# **Rethinking Efficient Lane Detection via Curve Modeling**

Zhengyang Feng<sup>1\*</sup>, Shaohua Guo<sup>1</sup><sup>\*</sup>, Xin Tan<sup>2,1</sup>, Ke Xu<sup>3</sup>, Min Wang<sup>4</sup>, Lizhuang Ma<sup>1,2,5†</sup> <sup>1</sup>Shanghai Jiao Tong University <sup>2</sup>East China Normal University <sup>3</sup>City University of Hong Kong <sup>4</sup>SenseTime Research <sup>5</sup>MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University

zyfeng97@sjtu.edu.cn;guoshaohua@sjtu.edu.cn;tanxin2017@sjtu.edu.cn;kkangwing@gmail.com; wangmin@sensetime.com;ma-lz@cs.sjtu.edu.cn

### Abstract

This paper presents a novel parametric curve-based method for lane detection in RGB images. Unlike stateof-the-art segmentation-based and point detection-based methods that typically require heuristics to either decode predictions or formulate a large sum of anchors, the curvebased methods can learn holistic lane representations naturally. To handle the optimization difficulties of existing polynomial curve methods, we propose to exploit the parametric Bézier curve due to its ease of computation, stability, and high freedom degrees of transformations. In addition, we propose the deformable convolution-based feature flip fusion, for exploiting the symmetry properties of lanes in driving scenes. The proposed method achieves a new state-ofthe-art performance on the popular LLAMAS benchmark. It also achieves favorable accuracy on the TuSimple and CULane datasets, while retaining both low latency (>150 FPS) and small model size (<10M). Our method can serve as a new baseline, to shed the light on the parametric curves modeling for lane detection. Codes of our model and PytorchAutoDrive: a unified framework for self-driving perception, are available at: https://github.com/ voldemortX/pytorch-auto-drive.

### 1. Introduction

Lane detection is a fundamental task in autonomous driving systems, which supports the decision-making of lanekeeping, centering and changing, *etc.* Previous lane detection methods [2, 12] typically rely on expensive sensors such as LIDAR. Advanced by the rapid development of deep learning techniques, many works [19, 21, 22, 33, 41] are proposed to detect lane lines from RGB inputs captured by commercial front-mounted cameras.



Figure 1. Lane detection strategies. Segmentation-based and point detection-based representations are local and indirect. The abstract coefficients (a, b, c, d) used in polynomial curve are hard to optimize. The cubic Bézier curve is defined by 4 actually existing control points, which roughly fit line shape and wrap the lane line in its convex hull (dashed red lines). Best viewed in color.

Deep lane detection methods can be classified into three categories, i.e., segmentation-based, point detection-based, and curve-based methods (Figure 1). Among them, by relying on classic segmentation [5] and object detection [28] networks, the segmentation-based and point detectionbased methods typically achieve state-of-the-art lane detection performance. The segmentation-based methods [21,22, 41] exploit the foreground texture cues to segment the lane pixels and decode these pixels into line instances via heuristics. The point detection-based methods [15, 33, 39] typically adopt the R-CNN framework [9, 28], and detect lane lines by detecting a dense series of points (e.g., every 10 pixels in the vertical axis). Both kinds of approaches represent lane lines via indirect proxies (i.e., segmentation maps and points). To handle the learning of holistic lane lines, under cases of occlusions or adverse weather/illumination conditions, they have to rely on low-efficiency designs, such as recurrent feature aggregation (too heavy for this realtime task) [22, 41], or a large number of heuristic anchors (> 1000, which may be biased to dataset statistics) [33].

On the other hand, there are only a few methods [19, 32] proposed to model the lane lines as holistic curves (typically the polynomial curves, *e.g.*,  $x = ay^3 + by^2 + cy + d$ ).

<sup>\*</sup>Equal Contribution.

<sup>&</sup>lt;sup>†</sup>Lizhuang Ma is a member of Qing Yuan Research Institute, Shanghai Jiao Tong University.

While we expect the holistic curve to be a concise and elegant way to model the geometric properties of lane line, the abstract polynomial coefficients are difficult to learn. Previous studies show that their performance lag behind the well-designed segmentation-based and point detectionbased methods by a large margin (up to 8% gap to state-ofthe-art methods on the CULane [22] dataset). *In this paper, we aim to answer the question of whether it is possible to build a state-of-the-art curve-based lane detector.* 

We observe that the classic cubic Bézier curves, with sufficient freedom degrees of parameterizing the deformations of lane lines in driving scenes, remain low computation complexity and high stability. This inspires us to propose to model the thin and long geometric shape properties of lane lines via Bézier curves. The ease of optimization from on-image Bézier control points enables the network to be end-to-end learnable with the bipartite matching loss [38], using a sparse set of lane proposals from simple column-wise Pooling (e.g., 50 proposals on the CU-Lane dataset [22]), without any post-processing steps such as the Non-Maximum Suppression (NMS), or hand-crafted heuristics such as anchors, hence leads to high speed and small model size. In addition, we observe that lane lines appear symmetrically from a front-mounted camera (e.g., between ego lane lines, or immediate left and right lanes). To model this global structure of driving scenes, we further propose the feature flip fusion, to aggregate the feature map with its horizontally flipped version, to strengthen such coexistences. We base our design of feature flip fusion on the deformable convolution [42], for aligning the imperfect symmetries caused by, e.g., rotated camera, changing lane, non-paired lines. We conduct extensive experiments to analyze the properties of our method and show that it performs favorably against state-of-the-art lane detectors on three popular benchmark datasets. Our main contributions are summarized as follows:

- We propose a novel Bézier curve-based deep lane detector, which can model the geometric shapes of lane lines effectively, and be naturally robust to adverse driving conditions.
- We propose a novel deformable convolution-based feature flip fusion module, to exploit the symmetry property of lanes observed from front-view cameras.
- We show that our method is fast, light-weight, and accurate through extensive experiments on three popular lane detection datasets. Specifically, our method outperforms all existing methods on the LLAMAS benchmark [3], with the light-weight ResNet-34 backbone.

# 2. Related Work

Segmentation-based Lane Detection. These methods represent lanes as per-pixel segmentation. SCNN [22] formu-

lates lane detection as multi-class semantic segmentation and is the basis of the 1st-place solution in TuSimple challenge [1]. It's core spatial CNN module recurrently aggregates spatial information to complete the discontinuous segmentation predictions, which then requires heuristic postprocessing to decode the segmentation map. Hence, it has a high latency, and only struggles to be real-time after an optimization of Zheng et al. [41]. Others explore knowledge distillation [13] or generative modeling [8], but their performance merely surpasses the seminal SCNN. Moreover, these methods typically assume a fixed number (e.g., 4) of lines. LaneNet [21] leverages an instance segmentation pipeline to deal with a variable number of lines, but it requires post-inference clustering to generate line instances. Some methods leverage row-wise classification [26, 40], which is a customized down-sampling of per-pixel segmentation so that they still require post-processing. Qin et al. [26] propose to trade performance for low latency, but their use of fully-connected layers results in large model size.

In short, segmentation-based methods all require heavy post-processing due to the misalignment of representations. They also suffer from the locality of segmentation task, so that they tend to perform worse under occlusions or extreme lighting conditions.

**Point Detection-based Lane Detection.** The success of object detection methods drives researchers to formulate lane detection as to detect lanes as a series of points (*e.g.*, every 10 pixels in the vertical axis). Line-CNN [15] adapts classic Faster R-CNN [28] as a one-stage lane line detector, but it has a low inference speed (<30 FPS). Later, LaneATT [33] adopts a more general one-stage detection approach that achieves superior performance.

However, these methods have to design heuristic lane anchors, which highly depend on dataset statistics, and require the Non-Maximum Suppression (NMS) as post-processing. On the contrary, we represent lane lines as curves with a fully end-to-end pipeline (anchor-free, NMS-free).

**Curve-based Lane Detection.** The pioneering work [37] proposes a differentiable least squares fitting module to fit a polynomial curve (*e.g.*,  $x = ay^3 + by^2 + cy + d$ ) to points predicted by a deep neural network. The PolyLaneNet [32] then directly learns to predict the polynomial coefficients with simple fully-connected layers. Recently, LSTR [19] uses transformer blocks to predict polynomials in an end-to-end fasion based on the DETR [4].

