracy remarkably in certain classes like fence, traffic light, terrain, truck, bus, train, and small obstacle, which is benefited from the depth complementary information. Figure 5 shows some examples from the validation set of *Cityscapes* and *Lost and Found*, which demonstrates the excellent segmentation accuracy of our RFNet in various scenarios with or without small obstacles.

To explore how the proposed RFNet improves precision in different depth ranges, we perform analysis on mean IoU and the IoU of small obstacle in different depth ranges for RFNet and SwiftNet. We calculate the depth value of each pixel from disparity value. The maximum depth value is set to 100 and limited by the quality of disparity image, while all the unmatched pixels are set to 100. Bar graph 6 shows that RFNet performs better in all depth ranges in the case of mean IoU of 20 classes. Specifically, RFNet boosts the accuracy of unexpected small obstacle recognition in close and middle ranges remarkably. This is reasonable because disparity images derived from semi-global matching algorithm have higher accuracy at close range, and contribute more to the prediction than pixels with greater depth values.

In Table III we also compare our RFNet with other state-of-the-art networks on the *Cityscapes* validation set. The column of speed reports the inference speed of a full resolution image ( $2048 \times 1024$ ) on a single RTX 2080Ti. Specifically, ERF-PSPNet and SwiftNet are implemented on the same hardware. Compared to mainstream RGB semantic segmentation networks, our RFNet achieves better results while maintaining a real-time performance, which proves that exploiting depth information helps improving accuracy. In the table, we also list some other RGB-D fusion networks: LDFNet [24] and RFBNet based on ERFNet [25]. Our RFNet is both more accurate and faster than these multimodal networks. Overall, rare multi-modal semantic segmentation methods meet the real-time prediction speed, while our method achieves the highest accuracy on the validation set of *Cityscapes* to the best of our knowledge that meets both demands including real-time inference, highly qualified accuracy, and capacity to leverage complementary features in cross-modal imagery.

**Qualitative Performance Study.** We present the qualitative examples in Figure 8. In this paper, the main purpose of exploiting depth information is to enhance the segmentation accuracy in classes which are difficult for the RGB method, including small obstacles. The appearance and surface texture of the small obstacles are not fixed, and it is easy to be confused with graffiti, manhole covers, zebra crossings on the road. Comparatively, in the depth map where texture is ignored, the contour of small obstacles is clear. Graffiti, manhole covers are flat, making it part of the road surface in the depth map. All these features of depth maps enable to

TABLE I  
PERFORMANCE OF RFNET ON THE CITYSCAPES AND LOST AND FOUND VALIDATION SET WITH DIFFERENT DESIGN CHOICES.

Method	RGB-D Fusion	Dual-branch	Concatenation	Element-wise add	mIoU(%)	Params
Single RGB					69.20%	12.17M
RGB-D-Stack	✓				65.20%	12.17M
RGB-D-Fusion (concatenation)	✓	✓	✓		68.67%	25.08M
RGB-RGB-Fusion (element-wise add)		✓		✓	69.37%	23.69M
RFNet (without multi-dataset training strategy)	✓	✓		✓	53.83%	23.69M
RFNet	✓	✓		✓	72.22%	23.69M

TABLE II  
PER-CLASS IOU(%) RESULTS OF THREE NETWORKS ON THE BLENDED VALIDATION SET OF CITYSCAPES AND LOST AND FOUND DATASET. LIST OF CLASSES(FROM LEFT TO RIGHT): ROAD, SIDEWALK, BUILDING, WALL, FENCE, POLE, TRAFFIC LIGHT, TRAFFIC SIGN, VEGETATION, TERRAIN, SKY, PEDESTRIAN, RIDER, CAR, TRUCK, BUS, TRAIN, MOTORBIKE, BICYCLE AND SMALL OBSTACLE.

