

Semi-Local 3D Lane Detection and Uncertainty Estimation

Netalee Efrat, Max Bluvstein, Noa Garnett, Dan Levi
Shaul Oron, Bat El Shlomo

General Motors
Advanced Technical Center Israel

Abstract. We propose a novel camera-based DNN method for 3D lane detection with uncertainty estimation. Our method is based on a semi-local, BEV, tile representation that breaks down lanes into simple lane segments. It combines learning a parametric model for the segments along with a deep feature embedding that is then used to cluster segment together into full lanes. This combination allows our method to generalize to complex lane topologies, curvatures and surface geometries. Additionally, our method is the first to output a learning based uncertainty estimation for the lane detection task. The efficacy of our method is demonstrated in extensive experiments achieving state-of-the-art results for camera-based 3D lane detection, while also showing our ability to generalize to complex topologies, curvatures and road geometries as well as to different cameras. We also demonstrate how our uncertainty estimation aligns with the empirical error statistics indicating that it is well calibrated and truly reflects the detection noise.

Keywords: autonomous driving, lane detection, uncertainty prediction

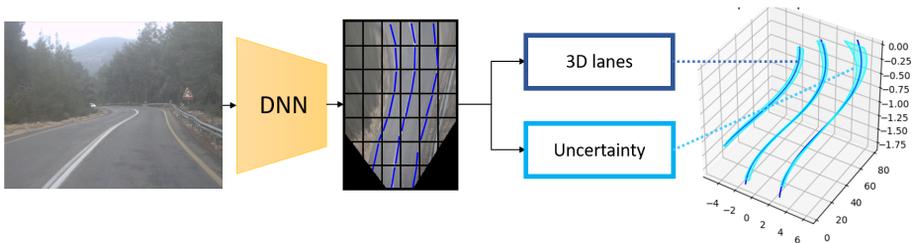


Fig. 1. Illustration of our camera-based 3D lane detection with uncertainty estimation network. Our method works in Bird Eye View perspective, which is rasterized to a coarse tiles grid. We output parametric 3D curve representations for all tiles, which are then processed to form entire 3D lane curves together with detection uncertainty estimates. See text for more details.

1 Introduction

Camera-based lane detection plays a central role in the semantic understanding of the world around a vehicle and can be used for many tasks including trajectory planning, lane keeping, vehicle localization and map generation.

Driving related applications in general and specifically autonomous driving require 3D lane detection with uncertainty estimation that can generalize well to all kinds of lane topologies (e.g. splits, merges, etc.), curvatures, and complex road surface geometries. In addition, as autonomous driving is a safety-critical system, it depends upon reliable estimation of its detection noise, in our case, lane position uncertainty. This uncertainty allows downstream modules like localization and planning to be robust to errors by weighing in the uncertainty when using lane information.

Despite the need for 3D lane detection, most existing methods [1–5] focus on 2D lane detection in the image plane. Additionally, existing methods typically support a limited number of lane topologies, mainly lane lines that are parallel to the vehicle direction of travel. Important topologies required for many driving scenarios such as splits, merges and intersection are for the most part not supported and disregarded. Another aspect not addressed in previous work is ensuring the lane detection system provides uncertainty estimates for its outputs. Recent work for object detection [6–8] suggest novel learning based methods for objects uncertainty estimation. However, to the best of our knowledge, none address a learning based solution for lane detection uncertainty.

In this work, we introduce a novel 3D lane representation and detection framework that is capable of detecting lanes, together with position uncertainty, for any arbitrary topology including splits, merges and lanes perpendicular to the vehicle travel direction. Our method generalizes to different road surface geometries and curvatures, as well as to different cameras.

Key to our solution is a compact semi-local representation that is able to capture local topology-invariant lane structures and road surface geometries. Our lane detection is done in Bird’s Eye View (BEV) which is divided into a regular grid of non-overlapping coarse tiles, as illustrated in Fig. 1. We assume lane segments passing through the tiles are simple and can be represented by a low dimensional parametric model. Specifically, each tile holds a line segment parameterized by an offset from the tile center, an orientation and a height offset from the BEV plane. This semi-local tile representation lies on the continuum between global representation (entire lane) to a local one (pixel level). Each tile output is more informative than a single pixel in a segmentation based solution as it is able to reason on the local lane structure but it is not as constrained as the global solution which has to capture together the complexity of the entire lane topology, curvature and surface geometry.