Curve is a holistic representation of lane line, which naturally eliminates occlusions, requires no post-processing, and can predict a variable number of lines. However, their performance on large and challenging datasets (*e.g.*, CU-Lane [22] and LLAMAS [3]) still lag behind methods of other categories. They also suffer from slow convergence (over 2000 training epochs on TuSimple), high latency architecture (*e.g.*, LSTR [19] uses transformer blocks which

n	Bézier	Polynomial
2nd	0.653	0.945
3rd	0.471	0.558
4th	0.315	0.330

Table 1. Comparison of *n*-order Bézier curves and polynomials  $(x = \sum_{i=0}^{n} a_i y^i)$  on TuSimple [1] *test* set (**lower is better**). Since the official metrics are too lose to show any meaningful difference, we use the fine-grained LPD metric following [32].

are difficult to optimize for low latency). We attribute their failure to the difficult-to-optimize and abstract polynomial coefficients. We propose to use the parametric Bézier curve, which is defined by actual control points on the image co-ordinate system<sup>1</sup>, to address these problems.

**Bézier curve in Deep Learning.** To our knowledge, the only known successful application of Bézier curves in deep learning is the ABCNet [20], which uses cubic Bézier curve for text spotting. However, their method cannot be directly used for our tasks. First, this method still uses NMS so that it cannot be end-to-end. We show in our work that NMS is not necessary so that our method can be an end-to-end solution. Second, it calculates  $L_1$  loss directly on the sparse Bézier control points, which results in difficulties of optimization. We address this problem in our work by leveraging a fine-grained sampling loss. In addition, we propose the feature flip fusion module, which is specifically designed for the lane detection task.

## 3. BézierLaneNet

#### 3.1. Overview

**Preliminaries on Bézier Curve.** The Bézier curve's formulation is shown in Equation (1), which is a parametric curve defined by n + 1 control points:

$$\mathcal{B}(t) = \sum_{i=0}^{n} b_{i,n}(t) \mathcal{P}_i, \ 0 \le t \le 1,$$

$$(1)$$

where  $\mathcal{P}_i$  is the i - th control point,  $b_{i,n}$  are Bernstein basis polynomials of degree n:

$$b_{i,n} = C_n^i t^i (1-t)^{n-i}, \ i = 0, ..., n.$$
 (2)

We use the classic cubic Bézier curve (n = 3), which is empirically found sufficient for modeling lane lines. It shows better ground truth fitting ability than 3rd order polynomial (Table 1), which is the base function for previous curve-based methods [19, 32]. Higher-order curves do not bring substantial gains while the high degrees of freedom leads to instability. All coordinates for points discussed here are relative to the image size (*i.e.*, mostly in range [0, 1]).



Figure 2. Pipeline. Feature from a typical encoder (*e.g.*, ResNet) is strengthened by feature flip fusion, then pooled to 1D and two 1D convolution layers are applied. At last the network predicts Bézier curves through a classification branch and a regression branch.

The Proposed Architecture. The overall model architecture is shown in Figure 2. Specifically, we use layer-3 feature of ResNets [11] as backbone following RESA [41], but we replace the dilation inside the backbone network by two dilated blocks outside with dilation rates [4, 8] [6]. This strikes a better speed-accuracy trade-off for our method, which leaves a  $16 \times$  down-sampled feature map with a larger receptive field. We then add the feature flip fusion module (Section 3.2) to aggregate opposite lane features. The enriched feature map  $(C \times \frac{H}{16} \times \frac{W}{16})$  is then pooled to  $(C \times \frac{W}{16})$  by average pooling, resulting in  $\frac{W}{16}$  proposals (50) for CULane [22]). Two  $1 \times 3$  1D convolutions are used to transform the pooled features, while also conveniently modeling interactions between nearby lane proposals, guiding the network to learn a substitute for the non-maximal suppression (NMS) function. Lastly, the final prediction is obtained by the classification and regression branches (each is only one  $1 \times 1$  1D convolution). The outputs are  $\frac{W}{16} \times 8$ for regression of 4 control points, and  $\frac{W}{16} \times 1$  for existence of lane line object.

#### **3.2. Feature Flip Fusion**

By modeling lane lines as holistic curves, we focus on the geometric properties of individual lane lines (e.g., thin, long, and continuous). Now we consider the global structure of lanes from a front-mounted camera view in driving scenes. Roads have equally spaced lane lines, which appear symmetrical and this property is worth modeling. For instance, the existence of left ego lane line should very likely indicate its right counterpart, the structure of immediate left lane could help describe the immediate right lane, etc.

To exploit this property, we fuse the feature map with its horizontally flipped version (Figure 3). Specifically, two separate convolution and normalization layers transform each feature map, they are then added together before a ReLU activation. With this module, we expect the model to base its predictions on both feature maps.

To account for the slight misalignment of camera captured image (e.g., rotated, turning, non-paired), we apply

<sup>&</sup>lt;sup>1</sup>Actually control points of Bézier curves can be outside the image, but statistically that rarely happens in autonomous driving scenes.



Figure 3. Feature flip fusion. Alignment is achieved by calculating deformable convolution offsets, conditioned on both the flipped and original feature map. Best viewed in color.

deformable convolution [42] with kernel size  $3 \times 3$  for the flipped feature map while learning the offsets conditioned on the original feature map for feature alignment.

We add an auxiliary binary segmentation branch (to segment lane line or non-lane line areas, which would be removed after training) to the ResNet backbone, and we expect it to enforce the learning of spatial details. Interestingly, we find this auxiliary branch improves the performance only when it works with the feature fusion. This is because the localization of the segmentation task may provide a more spatially-accurate feature map, which in turn supports accurate fusion between the flipped features.

Visualizations are shown in Figure 4, from which we can see that the flipped feature does correct the error caused by the asymmetry introduced by the car (Figure 4(a)).

## 3.3. End-to-end Fit of a Bézier Curve

Distances Between Bézier Curves. The key to learning Bézier curves is to define a good distance metric measuring the distances between the ground truth curve and prediction. Naively, one can directly calculate the mean  $L_1$ distance between Bézier curve control points, as in ABC-Net [20]. However, as shown in Figure 5(a), a large  $L_1$  error in curvature control points can demonstrate a very small visual distance between Bézier curves, especially on small or medium curvatures (which is often the case for lane lines). Since Bézier curves are parameterized by  $t \in [0, 1]$ , we propose the more reasonable sampling loss for Bézier curves (Figure 5(b)), by sampling curves at a uniformly spaced set of t values (T), which means equal curve length between adjacent sample points. The t values can be further transformed by a re-parameterization function f(t). Specifically, given Bézier curves  $\mathcal{B}(t)$ ,  $\hat{\mathcal{B}}(t)$ , the sampling loss  $\mathcal{L}_{reg}$  is:

$$\mathcal{L}_{reg} = \frac{1}{n} \sum_{t \in T} ||\mathcal{B}(f(t)) - \hat{\mathcal{B}}(f(t))||_1, \quad (3)$$

where n is the total number of sampled points and is set to 100. We empirically find f(t) = t works well. This simple yet effective loss formulation makes our model easy to converge and less sensitive to hyper-parameters that typically involved in other curved-based or point detection-





Figure 4. Grad-CAM [31] visualization on the last layer of ResNet backbone. (a) Our model can infer existence of an ill-marked lane line, from clear markings and cars around the opposite line. Note that the car is deviated to the left, this scene was not captured with perfect symmetry. (b) When entire road lacks clear marking, both sides are used for a better prediction. Best viewed in color.

based methods, *e.g.*, loss weighting for endpoints loss [19] and line length loss [33] (see Figure 5(b,c)).

**Bézier Ground Truth Generation.** Since lane datasets are currently annotated by on-line key points, we need the Bézier control points for the above sampling loss. Given the annotated points  $\{(k_{x_i}, k_{y_i})\}_{i=1}^m$  on one lane line, where  $(k_{x_i}, k_{y_i})$  denotes the 2D-coordinates of the *i*-th point. Our goal is to obtain control points  $\{\mathcal{P}_i(x_i, y_i)\}_{i=1}^n$ . Similar to [20], we use least squares fitting to solve this equation:

$$\begin{bmatrix} b_{0,n}(t_0) \cdots b_{n,n}(t_0) \\ b_{0,n}(t_1) \cdots b_{n,n}(t_1) \\ \vdots & \ddots & \vdots \\ b_{0,n}(t_m) \cdots b_{n,n}(t_m) \end{bmatrix} \begin{bmatrix} \mathcal{P}_0 \\ \mathcal{P}_1 \\ \vdots \\ \mathcal{P}_n \end{bmatrix} = \begin{bmatrix} k_{x_0} & k_{y_0} \\ k_{x_1} & k_{y_1} \\ \vdots & \vdots \\ k_{x_m} & k_{y_m} \end{bmatrix}, \quad (4)$$

 $\{t_i\}_{i=0}^m \in [0, 1]$ . Considering m >> n and the linear independence of b, it can be efficiently solved by a pseudoinverse of the leftmost matrix (mostly full column rank). Different from [20], we do not set  $\mathcal{P}_0$  and  $\mathcal{P}_n$  to original endpoints, which leads to better quality labels.

