

# LIDAR-Camera Fusion for Road Detection Using Fully Convolutional Neural Networks

Luca Caltagirone\*, Mauro Bellone, Lennart Svensson, Mattias Wahde

{luca.caltagirone, mauro.bellone, lennart.svensson, mattias.wahde}@chalmers.se

**Abstract**—In this work, a deep learning approach has been developed to carry out road detection by fusing LIDAR point clouds and camera images. An unstructured and sparse point cloud is first projected onto the camera image plane and then upsampled to obtain a set of dense 2D images encoding spatial information. Several fully convolutional neural networks (FCNs) are then trained to carry out road detection, either by using data from a single sensor, or by using three fusion strategies: early, late, and the newly proposed *cross fusion*. Whereas in the former two fusion approaches, the integration of multimodal information is carried out at a predefined depth level, the cross fusion FCN is designed to directly learn from data where to integrate information; this is accomplished by using trainable cross connections between the LIDAR and the camera processing branches.

To further highlight the benefits of using a multimodal system for road detection, a data set consisting of visually challenging scenes was extracted from driving sequences of the KITTI raw data set. It was then demonstrated that, as expected, a purely camera-based FCN severely underperforms on this data set. A multimodal system, on the other hand, is still able to provide high accuracy. Finally, the proposed cross fusion FCN was evaluated on the KITTI road benchmark where it achieved excellent performance, with a MaxF score of 96.03%, ranking it among the top-performing approaches.

## I. INTRODUCTION

Road detection is an important task that needs to be solved accurately and robustly in order to achieve higher automation levels. Knowing what regions of the road surface are available for driving is in fact a crucial prerequisite for carrying out safe trajectory planning and decision-making. Although some automated driving vehicles are already available on the market, the recent crash of a Tesla car controlled by its autopilot system highlighted that further research and testing are very much necessary. In that case, it was pointed out that a possible reason for the crash was that the autopilot system misinterpreted the trailer of a truck as free road due to unfavourable lighting conditions [1], [2].

Current approaches for road detection use either cameras or LIDAR sensors. Cameras can work at high frame-rate, and provide dense information over a long range under good illumination and fair weather. However, being passive sensors, they are strongly affected by the level of illumination. A passive sensor is able to receive a specific amount of energy from the environment, light waves in the case of cameras,

and transform it into a quantitative measure (image). Clearly, the process depends on the amplitude and frequency of the light waves, influencing the overall result, while a reliable system should be invariant with respect to changes in illumination [3]. LIDARs sense the environment by using their own emitted pulses of laser light and therefore they are only marginally affected by the external lighting conditions. Furthermore, they provide accurate distance measurements. However, they have a limited range, typically between 10 and 100 meters, and provide sparse data.

Based on this description of benefits and drawbacks of these two sensor types, it is easy to see that using both might provide an improved overall reliability. Inspired by this consideration, the work presented here investigates how LIDAR point clouds and camera images can be integrated for carrying out road segmentation. The choice to use a fully convolutional neural network (FCN) for LIDAR-camera fusion is motivated by the impressive success obtained by deep learning algorithms in recent years in the fields of computer vision and pattern recognition [4].

In summary, this work makes the following two main contributions: (i) A novel LIDAR-camera fusion FCN that outperforms established approaches found in the literature and achieves state-of-the-art performance on the KITTI road benchmark; (ii) a data set of visually challenging scenes extracted from KITTI driving sequences that can be used to further highlight the benefits of combining LIDAR data and camera images for carrying out road segmentation.

The remainder of the paper is structured as follows: Sect. II gives a brief overview of related approaches that deal with the problems of road detection or sensor fusion. The FCN base architecture and the fusion strategies are presented in Sect. III. Section IV describes the procedure to transform a sparse 3D point cloud into a set of dense 2D images. The experimental results and discussion are reported in Sect. V which is followed, in Sect. VI, by a summary and the conclusions.

## II. RELATED WORK

The study of road detection can be tracked back a few decades; in [5] and [6], Broggi *et al.* already presented an algorithm for the binarization, classification, and interpretation of visual images for road detection. However, recent advances in sensor development and hardware computation capacity has made possible the use of high-accuracy methods. Nowadays, the large majority of state-of-the-art algorithms for road detection use, to different extent, machine learning

\*Corresponding author. Luca Caltagirone, Mauro Bellone, and Mattias Wahde are with the Adaptive Systems Research Group, Department of Mechanics and Maritime Sciences, Chalmers University of Technology, Gothenburg, Sweden. Lennart Svensson is with the Department of Electrical Engineering, also at Chalmers University of Technology.

