

SwiftLane: Towards Fast and Efficient Lane Detection

Oshada Jayasinghe, Damith Anhettigama, Sahan Hemachandra, Shenali Kariyawasam,
Ranga Rodrigo and Peshala Jayasekara

Department of Electronic and Telecommunication Engineering,
University of Moratuwa, Sri Lanka

Email: oshadajayasinghe@gmail.com, damithkawshan@gmail.com, sahanhemachandra@gmail.com,
shenali1997@gmail.com, ranga@uom.lk, peshala@uom.lk

Abstract—Recent work done on lane detection has been able to detect lanes accurately in complex scenarios, yet many fail to deliver real-time performance specifically with limited computational resources. In this work, we propose SwiftLane: a simple and light-weight, end-to-end deep learning based framework, coupled with the row-wise classification formulation for fast and efficient lane detection. This framework is supplemented with a false positive suppression algorithm and a curve fitting technique to further increase the accuracy. Our method achieves an inference speed of 411 frames per second, surpassing state-of-the-art in terms of speed while achieving comparable results in terms of accuracy on the popular CULane benchmark dataset. In addition, our proposed framework together with TensorRT optimization facilitates real-time lane detection on a Nvidia Jetson AGX Xavier as an embedded system while achieving a high inference speed of 56 frames per second.

Index Terms—lane detection, deep learning, convolutional neural network, row-wise classification, embedded system

I. INTRODUCTION

Lane detection is a pivotal element in driver assistance systems and autonomous vehicles as lane marker information is essential in maneuvering the vehicle safely on roads. Detecting lanes in real-world scenarios is a challenging task due to adverse weather, lighting conditions and occlusions. As the computational budget available for lane detection in the aforementioned systems is limited, a light-weight, fast and accurate lane detection system is crucial.

Recent lane detection approaches fall into two broad classes: semantic segmentation based methods and row-wise classification based methods. While semantic segmentation based methods [1]–[3] provide competitive results in terms of accuracy, a common drawback is the reduced speed due to per-pixel classification and large backbones. On the other hand, row-wise classification based methods [4], [5] focus on improving speed and obtaining real-time performance. However, the inherent limitation of a grid-based representation in row-wise classification methods and the bias towards overfitting due to the similar structure of lanes in the training set may result in

reduced accuracy, highlighting the speed-accuracy trade-off in lane detection models.

In this work, we propose a simple, light-weight, end-to-end deep learning based lane detection framework with a smaller backbone and a lesser number of multiply-accumulate operations (MACs) following the row-wise classification approach. The inference speed is significantly increased by reducing the computational complexity, and the light-weight network architecture is less prone to overfitting. Moreover, we also introduce a false positive suppression algorithm based on the length of the lane segment and the Pearson correlation coefficient, and a second-order polynomial fitting method as post-processing techniques to improve the overall accuracy of the system. Comprehensive experimental results are shown on the CULane [1] benchmark dataset, accompanied by a comparison of our results with other state-of-the-art approaches. An ablation study shows how each of the proposed methods contributes to the speed and the accuracy.

Furthermore, we deploy our lane detection framework on a Nvidia Jetson AGX Xavier integrated with Robot Operating System (ROS) [6] to demonstrate the capability of our light-weight network architecture to perform real-time lane detection in an embedded system. The trained model is optimized and quantized using TensorRT for increasing the inference speed. We also provide qualitative results for locally captured street view images to showcase how well our model generalizes for the task of lane detection.

In summary, our contributions are as follows: we introduce a novel, light-weight, end-to-end deep learning architecture supplemented with two effective post-processing techniques for fast and efficient lane detection. Our proposed method drastically improves the inference speed, reaching 411 frames per second (FPS) to surpass state-of-the-art while achieving comparable accuracy. We further optimize the trained model using TensorRT and implement it on an embedded system in the ROS ecosystem. The overall system achieves an inference speed of 56 FPS, demonstrating the capability of our method to perform real-time lane detection.

II. RELATED WORK

Initially, lane detection research mainly focused on classical image processing algorithms, such as using basic hand-