Label and Prediction Matching. After obtaining the ground truth, in training, we perform a one-to-one assignment between G labels and N predictions (G < N) using optimal bipartite matching, to attain a fully end-toend pipeline. Following Wang *et al.* [38], we find a Gpermutation of N predictions  $\pi \in \Pi_G^N$  that formulates the best bipartite matching:

$$\hat{\pi} = \operatorname*{arg\,max}_{\pi \in \Pi_G^N} \sum_{i}^G Q_{i,\pi(i)},\tag{5}$$

$$Q_{i,\pi(i)} = \left(\hat{p}_{\pi(i)}\right)^{1-\alpha} \cdot \left(1 - L_1(b_i, \hat{b}_{\pi(i)})\right)^{\alpha}, \quad (6)$$



Figure 5. Lane loss functions. (a) The  $L_1$  distance of control points is not highly correlated with the actual distance between curves. (b) The proposed sampling loss is one unified distance metric by *t*-sampling. (c) Typical loss for polynomial regression [19], at least 3 separate losses are required: *y*-sampling loss, *y* start point loss, *y* end point loss.

where  $Q_{i,\pi(i)} \in [0, 1]$  represents matching quality of the *i*th label with the  $\pi(i)$ -th prediction, based on  $L_1$  distance between curves  $b_i, \hat{b}_{\pi(i)}$  (sampling loss) and class score  $\hat{p}_{\pi(i)}$ .  $\alpha$  is set to 0.8 by default. The above equations can be efficiently solved by the well-known Hungarian algorithm.

Wang *et al.* [38] also use a spatial prior that restricts the matched prediction to a spatial neighborhood of the label (object center distance, the *centerness* prior in FCOS [35]). However, since lots of lanes are long lines with a large slope, this centerness prior is not useful. See Appendix E for more investigations on matching priors.

**Overall Loss.** Other than Bézier curve sampling loss, there is also the classification loss  $\mathcal{L}_c$  for the lane object classification (existence) branch. Since the imbalance between positive and negative examples is not as severe in lane detection as in object detection, instead of the focal loss [16], we use the simple weighted binary cross-entropy loss:

$$\mathcal{L}_{cls} = -(y\log(p) + w(1-y)\log(1-p)), \quad (7)$$

where w is the weighting for negative samples, which is set to 0.4 in all experiments. The loss  $\mathcal{L}_{seg}$  for the binary segmentation branch (Section 3.2) takes the same format.

The overall loss is a weighted sum of all three losses:

$$\mathcal{L} = \lambda_1 \mathcal{L}_{reg} + \lambda_2 \mathcal{L}_{cls} + \lambda_3 \mathcal{L}_{seg},\tag{8}$$

where  $\lambda_1, \lambda_2, \lambda_3$  are set to 1, 0.1, 0.75, respectively.

#### 4. Experiments

#### 4.1. Datasets

To evaluate the proposed method, we conduct experiments on three well-known datasets: TuSimple [1], CU-Lane [22] and LLAMAS [3]. TuSimple dataset was collected on highways with high-quality images, under fair

Dataset	Train	Val	Test	Resolution	#Lines
TuSimple [1]	3268	358	2782	$720\times1280$	$\leq 5$
CULane [22]	88880	9675	34680	$590 \times 1640$	$\leq 4$
LLAMAS [3]	58269	20844	20929	$717\times1276$	$\leq 4^*$

Table 2. Details of datasets. \*Number of lines in LLAMAS dataset is more than 4, but official metric only evaluates 4 lines.

weather conditions. CULane dataset contains more complex urban driving scenarios, including shades, extreme illuminations, and road congestion. LLAMAS is a newly formed large-scale dataset, it is the only lane detection benchmark without public *test* set labels. Details of these datasets can be found in Table 2.

#### 4.2. Evalutaion Metics

For CULane [22] and LLAMAS [3], the official metric is F1 score from [22]:

$$F1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall},$$
(9)

where Precision =  $\frac{TP}{TP+FP}$  and Recall =  $\frac{TP}{TP+FN}$ . Lines are assumed to be 30 pixels wide, prediction and ground truth lines with pixel IoU over 0.5 are considered a match.

For TuSimple [1] dataset, the official metrics include Accuracy, false positive rate (FPR), and false negative rate (FNR). Accuracy is computed as  $\frac{N_{pred}}{N_{gt}}$ , where  $N_{pred}$  is the number of correctly predicted on-line points and  $N_{gt}$  is the number of ground truth on-line points.

### 4.3. Implementation Details

Fair Comparison. To fairly compare among different stateof-the-art methods, we re-implement representative methods [19, 22, 41] in a unified PyTorch framework. We Also provide a semantic segmentation baseline [5] originally proposed in [22]. All our implementations do **not** use *val* set in training, and tune hyper-parameters **only** on *val* set. Some methods with reliable open-source codes are reported from their own codes [26, 32, 33]. For platform sensitive metric Frames-Per-Second (FPS), we re-evaluate all reported methods on the same RTX 2080 Ti platform. More details for implementations and FPS tests are in Appendices A to C.

**Training.** We train 400, 36, 20 epochs for TuSimple, CU-Lane, and LLAMAS, respectively (training of our model takes only 12 GPU hours on a single RTX 2080 Ti), and the input resolution is  $288 \times 800$  for CULane [22] and  $360 \times 640$ for others, following common practice. Other than these, all hyper-parameters are tuned on CULane [22] *val* set and remain the same for our method across datasets. We use Adam optimizer with learning rate  $6 \times 10^{-4}$ , weight decay  $1 \times 10^{-4}$ , batch size 20, Cosine Annealing learning rate schedule as in [33]. Data augmentation includes random affine transforms, random horizontal flip, and color jitter.

	CULane [22]								Ti	uSimple [	1]					
Method	Ep.	Total	Normal	Crowd	Night	No line	Shadow	Arrow	Dazzle light	Curve	$\mathbf{Cross}\downarrow$	train+val	Ep.	Acc.	$\mathbf{FPR}\downarrow$	$\mathbf{FNR}\downarrow$
Segmentation-based																
Baseline (ResNet-18)*	12	65.30	85.45	62.63	61.04	33.88	51.72	78.15	53.05	59.70	1915		50	94.25	0.088	0.089
Baseline (ResNet-34)*	12	69.92	89.46	66.66	65.38	40.43	62.17	83.18	58.51	63.00	1713		50	95.31	0.064	0.062
Baseline (ResNet-101)*	12	71.37	90.11	67.89	67.01	43.10	70.56	85.09	61.77	65.47	1883		50	95.19	0.062	0.062
SCNN (ResNet-18) [22]*	12	72.19	90.98	70.17	66.54	43.12	66.31	85.62	62.20	65.58	1808		50	94.77	0.075	0.074
SCNN (ResNet-34) [22]*	12	72.70	91.06	70.41	67.75	44.64	68.98	86.50	61.57	65.75	2017		50	95.25	0.063	0.063
SCNN (ResNet-101) [22]*	12	73.58	91.10	71.43	68.53	46.39	72.61	86.87	61.95	67.01	1720		50	95.69	0.052	0.050
UFLD (ResNet-18) [26]**	50	68.4	87.7	66.0	62.1	40.2	62.8	81.0	58.4	57.9	1743	-	_	-	_	_
UFLD (ResNet-34) [26]**	50	72.3	90.7	70.2	66.7	44.4	69.3	85.7	59.5	69.5	2037	-	_	-	_	_
RESA (ResNet-18) [41]*	12	72.90	91.23	70.57	67.16	45.24	68.01	86.56	64.32	66.19	1679		50	95.24	0.069	0.057
RESA (ResNet-34) [41]*	12	73.66	91.31	71.80	67.54	46.57	72.74	86.94	64.46	67.31	1701		50	95.15	0.069	0.059
RESA (ResNet-101) [41]*	12	<b>74.04</b>	91.45	71.51	69.01	46.54	75.83	87.75	63.90	68.24	1522		50	95.56	0.058	0.051
Point detection-based																
FastDraw (ResNet-18) [25]	-	_	_	_	_	_	_	_	_	_	_	$\checkmark$	7	94.9	0.061	0.047
CurveLanes-NAS-S [39]	12	71.4	88.3	68.6	66.2	47.9	68.0	82.5	63.2	66.0	2817	-	_	-	_	_
CurveLanes-NAS-M [39]	12	73.5	90.2	70.5	68.2	48.8	69.3	85.7	65.9	67.5	2359	-	_	_	_	_
CurveLanes-NAS-L [39]	12	74.8	90.7	72.3	68.9	49.4	70.1	85.8	67.7	68.4	1746	-	_	-	_	_
LaneATT (ResNet-18) [33]**	15	74.88	90.98	72.78	68.61	48.23	69.68	85.44	65.43	63.18	1163	$\checkmark$	100	95.57	0.036	0.030
LaneATT (ResNet-34) [33]**	15	76.42	91.94	74.76	70.32	49.17	77.68	88.14	65.92	68.07	1323	$\checkmark$	100	95.63	0.035	0.029
LaneATT (ResNet-122) [33]**	15	76.79	91.50	76.04	70.43	50.29	75.96	86.16	68.99	63.99	1265	$\checkmark$	100	96.10	0.056	0.022
Curve-based																
PolyLaneNet (EfficientNet-B0) [32]**	-	_	_	_	_	_	_	_	_	_	_	$\checkmark$	2695	93.36	0.094	0.093
LSTR (ResNet-18, 1×) [19]*	_	_	_	_	_	_	_	_	_	_	_		2000	95.06	0.049	0.042
LSTR (ResNet-18, 2×) [19]*	150	68.72	86.78	67.34	59.92	40.10	59.82	78.66	56.63	56.64	1166	-	_	_	_	_
BézierLaneNet (ResNet-18)	36	73.67	90.22	71.55	68.70	45.30	70.91	84.09	62.49	58.98	996		400	95.41	0.053	0.046
BézierLaneNet (ResNet-34)	36	75.57	91.59	73.20	69.90	48.05	76.74	87.16	69.20	62.45	888		400	95.65	0.051	0.039