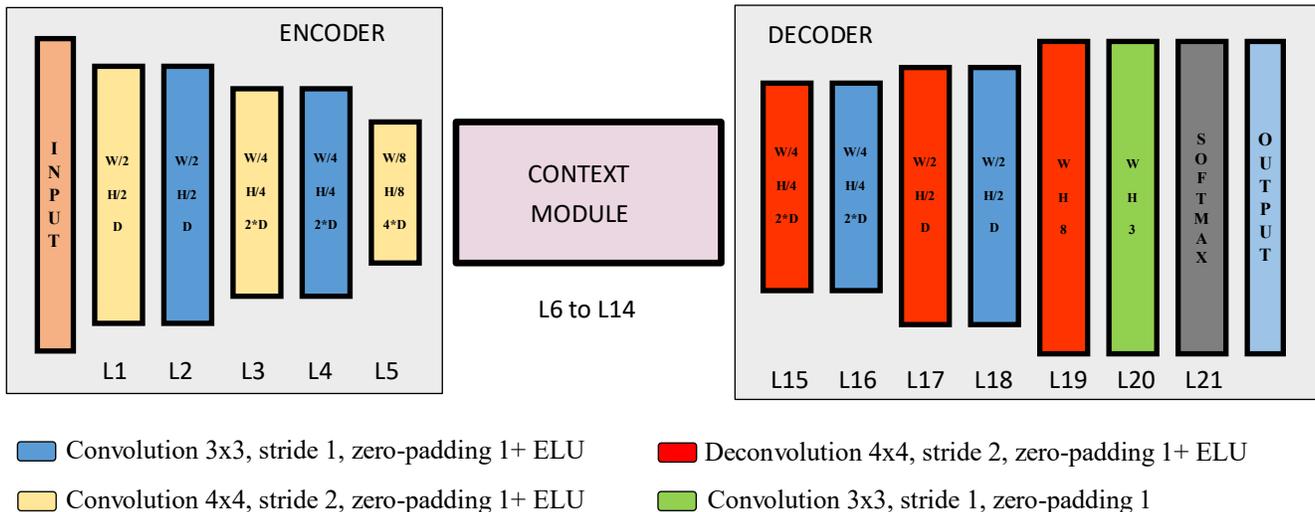


Fig. 1. A schematic illustration of the proposed base FCN architecture which consists of 21 layers. W represents the width, H denotes the height, and D is the number of feature maps in the first layer which was set to 32. The FCN uses the exponential linear unit (ELU) activation function after each convolutional layer. See Table I for details about the context module architecture.

techniques. Teichmann *et al.* [7], for example, trained a convolutional neural network (CNN) to simultaneously solve the tasks of road segmentation and vehicle detection in monocular camera images. Chen *et al.* [8] developed a deep neural network within a Bayesian framework to jointly estimate the road surface and its boundaries. In [9], LIDAR point clouds are transformed into 2D top-view images that are then used as input for an FCN to carry out road segmentation. Shinzato *et al.* [10] instead projected the point clouds onto the camera image plane and then used a graph-based approach to discriminate between obstacle and obstacle-free space. Some methods, such as those found in [11] and [12], tackled road detection by performing LIDAR-camera fusion within the framework of conditional random fields (CRFs).

Eitel *et al.* [13] proposed to carry out objection recognition by fusing depth maps and color images with a CNN. In [14], LIDAR point clouds were transformed into their HHA (horizontal disparity, height above the ground, and angle) representation [15] and then combined with RGB images using a variety of CNN fusion strategies for performing pedestrian detection. More recently, Asvadi *et al.* [16] developed a system for vehicle detection that integrates LIDAR and color camera data within a deep learning framework.

Investigating another line of research, in [17] a support vector machine (SVM) to carry out road detection on 3D cloud data in challenging scenarios. Using SVM, Zhou *et al.* [18], built a road detection algorithm enabling on-line learning, meaning that this method is able to update the training data, thus reducing the probability of misclassification. Moreover, in more recent research the task of road detection has been extended to challenging scenarios such as slippery roads and adverse weather [19].

