

Lane Detection and Classification using Cascaded CNNs*

Fabio Pizzati¹, Marco Allodi², Alejandro Barrera³, and Fernando García³

¹ University of Bologna, Viale Risorgimento 2, 40136 Bologna BO, Italy
fabio.pizzati2@unibo.it

² University of Parma, Via delle Scienze, 181 / a, 43124 Parma PR, Italy
marco.allodi1@studenti.unipr.it

³ Universidad Carlos III de Madrid, Av. de la Universidad, 30, 28911 Leganés, Madrid, Spain
alebarre@pa.uc3m.es, fegarcia@ing.uc3m.es

Abstract. Lane detection is extremely important for autonomous vehicles. For this reason, many approaches use lane boundary information to locate the vehicle inside the street, or to integrate GPS-based localization. As many other computer vision based tasks, convolutional neural networks (CNNs) represent the state-of-the-art technology to identify lane boundaries. However, the position of the lane boundaries w.r.t. the vehicle may not suffice for a reliable positioning, as for path planning or localization information regarding lane types may also be needed. In this work, we present an end-to-end system for lane boundary identification, clustering and classification, based on two cascaded neural networks, that runs in real-time. To build the system, 14336 lane boundaries instances of the TuSimple dataset for lane detection have been labelled using 8 different classes. Our dataset and the code for inference are available online.

Keywords: Lane boundary detection · Lane boundary classification · Deep learning

1 Introduction

In autonomous driving, a deep understanding of the surrounding environment is vital to safely drive the vehicle. For this reason, a precise interpretation of the visual signals is necessary to identify all the components essential for navigation. One of them of particular importance is lane boundary position, which is needed to avoid collisions with other vehicles, and to localize the vehicle inside the street. Besides easing localization, lane detection is employed in many ADAS for lane departure warning and lane keeping assist. As many others computer vision based tasks, lanes boundaries detection accuracy has been significantly improved after the introduction of deep learning. The majority of recent systems, indeed, use convolutional neural networks to process sensorial data and

* The GPU used has been donated by the NVIDIA corporation.

infer high-level information relative to lanes. Some of them process LiDAR data to exploit differences in lane markings reflectivity [1,2]. However, LiDARs are extremely expensive, thus not always available on a vehicle. On the other hand, cameras are cheaper, and they make it possible to exploit chromatic differences on the road surface. Among the deep learning based lane detection approaches that rely solely on visual data, there is significant interest in lane marking detection [3,4,5,6,7]. In [4], a modified version of Faster R-CNN [8] is used to identify road patches that belong to lane markings. Many of those patches are then joined to obtain a complete representation of the marking. However, the system runs at approximately 4 frames per seconds, so it is not suitable for real-time elaboration, that is often a hard requirement for high-speed driving. In other works [6,7,9] lane markings detection and classification is achieved using fully-convolutional networks [10], and this enables more complex path planning tasks, where lane changes could be considered. Nonetheless, a comprehensive understanding of lane boundaries may be needed for path planning, so it may be necessary to join the detected markings with post processing algorithms. An alternative approach is to directly identify the boundaries, in order to reduce post-processing times. In [11,12,13,14], fully-convolutional networks are used to obtain a pixelwise representation of lane boundaries. A slightly different approach is proposed in [15], where a CNN is used to estimate polylines points, in order to solve fragmentation issues that often occur in segmentation networks. They classify the obtained boundaries, but only in terms of position w.r.t. the ego vehicle, so no information regarding the lane boundary type (e.g. dashed, continuous) is extracted. Similarly to [15], Ghafoorian et al. [16] exploit adversarial training to reduce fragmentation. In [17], an end-to-end approach is proposed, where lane boundary parameters are directly estimated by a CNN. As lane boundaries position, also lane boundaries types could be exploited to achieve a high grade of scene understanding. For example, knowing if a lane is dashed is indispensable for a lane change. Nonetheless, there is little interest in literature for simultaneous lane boundary identification and classification using deep learning. This could be caused by the lack of datasets that contains both information. For this reason, we extended a lane detection dataset with lane class annotations. Then, we developed a novel approach, based on the concatenation of multiple neural networks, that is used to perform lane boundary instance segmentation and classification, in an end-to-end deep learning based fashion. Our system satisfy real-time constraints on a NVIDIA Titan Xp GPU. Code for inference and pretrained models are available online.