Network	Roa	Sid	Bui	Wal	Fen	Pol	TLi	TSi	Veg	Ter	Sky	Ped	Rid	Car	Tru	Bus	Tra	Mot	Bic	SOB	mIoU
ERF-PSPNet	89.3	65.1	82.1	43.4	39.8	46.9	48.2	46.1	86.4	45.2	88.4	68.5	57.2	90.1	55.6	62.1	65.6	56.9	64.8	60.3	63.1
SwiftNet	95.7	60.6	89.1	50.9	53.6	56.9	61.1	71.4	90.7	55.0	92.2	75.2	58.5	92.7	65.3	81.3	70.0	56.2	72.1	62.8	70.6
Our RFNet	96.0	60.6	90.8	50.2	59.9	60.0	62.6	72.8	91.1	57.3	92.5	76.1	57.9	93.3	73.8	82.3	73.2	54.0	72.7	67.9	72.2

TABLE III  
COMPARISON OF SEMANTIC SEGMENTATION METHODS ON THE VALIDATION SET OF CITYSCAPES.

Network	Multimodal	mIoU(%)	Speed (FPS)
FCN8s [1]	✗	65.3%	2.0 *
DeepLabV2-CRF [2]	✗	70.4%	n/a
ENet [40]	✗	58.3%	76.9 *
ERFNet [26]	✗	65.8%	20.8
ERF-PSPNet [33]	✗	64.1%	20.4
SwiftNet [31]	✗	72.0%	41.0
VGG-D (ScaleInvariant) [23]	✓	64.4%	n/a
LDNet [24]	✓	68.5%	18.4
GoogLeNet (NiN-2) [17]	✓	69.1%	n/a
RFBNet (ERFNetEnc) [25]	✓	72.0%	n/a
RFNet (Ours)	✓	72.5%	22.2

\* Speed on half resolution images.

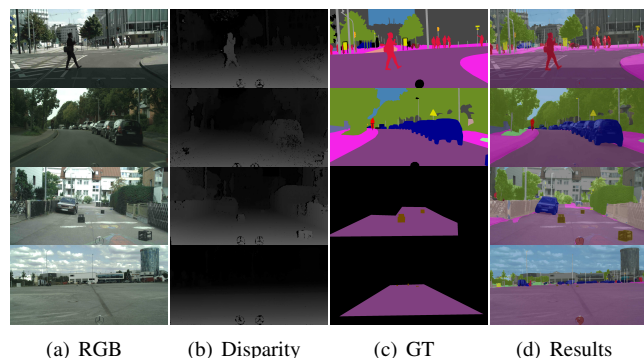


Fig. 5. Predictions with additional unexpected obstacle class from RFNet.

reduce the chance of false alarm in detecting small obstacles. We compare our RFNet with SwiftNet [31], which has a similar network architecture with RFNet, where Figure 8 shows representative contrast results from the two networks. As a purely RGB-based method, SwiftNet fails to predict some small obstacles on the road and predicts manhole cover as small obstacle. RFNet correctly detect small obstacles and classifies the manhole as part of the road, which demonstrates the superiority of our method for safety-critical road sensing. RFNet also performs better in large-scale objects like bus and truck because contours of these classes are much clearer in

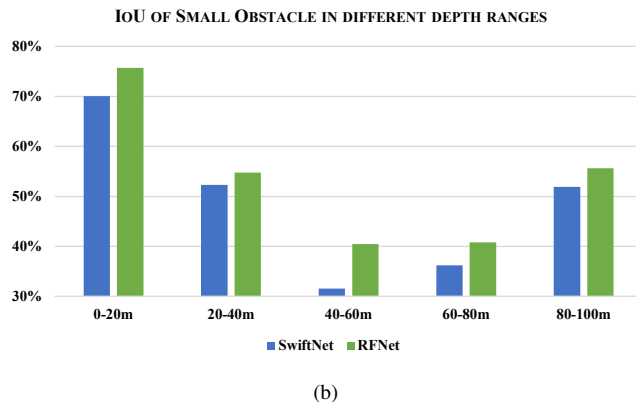
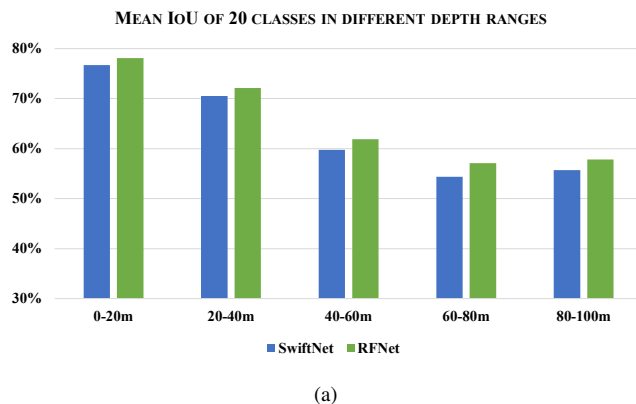
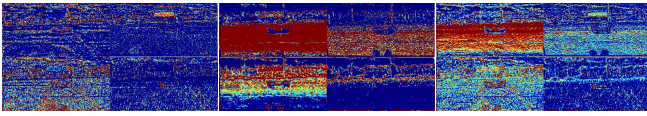