### III. NETWORK ARCHITECTURES

#### A. Base FCN

The base neural network used in this work consists of a fully convolutional encoder-decoder that also contains an intermediate context module. This type of architecture has been successfully used in previous publications, such as [9] and [20], and it is illustrated in Fig. 1. The encoder consists of 5 convolutional layers:  $4 \times 4$  convolutions with stride 2 are used in order to downsample the input tensors thus reducing memory requirements. The context module consists of 9 convolutional layers with  $3 \times 3$  kernels and exponentially growing dilation [21]. This makes it possible to quickly grow the network's receptive field while limiting the number of layers. A large receptive field is beneficial for aggregating information within a large region of the input. More details about the context module are provided in Table I. The decoder contains 6 convolutional layers and its purpose is to further process and upsample the input tensors. Upsampling is achieved by using 3 strided convolutional layers with  $4 \times 4$  kernels and stride 2. Each convolutional layer is followed by an *exponential linear unit* (ELU) layer [22] which implements the following function:

$$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ e^x - 1 & \text{otherwise} \end{cases}$$

For regularization, spatial dropout layers, with dropout probability  $p = 0.25$ , have been added after each convolutional layer within the context module. This means that, during training, each feature map of a given convolutional layer has a probability  $p$  of being set to zero. This technique was shown to be more effective [23] than the original dropout implementation [24] for improving generalization performance.

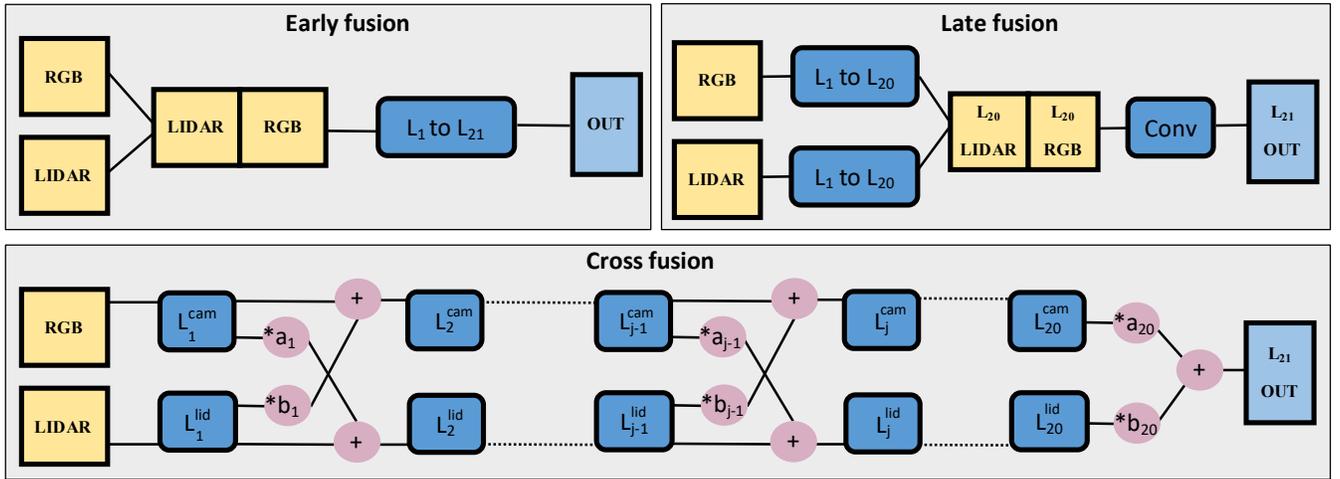


Fig. 2. Fusion strategies considered in this work. 1) Early fusion. In this case, the input camera and LIDAR images are concatenated in the depth dimension thus producing a tensor of size  $6 \times H \times W$ . This input tensor is then processed using the base FCN described in Sect. III-A. 2) Late fusion. Two parallel streams process the LIDAR and RGB images independently until layer 20. The outputs of  $L_{20}$  are then concatenated in the depth dimension and finally fed through a convolutional layer that carries out high-level information fusion. 3) Cross fusion. Also in this case there are two processing branches that, however, are connected by trainable scalar cross connections,  $a_j$  and  $b_j$  with  $j \in \{1, \dots, 20\}$ . The inputs of each layer, at a given depth, are calculated according to the illustrated computational operations.

TABLE I

CONTEXT MODULE ARCHITECTURE. THE CONTEXT MODULE CONSISTS OF 9 CONVOLUTIONAL LAYERS WITH EXPONENTIALLY GROWING DILATION FACTOR. EACH CONVOLUTIONAL LAYER IS FOLLOWED BY A SPATIAL DROPOUT LAYER WITH  $p = 0.25$ . ZERO-PADDING IS APPLIED THROUGHOUT THE CONTEXT MODULE IN ORDER TO PRESERVE THE WIDTH AND HEIGHT OF THE FEATURE MAPS.