2 Method

Our method is composed by two main sections, as presented in figure 1. As a first step, we train a CNN for lane boundary instance segmentation. Then, we extract a descriptor for each detected lane boundary and process it with a second CNN.

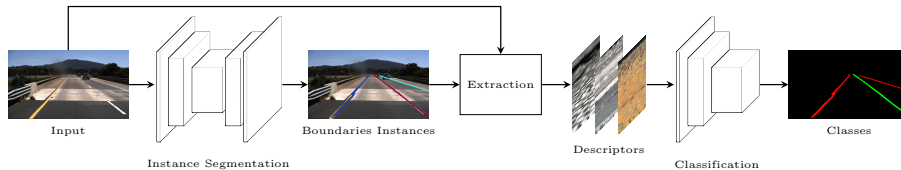


Fig. 1: System overview

2.1 Instance Segmentation

As discussed in section 1, several state-of-the-art approaches employ pixelwise classifications in order to differentiate pixels belonging to lane boundaries and background. In our case, different approaches are possible, so several design guidelines have been defined. First of all, we train the CNN to recognize lane boundaries, rather than lane markings. Doing this, it is indeed avoidable to group different lane markings in a lane boundary, considerably saving processing times. For similar reasons, we perform instance segmentation on lane boundaries instead of semantic segmentation. In this way, it is possible to distinguish different lane boundaries without relying on clustering algorithms. The state-of-the-art network for instance segmentation is Mask R-CNN [18]. However, two-step networks like Mask R-CNN are typically not amenable to use in real-time application, as they are typically slower than single-step ones. Furthermore, they are usually employed to detect objects easily enclosable in rectangular boxes. Being lane boundaries appearance heavily influenced by the perspective effects, other kinds of architectures should be preferred. Taking into account the previous assumptions, a fully-convolutional network has been trained.

We choose ERFNet [19] as our baseline model, as it is the model that, at the time of writing, has the best performances among real-time networks in the Cityscapes Dataset benchmark for semantic segmentation. As in all deep-learning based approaches, large amounts of images and annotations are crucial to achieve correct predictions and to avoid overfitting. For this reason, the TuSimple dataset for lane detection has been used. It is composed by 6408 1280×720 images, divided in 3626 for training, and 2782 for testing. 410 images extracted from the training set have been used as validation set during training. The main peculiarity in the TuSimple dataset is that entire lane boundaries are annotated, rather than lane markings. This makes the TuSimple dataset ideal for our needs. The lane boundaries are represented as polylines. In order to avoid clustering via post-processing, it is possible to make use of the loss function presented in [20], that is based on the Kullback-Leibler divergence minimization between the output probability distributions associated to pixels belonging to the same lane boundary instance. Please note that we do not address an unlimited number of possible instances for lane boundaries, as we decided to detect only the ego-lane boundaries, and the ones of the lanes on the sides of the ego-lane. Considering that two boundaries are shared for different lanes, we set a fixed maximum number of detected boundaries to 4. However, directly training

the network with [20] ultimately leads to gradient explosion and loss divergence. For this reason, the curriculum learning strategy [21] has been used to achieve convergence. In fact, in a first step a binary cross entropy loss has been used to train the network to distinguish between points belonging to a generic lane boundary and background. The resulting model is fine-tuned using [20] as loss function. The network has been trained using the images in the dataset resized at 512×256 resolution, for 150 epochs. This resolution led to satisfying results, while keeping the computational cost low. In order to represent the ground truth data as images, the polylines in the dataset have been projected with a fixed width of 5px on semantic maps of size 512×256 . We use the Adam optimizer, with learning rate $5 \cdot 10^{-4}$, and polynomial learning rate decay with exponent set to 0.9.