Our representation breaks down lane curves into multiple lane segments but does not explicitly capture any relation between them. Adjacent tiles will have overlapping receptive fields, and thus correlated results, but the fact that several tiles represent the same lane entity is not captured. In order to generate full lane curves we learn an embedding for each tile which is globally consistent across

the lane. This enables clustering small lane segments into full curves. As we show in our experiments (Sec. 5), the combination between the semi-local tile representation and the embedding based clustering allows the network to output full 3D lane curves of any topology or surface geometry.

Another key component of our method is its ability to provide a noise estimate for the detected lane positions. This uncertainty estimation is achieved by modeling the network outputs as Gaussian distributions and estimating their mean and variance values. This is done for each lane segment parameter and then combined together to produce the final *Covariance* matrix for each lane point. Unlike the segment parameters that can be learned locally across tiles, the empirical errors required for training the uncertainty depend on all the tiles composing an entire lane and have to be reasoned globally as will be further explained in Sec. 3.3 and shown in our experiments. To the best of our knowledge, this is the first learning based uncertainty estimation method for lane detection.

We run extensive experiments, using three datasets, that show our method improves the average precision (*AP*) over the current 3D camera-based state-of-the-art 3D-LaneNet [9] by large margins. We demonstrate qualitatively and quantitatively the efficacy of our learning based clustering, and our method generalization to new lane curvatures and surface geometries as well as new cameras and unseen data. Finally, we present our learning based lane position uncertainty results, and show that it can properly capture the statistics of the actual error of our predicted lanes.

To summarize, the main contribution of our work is twofold: (a) We present a novel 3D semi-local lane representation and detection framework that generalizes to arbitrary lane topologies, curvatures and road surface geometries as well as different camera setups. (b) We propose the first learning based method to provide position uncertainty estimation for the lane detection task.

2 Related work

2D lane detection Most existing lane detection methods focus on lane detection in the image plane and are mostly limited to parallel lane topologies. The literature is vast and includes methods performing 2D lane detection by using self attention [3], employing GANs [4], using new convolution layers [2], exploit vanishing points to guide the training [5] or use differentiable least-squares fitting [10]. Most related to ours is the method of [11] that uses a grid based representation in the image plane, with a line parametrizations and density based spatial clustering for highway lane detection. Our approach uses BEV and a different parametrization than [11] and performs 3D lane detection. We also use learning based clustering as well as output uncertainty estimates for detected lanes. Another work related to ours is [1] that uses learned embedding to perform lane clustering. While [1] perform segmentation at the image pixel level, we cluster the lane segments in BEV on the semi-local tile scale, which is far less computationally expensive.

3D lane detection Detecting lanes in 3D is a challenging task drawing increasing attention in recent years. 3D lane detection methods can be roughly divided to LiDAR-based methods, camera-based and hybrid-methods using both like in [12]. In that work a CNN uses LiDAR to estimate road surface height and then re-projects the camera to BEV accordingly. The network doesn't detect lane instances end-to-end, but rather outputs a dense detection probability map that needs to be further processed and clustered. More related to ours are camera-based methods. The DeepLanes method [13] uses a BEV representation but works with top-viewing cameras that only detect lanes in the immediate surrounding of the vehicle without providing height information. Another work closely related to ours is 3D-LaneNet [9] which also performs camera-based 3D lane detection using BEV. However, unlike our semi-local representation they use a global description that relies on strong assumptions regarding lane geometry, and therefore are unable to detect lanes which are not roughly parallel to ego vehicle direction, lanes starting further ahead, and other non-trivial topologies as will be shown in our experiments (Sec. 5).

Uncertainty estimation Despite its importance, uncertainty estimation for lane detection is not addressed in the literature. We therefore review work done on uncertainty estimation for object detection and classification. To estimate the prediction uncertainty during inference, the machine learning module should output a full distribution over the target domain. Among the available approaches are Bayesian neural networks [14, 15], ensembles [16] and outputting a parametric distribution *directly* [8, 17]. In addition, since uncertainty estimates rely on observed errors on the *training set*, they are often underestimated on the test set and require post training re-calibration. In the context of on-road perception, several works estimate uncertainty in object localization [8, 18] but as far as we know, no previous work applies such techniques to lane detection. In this work we follow [8] which provides a practical solution to both uncertainty estimation and re-calibration. Moreover, as further discusses in Sec. 3.3, uncertainty estimation for general curves requires additional reasoning to that of object localization, such that properly reflects the error of each locally estimated segment with its associated global lane entity.