Layer	6	7	8	9	10	11	12	13	14
Dilation H	1	1	1	2	4	8	16	1	-
Receptive field H	3	5	7	11	19	35	67	69	69
Dilation W	1	1	2	4	8	16	32	1	-
Receptive field W	3	5	9	17	33	65	129	131	131
# Feature maps	128	128	128	128	128	128	128	128	128
Filter size	3	3	3	3	3	3	3	3	1

### B. Early and late fusion

This work addresses the task of integrating information acquired with two sensors, an RGB camera and a rotating LIDAR. As will be explained in detail in Sect. IV, the LIDAR point clouds are transformed into a set of 2D images (in the following denoted as ZYX) that have the same spatial size as the camera images. Given this setup, the integration of camera and LIDAR data can be carried out in a straightforward manner using well-known CNN fusion strategies such as *early* and *late fusion* (see, for example, [13] and [14]).

In the early fusion approach, the input LIDAR and camera tensors are simply concatenated in the depth dimension thus producing a tensor with 6 channels (RGBZYX). This tensor then becomes the input for the base FCN described in Sect. III-A which has to learn, from the very beginning, features that combine both sensing modalities; in this case, fusion happens at a very low abstraction level. A graphical illustration of this strategy is presented in Panel 1 of Fig. 2.

At the other side of the spectrum is the late fusion. Here, the integration of LIDAR and camera information is carried out at the very end of two independent processing branches, as illustrated in Panel 2 of Fig. 2. In this case, fusion happens at the decision level.

A drawback of those approaches is that the developer has to manually decide at which stage the fusion should be done. Here, instead, a novel fusion strategy (*cross fusion*) has been introduced such that the FCN can learn from the data itself, during the training process, where fusion is necessary and to what extent.

### C. Cross fusion

The approach proposed in this work is represented in Panel 3 of Fig. 2. The rationale behind this strategy is to allow the FCN to integrate information at any processing depth instead of limiting it to a single level, as was the case in the previously mentioned methods. For example, the input tensors at depth  $j$ , denoted as  $I_j^{\text{Cam}}$  and  $I_j^{\text{Lid}}$ , that are fed to layers  $L_j^{\text{Cam}}$  and  $L_j^{\text{Lid}}$ , respectively, are given by the following expressions:

$$I_j^{\text{Lid}} = L_{j-1}^{\text{Lid}} + a_{j-1} L_{j-1}^{\text{Cam}} \quad (1)$$

$$I_j^{\text{Cam}} = L_{j-1}^{\text{Cam}} + b_{j-1} L_{j-1}^{\text{Lid}} \quad (2)$$

where  $a_j, b_j \in \mathbf{R}$  with  $j \in \{1, \dots, 20\}$  are trainable *cross fusion parameters*. The cross fusion parameters are initialized to zero which corresponds to the case of no information flow between the two processing branches. Afterwards, during training, these parameters are adjusted automatically in order to integrate the two information modalities.

#### IV. DATA PREPROCESSING

In this work, each LIDAR point cloud is converted to a set of three 2D images that make it straightforward to establish correspondences between color intensities and 3D information. Structuring a 3D point cloud in this manner is also convenient for the purpose of using the CNN machinery originally developed for processing camera images.

A point cloud acquired with a Velodyne HDL64 consists of approximately 100k points where each point  $p$  is specified by its spatial coordinates in the LIDAR coordinate system, that is  $p = [x, y, z, 1]^T$ . Given the LIDAR-camera transformation matrix  $\mathbf{T}$ , the rectification matrix  $\mathbf{R}$ , and the camera projection matrix  $\mathbf{P}$ , one can calculate the column position,  $u$ , and the row position,  $v$ , where the projection of  $p$  intersects the camera plane:

$$\lambda [u, v, 1]^T = \mathbf{P R T} p \quad (3)$$

where  $\lambda$  is a scaling factor that is determined by solving System (3). The above transformation is applied to every point in the point cloud, while discarding points such that  $\lambda < 0$  or when  $[u, v]$  falls outside the image.

Whereas an RGB image contains information about the red, green, and blue intensities of each pixel, the above procedure generates three images, X, Y, and Z where each pixel contains the  $x$ ,  $y$ , and  $z$  coordinates of the 3D point that is projected into it. An important difference between camera and LIDAR is that RGB images have valid values for each pixel, whereas, in the LIDAR images, many pixels are set to a default zero value because none of the laser beams hit the corresponding regions. For this reason, it is common practice [14], [25], [26] to upsample the LIDAR images before processing them with machine learning algorithms. This work makes use of the approach introduced by Premebida *et al.* [25] to accomplish that. Figure 3 shows an example of dense LIDAR images obtained by applying this procedure.