2.2 Classification

To the best of our knowledge, there are currently no publicly available datasets where entire lane boundaries with class-related information are annotated. For this reason, all the lanes in the TuSimple dataset have been manually classified using 8 different classes: *single white continuous*, *double white continuous*, *single yellow continuous*, *double yellow continuous*, *dashed*, *double-dashed*, *Botts' dots* and *unknown*. The obtained annotations are available online at <https://github.com/fabvio/TuSimple-lane-classes>.

Associating a class to each lane boundary detected could be addressed in several different ways. One possibility that has been considered in an early stage of the development was branching the instance segmentation network, to perform a pixel-wise classification of lane boundaries with dedicated convolutional layers, then fuse the outputs of the two branches. This approach has been discharged as it is memory intensive, because it requires two decoders in the same network, and it may generate inconsistencies between the detection of the two branches, that should be solved using post-processing algorithms. For example, there could be pixels classified as background from the instance segmentation branch, but classified as lane from the other branch. For this reason, we perform a classification for each lane boundary with another CNN, associating the detected boundaries to the ground truth. A problem with this approach is that each lane boundary is constituted by a different number of points in the input image. For this reason, it is difficult to extract a representation of them that is position-independent w.r.t. the ego vehicle. This may be essential to achieve a correct classification. Thus, we extract a descriptor for each boundary, sampling a fixed number of points from the input image which belong to the detected lane boundary. The points extracted in this way are then ordered following their index in the original image, and arranged in squared images, that are processed by the second neural network. In this way, a spatially normalized compact representation of lane boundaries is obtained, while preserving information given by visual clues such as lane markings. Furthermore, using this approach we are able to perform lane boundary instance segmentation and classification with only two inferences, in an end-to-end fashion. In fact, the descriptors of different lane

boundaries detected could be grouped in batches and classified simultaneously. Examples of descriptors are shown in image 2.

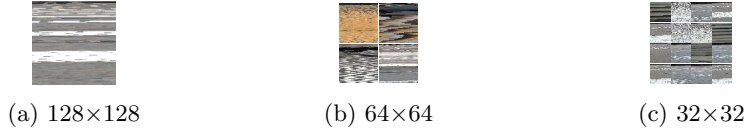


Fig. 2: Descriptors of different sizes

The architecture we use for this task is derived from H-Net [11]. A detailed description of its structure is given in image 3. We trained this network separately from the first one. To do that, the TuSimple dataset has been processed by the instance segmentation network. Each detected lane boundary is then compared with the ground truth, and it is associated to the corresponding class if the average distance between the detected points and the ground truth is under a threshold. This is needed to filter false positives generated by the first network. In fact, only lane boundaries that are effectively in the training set have a ground truth class, while others detected by the CNN should be excluded. As a result, we obtain a set of $\{descriptor, class\}$ objects that can be used to train the classification neural network. This has been trained with the same hyperparameters of the instance segmentation network. Examples of extracted descriptors are shown in image 2. Code for inference, descriptor extraction and pretrained models are publicly available at <https://github.com/fabvio/Cascade-LD>.

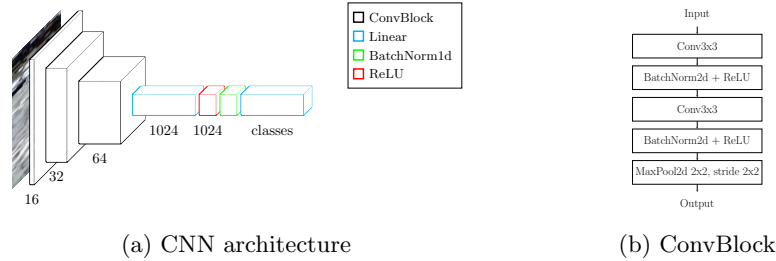


Fig. 3: Classification Network. Output channels are listed below each layer.