Fig. 6. Mean IoU of 20 classes and IoU of Small Obstacle from SwiftNet and RFNet respectively. RFNet improves precision in all depth ranges, especially in close and middle ranges.

depth maps compared to RGB images.

Furthermore, for the input image from Figure 1, Figure 7 shows the feature maps after the second block from RFNet, in which the first two are feature maps from RGB and depth branch respectively, and the merged feature maps are from the output part of the AFC module. As it can be clearly seen, compared to RGB feature maps, small obstacle is much more clear and manhole cover disappears in depth feature maps,



(a) RGB feature maps (b) Depth feature maps (c) Merged feature maps

Fig. 7. Visualization of feature maps from the second block of RFNet.

while the feature map after AFC module takes the advantages of both branches. In summary, the AFC module enables RFNet to effectively exploit the depth features in a complementary way, improving the accuracy of obstacle detection evidenced by both numerical and qualitative results.

## V. CONCLUSION

In this study, we propose RFNet, a real-time fusion network for RGB-D semantic segmentation on road-driving images. With the designed AFC module, RFNet exploits complementary depth information effectively and significantly improves the accuracy over purely RGB-based methods. With the presented multi-source training strategy, RFNet can also detect unexpected small obstacles, enriching the recognizable classes required to face the real world with unforeseen hazards. More importantly, RFNet operates at 22Hz with full resolution *Cityscapes* images and 41.6Hz with half resolution on a single Nvidia GTX2080Ti GPU, which makes it ideally suitable for autonomous driving applications. Our RFNet outperforms state-of-the-art RGB-D fusion methods in terms of accuracy and speed. In the future, we plan to further streamline RFNet and deploy it to portable TPU devices with robustness augmented. The source code of our RFNet is available at <https://github.com/AHupuJR/RFNet>.