#### V. EXPERIMENTS AND DISCUSSION

##### A. Data set

In this work, five different FCNs were considered: ZYX, RGB, Early fusion, Late fusion, and Cross fusion. ZYX denotes the base FCN (see Sect. III-A) trained only on LIDAR images. Similarly, RGB is the base FCN trained only on camera images. Early, Late, and Cross fusion are the FCNs implementing the homonymous fusion strategy (see Sect. III-B). Each FCN was trained using exclusively the KITTI road data set which consists of 289 training images and 290 test images taken over several days in city, rural, and highway settings. It is important to mention that most of the training examples were captured in rather ideal weather and lighting conditions, something that might obscure the benefits of combining camera images with additional sensing modalities. For this reason, as will be described in Sect. V-C, an additional data set of more challenging scenes was included for performance evaluation. Table II provides further information regarding the data set splits. Given that the RGB images had different sizes due to the rectification

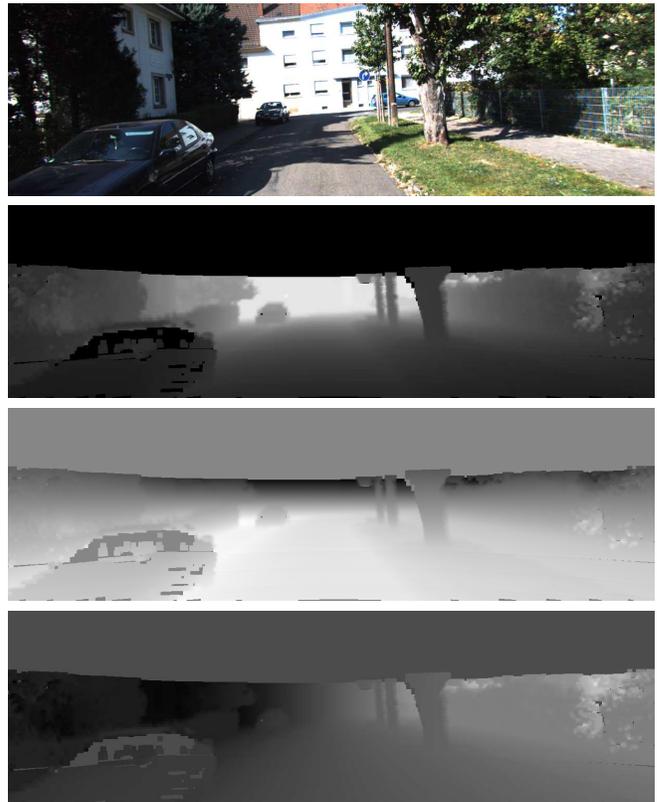


Fig. 3. Dense LIDAR images obtained by projecting the point cloud onto the camera plane and then applying the upsampling procedure described in [25]. From top to bottom: RGB image, Z channel, Y channel, and X channel. The gray-scale intensities are proportional to the numerical values of the corresponding quantities.

procedure, zero padding was applied to ensure that each training example had the same size of  $384 \times 1248$  pixels.

TABLE II

KITTI ROAD DATA SET: SIZE AND NUMBER OF IMAGES FOR EACH CATEGORY AND SPLIT. THE CHALLENGING SET WAS CREATED BY THE AUTHORS AND IT IS NOT PART OF THE STANDARD KITTI ROAD DATA SET (SEE SECT. V-D).

Category	Train	Validation	Test
urban marked	78	17	96
urban multiple marked	80	16	94
urban unmarked	81	17	100
challenging	-	33	-

##### B. Training procedure

Training was carried out for  $N = 100k$  iterations using the Adam optimization algorithm [27]. The performance of the FCN on the validation set was computed every 1000 iterations and the network's weights were saved if any improvement occurred. The learning rate  $\eta$  was decayed using the poly learning policy [28] implemented as:

$$\eta(i) = \eta_0 \left(1 - \frac{i}{N}\right)^\alpha, \quad (4)$$

where  $i$  denotes the current iteration number,  $\eta_0$  is the starting learning rate which was set to 0.0005, and  $\alpha = 0.9$ . The batch size was set to 1. Given the small size of the data set, data augmentation was also carried out by applying random rotations in the range  $[-20^\circ, 20^\circ]$  about the center of the images. The FCNs were implemented in PyTorch and trained using an Nvidia GTX1080 GPUs. The evaluation measures used in the following comparisons are the pixel-wise maximum F-measure (MaxF), precision (PRE), recall (REC), and average precision (AP) [29].

TABLE III  
PERFORMANCE COMPARISON OF SINGLE MODALITY AND FUSION FCNS EVALUATED ON THE VALIDATION SET.