3 Lane detection and uncertainty estimation

We now describe our 3D lane detection and uncertainty estimation framework. A schematic overview appears in Fig. 2. We first present our semi-local tile representation and lane segment parameterization (Sec. 3.1) followed by how lane segments are clustered together using a learned embedding (Sec. 3.2). Next, we discuss how uncertainty is estimated and calibrated (Sec. 3.3) and finally how the lane structure is inferred from the network's output (Sec. 3.4).

3.1 Learning 3D lane segments with Semi-local tile representation

Lane curves have many different global topologies and lie on road surfaces with complex geometries. This makes reasoning for entire 3D lane curves a very chal-

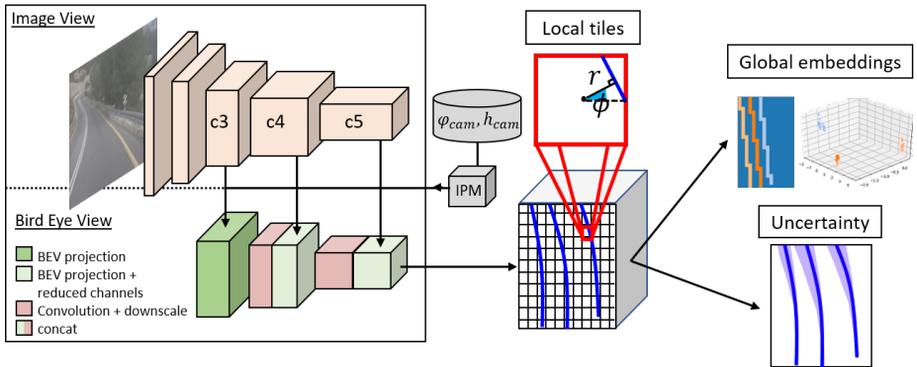


Fig. 2. Method overview. Our network is comprised of two processing pipelines: in image view (top) and in Bird Eye View (bottom). The image view encoder is composed of resnet blocks each one multiplying the number of channels. The BEV backbone is comprised of projected image view feature maps which are concatenated with the convoluted projected feature map from the former block. The final decimated BEV feature map is the input to the lane prediction head which outputs local lane segments, global embedding for clustering the segments to entire lanes, and lane point position uncertainty which relies both on the local tiles and on the entire lane curves.

lensing task. Our key observation is that despite this global complexity, on a local level, lane segments can be represented by low dimensional parametric models. Taking advantage of this observation we propose a semi-local representation that allows our network to learn local lane segments thus generalizes well to unseen lane topologies, curvatures and surface geometries.

The input to our network is a single camera image. We adopted the dual pathway backbone proposed by Garnett et al. [9] which uses an encoder and an Inverse Perspective Mapping (IPM) module to project feature maps to Bird Eye View (BEV). The projection applies a homography, defined by camera pitch angle φ_{cam} and height h_{cam} , that maps the image plane to the road plane (see Fig. 3). The final decimated BEV feature map is spatially divided into a grid $G_{W \times H}$ comprised of $W \times H$ non-overlapping tiles. Similar to [9], the projection ensures each pixel in the BEV feature map corresponds to a predefined position on the road, independent of camera intrinsics and pose.

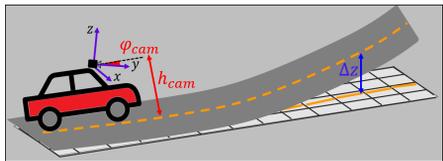


Fig. 3. The road projection plane is defined according to the camera mounting pitch angle φ_{cam} and height h_{cam} , hence our representation is invariant to the camera extrinsics. We represent the GT lanes in full 3D relatively to that plane.

We assume that through each tile $g_{ij} \in G_{W \times H}$ can pass a single line segment which can be approximated by a straight line. Specifically, the network regresses, per each tile g_{ij} , three parameters: lateral offset distance relative to tile center

\tilde{r}_{ij} , line angle $\tilde{\phi}_{ij}$, (see Local tiles in Fig. 2) and height offset $\tilde{\Delta z}_{ij}$ (see Fig. 3). In addition to these line parameters, the network also predicts a binary classification score \tilde{c}_{ij} indicating the probability that a lane intersects a particular tile. GT regression targets for the offsets and angles are calculated by approximating the lane segments intersecting the tiles to straight lines using the GT lane points after they were projected to the road plane (Fig. 3).