## 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, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.
- [3] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.
- [4] B. C. Tefft, "The prevalence of motor vehicle crashes involving road debris, united states, 2011-2014." *Age (years)*, vol. 20, no. 5.7, pp. 10–1, 2016.
- [5] P. Pinggera, S. Ramos, S. Gehrig, U. Franke, C. Rother, and R. Mester, "Lost and found: detecting small road hazards for self-driving vehicles," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 1099–1106.
- [6] K. Gupta, S. A. Javed, V. Gandhi, and K. M. Krishna, "Mergenet: A deep net architecture for small obstacle discovery," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 1–7.
- [7] X. Hu, K. Yang, L. Fei, and K. Wang, "Acnet: Attention based network to exploit complementary features for rgb-d semantic segmentation," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 1440–1444.
- [8] S. Ramos, S. Gehrig, P. Pinggera, U. Franke, and C. Rother, "Detecting unexpected obstacles for self-driving cars: Fusing deep learning and geometric modeling," in *2017 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2017, pp. 1025–1032.
- [9] C. Hazirbas, L. Ma, C. Domokos, and D. Cremers, "Fusenet: Incorporating depth into semantic segmentation via fusion-based cnn architecture," in *Asian conference on computer vision*. Springer, 2016, pp. 213–228.
- [10] J. Jiang, L. Zheng, F. Luo, and Z. Zhang, "Rednet: Residual encoder-decoder network for indoor rgb-d semantic segmentation," *arXiv preprint arXiv:1806.01054*, 2018.
- [11] L. Sun, K. Wang, K. Yang, and K. Xiang, "See clearer at night: towards robust nighttime semantic segmentation through day-night image conversion," in *Artificial Intelligence and Machine Learning in Defense Applications*, vol. 11169. International Society for Optics and Photonics, 2019, p. 111690A.
- [12] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [13] P. Meletis and G. Dubbelman, "Training of convolutional networks on multiple heterogeneous datasets for street scene semantic segmentation," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 1045–1050.
- [14] C. Couprie, C. Farabet, L. Najman, and Y. LeCun, "Indoor semantic segmentation using depth information," *arXiv preprint arXiv:1301.3572*, 2013.
- [15] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Y. Ng, "Multimodal deep learning," in *Proceedings of the 28th international conference on machine learning (ICML-11)*, 2011, pp. 689–696.
- [16] W. Wang and U. Neumann, "Depth-aware CNN for RGB-D segmentation," *CoRR*, vol. abs/1803.06791, 2018. [Online]. Available: <http://arxiv.org/abs/1803.06791>
- [17] L. Schneider, M. Jasch, B. Fröhlich, T. Weber, U. Franke, M. Pollefeys, and M. Rätzsch, "Multimodal neural networks: Rgb-d for semantic segmentation and object detection," in *Scandinavian conference on image analysis*. Springer, 2017, pp. 98–109.
- [18] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 2015, pp. 1–9.
- [19] S. Gupta, R. Girshick, P. Arbeláez, and J. Malik, "Learning rich features from rgb-d images for object detection and segmentation," in *European conference on computer vision*. Springer, 2014, pp. 345–360.
- [20] S.-J. Park, K.-S. Hong, and S. Lee, "Rdfnet: Rgb-d multi-level residual feature fusion for indoor semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 4980–4989.
- [21] L. Li, B. Qian, J. Lian, W. Zheng, and Y. Zhou, "Traffic scene segmentation based on rgb-d image and deep learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 5, pp. 1664–1669, 2017.
- [22] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *Computer Science*, 2014.
- [23] I. Krešo, D. Čaušević, J. Krapac, and S. Šegvić, "Convolutional scale invariance for semantic segmentation," in *German Conference on Pattern Recognition*. Springer, 2016, pp. 64–75.
- [24] S.-W. Hung, S.-Y. Lo, and H.-M. Hang, "Incorporating luminance, depth and color information by a fusion-based network for semantic segmentation," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 2374–2378.
- [25] L. Deng, M. Yang, T. Li, Y. He, and C. Wang, "Rfbnet: deep multimodal networks with residual fusion blocks for rgb-d semantic segmentation," *arXiv preprint arXiv:1907.00135*, 2019.
- [26] E. Romera, J. M. Alvarez, L. M. Bergasa, and R. Arroyo, "Erfnet: Efficient residual factorized convnet for real-time semantic segmentation," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 1, pp. 263–272, 2017.
- [27] D. Pfeiffer and U. Franke, "Towards a global optimal multi-layer stixel representation of dense 3d data," in *BMVC*, vol. 11, 2011, pp. 51–1.
- [28] R. Manduchi, A. Castano, A. Talukder, and L. Matthies, "Obstacle detection and terrain classification for autonomous off-road navigation," *Autonomous robots*, vol. 18, no. 1, pp. 81–102, 2005.
- [29] A. Broggi, M. Buzzoni, M. Felisa, and P. Zani, "Stereo obstacle detection in challenging environments: the viac experience," in *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2011, pp. 1599–1604.
- [30] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *European conference on computer vision*. Springer, 2016, pp. 630–645.
- [31] M. Orsic, I. Kreso, P. Bevandic, and S. Segvic, "In defense of pre-trained imagenet architectures for real-time semantic segmentation of road-