Table 3. Results on *test* set of CULane [22] and TuSimple [1]. \*reproduced results in our code framework, best performance from three random runs. \*\*reported from reliable open-source codes from the authors.

**Testing.** No post-processing is required for curve methods. Standard Gaussian blur and row selection post-processing is applied to segmentation methods. NMS is used for LaneATT [33], while we remove its post-inference B-Spline interpolation in CULane [22], to align with our framework.

### 4.4. Comparisons

**Overview.** Experimental results are shown in Tables 3 and 4. TuSimple [1] is a small dataset that features clearweather highway scenes and has a relatively easy metric, most methods thrive in this dataset. Thus, we mainly focus on the other two large-scale datasets [3, 22], where there is still a rather clear difference between methods. For highperformance methods (> 70% F1 on CULane [22]), we also show efficiency metrics (FPS, Parameter count) in Table 5. Comparison with Curve-based Methods. As shown in Tables 3 and 4, in all datasets, BézierLaneNet outperforms previous curve-based methods [19, 32] by a clear margin, advances the state-of-the-art of curve-based methods by 6.85% on CULane [22] and 6.77% on LLAMAS [3]. Thanks to our fully convolutional and fully end-toend pipeline, BézierLaneNet runs over 2× faster than LSTR [19]. LSTR has a speed bottleneck from transformer architecture, the  $1 \times$  and  $2 \times$  model have FPS 98 and 97, respectively<sup>2</sup>. While curves are difficult to learn, our method converges  $4-5 \times$  faster than LSTR. For the first time, an elegant curve-based method can challenge well-designed segmentation methods or point detection methods on these datasets

while showing a favorable trade-off, with an acceptable convergence time.

**Comparison with Segmentation-based Methods.** These methods tend to have a low speed due to recurrent feature aggregation [22, 41], and the use of high-resolution feature map [5, 22, 41]. BézierLaneNet outperforms them in both speed and accuracy. Our small models even compare favorably against RESA [41] and SCNN [22] with large ResNet-101 backbone, surpassing them in CULane [22] with a clear margin (1 ~ 2%). On LLAMAS [3], where the dataset restricts testing on 4 center lines, the segmentation approach shows strong performance (Table 4). Nevertheless, our ResNet-34 model still outperforms SCNN by 0.92%.

UFLD [26] reformulates segmentation to row-wise classification on a down-sampled feature map to achieve fast speed, at the cost of accuracy. Compared to us, UFLD (ResNet-34) is 0.9% lower on CULane Normal, while 7.4%, 3.0%, 3.2% worse on Shadow, Crowd, Night, respectively. Overall, our method with the same backbones outperforms UFLD by  $3 \sim 5\%$ , while being faster on ResNet-34. Besides, UFLD uses large fully-connected layers to optimize latency, which causes a huge model size (the largest in Table 5).

A drawback for all segmentation methods is the weaker performance on **Dazzle Light**. Per-pixel (or per-pixel grid for UFLD [26]) segmentation methods may rely on information from local textures, which is destroyed by extreme exposure to light. While our method predicts lane lines as holistic curves, hence robust to changes in local textures.

**Comparison with Point Detection-based Methods.** Xu *et al.* [39] finds a series of point detection-based models with

 $<sup>^{2}</sup>$ The original 420 FPS report from LSTR paper [19], is throughput with batch size 16, detailed discussions in Appendix A.

	LLAMAS [3]				
Method	Ep.	F1	Precision	Recall	
Segmentation-based					
Baseline (ResNet-34)*	18	93.43	92.61	94.27	
SCNN (ResNet-34) [22]*	18	94.25	94.11	94.39	
Point detection-based					
LaneATT (ResNet-18) [33]**	- 15	93.46	96.92	90.24	
LaneATT (ResNet-34) [33]**	15	93.74	96.79	90.88	
LaneATT (ResNet-122) [33]**	15	93.54	96.82	90.47	
Curve-based					
PolyLaneNet (EfficientNet-B0) [32]**	- 75	88.40	88.87	87.93	
BézierLaneNet (ResNet-18)	20	94.91	95.71	94.13	
BézierLaneNet (ResNet-34)	20	95.17	95.89	94.46	

Table 4. Results from LLAMAS [3] test server.

neural architecture search techniques called CurveLanes-NAS. Despite its complex pipeline and extensive architecture search for the best accuracy-FLOPs trade-off, our simple ResNet-34 backbone model (29.9 GFLOPs) still surpasses its large model (86.5 GFLOPs) by 0.8% on CULane. CurveLanes-NAS also performs worse under occlusions, a similar drawback as the segmentation methods without recurrent feature fusion [5, 26]. As shown in Table 3, with similar model capacity compared to our ResNet-34 model, CurveLanes-NAS-M (35.7 GFLOPs) is 1.4% worse on **Normal** scenes, but the gap on **Shadow** and **Crowd** are 7.4% and 2.7%.

Recently, LaneATT [33] achieves higher performance with a point detection network. However, their design is not fully end-to-end (requires Non-Maximal Suppression (NMS)), based on heuristic anchors (>1000), which are calculated directly from the dataset's statistics, thus may systematically pose difficulties in generalization. Still, with ResNet-34, our method outperforms LaneATT on the LLA-MAS [3] test server (1.43%), with a significantly higher recall (3.58%). We also achieve comparable performance to LaneATT on TuSimple [1] using only the train set, and only  $\sim 1\%$  worse on CULane. Our method performs significantly better in **Dazzle Light** (3.3% better), comparably in **Night** (0.4% lower). It also has a lower False Positive (FP) rate on Crossroad scenes (Cross), even though LaneATT shows an extremely low-FP characteristic (large Precision-Recall gap in Table 4). Methods that rely on heuristic anchors [33] or heuristic decoding process [22, 26, 39, 41] tend to have more false predictions in this scene. Moreover, the NMS is a sequential process that could have unstable runtime in real-world applications. Even when NMS was not evaluated on real inputs, our models are 29%, 28% faster, have  $2.9 \times$ ,  $2.3 \times$  fewer parameters, compared to LaneATT on ResNet-18 and ResNet-34 backbones, respectively.

To summarize, previous curve-based methods (Poly-LaneNet [32], LSTR [19]) have significantly worse performance. Fast methods trades either accuracy (UFLD [26]) or model size (UFLD [26], LaneATT [33]) for speed. Accurate

Method	FPS ↑	Params (M) $\downarrow$
Segmentation-based (ignored post-processing time)		
Baseline (ResNet-101)	27	43.56
SCNN (ResNet-18) [22]	21	12.63
SCNN (ResNet-34) [22]	21	22.74
SCNN (ResNet-101) [22]	14	44.15
UFLD (ResNet-34) [26]	144	71.58
RESA (ResNet-18) [41]	68	6.61
RESA (ResNet-34) [41]	54	11.99
RESA (ResNet-101) [41]	25	31.46
Point detection-based (ignored NMS time in real images)		
LaneATT (ResNet-18) [33]	165	12.02
LaneATT (ResNet-34) [33]	117	22.13
LaneATT (ResNet-122) [33]	26	8.55
Curve-based (entirely end-to-end)		
BézierLaneNet (ResNet-18)	213	4.10
BézierLaneNet (ResNet-34)	150	9.49

Table 5. FPS (*image/s*) and model size. All FPS results are tested with  $360 \times 640$  random inputs on the same platform. Here only shows models with > 70% CULane [22] F1 score.

methods either discards the end-to-end pipeline (LaneATT [33]), or entirely fails the real-time requirement (SCNN [22], RESA [41]). While our BézierLaneNet is fully end-to-end, fast (>150 FPS), light-weight (<10 million parameters) and maintains consistent high accuracy across datasets.