Fusion strategy	# param.	MaxF [%]	PRE [%]	REC [%]
ZYX	1623395	94.96	94.05	95.89
RGB	1623395	96.00	96.16	95.84
Early fusion	1624931	95.41	94.62	96.21
Late fusion	3246787	96.06	95.97	96.15
Cross fusion	3246830	<b>96.25</b>	<b>96.17</b>	<b>96.34</b>

### C. Comparison of fusion strategies

The first experiment involved a performance comparison of single modality and fusion FCNs. Table III reports the results obtained on the validation set. As can be seen, the overall best performance was achieved by the cross fusion network with a MaxF score of 96.25%. This is followed by the late fusion FCN that obtained a MaxF score of 96.06% and then the single modality RGB-FCN at 96.00%. The worst performance was obtained by the FCN that only had access to the LIDAR images resulting in a MaxF score of 94.82%. This suggests that in scenarios presenting good lighting conditions, camera images are more informative than LIDAR point clouds for the task of road detection.

### D. Challenging scenarios

As was mentioned in Sect. V-A, the KITTI road data set mostly consists of examples captured in rather ideal lighting and weather conditions. In those situations, camera images are already, by themselves, quite informative and provide rich discriminative clues for carrying our accurate road detection. This is in part confirmed by noticing that most state-of-the-art algorithms in the KITTI road benchmark are purely camera-based. For this reason, by limiting the evaluation exclusively to the KITTI road data set, it might be difficult to fully appreciate the benefits provided by combining RGB cameras with other sensors, such as LIDARs.

TABLE IV  
PERFORMANCE COMPARISON OF SINGLE MODALITY AND FUSION FCNS EVALUATED ON THE CHALLENGING SET.

Fusion strategy	# params	MaxF [%]	PRE [%]	REC [%]
ZYX	1623395	95.21	93.40	97.09
RGB	1623395	91.81	89.18	94.61
Early fusion	1624931	95.44	93.54	97.42
Late fusion	3246787	95.24	92.73	97.09
Cross fusion	3246830	<b>96.02</b>	<b>94.39</b>	<b>97.70</b>

With this consideration in mind, an additional set<sup>1</sup> consisting of 33 images was extracted from the driving sequences of the KITTI raw data set [30] by looking for scenes that appeared particularly challenging for road segmentation using only the camera sensor, specifically images that presented shadows, strong light reflections, or peculiar lighting conditions affecting the appearance of the road surface. Four such examples and their ground truth annotations are shown in the top two rows of Fig. 4. The networks trained in Sect. V-C were also evaluated on this challenging set and their performance is reported in Table IV. As can be noticed, also in this case the cross fusion FCN performed best by achieving a MaxF score of 96.02%, once more supporting the previous finding that this fusion strategy is more efficient at integrating multimodal information. The RGB-FCN, on the other hand, significantly underperformed, obtaining a MaxF score of 91.81%. The fusion FCNs and the single modality ZYX-FCN all achieved MaxF scores above 95%. These results support the intuitive assumption that combining camera images with the spatial information acquired by a LIDAR sensor is indeed beneficial for carrying out road detection in more challenging illumination conditions. Figure 4 shows some examples of road segmentations obtained with the RGB-FCN (third row) and the cross fusion FCN (fourth row) that qualitatively illustrate the above remark.

TABLE V  
KITTI ROAD BENCHMARK RESULTS (IN %) ON THE URBAN ROAD CATEGORY. ONLY RESULTS OF PUBLISHED METHODS ARE REPORTED.

Method	MaxF	AP	PRE	REC	Time (s)
<b>LidCamNet (our)</b>	<b>96.03</b>	<b>93.93</b>	<b>96.23</b>	95.83	0.15
RBNNet [8]	94.97	91.49	94.94	95.01	0.18
StixelNet II [31]	94.88	87.75	92.97	<b>96.87</b>	1.2
MultiNet [7]	94.88	93.71	94.84	94.91	0.17
LoDNN [9]	94.07	92.03	92.81	95.37	<b>0.018</b>
DEEP-DIG [32]	93.98	93.65	94.26	93.69	0.14
Up-Conv-Poly [28]	93.83	90.47	94.00	93.67	0.08

### E. KITTI road benchmark

The cross fusion FCN was also evaluated on the KITTI road benchmark test set. Its performance on the *urban road category* is reported in Table V together with the results obtained by other state-of-the-art approaches. At the time of submission, the proposed system was among the best methods in the benchmark. Some examples of road detections on the test set are shown in Fig. 5. The results on individual categories are reported in Table VI. Additional evaluation metrics and further examples of detections can be found at the KITTI road benchmark website<sup>2</sup>: The proposed system is called LidCamNet which stands for LIDAR-Camera network. Lastly, several videos of road segmentations on full driving sequences are available at this link <https://goo.gl/1oLcmz>.