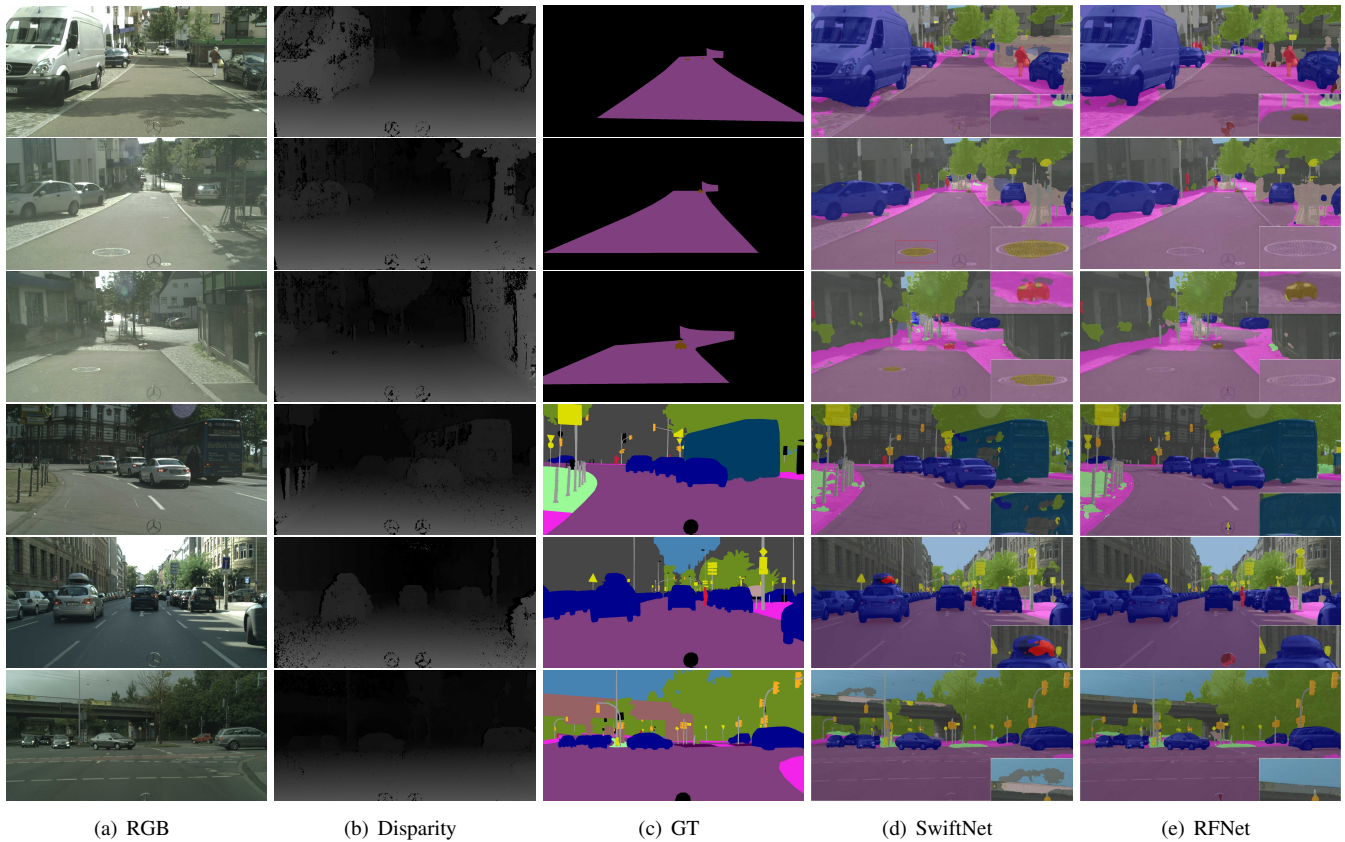


Fig. 8. Qualitative semantic segmentation results from SwiftNet and the proposed RFNet that exploits both RGB and depth information.

- driving images,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 607–12 616.
- [32] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [33] K. Yang, K. Wang, L. M. Bergasa, E. Romera, W. Hu, D. Sun, J. Sun, R. Cheng, T. Chen, and E. López, “Unifying terrain awareness for the visually impaired through real-time semantic segmentation,” *Sensors*, vol. 18, no. 5, p. 1506, 2018.
- [34] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.
- [35] H. Hirschmuller, “Accurate and efficient stereo processing by semi-global matching and mutual information,” in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, vol. 2. IEEE, 2005, pp. 807–814.
- [36] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [37] I. Loshchilov and F. Hutter, “Sgdr: Stochastic gradient descent with warm restarts,” *arXiv preprint arXiv:1608.03983*, 2016.
- [38] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [39] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1026–1034.
- [40] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, “Enet: A deep neural network architecture for real-time semantic segmentation,” *arXiv preprint arXiv:1606.02147*, 2016.