Position and z offsets are trained using an $L1$ loss:

$$\mathcal{L}_{ij}^{Offsets} = \|\tilde{r}_{ij} - r_{ij}\|_1 + \|\tilde{\Delta z}_{ij} - \Delta z_{ij}\|_1 \quad (1)$$

Predicting the line angle $\tilde{\phi}_{ij}$ is done using the hybrid classification-regression framework of [19] in which we classify the angle ϕ (omitting tile indexing for brevity) to be in one of N_α bins, centered at $\alpha = \{\frac{2\pi}{N_\alpha} \cdot i\}_{i=1}^{N_\alpha}$. In addition, we regress a vector Δ^α , corresponding to the residual offset relative to each bin center. Our angle bin estimation is optimized using a soft multi-label objective, and the GT probabilities are calculated as $p^\alpha = [1 - |\frac{2\pi}{N_\alpha} \cdot i - \phi| / \frac{2\pi}{N_\alpha}]_+$. The GT offsets Δ^α are the difference between the GT angle and the bin centers, and their training is supervised on the GT angle bin and adjacent bins to ensure that the delta offset can account for erroneous bin class prediction.

The angle loss is the sum of the classification and offset regression losses:

$$\mathcal{L}_{ij}^{angle} = \sum_{\alpha=1}^{N_\alpha} [p_{ij}^\alpha \cdot \log \tilde{p}_{ij}^\alpha + (1 - p_{ij}^\alpha) \cdot \log (1 - \tilde{p}_{ij}^\alpha) + \delta_{ij}^\alpha \cdot \|\tilde{\Delta}_{ij}^\alpha - \Delta_{ij}^\alpha\|_1] \quad (2)$$

where δ_{ij}^α is the indicator function masking the relevant bins for the offset learning.

The lane tile probability \tilde{c}_{ij} is trained using a binary cross entropy loss:

$$\mathcal{L}_{ij}^{score} = c_{ij} \cdot \log \tilde{c}_{ij} + (1 - c_{ij}) \cdot \log (1 - \tilde{c}_{ij}) \quad (3)$$

Finally, the overall tile loss is the sum over all the tiles in the BEV grid:

$$\mathcal{L}^{tiles} = \sum_{i,j \in W \times H} (\mathcal{L}_{ij}^{score} + c_{ij} \cdot \mathcal{L}_{ij}^{angle} + c_{ij} \cdot \mathcal{L}_{ij}^{offsets}) \quad (4)$$

3.2 Global embedding for lane curve clustering

In order to provide complete lane curves we need to cluster together multiple lane segments into complete lane entities. To this end, we learn an embedding vector f_{ij} for each tile such that vectors representing tiles belonging to the same lane would reside close in embedded space while vectors representing tiles of different lanes would reside far apart. For this we adopted the approach of [1,20], and use a discriminative push-pull loss. Unlike previous work, we use the discriminative loss on the decimated tiles grid, which requires far less computations than operating at the pixel level.

The discriminative push-pull loss is a combination of two losses:

$$\mathcal{L}^{embedding} = \mathcal{L}^{pull} + \mathcal{L}^{push} \quad (5)$$

A pull loss aimed at pulling the embeddings of the same lane tiles closer together:

$$\mathcal{L}^{pull} = \frac{1}{C} \sum_{c=1}^C \frac{1}{N_c} \sum_{ij \in W \times H} [\delta_{ij}^c \cdot \|\mu_c - f_{ij}\| - \Delta_{pull}]_+^2 \quad (6)$$

and a push loss aimed at pushing the embedding of tiles belonging to different lanes farther apart:

$$\mathcal{L}^{push} = \frac{1}{C(C-1)} \sum_{c_A=1}^C \sum_{c_B=1, c_B \neq c_A}^C [\Delta_{push} - \|\mu_{c_A} - \mu_{c_B}\|]_+^2 \quad (7)$$

where C is the number of lanes (can vary), N_c is the number of tiles belonging to lane c , δ_{ij}^c indicates if tile i, j belongs to lane c , $\mu_c = \frac{1}{N_c} \sum_{ij \in W \times H} \delta_{ij}^c \cdot f_{ij}$ is the average of f_{ij} belonging to lane c , Δ_{pull} constraints the maximal inter-cluster distance and Δ_{push} is the intra-cluster minimal required distance.

Given the learnt feature embedding we can use a simple clustering algorithm to extract the tiles that belong to individual lanes. We adopted the clustering methodology from Neven et al. [1] which uses mean-shift to find the clusters centers and set a threshold around each center to get the cluster members. We set the threshold to $\frac{\Delta_{push}}{2}$.