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

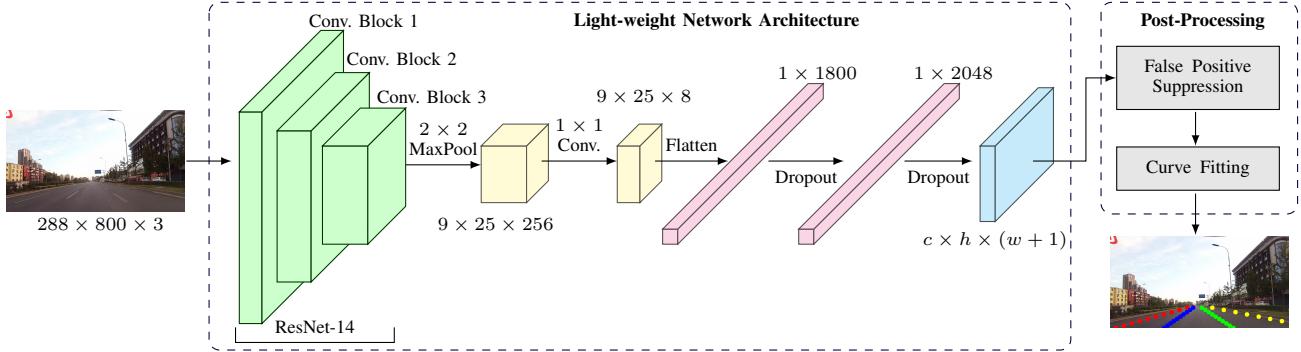


Fig. 1: Proposed model architecture. ResNet-14 backbone generates feature maps from the input image. A 2×2 max pooling layer and a 1×1 convolutional layer are used to reduce the spatial dimensions and the number of channels. Resulting feature maps are flattened and passed through two fully-connected layers with dropout layers in between. The model predictions are fed through false positive suppression and curve fitting modules to obtain the lane output.

crafted features [7]–[9], color-based approaches [10], [11], and traditional feature extraction methods with machine learning algorithms such as decision trees and support vector machines [12], [13]. Although these methods are computationally less expensive, the performance is poor in complex scenarios with occlusions, shadows and different lighting conditions.

Recent deep learning based approaches outperform classical methods and can be further divided into two broad classes: semantic segmentation based methods and row-wise classification based methods. In semantic segmentation based methods [1]–[3], classification is done on a per-pixel basis by classifying each pixel as lane or background. A special convolution method known as slice-by-slice convolution is proposed in SCNN [1], which enables information propagation within the same layer to improve the detection of long thin structures such as lanes. CurveLane-NAS [2] focuses on capturing long-range contextual information and short-range curved trajectory information using a lane-sensitive neural architecture search framework. Attention maps extracted from different layers of a trained model which contain important contextual information are used as distillation targets for the lower layers in SAD [3]. The pixel-wise computation in semantic segmentation based approaches increases the computational complexity and reduces the inference speed drastically.

Row-wise classification based methods [4], [5] have been able to progress towards real-time lane detection by addressing the computational complexity problem. In these approaches, the input image is divided into a grid and for each row, the model outputs the probability of each cell belonging to a lane. This approach is first introduced in E2E-LMD [4] by converting the output of the segmentation backbone to a row-wise representation using a special module called horizontal reduction module. The no-visual-clue problem in lane detection is addressed in UltraFast [5] using a low-cost, row-wise classification based network, which utilizes global and structural information. Although their approach achieves state-of-the-art speed of 322.5 FPS, the accuracy is low when compared with other methods.

Almost all of the above mentioned algorithms have been

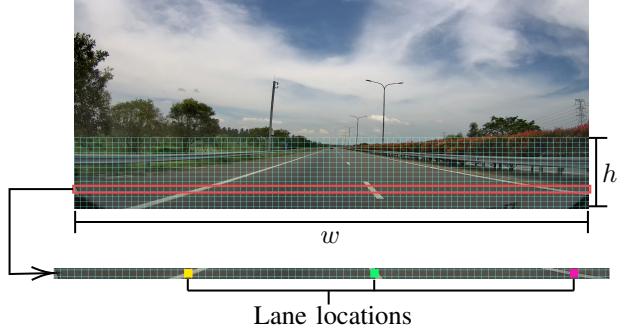


Fig. 2: Lane Representation. The region comprising lanes is divided into a pre-defined number of row anchors (h) and gridding cells (w).

implemented in high-end computational platforms and implementation of lane detectors in embedded systems is comparatively a less researched area. A lane detection algorithm optimized for PXA255 embedded device has been introduced by [14] which achieves a frame rate of 13 FPS. PathMark [15] is another lane detection algorithm running at 13 FPS in a TI-OMAP4430 based embedded system. A Nvidia Jetson-TK1 board has been used in [16] for implementing a real-time lane detection and departure warning system at 44 FPS. In [17], a lane detection and modeling pipeline has been presented for embedded platforms which delivers real-time performance in a Jetson-TX2 embedded device. All of these approaches rely on classical image processing based techniques and do not perform well in complex scenarios when compared with deep learning based approaches.

III. METHODOLOGY

In this section, we present the lane representation mechanism, a detailed explanation of our model architecture and the algorithms used to further increase the model accuracy.