#### 4.5. Analysis

Although we develop our method by tuning on the *val* set, we re-run ablation studies with ResNet-34 backbone (including our full method) and report performance on the CULane *test* set for clear comparison.

Curve representation	F1
Cubic Bézier curve baseline 3rd Polynomial baseline	$68.89 \\ 1.49$
BézierLaneNet 3rd Polynomial from BézierLaneNet	$75.41 \\ 5.01$

Table 6. Curve representations. Baselines directly predict curve coefficients without feature flip fusion.

**Importance of Parametric Bézier Curve.** We first replace the Bézier curve prediction with a 3rd order polynomial, adding auxiliary losses for start and end points. As shown in Table 6, polynomials catastrophically fail to converge in our fully convolutional network, even when trained with 150 epochs (details in Appendix B.8). Then we consider modifying the LSTR [19] to predict cubic Bézier curves, the performance is similar to predicting polynomials. We conclude that heavy MLP may be necessary to learn polynomials [19, 32], while predicting Bézier control points from position-aware CNN is the best choice. The transformerbased LSTR decoder destroys the fine spatial information, suppresses the advancement of curve function.

Feature Flip Fusion Design. As shown in Table 7, feature flip fusion brings 4.07% improvement. We also find that the auxiliary segmentation loss can regularize and increase

СР	SP	Flip	Deform	Seg	F1
$\checkmark$					63.74
	$\checkmark$				68.89
	$\checkmark$			$\checkmark$	65.82
	$\checkmark$	$\checkmark$			70.28
	$\checkmark$	$\checkmark$	$\checkmark$		72.96
	$\checkmark$	$\checkmark$		$\checkmark$	73.97
	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	75.41

Table 7. Ablations. **CP**: Control point loss [20]. **SP**: The proposed sampling loss. **Flip**: The feature flip fusion module. **Deform**: Employ the deformable convolution in feature flip fusion. **Seg**: Auxiliary segmentation loss.

the performance further, by 2.45%. It is worth noting that auxiliary loss only works with feature fusion, it can lead to degenerated results when directly applied on the baseline (-3.07%). A standard  $3 \times 3$  convolution performs worse than deformable convolution, by 2.68% and 1.44%, before and after adding the auxiliary segmentation loss, respectively. We attribute this to the effects of feature alignment. **Bézier Curve Fitting Loss.** As shown in Table 7, replacing the sampling loss by direct loss on control points lead to inferior performance (-5.15%) in the baseline setup). Inspired by the success of IoU loss in object detection. We also implemented a IoU loss (formulas in Appendix D) for the convex hull of Bézier control points. However, the convex hull of close-to-straight lane lines are too small, the IoU loss is numerically unstable, thus failing to facilitate the sampling loss.

Model	Aug	F1
LSTR (ResNet-18, 2×) [19] LSTR (ResNet-18, 2×) [19]	$\checkmark$	$\begin{array}{c} 68.72 \\ 39.77(-28.95) \end{array}$
BézierLaneNet (ResNet-34) BézierLaneNet (ResNet-34)	$\checkmark$	$75.41 \\ 55.11(-20.30)$

Table 8. A	Augmentation	ablations.	Ang:	Strong	data aug	mentation.
14010 0. 1	iaginontation	aoranomo.		ouons	autu uu_	incincution.
	0					

**Importance of Strong Data Augmentation.** Strong data augmentation is defined by a series of affine transforms and color distortions, the exact policy may slightly vary for different methods. For instance, we use random affine transform, random horizontal flip, and color jitter. LSTR [19] also uses random lighting. Default augmentation includes only a small rotation (3 degrees). As shown in Table 8, strong augmentation is essential to avoid over-fitting for curve-based methods.

For segmentation-based methods [5, 22, 41], we fast validated strong augmentation on the smaller TuSimple [1] dataset. All shows a  $1 \sim 2\%$  degradation. This suggests that they may be robust due to per-pixel prediction and heuristic post-processing. But they highly rely on learning the distribution of local features such as texture, which could become confusing by strong augmentation.

#### 4.6. Limitations and Discussions

Curves are indeed a natural representation of lane lines. However, their elegance in modeling inevitably brings a drawback. It is difficult for the curvature coefficients to generalize when the data distribution is highly biased (almost all lane lines are straight lines in CULane). Our Bézier curve approach has already alleviated this problem to some extent and has achieved an acceptable performance (62.45) in CULane **Curve**. On datasets such as TuSimple and LLA-MAS [1,3], where the curvature distribution is fair enough for learning, our method achieves even better performance. To handle broader corner cases, *e.g.*, sharp turns, blockages and bad weather, datasets such as [30, 34, 39] may be useful.

The feature flip fusion is specifically designed for a front-mounted camera, which is the typical use case of deep lane detectors. Nevertheless, there is still a strong inductive bias by assuming scene symmetry. In future work, it would be interesting to find a replacement for this module, to achieve better generalization and to remove the deformable convolution operation, which poses difficulty for effective integration into edge devices such as Jetson.

More discussions in Appendix G.

### **5.** Conclusions

In this paper, we have proposed BézierLaneNet: a novel fully end-to-end lane detector based on parametric Bézier curves. The on-image Bézier curves are easy to optimize and naturally model the continuous property of lane lines, without heavy designs such as recurrent feature aggregation or heuristic anchors. Besides, a feature flip fusion module is proposed. It efficiently models the symmetry property of the driving scene, while also being robust to slight asymmetries by using deformable convolution. The proposed model has achieved favorable performance on three datasets, defeating all existing methods on the popular LLAMAS benchmark. It is also both fast (>150 FPS) and light-weight (<10 million parameters).

Acknowledgements. This work has been sponsored by National Key Research and Development Program of China (2019YFC1521104), National Natural Science Foundation of China (61972157, 72192821), Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102), Shanghai Science and Technology Commission (21511101200), Art major project of National Social Science Fund (I8ZD22), and SenseTime Collaborative Research Grant. We thank Jiaping Oin for guidance on road design and geometry, Yuchen Gong and Pan Chen for helping with CAM visualizations, Zhijun Gong, Jiachen Xu and Jingyu Gong for insightful discussions about math, Fengqi Liu for providing GPUs, Lucas Tabelini for cooperation in evaluating [32, 33], and the CVPR reviewers for constructive comments.

## References

- [1] TuSimple benchmark. https://github.com/ TuSimple/tusimple-benchmark, 2017. 2, 3, 5, 6, 7, 8, 12, 14, 15
- [2] Claudine Badue, Rânik Guidolini, Raphael Vivacqua Carneiro, Pedro Azevedo, Vinicius B Cardoso, Avelino Forechi, Luan Jesus, Rodrigo Berriel, Thiago M Paixao, Filipe Mutz, et al. Self-driving cars: A survey. *Expert Systems with Applications*, 2021. 1
- [3] Karsten Behrendt and Ryan Soussan. Unsupervised labeled lane markers using maps. In *ICCV*, 2019. 2, 5, 6, 7, 8, 14, 15
- [4] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-toend object detection with transformers. In ECCV, 2020. 2
- [5] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L. Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. 1, 5, 6, 7, 8, 11
- [6] Qiang Chen, Yingming Wang, Tong Yang, Xiangyu Zhang, Jian Cheng, and Jian Sun. You only look one-level feature. In CVPR, 2021. 3
- [7] MOT Highway Department and Highway Engineering Committee under China Association for Engineering Construction Standardization. *Technical Standard of Highway Engineering*. 2004. 14
- [8] Mohsen Ghafoorian, Cedric Nugteren, Nóra Baka, Olaf Booij, and Michael Hofmann. El-gan: Embedding loss driven generative adversarial networks for lane detection. In ECCV Workshops, 2018. 2
- [9] Ross Girshick. Fast r-cnn. In ICCV, 2015. 1
- [10] Alexandru Gurghian, Tejaswi Koduri, Smita V Bailur, Kyle J Carey, and Vidya N Murali. Deeplanes: End-to-end lane position estimation using deep neural networksa. In CVPR Workshops, 2016. 14
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In CVPR, 2016. 3
- [12] Aharon Bar Hillel, Ronen Lerner, Dan Levi, and Guy Raz. Recent progress in road and lane detection: a survey. *MVA*, 2014.
- [13] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning lightweight lane detection cnns by self attention distillation. In *ICCV*, 2019. 2
- [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, 2012. 12
- [15] Xiang Li, Jun Li, Xiaolin Hu, and Jian Yang. Line-cnn: End-to-end traffic line detection with line proposal unit. *ITS*, 2019. 1, 2
- [16] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 5
- [17] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In ECCV. Springer, 2014. 12