3 Results

In order to validate our method, we evaluate the performances of both networks separately. For lane boundary instance segmentation, the evaluation formula for the TuSimple benchmark is presented in equation 1. In it, C_i and S_i are

the number of correctly detected points and ground truth points in image i , respectively. A point is defined as correctly detected if it has a distance w.r.t. a ground truth point under 20 pixels. Additionally, false positive and false negative lane boundaries are evaluated. Given that our detected lane boundaries have width over 1 pixel, we average the x coordinates of the detected pixels for a given row, to obtain a single value. In 2, F_{pred} refers to the number of erroneously detected lanes, while N_{pred} is the total number of detected lanes. In 3, M_{pred} is the total number of unidentified lanes, and N_{gt} is the total number of lane boundaries annotated.

$$accuracy = \frac{\sum_i C_i}{\sum_i S_i} \quad (1) \quad FP = \frac{F_{pred}}{N_{pred}} \quad (2) \quad FN = \frac{M_{pred}}{N_{gt}} \quad (3)$$

We compare our instance segmentation network with the top-three approaches in the TuSimple benchmark for lane detection. We do not evaluate lanes that are composed than less of three points, in order to filter false positives. Results are presented in table 1. Inference times are evaluated on 512×256 images. Our network is slightly less accurate than the others. However, taking into account the computational times reduction, we found this tradeoff acceptable.

Table 1: TuSimple Lane detection metrics results and comparison.

Method	Accuracy	FP	FN	FPS
Xingang Pan [12]	96.53	6.17	1.80	5.31
Yen-Chang Hsu [20]	96.50	8.51	2.69	55.55
Davy Neven [11]	96.40	23.65	2.76	52.63
Ours	95.24	11.97	6.20	58.93

For classification, two different experiments are performed. In a first phase, we train the network to distinguish between two different classes: *dashed* and *continuous*. To do that, the *single white continuous*, *double white continuous*, *single yellow continuous*, *double yellow continuous* classes are mapped to the *continuous* class. On the other hand, *dashed*, *Botts' dots* and *double-dashed* are equally labelled as *dashed*. *Unknown* descriptors are ignored. In this way, it is possible to distinguish between lane boundaries that may or may not be crossed. In the second experiment, we treat the *double-dashed* class as independent. Doing this, we could identify also the boundaries of lanes that may be crossed only in specific conditions, as highway entry or exit. An ablation study regarding the descriptor size has been performed. We evaluate classification performances on the validation set, as the test set labels for lane boundaries are not publicly available. Results are reported in table 2. Inference times for the classification network are around 1ms.

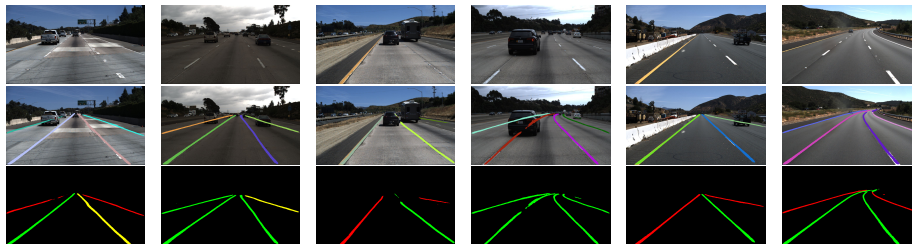


Fig. 4: Qualitative results on the test set. From top to bottom: original image, instance segmentation, classification. For instance segmentation, different colors represent different boundaries. For classification, green represents dashed lanes, yellow double-dashed, red continuous.

Table 2: Ablation study on descriptor size.

Descriptor size	Acc. (two classes)	Acc. (three classes)
256×256	0.9698	0.9600
128×128	0.9596	0.9600
64×64	0.9519	0.9443
32×32	0.9527	0.9436
16×16	0.9359	0.9203

As it is visible, it is possible to achieve better performances increasing the descriptor spatial resolution. However, this leads to a major occupation of GPU RAM. On the other hand, our results demonstrate that it is possible to achieve satisfying accuracies with only 256 points.

4 Conclusions

In this work, we presented a novel approach to lane boundary identification and classification in a end-to-end deep learning fashion. With our method, it is possible to achieve high accuracy in both tasks, in real-time. We formalized a descriptor extraction strategy that is useful when it is needed to combine instance segmentation and classification without relying on two-step detection networks. Furthermore, we performed an ablation study on the descriptor size, in order to define the tradeoff between detection accuracy and needed GPU RAM.