3.3 Uncertainty estimation

We now explain how we estimate the noise (uncertainty) of our lane detector. As it is a statistical property, its estimation is done by casting the tile prediction problem as a distribution estimation task. This means that we formulate each one of the lane segment parameters (omitting tile indexing for brevity) $y \in \{r, \phi, \Delta z\}$ as a Gaussian distribution such that

$$y|x, \theta \sim \mathcal{N}(\mu(x, \theta), \sigma^2(x, \theta)) \quad (8)$$

where x is the network input and θ the network parameters. The mean values of the above distributions are the predicted values for the tile parameters i.e. $\mu_r(x, \theta) = \tilde{r}$, $\mu_\phi(x, \theta) = \tilde{\phi}$, $\mu_{\Delta z}(x, \theta) = \tilde{\Delta z}$, estimated using the methodology described in Sec. 3.1. In this section we focus on the estimation of the variances, given the predicted mean values, as a second training stage. The variances $\sigma^2(x, \theta)$ are estimated through the optimization of the Negative Log Likelihood (NLL) which is a standard measure for probabilistic models quality [21, 22]

$$NLL = -\log p(y|x, \theta) = \frac{1}{2} \log \sigma^2(x, \theta) + \frac{(y - \mu(x, \theta))^2}{2\sigma^2(x, \theta)} + const. \quad (9)$$

where y is the GT value and $(y - \mu(x, \theta))^2$ is the empirical Squared Error (SE). Proper evaluation of the empirical SE is therefore a key component in quantifying the uncertainty. Measuring the SE for the predicted tiles is not trivial because it is not obvious which is the corresponding GT segment the error should be calculated with respect to. This correspondence depends not only on the predicted tile itself, but on the entire curve it belongs to. Therefore, we propose measuring the SE in a global context of full lane curves as illustrated in Fig. 4b. Alternative solution for measuring the SE would be to simply calculate the error on the same tiles supervised for the tile parameters prediction, i.e., the tiles for which $c_{ij} = 1$ (Fig. 4a). Using this solution, the errors originate only from the semi-local tile context and essentially bounded to the tile size thus would generate a skewed statistic. Obtaining the global-context error values requires that we first cluster lane segments together into full lane curves and then associate them to GT lanes. Once this is performed we can find, for each predicted lane segment, a corresponding GT segment *in the context of the full associated GT lane* that can now be far from its semi-local context. This makes the uncertainty training a multi-stage process that requires first inferring lane segment parameters on the tile level $\mu(x, \theta)$, then clustering tiles together into full lanes, associating these lanes to GT lanes, computing the SE of the lane segments with respect to the associated GT segments, and finally using these SE values to compute the NLL loss and update the network parameters.

Despite the supervision involved in the training process, deep neural networks tend to produce over confident uncertainty predictions (high probabilities, and small variances) [8, 23]. It is therefore useful to further calibrate¹ the uncertainty estimation of the network after it is trained. For this we adopt the Temperature Scaling [8, 23] procedure and use it to calibrate the tiles classification scores and regressed offsets and angles uncertainties. Temperature scaling uses a single scalar parameter T per parameter ($r, \phi, \Delta z$) that multiplies the estimated variance such that $\sigma_T^2(x, \theta) = T \cdot \sigma^2(x, \theta)$. During the calibration training the NLL is optimized with $\sigma_T^2(x, \theta)$ but only T is updated. The training is performed on a different train set and the parameter T adjusts the learnt variances to better capture the statistics on the new dataset.

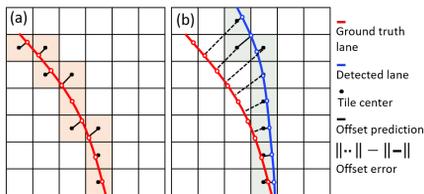


Fig. 4. Error calculation for uncertainty prediction supervision. (a) The orange tiles are occupied by the ground truth segments. On these tiles we supervise the tiles parameters prediction, and they can be easily extended to output those parameters uncertainty based on the within tile error. (b) By calculating the error while considering the actual network predictions (green tiles), together with the entire curve, we can better quantify the detector uncertainty.

¹ Calibrated uncertainty is broadly defined as having the variance reflect the actual MSE error of the predictions in case of regression, and the probability output to match the actual accuracy in case of classification.

3.4 Final output

Our 3D lane detection module outputs lanes represented as sets of 3D directed lane points, where each lane point has an associated 3D covariance matrix indicating its estimated uncertainty, and local direction vector based on $\tilde{\phi}_{ij}$.