- [18] Lizhe Liu, Xiaohao Chen, Siyu Zhu, and Ping Tan. Condlanenet: a top-to-down lane detection framework based on conditional convolution. In *ICCV*, 2021. 14
- [19] Ruijin Liu, Zejian Yuan, Tie Liu, and Zhiliang Xiong. Endto-end lane shape prediction with transformers. In WACV, 2021. 1, 2, 3, 4, 5, 6, 7, 8, 11, 12
- [20] Yuliang Liu, Hao Chen, Chunhua Shen, Tong He, Lianwen Jin, and Liangwei Wang. Abcnet: Real-time scene text spotting with adaptive bezier-curve network. In *CVPR*, 2020. 3, 4,8
- [21] Davy Neven, Bert De Brabandere, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Towards end-to-end lane detection: an instance segmentation approach. In *IV*, 2018. 1, 2
- [22] Xingang Pan, Jianping Shi, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Spatial as deep: Spatial cnn for traffic scene understanding. In AAAI, 2018. 1, 2, 3, 5, 6, 7, 8, 11, 14, 15
- [23] Tim A Pastva. Bezier curve fitting. Technical report, NAVAL POSTGRADUATE SCHOOL MONTEREY CA, 1998. 13
- [24] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 11, 14
- [25] Jonah Philion. Fastdraw: Addressing the long tail of lane detection by adapting a sequential prediction network. In *CVPR*, 2019. 6
- [26] Zequn Qin, Huanyu Wang, and Xi Li. Ultra fast structureaware deep lane detection. In ECCV, 2020. 2, 5, 6, 7, 11
- [27] Zhan Qu, Huan Jin, Yang Zhou, Zhen Yang, and Wei Zhang. Focus on local: Detecting lane marker from bottom up via key point. In CVPR, 2021. 14
- [28] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *NeurIPS*, 2015. 1, 2
- [29] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *CVPR*, 2019. 14
- [30] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Semantic foggy scene understanding with synthetic data. *IJCV*, 2018. 8
- [31] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In CVPR, 2017. 4
- [32] Lucas Tabelini, Rodrigo Berriel, Thiago M Paixao, Claudine Badue, Alberto F De Souza, and Thiago Oliveira-Santos. Polylanenet: Lane estimation via deep polynomial regression. In *ICPR*, 2020. 1, 2, 3, 5, 6, 7, 8, 12
- [33] Lucas Tabelini, Rodrigo Berriel, Thiago M Paixao, Claudine Badue, Alberto F De Souza, and Thiago Oliveira-Santos. Keep your eyes on the lane: Real-time attention-guided lane detection. In *CVPR*, 2021. 1, 2, 4, 5, 6, 7, 8, 11, 12
- [34] Xin Tan, Ke Xu, Ying Cao, Yiheng Zhang, Lizhuang Ma, and Rynson W. H. Lau. Night-time scene parsing with a large real dataset. *TIP*, 2021. 8

- [35] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *ICCV*, 2019. 5
- [36] Federal Highway Administration under United States Department of Transportation. Standard Specifications for Construction of Roads and Bridges on Federal Highway Projects. 2014. 14
- [37] Wouter Van Gansbeke, Bert De Brabandere, Davy Neven, Marc Proesmans, and Luc Van Gool. End-to-end lane detection through differentiable least-squares fitting. In *ICCV Workshops*, 2019. 2
- [38] Jianfeng Wang, Lin Song, Zeming Li, Hongbin Sun, Jian Sun, and Nanning Zheng. End-to-end object detection with fully convolutional network. In *CVPR*, 2021. 2, 4, 5
- [39] Hang Xu, Shaoju Wang, Xinyue Cai, Wei Zhang, Xiaodan Liang, and Zhenguo Li. Curvelane-nas: Unifying lanesensitive architecture search and adaptive point blending. In *ECCV*, 2020. 1, 6, 7, 8
- [40] Seungwoo Yoo, Hee Seok Lee, Heesoo Myeong, Sungrack Yun, Hyoungwoo Park, Janghoon Cho, and Duck Hoon Kim. End-to-end lane marker detection via row-wise classification. In *CVPR Workshops*, 2020. 2
- [41] Tu Zheng, Hao Fang, Yi Zhang, Wenjian Tang, Zheng Yang, Haifeng Liu, and Deng Cai. Resa: Recurrent feature-shift aggregator for lane detection. In AAAI, 2021. 1, 2, 3, 5, 6, 7, 8, 11
- [42] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *CVPR*, 2019. 2, 4

**Appendix Overview.** The Appendix is organized as follows: Appendix **A** describes the FPS test protocol and environments; Appendix **B** introduces implementation details for each compared method (including ours in Appendix **B**.8); Appendix **C** provides implementation details for Bézier curves, including sampling, ground truth generation and transforms; Appendix **D** formulates the IoU loss for Bézier curves and discusses why it failed; Appendix **E** explores matching priors other than the *centerness* prior; Appendix **F** shows extra ablation studies on datasets other than CULane [22], to verify the generalization of feature flip fusion; Appendix **G** discusses limitations and recognizes new progress in the field; Appendix **H** presents qualitative results from our method, visualized on three datasets.

## **A. FPS Test Protocol**

Let one Frames-Per-Second (FPS) test trial be the average runtime of 100 consecutive model inference with its Py-Torch [24] implementation, without calculating gradients. The input is a 3x360x640 random Tensor (some use all 1 [33], which does not have impact on speed). Note that all methods do **not** use optimization from packages like TensorRT. We wait for all CUDA kernels to finish before counting the whole runtime. We use Python *time.perf\_counter()* since it is more precise than *time.time()*. For all methods, the FPS is reported as the best result from 3 trials.

Before each test trial, at least 10 forward pass is conducted as warm-up of the device. For each new method to be tested, we keep running warm-up trials of a recorded method until the recorded FPS is reached again, so we can guarantee a similar peak machine condition as before.

**Evaluation Environment.** The evaluation platform is a 2080 Ti GPU (standard frequency), on a Intel Xeon-E3 CPU server, with CUDA 10.2, CuDNN 7.6.5, PyTorch 1.6.0. FPS is a platform-sensitive metric, depending on GPU frequency, condition, bus bandwidth, software versions, *etc.* Also using 2080 Ti, Tabelini *et al.* [33] can achieve a better peak performance for all methods. Thus we use the same platform for all FPS tests, to provide fair comparisons.

**Remark.** Note that **FPS** (image/s) is different from throughput (image/s). Since FPS restricts batch size to 1, which better simulates the real-time application scenario. While throughput considers a batch size more than 1. LSTR [19] reported a 420 FPS for its fastest model, which is actually throughput with batch size 16. Our re-tested FPS is 98.

## **B.** Specifications for Compared Methods

### **B.1. Segmentation Baseline**

The segmentation baseline is based on DeeplabV1 [5], originally proposed in SCNN [22]. It is essentially the orig-

inal DeeplabV1 without CRF, lanes are considered as different classes, and a separate lane existence branch (a series of convolution, pooling and MLP) is used to facilitate lane post-processing. We optimized its training and testing scheme based on recent advances [41]. Re-implemented in our codebase, it attains higher performance than what recent papers usually report.

**Post-processing.** First, the existence of a lane is determined by the lane existence branch. Then, the predicted per-pixel probability map is interpolated to the input image size. After that, a  $9 \times 9$  Gaussian blur is applied to smooth the predictions. Finally, for each existing lane class, the smoothed probability map is traversed by pre-defined Y coordinates (quantized), and corresponding X coordinates are recorded by the maximum probability position on the row (provided it passes a fixed threshold). Lanes with less than two qualified points are simply discarded.

**Data Augmentation.** We use a simple random rotation with small angles (3 degrees), then resize to input resolution.

### **B.2. SCNN**

Our SCNN [22] is re-implemented from the Torch7 version of the official code. Advised by the authors, we added an initialization trick for the spatial CNN layers, and learning rate warm-up, to prevent gradient explosion caused by recurrent feature aggregation. Thus, we can safely adjust the learning rate. Our improved SCNN achieves significantly better performance than the original one.

Some may find reports of 96.53 accuracy of SCNN on TuSimple. However, that was a competition entry trained with external data. We report SCNN with ResNet backbones, trained with the same data as other re-implemented methods in our codebase.

**Post-processing.** Same as Appendix B.1.

Data Augmentation. Same as Appendix B.1.