<sup>1</sup>The challenging data set can be found at <https://goo.gl/Z5amjQ>

<sup>2</sup><https://goo.gl/QNveL1>

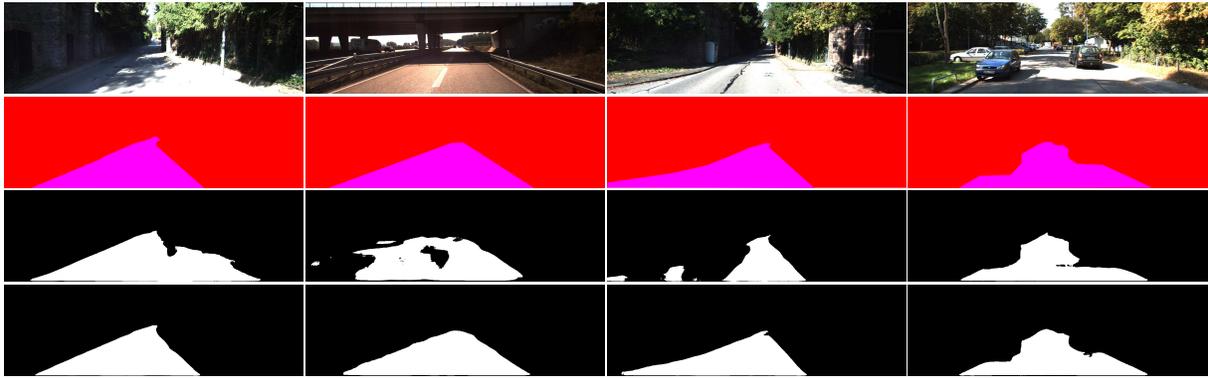


Fig. 4. (Top row) Some examples of camera images captured in difficult lighting conditions and included in the challenging set. (Second row) Corresponding ground truth annotations: The road is depicted as violet, whereas not-road is red. (Third row) Road segmentations obtained by the RGB-FCN. (Fourth row) Road segmentations generated by the cross fusion FCN.



Fig. 5. Examples of road segmentations in scenes from the KITTI test set. Correct road classifications are green. Red pixels correspond to false negatives, whereas blue pixels denote false positives.

TABLE VI  
KITTI ROAD BENCHMARK RESULTS (IN %) ON THE INDIVIDUAL CATEGORIES. FPR = FALSE POSITIVE RATE. FNR = FALSE NEGATIVE RATE.

Benchmark	MaxF	AP	PRE	REC	FPR	FNR
UM_ROAD	95.62	93.54	95.77	95.48	1.92	4.52
UMM_ROAD	97.08	95.51	97.28	96.88	2.98	3.12
UU_ROAD	94.54	92.74	94.64	94.45	1.74	5.55
URBAN_ROAD	96.03	93.93	96.23	95.83	2.07	4.17

## VI. CONCLUSION

In this paper, a novel fusion FCN has been developed to integrate camera images and LIDAR point clouds for carrying out road detection. Whereas other established fusion strategies found in the literature, such as early and late fusion, are designed to carry out information fusion at a single predefined processing depth, the proposed system incorporates trainable cross connections between the LIDAR and the camera processing branches, in all layers. These connections are initialized to zero, which corresponds to the case of no fusion, and are adjusted during training in order to find a suitable integration level.

The cross fusion FCN performed best among the single modality and fusion networks considered in this work. Its performance was also evaluated on the KITTI road benchmark where it achieved excellent results, with a MaxF score of 96.03% in the urban category, and it is currently among the top-performing algorithms.

An additional data set, consisting of visually challenging examples, was also considered to further highlight the benefits provided by using multiple sensors for carrying out road detection. It was shown that a camera-based FCN that performs quite well in good lighting conditions, will likely underperform in less forgiving situations, whereas a multimodal system that can leverage the information obtained with a different sensing mechanism can provide more robust and accurate segmentations in a wider spectrum of external conditions.

## ACKNOWLEDGMENT

The authors gratefully acknowledge financial support from Vinnova/FFI.