Lane segments are clustered together as explained in Sec. 3.2. Each lane segment (tile) contributes a point to the lane point set. We begin by thresholding tile score to output only tiles with lanes. We then convert the offsets ($\tilde{r}_{ij}, \tilde{\Delta z}_{ij}$) and angles ($\tilde{\phi}_{ij}$) to points by converting them from Polar to Cartesian coordinates. Finally, we transform the points from the BEV plane to the camera coordinate frame by subtracting h_{cam} and rotating by $-\varphi_{cam}$ (see Fig. 3).

$$\begin{bmatrix} \tilde{x}_{ij} \\ \tilde{y}_{ij} \\ \tilde{z}_{ij} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(\varphi_{cam}) & \sin(\varphi_{cam}) \\ 0 & -\sin(\varphi_{cam}) & \cos(\varphi_{cam}) \end{bmatrix} \cdot \begin{bmatrix} \tilde{r}_{ij} \cdot \cos(\tilde{\phi}_{ij}) \\ \tilde{r}_{ij} \cdot \sin(\tilde{\phi}_{ij}) \\ \tilde{\Delta z}_{ij} - h_{cam} \end{bmatrix} \quad (10)$$

Although the network learns to predict the variances of offsets and angles independently, we output a full covariance matrix for each predicted lane point in Cartesian coordinates. This is done by transforming the offsets and angle covariance matrix $Cov_{ij}^{pol} = \text{diag}[\tilde{\sigma}_r^2, \tilde{\sigma}_\phi^2, \tilde{\sigma}_{\Delta z}^2]_{ij}$ to Cartesian space:

$$Cov_{ij}^{cart} = J_{ij} \cdot Cov_{ij}^{pol} \cdot J_{ij}^T \quad (11)$$

Where J_{ij} is the Jacobian matrix of the Polar to Cartesian conversion in Eq. (10) approximated for segment i, j .

4 Experimental Setup

We study the performance of our 3D lane detection framework using several 3D-lane datasets, comparing it to [9] which is the current state-of-the-art (SOTA) camera-based 3D-lane detector. We demonstrate the method ability to detect difficult lane topologies and generalize to complex surface geometries and different cameras setups. Finally, we show the accuracy of our uncertainty estimation.

Datasets Evaluation is done using two 3D-lane datasets. The first is *synthetic-3D-lanes* [9] containing synthetic images of complex road geometries with 3D ground truth lane annotations. The second, is a dataset we collected and annotated² referred to as *3D-lanes*. This dataset contains 327K images from 19 distinct recordings (different geographical locations at different times) taken at 20 fps. The data is split such that train and test sets have different distributions. Specifically, the train set, 298K images, is comprised mostly of highway scenarios while the test set is comprised of a rural scenario with complex curvatures and road surface geometries, taken at a geographic location not in the train set. To reduce temporal correlation we sampled every 30th frames giving us a test set

² Annotation protocol is similar to that used in [9].

of 1000 images. We also set aside 5K images for uncertainty calibration training. Example images from the train and test sets can be seen in Fig. 6.

To quantitatively demonstrate our methods ability to generalize to new cameras and scenes we use the tuSimple 2D benchmark [24]. Additional qualitative evaluation using a new camera setup is also shown.

Evaluation We adopt the AP metric commonly used in object detection [25, 26] that averages the area under ROC curves, generated with different IOU thresholds. Unlike object detection, where intersection and union for bounding boxes are well defined, intersection and union for 3D curves are not. Similarly to [27] we define the intersection for curves as the length of the curve sections that are closer than a threshold to the GT curve, and the union as the length of the longer curve out of the two: detected and GT curves. However, unlike [26], we calculate the True Positive (TP), False Positive (FP) and miss detections, for every lane curve, not for every lane point. This gives a better estimate of the number of lanes properly detected, regardless of their length, distance, or topology (merges, splits or intersections). For example, in per-point metrics, detecting half of the points out of two lanes, or detecting only one lane out of the two, would get the same score, whereas per-lane metrics will give these two cases different scores. Note that in order to determine if a certain lane section intersects or not, we have to define a distance threshold. This is a heuristic that does not exist when defining IOU in object detection. To this end we add another set of distance accuracy metrics to account for the location error of each detected lane point with respect to its associated GT curve. We divided the entire dataset to lane points in the near range (0-30m) and far range (30-80m) and calculate the mean absolute lateral error for every range.

Implementation details We use the dual-pathway architecture [9] with a ResNet34 [28] backbone. Our BEV projection covers 20.4m x 80m divided in the last decimated feature map to our tile grid $G_{W \times H}$ with $W = 16, H = 26$ such that each tile represents $1.28m \times 3m$ of road surface. We found that predicting the camera angle φ_{cam} and height h_{cam} gave negligible boost in performance compared to using the fixed mounting parameters on *3D-lanes*, however, on *synthetic-3D-lanes* we followed [9] methodology and trained the network to output φ_{cam} and h_{cam} as well.