A. Lane Representation

We address the lane detection task as a row-wise classification problem following the formulation introduced by [5].

The region of the image which contains lanes is divided into a pre-defined number of row anchors (h) and each row anchor is divided into a pre-defined number of gridding cells (w) as shown in Fig. 2. The number of lanes (c) is pre-defined, and for each lane, the lane locations are represented by a $h \times w$ grid. An additional cell is attached to the end of each row anchor to indicate the absence of a particular lane in that row anchor.

B. Model Architecture

We propose a simple end-to-end light-weight convolutional neural network based model architecture for the lane detection task as shown in Fig. 1. The first stage of the proposed model is the backbone which extracts features from the input image. As the backbone we use “ResNet-14” which is obtained by dropping the last four convolutional layers of ResNet-18 [18] to increase the speed by reducing the computational complexity.

The output of the backbone is a feature representation of the image which would then be fed into a 2×2 max pooling layer for dimensionality reduction in the spatial dimensions. For dimensionality reduction in the channel dimension a 1×1 convolution layer is applied. This output is flattened to obtain a one-dimensional tensor which is then passed through two fully connected layers to obtain the output tensor. Dropout layers are implemented in between to further prevent the network from overfitting.

The output tensor represents the score of each gridding cell (including the no lane cell) belonging to each lane in each row anchor. $S_{i,j,k}$ represents the score of k^{th} gridding cell in j^{th} row anchor belonging to i^{th} lane which can be obtained by,

$$S_{i,j,k} = f(X), \text{ s.t. } i \in [1, c], j \in [1, h], k \in [1, w + 1] \quad (1)$$

Here, f , X , c , h and w stands for the classification model, the input image, the number of lanes, the number of row anchors and the number of gridding cells, respectively. The lane points can then be extracted by choosing the gridding cell with the highest score in each row anchor for each lane. If the last gridding cell is not the cell with the highest score, the location of i^{th} lane in j^{th} row anchor is given by,

$$Loc_{i,j} = \operatorname{argmax}_k (S_{i,j,k}), \text{ s.t. } k \in [1, w] \quad (2)$$

Having the highest score in the last gridding cell implies that the considered lane is not present in the selected row anchor. For training the model, we define the classification loss as the negative log likelihood loss which is given by,

$$L_{cls} = \sum_{i=1}^c \sum_{j=1}^h -\alpha_{i,j,T_{i,j}} \cdot \log(P_{i,j,T_{i,j}}) \quad (3)$$

Here, $T_{i,j}$ denotes the correct location (gridding cell) of i^{th} lane in j^{th} row anchor as per the ground truth and $P_{i,j,k}$ denotes the probability of k^{th} gridding cell in j^{th} row anchor belonging to i^{th} lane which can be obtained by,

$$P_{i,j,k} = \operatorname{softmax}(S_{i,j,k}) \quad (4)$$

$\alpha_{i,j,k}$ is the modulating factor for the focal loss adjustment as mentioned in [19].

$$\alpha_{i,j,k} = (1 - P_{i,j,k})^\gamma \quad (5)$$

C. False Positive Suppression

We propose two post processing techniques to reduce false detections in the model output. First, we remove all instances of small lane segments which have a less number of detected lane points than a threshold value. Second, we remove all instances of lanes which have a considerable deviation from a straight line. Pearson correlation coefficient measures the linear correlation between two variables which is given by (6), where x_i and y_i are the sample data points of x and y variables and \bar{x} and \bar{y} are the respective means.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2(y_i - \bar{y})^2}} \quad (6)$$

In our case, Pearson correlation coefficient of row anchors and gridding cells of an identified lane segment is used to measure how well the lane points can be represented using a straight line. Since majority of the lanes have a slight deviation from a straight line, the Pearson correlation coefficient should be close to one in magnitude. Therefore, we remove all instances of lanes which have a Pearson correlation coefficient below a threshold value.

D. Curve Fitting

In most of the scenarios, lanes are straight lines or curve segments with small curvature values. Therefore, lanes can be approximated to a greater extent by second-order polynomials. Since we use a finite number of gridding cells, lanes in the model output are represented in the discrete domain. Second-order polynomial fitting can be used to replace these discrete gridding cell numbers by continuous values which results in smooth lane segments.

IV. EXPERIMENTS

In this section, we present the details about the dataset used to evaluate our model, the training process and a detailed description on the embedded system implementation for real-time applications.

A. Dataset Description

For the training and quantitative evaluation of our model, we use the publicly available CULane [1] benchmark dataset which is one of the largest lane detection datasets with 133,235 total frames having a resolution of 1640×590 . The dataset is divided into the train set, the validation set and the test set which comprises 88,880 frames, 9,675 frames and 34,680 frames, respectively. The dataset covers several complex scenarios and the test images are divided into 9 categories: Normal, Crowded, Dazzle light, Shadow, No line, Arrow, Curve, Crossroad and Night.