References

1. Caltagirone, L., Scheidegger, S., Svensson, L., Wahde, M.: Fast lidar-based road detection using fully convolutional neural networks. In: 2017 IEEE Intelligent Vehicles Symposium (IV). (2017)

2. Bai, M., Mattyus, G., Homayounfar, N., Wang, S., Lakshmikanth, S.K., Urtasun, R.: Deep multi-sensor lane detection. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). (2018)
3. Chen, P., Lo, S., Hang, H., Chan, S., Lin, J.: Efficient road lane marking detection with deep learning. CoRR [abs/1809.03994](https://arxiv.org/abs/1809.03994) (2018)
4. Tian, Y., Gelernter, J., Wang, X., Chen, W., Gao, J., Zhang, Y., Li, X.: Lane marking detection via deep convolutional neural network. *Neurocomputing* (2018)
5. Li, J., Mei, X., Prokhorov, D., Tao, D.: Deep neural network for structural prediction and lane detection in traffic scene. *IEEE transactions on neural networks and learning systems* **28**(3) (2017)
6. Lee, S., Kim, J.S., Yoon, J.S., Shin, S., Bailo, O., Kim, N., Lee, T.H., Hong, H.S., Han, S.H., Kweon, I.S.: Vpgnet: Vanishing point guided network for lane and road marking detection and recognition. *ICCV* (2017)
7. Zang, J., Zhou, W., Zhang, G., Duan, Z.: Traffic lane detection using fully convolutional neural network. In: APSIPA ASC, IEEE (2018)
8. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. In: *Advances in NIPS*. (2015)
9. John, V., Liu, Z., Mita, S., Guo, C., Kidono, K.: Real-time road surface and semantic lane estimation using deep features. *Signal, Image and Video Processing* **12**(6) (Sep 2018) 1133–1140
10. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: *CVPR*. (2015) 3431–3440
11. Neven, D., De Brabandere, B., Georgoulis, S., Proesmans, M., Van Gool, L.: Towards end-to-end lane detection: an instance segmentation approach. In: 2018 IEEE Intelligent Vehicles Symposium (IV), IEEE (2018) 286–291
12. Pan, X., Shi, J., Luo, P., Wang, X., Tang, X.: Spatial as deep: Spatial cnn for traffic scene understanding. In: 32nd AAAI Conference on Artificial Intelligence. (2018)
13. Zhang, J., Xu, Y., Ni, B., Duan, Z.: Geometric constrained joint lane segmentation and lane boundary detection. In: *ECCV*. (2018) 486–502
14. Kim, J., Park, C.: End-to-end ego lane estimation based on sequential transfer learning for self-driving cars. In: *Proceedings of the IEEE CVPR Workshops*. (2017)
15. Chougule, S., Koznek, N., Ismail, A., Adam, G., Narayan, V., Schulze, M. In: *Reliable Multilane Detection and Classification by Utilizing CNN as a Regression Network: Munich, Germany, September 8-14, 2018, Proceedings, Part V*
16. Ghafoorian, M., Nugteren, C., Baka, N., Booij, O., Hofmann, M.: El-gan: embedding loss driven generative adversarial networks for lane detection. In: *ECCV*. (2018)
17. De Brabandere, B., Van Gansbeke, W., Neven, D., Proesmans, M., Van Gool, L.: End-to-end lane detection through differentiable least-squares fitting. *arXiv preprint arXiv:1902.00293* (2019)
18. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: *Proceedings of the IEEE international conference on computer vision*. (2017)
19. Romera, E., Alvarez, J.M., Bergasa, L.M., Arroyo, R.: Erfnet: Efficient residual factorized convnet for real-time semantic segmentation. *IEEE Transactions on ITS*
20. Hsu, Y.C., Xu, Z., Kira, Z., Huang, J.: Learning to cluster for proposal-free instance segmentation. (07 2018) 1–8
21. Bengio, Y., Louradour, J., Collobert, R., Weston, J.: Curriculum learning. In: *Proceedings of the 26th annual international conference on machine learning*. (2009)