### **B.3. RESA**

Our RESA [41] is implemented based on its published paper. A main difference to the official code release is we do not cutout no-lane areas (in each dataset, there is a certain height range for lane annotation). Because that trick is dataset specific and not generalizable, we do not use that for all compared methods. Other differences are all validated to have better performance than the official code, at least on the CULane *val* set.

Post-processing. Same as Appendix B.1.

**Data Augmentation.** Same as Appendix B.1. The original RESA paper [41] also apply random horizontal flip, which was found ineffective in our re-implementation.

#### **B.4. UFLD**

Ultra Fast Lane Detection (UFLD) [26] is reported from their paper and open-source code. Since TuSimple FP and FN information is not in the paper, and training from source code leads to very high FP rate (almost 20%), we did not report their performance on this dataset. We adjusted its profiling scripts to calculate number of parameters and FPS in our standard.

**Post-processing.** Since this method uses gridding cells (each cell is equivalent to several pixels in a segmentation probability map), each point's X coordinate is calculated as the expectation of locations (cells from the same row), i.e. a weighted average by probability. Differently from segmentation post-processing, it is possible to be efficiently implemented.

**Data Augmentation.** Augmentations include random rotation and some form of random translation.

### **B.5.** PolyLaneNet

PolyLaneNet [32] is reported from their paper and opensource code. We added a profiling script to calculate number of parameters and FPS in our standard.

**Post-processing.** This method requires no post-processing. **Data Augmentation.** Augmentations include large random rotation (10 degrees), random horizontal flip and random crop. They are applied with a probability of  $\frac{10}{11}$ .

## **B.6.** LaneATT

LaneATT [33] is reported from their paper and opensource code. We adjusted its profiling scripts to calculate parameters and FPS in our standard.

**Post-processing.** Non-Maximal Suppression (NMS) is implemented by a customized CUDA kernel. An extra interpolation of lanes by B-Spline is removed both in testing and profiling, since it is slowly executed on CPU and provides little improvement ( $\sim 0.2\%$  on CULane).

**Data Augmentation.** LaneATT uses random affine transforms including scale, translation and rotation. While it also uses random horizontal flip.

**Followup.** We did not have time to validate the reimplementation of LaneATT in our codebase, prior the submission deadline. Therefore, the LaneATT performance is still reported from the official code. Our re-implementation indicates that all LaneATT results are reproducible except for the ResNet-34 backbone on CULane, which is slightly outside the standard deviation range, but still reasonable.

## **B.7. LSTR**

LSTR [19] is re-implemented in our codebase. All ResNet backbone methods start from ImageNet [14] pretraining. While LSTR [19] use 256 channels ResNet-18 for CULane (2×), 128 channels for other datasets (1×), which makes it impossible to use off-the-shelf pre-trained ResNets. Although whether ImageNet pre-training helps lane detection is still an open question. Our reported performance of LSTR on CULane, is the first documented report of LSTR on this dataset. With tuning of hyper-parameters (learning rate, epochs, prediction threshold), bug fix (the original classification branch has 3 output channels, which should be 2), we achieve 4% better performance on CU-Lane than the authors' trial. Specifically, we use learning rate  $2.5 \times 10^{-4}$  with batch size 20. 150 and 2000 epochs, 0.95 and 0.5 prediction thresholds, for CULane and TuSimple. The lower threshold in TuSimple is due to the official test metric, which significantly favors a high recall. However, for real-world applications, a high recall leads to high False Positive rate, which is undesired.

We divide the curve loss weighting by 10 with our LSTR-Beizer ablation, since there were 100 sample points with both X and Y coordinates to fit, that is a loss scale about 10 times the original loss (LSTR loss takes summation of point L1 distances instead of average). This modulation achieves a similar loss landscape to original LSTR.

**Post-processing.** This method requires no post-processing. **Data Augmentation.** Data augmentation includes Poly-LaneNet's (Appendix B.5), then appends random color distortions (brightness, contrast, saturation, hue) and random lighting by a light source calculated from the COCO dataset [17]. That is by far the most complex data augmentation pipeline in this research field, we have validated that all components of this pipeline helps LSTR training.

**Remark.** The polynomial coefficients of LSTR are unbounded, which leads to numerical instability (while the bipartite matching requires precision), and high failure rate of training. The failure rate of fp32 training on CULane is  $\sim 30\%$ . This is circumvented in BézierLaneNet, since our L1 loss can be bounded to [0, 1] without influence on learning (control points easily converges to on-image).

#### **B.8.** BézierLaneNet

BézierLaneNet is implemented in the same code framework where we re-implemented other methods. Same as LSTR, the default prediction threshold is set to 0.95, while 0.5 is used for TuSimple [1].

**Post-processing.** This method requires no post-processing. **Data Augmentation.** We use augmentations similar to LSTR (Appendix B.7). Concretely, we remove the random lighting from LSTR (to strictly avoid using knowledge from external data), and replace the PolyLaneNet  $\frac{10}{11}$  chance augmentations with random affine transforms and random horizontal flip, like LaneATT (Appendix B.6). The random affine parameters are: rotation (10 degrees), translation (maximum 50 pixels on X, 20 on Y), scale (maximum 20%).

**Polynomial Ablations.** For the polynomial ablations (Table 7), we modified the network to predict 6 coefficients for 3rd order Polynomial (4 curve coefficients and start/end Y coordinates). Extra *L*1 losses are added for the start/end Y coordinates similar to LSTR [19]. With extensive tryouts (ad-

justing learning rate, loss weightings, number of epochs), even at the full BézierLaneNet setup, with 150 epochs on CULane, the models still can not converge to a good enough solution. In other word, not precise enough to pass the CULane metric. The sampling loss on polynomial curves can only get to 0.02, which means  $0.02 \times 1640$  pixels = 32.8 pixels average X coordinate error on training set. CU-Lane requires a 0.5 IoU between curves, which are enlarged to 30 pixels wide, thus at least around 10 pixels average error is needed to get meaningful results. By loosen up the IoU requirement to 0.3, we can get F1 score 15.82 for "3rd Polynomial from BézierLaneNet". Although the reviewing committee suggested adding simple regularization for this ablation to converge, regretfully we failed to do this.

### **C. Bézier Curve Implementation Details**

Fast Sampling. The sampling of Bézier curves may seem tiresome due to the complex Bernstein basis polynomials. To fast sample a Bézier curve by a series of fixed t values, simply pre-compute the results from Bernstein basis polynomials, thus only one simple matrix multiplication is left. Remarks on GT Generation. The ground truth of Bézier curves are generated with least squares fitting, a common technique for polynomials. We use it for its simplicity and the fact that it already shows near-perfect lane line fitting ability (99.996 and 99.72 F1 score on CULane test and LLAMAS val, respectively). However, it is not an ideal algorithm for parametric curves. There is a whole research field for fitting Bézier curves better than least squares [23]. Bézier Curve Transform. Another implementation difficulty on Bézier curves is how to apply affine transform (for transforming ground truth curves in data augmentation). Mathematically, affine transform on the control points is equivalent to affine transform on the entire curve. However, translation or rotation can move control points out of the image. In this case, a cutting of Bézier curves is required. The classical De Casteljau's algorithm is used for cutting an on-image Bézier curve segment. Assume a continuous onimage segment, valid sample points with minimum boundary  $t = t_0$ , maximum boundary  $t = t_1$ . The formula to cut a cubic Bézier curve defined by control points  $\mathcal{P}_0, \mathcal{P}_1, \mathcal{P}_2, \mathcal{P}_3$ to its on-image segment  $\mathcal{P}'_0, \mathcal{P}'_1, \mathcal{P}'_2, \mathcal{P}'_3$ , is derived as:

$$\begin{aligned} \mathcal{P}'_{0} &= u_{0}u_{0}u_{0}\mathcal{P}_{0} + (t_{0}u_{0}u_{0} + u_{0}t_{0}u_{0} + u_{0}u_{0}t_{0})\mathcal{P}_{1} \\ &+ (t_{0}t_{0}u_{0} + u_{0}t_{0}t_{0} + t_{0}u_{0}t_{0})\mathcal{P}_{2} + t_{0}t_{0}t_{0}\mathcal{P}_{3}, \\ \mathcal{P}'_{1} &= u_{0}u_{0}u_{1}\mathcal{P}_{0} + (t_{0}u_{0}u_{1} + u_{0}t_{0}u_{1} + u_{0}u_{0}t_{1})\mathcal{P}_{1} \\ &+ (t_{0}t_{0}u_{1} + u_{0}t_{0}t_{1} + t_{0}u_{0}t_{1})\mathcal{P}_{2} + t_{0}t_{0}t_{1}\mathcal{P}_{3}, \\ \mathcal{P}'_{2} &= u_{0}u_{1}u_{1}\mathcal{P}_{0} + (t_{0}u_{1}u_{1} + u_{0}t_{1}u_{1} + u_{0}u_{1}t_{1})\mathcal{P}_{1} \\ &+ (t_{0}t_{1}u_{1} + u_{0}t_{1}t_{1} + t_{0}u_{1}t_{1})\mathcal{P}_{2} + t_{0}t_{1}t_{1}\mathcal{P}_{3}, \\ \mathcal{P}'_{3} &= u_{1}u_{1}u_{1}\mathcal{P}_{0} + (t_{1}u_{1}u_{1} + u_{1}t_{1}u_{1} + u_{1}u_{1}t_{1})\mathcal{P}_{1} \\ &+ (t_{1}t_{1}u_{1} + u_{1}t_{1}t_{1} + t_{1}u_{1}t_{1})\mathcal{P}_{2} + t_{1}t_{1}t_{1}\mathcal{P}_{3}, \end{aligned}$$
(10)