As the evaluation metric, F1-measure is used to compare the performance in the CULane benchmark. Each lane is represented by a 30-pixel-width line and each prediction which



Fig. 5: Visualization of lane detection result on locally captured images. The first six images show accurate detections while the last two show failure cases including false detections and undetected lanes.

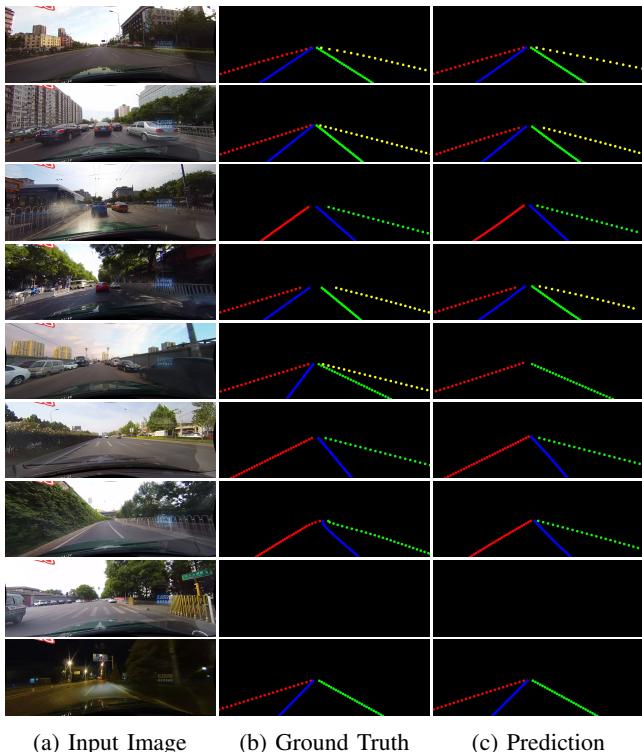


Fig. 6: Visualization of results on CULane. The nine rows represent the nine scenarios in CULane; Normal, Crowded, Dazzle light, Shadow, No line, Arrow, Curve, Crossroad and Night respectively.

faster inference speed, we use the FP16 quantized TensorRT engine for the lane detection task. The visualizer node marks the detected lane points in the current frame and publishes the resultant image to the *output_frame* topic. The RViz visualization tool is used to visualize the lane detections in real-time.

V. RESULTS

The performance of our method on the CULane benchmark dataset is compared against state-of-the-art lane detection

TABLE II: Performance on the embedded system

Model	F1-measure	Speed (FPS)
Pytorch Model	74.02	23
TensorRT Engine (FP32)	74.02	35
TensorRT Engine (FP16)	74.03	56

approaches in Table I. The number of false positives are displayed under the “Cross” category since there are no true positives in the ground truth for that category. The inference speed is measured by taking the average frames per second (FPS) value for 1000 runs including the forward pass of the model and the post-processing steps. The number of multiply-accumulate operations in billions is represented in the “GMACs” column. For a fairer comparison, we measured the speed of [5] under the same conditions as ours.

It can be observed that while being the fastest, our method achieves competitive results with other state-of-the-art methods in F1-measure. Our method also uses the least number of multiply-accumulate operations (MACs) which highlights the efficiency of our formulation. The low number of false positives in the “Cross” category validates the effectiveness of our false positive suppression technique. Compared to the segmentation based methods [1], [3], the inference speed improves substantially while providing better results at the same time. When compared with [5], which is the fastest among other approaches, our method achieves better results with a 6.6% increase in F1-measure. While we obtain comparable performance with [2] and [4], a direct comparison cannot be made in terms of the speed, as their inference speeds are not mentioned. Although [20], [21], [22] and [23] achieves on par or better results than our method, the low inference speeds of their best performing models act as a barrier for real-time implementation especially on resource constrained environments.

The performance of the Pytorch model and the generated FP32 and FP16 TensorRT engines on the Nvidia Jetson AGX Xavier are shown in Table II in terms of the F1-measure and speed. The inference speed is calculated as the average

- guided Lane Detection,” in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [23] Z. Qu, H. Jin, Y. Zhou, Z. Yang, and W. Zhang, “Focus on local: Detecting lane marker from bottom up via key point,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 14 122–14 130.
- [24] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, “On the importance of initialization and momentum in deep learning,” in *Proceedings of the 30th International Conference on Machine Learning*, 2013, pp. 1139–1147.
- [25] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Advances in Neural Information Processing Systems 32*, 2019, pp. 8024–8035.