## REFERENCES

- [1] A. Evan, "Fatal tesla self-driving car crash reminds us that robots aren't perfect," in *IEEE Spectrum, Tech. Rep., Jul. 2016*. [Online]. Available: <https://spectrum.ieee.org/cars-that-think/transportation/self-driving/fataltesla-autopilot-crash-reminds-us-that-robots-arent-perfect>, 2016.
- [2] M. R. Endsley, "Autonomous driving systems: A preliminary naturalistic study of the tesla model s," *Journal of Cognitive Engineering and Decision Making*, vol. 11, no. 3, pp. 225–238, 2017.
- [3] J. M. Alvarez, A. Lopez, and R. Baldrich, "Illuminant-invariant model-based road segmentation," in *2008 IEEE Intelligent Vehicles Symposium*, June 2008, pp. 1175–1180.
- [4] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [5] A. Broggi, "Robust real-time lane and road detection in critical shadow conditions," in *Proceedings of International Symposium on Computer Vision - ISCV*, Nov 1995, pp. 353–358.
- [6] M. Bertozzi and A. Broggi, "Gold: A parallel real-time stereo vision system for generic obstacle and lane detection," *IEEE Transactions on Image Processing*, vol. 7, no. 1, pp. 62–81, Jan 1998.
- [7] M. Teichmann, M. Weber, M. Zoellner, R. Cipolla, and R. Urtaun, "Multinet: Real-time joint semantic reasoning for autonomous driving," *arXiv preprint arXiv:1612.07695*, 2016.
- [8] 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.
- [9] 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.
- [10] P. Y. Shinzato, D. F. Wolf, and C. Stiller, "Road terrain detection: Avoiding common obstacle detection assumptions using sensor fusion," in *2014 IEEE Intelligent Vehicles Symposium Proceedings*. IEEE, 2014, pp. 687–692.
- [11] L. Xiao, B. Dai, D. Liu, T. Hu, and T. Wu, "Crf based road detection with multi-sensor fusion," in *2015 IEEE Intelligent Vehicles Symposium (IV)*, June 2015, pp. 192–198.
- [12] 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.
- [13] A. Eitel, J. T. Springenberg, L. Spinello, M. Riedmiller, and W. Burgard, "Multimodal deep learning for robust rgb-d object recognition," in *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*. IEEE, 2015, pp. 681–687.
- [14] J. Schlosser, C. K. Chow, and Z. Kira, "Fusing lidar and images for pedestrian detection using convolutional neural networks," in *Robotics and Automation (ICRA), 2016 IEEE International Conference on*. IEEE, 2016, pp. 2198–2205.
- [15] 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.
- [16] A. Asvadi, L. Garrote, C. Premebida, P. Peixoto, and U. J. Nunes, "Multimodal vehicle detection: fusing 3d-lidar and color camera data," *Pattern Recognition Letters*, 2017.
- [17] M. Bellone, G. Reina, L. Caltagirone, and M. Wahde, "Learning traversability from point clouds in challenging scenarios," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 1, pp. 296–305, Jan 2018.
- [18] S. Zhou, J. Gong, G. Xiong, H. Chen, and K. Iagnemma, "Road detection using support vector machine based on online learning and evaluation," in *2010 IEEE Intelligent Vehicles Symposium*, June 2010, pp. 256–261.
- [19] J. Zhao, H. Wu, and L. Chen, "Road surface state recognition based on svm optimization and image segmentation processing," *Journal of Advanced Transportation*, vol. 2017, 2017.
- [20] L. Caltagirone, M. Bellone, L. Svensson, and M. Wahde, "Lidar-based driving path generation using fully convolutional neural networks," in *International Conference on Intelligent Transportation Systems (ITSC), 2017 IEEE*. IEEE, 2017, pp. 573–578.
- [21] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," in *ICLR*, 2016.
- [22] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *CoRR*, vol. abs/1511.07289, 2015.
- [23] J. Tompson, R. Goroshin, A. Jain, Y. LeCun, and C. Bregler, "Efficient object localization using convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 648–656.
- [24] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [25] C. Premebida, J. Carreira, J. Batista, and U. Nunes, "Pedestrian detection combining rgb and dense lidar data," in *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*. IEEE, 2014, pp. 4112–4117.
- [26] R. Fernandes, C. Premebida, P. Peixoto, D. Wolf, and U. Nunes, "Road detection using high resolution lidar," in *2014 IEEE Vehicle Power and Propulsion Conference (VPPC)*, Oct 2014, pp. 1–6.
- [27] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [28] G. L. Oliveira, W. Burgard, and T. Brox, "Efficient deep models for monocular road segmentation," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct 2016, pp. 4885–4891.
- [29] T. Fawcett, "An introduction to roc analysis," *Pattern Recognition Letters*, vol. 27, no. 8, pp. 861 – 874, 2006, rOC Analysis in Pattern Recognition. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S016786550500303X>
- [30] A. Geiger, P. Lenz, C. Stiller, and R. Urtaun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [31] 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.
- [32] 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," *Workshop on Deep Learning for Autonomous Driving on IEEE 20th International Conference on Intelligent Transportation Systems*, 2017.