where  $u_0 = 1 - t_0$ ,  $u_1 = 1 - t_1$ . This formula can be efficiently implemented by matrix multiplication. The possibility of noncontinuous cubic Bézier segment on lane detection datasets is extremely low and thus ignored for simplicity. If it does happen, Equation (10) will not change the curve, while our network can also predict out-of-image control points, which still fit the on-image lane segments.

### **D.** IoU Loss for Bézier Curves

Here we briefly introduce how we formulated the IoU loss between Bézier curves. Before diving into the algorithm, there are two preliminaries.

- Polar sort: By anchoring on an arbitrary point inside the N-sided polygon with vertices  $c_i(x_i, y_i)_{i=1}^N$  (normally the mean coordinate between vertices  $c' = (\frac{1}{N} \sum_{i=1}^N x_i, \frac{1}{N} \sum_{i=1}^N y_i)$ ), vertices are sorted by its *atan2* angles. This will return a clockwise or counterclockwise polygon.
- Convex polygon area: A sorted convex polygon can be efficiently cut into consecutive triangles by simple indexing operations. The convex polygon area is the sum of these triangles. The area S of triangle  $((x_1, y_1), (x_2, y_2), (x_3, y_3))$  is:  $S = \frac{1}{2}|x_1(y_2 - y_3) + x_2(y_3 - y_1) + x_3(y_1 - y_2)|$ .

Assume we have two convex hulls from Bézier curves (there are a lot of convex hull algorithms). Now the IoU between Bézier curves are converted to IoU between convex polygons. Based on the simple fact that the intersection of convex polygons is still a convex polygon, after polar sorting all the convex hulls and determining the intersected polygon, we can easily formulate IoU calculations as a series of convex polygon area calculations. The difficulty lies in how to efficiently determine the intersection between convex polygon pairs.

Consider two intersected convex polygons, their intersection includes two types of vertices:

- Intersections: intersection points between edges.
- Insiders: vertices inside/on both polygons.

For Intersections, we first represent every polygon edge as the general line equation: ax + by = c. Then, for line  $a_1x + b_1y = c_1$  and line  $a_2x + b_2y = c_2$ , the intersection (x', y') is calculated by:

$$x' = (b_2c_1 - b_1c_2)/det y' = (a_1c_2 - a_2c_1)/det,$$
(11)

where  $det = a_1b_2 - a_2b_1$ . All (x', y') that is on the respective line segments are Intersections.

For Insiders, there is a certain definition:

**Def. 1** For a convex polygon, point P(x, y) on the same side of each edge is inside the polygon.

A sorted convex polygon is a series of edges (line segments defined by  $P_0(x_0, y_0), P_1(x_1, y_1)$ ), the equation to decide which side a point is to a line segment is as follows:

$$sign = (y - y_0)(x_1 - x_0) - (x - x_0)(y_1 - y_0).$$
 (12)

sign > 0 means P is on the right side, sign < 0 is the left side, and sign = 0 means P is on the line segment. Note that equality is not a stable operation for float computations. But there are simple ways to circumvent that in coding, which we will not elaborate here.

There are other ways to determine Intersections and Insiders, but the above formulas can be efficiently implemented with matrix operations and indexing, making it possible to quickly train networks with batched inputs.

Finally, after being able to compute convex polygon intersections and areas, the Generalized IoU loss (GIoU) is simply (as in [29]):

input : Two arbitrary convex shapes:  $A, B \subseteq \mathbb{S} \in \mathbb{R}^n$  output: GIoU

1. For A and B, find the smallest enclosing convex object C where  $C \subseteq \mathbb{S} \subset \mathbb{P}^n$ 

C, where  $C \subseteq \mathbb{S} \in \mathbb{R}^n$ 2.  $IoU = \frac{|A \cap B|}{|A \cup B|}$ 3.  $GIoU = IoU - \frac{|C \setminus (A \cup B)|}{|C|}$ 

Union is computed as  $A \cup B = A + B - A \cap B$ . The enclosing convex object C can be computed as the convex hull of two convex polygons, or upper-bounded by a enclosing rectangle. We implement the IoU computation purely in PyTorch [24], the runtime for our implementation is only about  $5 \times$  the runtime of rectangle IoU loss computation.

However, lane lines are mostly straight based on road design regulations [7, 36]. This leads to extremely small convex hull area for Bézier curves, thus introduces numerical instabilities in optimization. Although succeeded in a toy polygon fitting experiment, we currently failed to observe the loss's convergence to help learning on lane datasets.

### E. GT and Prediction Matching Prior



Figure 6. Logits activation statistics  $(1 \times \frac{W}{16})$  on CULane [22].

Instead of the *centerness* prior, we explore a local maximum prior, *i.e.*, restricts matched prediction to have a local

maximum classification logit. This prior can facilitate the model to understand the spatially sparse structure of lane lines. As shown in Figure 6, the learned feature activation for classification logits exhibits a similar structure as an actual driving scene.

# F. Extra Results

	TuSimple [1]	LLAMAS [3]
Bézier Baseline	93.36	95.27
+ Feature Flip Fusion	95.26(+1.90)	96.00(+0.73)

Table 9. Ablation study on TuSimple (*test* set Accuracy) and LLA-MAS (*val* set F1), before and after adding the Feature Flip Fusion module. Reported 3-times average with the ResNet-34 backbone, since ablations often are not stable enough on these datasets to exhibit a clear difference between methods.

#### **G.** Discussions

There exists a primitive application of lane detectors from lateral views to estimate the distance to the border of the drivable area [10], which contradicts the use of feature flip fusion. In this case, possibly a lower order Bézier curve baseline (with row-wise instead of column-wise pooling) would suffice. This is out of the focus of this paper.

**Recent Progress.** Recently, others have explored alternative lane representation or formulation methods that do not fully fit in the three categories (segmentation, point detection, curve). Instead of the popular top-down regime, [27] propose a bottom-up approach that focus on local details. [18] achieve state-of-the-art performance, but the complex conditional decoding of lane lines results in unstable runtime depending on the input image, which is not desirable for a real-time system.

## **H.** Qualitative Results

Qualitative results are shown in Figure 7, from our ResNet-34 backbone models. For each dataset, 4 results are shown in two rows: first row shows qualitative successful predictions; second row shows typical failure cases.

**TuSimple.** As shown in Figure 7(a), our model fits highway curves well, only slight errors are seen on the far side where image details are destroyed by projection. Our typical failure case is a high FP rate, mostly attributed to the use of low threshold (Appendix B.8). However, in the bottomright wide road scene, our FP prediction is actually a meaningful lane line that is ignored in center line annotations.

**CULane.** As shown in Figure 7(b), most lanes in this dataset are straight. Our model can make accurate predictions under heavy congestion (top-left) and shadows (top-right, shadow cast by trees). A typical failure case is inaccurate prediction under occlusion (second row), in these cases



(a) TuSimple [1].



(b) CULane [22].



(c) LLAMAS [3].

Figure 7. Qualitative results from BézierLaneNet (ResNet-34) on *val* sets. False Positives (FP) are marked by red, True Positives (TP) are marked by green, ground truth are drawn in blue. Blue lines that are barely visible are precisely covered by green lines. Bézier curve control points are marked with solid circles. Images are slightly resized for alignment. Best viewed in color, in  $2 \times$  scale.

one often cannot visually tell which one is better (ground truth or our FP prediction).

**LLAMAS.** As shown in Figure 7(c), our method performs accurate for clear straight-lines (top-left), and also good for large curvatures in a challenging scene almost entirely cov-

ered by shadow. In bottom-left image, our model fails in a low-illumination, tainted road. While in the other lowillumination scene (bottom-right), the unsupervised annotation from LIDAR and HD-map is misled by the white arrow (see the zigzag shape of the right-most blue line).