The network is trained with batch size 16 using ADAM optimizer, with initial lr of 1e-5 for 80K iterations which is then reduced to 1e-6 for another 50K iterations. We set Δ_{pull} and Δ_{push} (Eqs. 6, 7) to 0.1 and 3 respectively, and used a coarse 0.3 threshold on the output segment scores \tilde{c}_{ij} prior to the clustering. In the evaluation we set a coarse distance threshold for association of 1m, and measured the AP as the average of $AP_{IOU\%}$ at IOU thresholds 0.1 : 0.1 : 0.9.

5 Results

Synthetic-3D-lanes dataset We compared our method to 3D-LaneNet [9]. Results are presented in Table 1. It can be seen that our AP and AP_{50} are far

superior to those of 3D-LaneNet³ while showing comparable lateral error (for $IOU = 0.5$ and $recall = 0.75$). We believe the main reason is our semi-local representation that allows our method to support many different lane topologies such as short lanes, splits and merges that emerge only at a certain distance from the ego vehicle. This is evident in Fig. 5 showing examples where splits and short lanes are not detected by 3D-LaneNet but detected by our method.

Table 1. comparison on *synthetic-3D-lanes*

Method	AP AP_{50}		recall	Lateral error (cm)	
				0-30m	30-80m
3D-LaneNet [9]	0.74	0.79	0.75	9.3	23.9
Ours	0.9	0.95	0.75	9.7	26.7

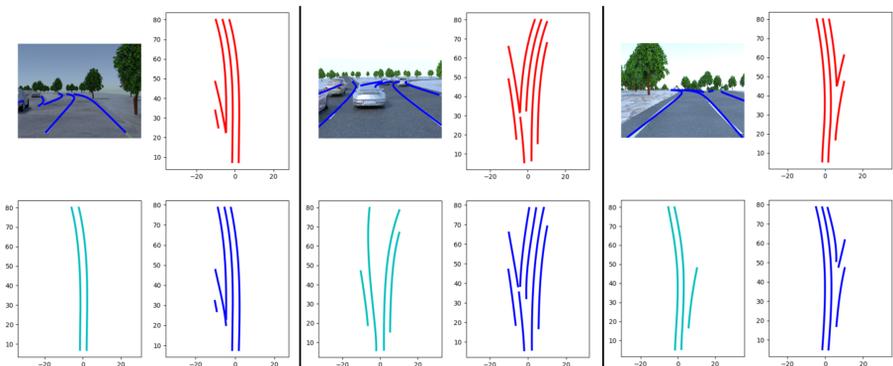


Fig. 5. Example results on *synthetic-3D-lanes*. Our detected lanes (blue), the ground truth (red) and 3d-LaneNet lanes (cyan). It is visible that our method is less constrained and detects all lanes in the scene.

Generalizing to new topologies and geometries We use the *3D-lanes* dataset to demonstrate our methods ability to generalize to new scenes with complex curvatures and surface geometries. Results comparing our method to [9], trained on the same train set, are summarized in Table 2. It can be seen that our method achieves better results improving by 9 points over 3D-LaneNet in overall AP as well as lowering the lateral error for the 3d lane points. This experiment is challenging compared to the *synthetic-3D-lanes* experiment in which train and test sets have the same distribution. In the case of *3D-lanes*, train and test sets have much different distribution with the test set exhibiting much more complex curvatures and surface geometries as also shown in the examples in Fig. 6. We believe our ability to generalize to this test set demonstrates the advantage of using the proposed semi-local tiles representation.

³ Results reported here differ than those in [9] since their evaluation disregards short lanes that start beyond 20m from the ego vehicle.

Table 2. comparison on *3D-lanes*

Method	AP	AP_{50}	AP_{90}	Recall	Lateral error (cm)	
					0-30m	30-80m
3D-LaneNet [9]	0.80	0.86	0.48	0.85	15.6	47.2
Ours - w/o global	0.84	0.94	0.43	0.85	14.5	45.5
Ours	0.89	0.95	0.60	0.85	14.1	44.7
Ours - w synthetic	0.9	0.95	0.59	0.85	12.9	36.3

To show the significance of our clustering approach using global feature embedding we compare it with a naïve clustering alternative (Table 2 ‘Ours - w/o global’). This alternative uses a simple greedy algorithm concatenating segments based on continuity and similarity heuristics. We find that in using naïve clustering we loose 5 points in overall AP and 17 in AP_{90} suggesting that detected lanes with greedy clustering are much shorter. In addition, we see that with feature embedding we obtain lower lateral error. This may suggest that feature embedding learning also helps predicting more accurate segments.

We also compare a model trained only on *3D-lanes* with a one trained on both *3D-lanes* and *synthetic-3D-lanes* (Table 2 ‘Ours -w synthetic’). We find that additional 3D training data of complex curvatures and geometries helps in the generalization despite it being synthetic and without using any domain adaptation techniques.

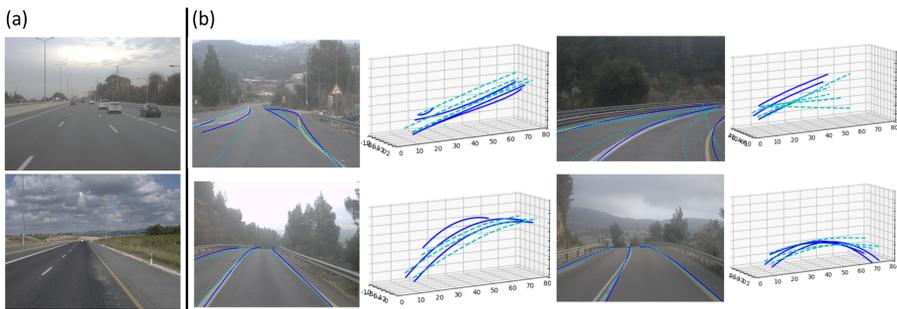


Fig. 6. Example results on *3D-lanes*. (a) Depicts examples from the training set, and (b) examples from the test set. It is clear that the surface geometries and curvatures appear in the test set are different from the train set. Our detected lanes are shown in blue and 3D-LaneNet lanes in cyan.

Generalization to new cameras We now examine our methods generalization to new unseen cameras. To this end, we first use the 2D lanes tuSimple dataset [24] and show that our network, trained on a different task (3D lane detection rather than 2D), and on different data (*synthetic-3D-lanes* and *3D-lanes*), can generalize to a new task on unseen cameras. For inference on tuSimple we use fixed camera angle φ_{cam} and height h_{cam} from [9] to projected the feature maps to BEV. The resulting tuSimple accuracy (acc) metric for this experiment, is 0.912. This result is surprisingly high given that our network is designed for

detecting lanes in 3D, and more importantly, was not trained on a single example from the tuSimple dataset (See Fig. 7a). Encouraged by this result, we next trained our network on the tuSimple dataset, lifting 2D lanes to 3D using a flat world assumption. When the flat world assumption is violated, the lifted 3D lanes no longer have the BEV properties of real 3D lanes, such as parallel lanes of constant distance between them, making this approach more challenging than solving the 2D detection problem directly. Once trained on tuSimple data the network reached an accuracy of 0.956 which is comparable to the 0.966 SOTA results of [3]. We also conduct a qualitative evaluation on another unseen camera using an *internal evaluation dataset* not used in training (See Fig. 7b), Here too, we find good generalization to new cameras and scenarios.

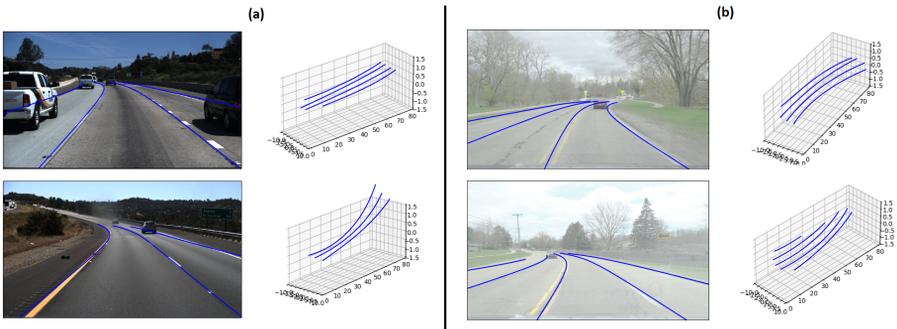


Fig. 7. generalization to new cameras: (a) depicts examples from *tuSimple* test set (b) examples from our *internal evaluation dataset*. Our method generalizes well to new cameras that weren't present in the training set without any need for adaptation.

Uncertainty estimation Our estimated uncertainty (variance) tries to quantify the noise in our detection model. That is, the empirical error between detected lane points and GT. Such noise is statistical by nature, thus can not be evaluate for a single sample (lane point). We do however, expect that the estimated uncertainty *on average* would reflect the average empirical error. Therefore, in order to evaluate the estimated uncertainty we first divide it to bins, and compare the Root Mean estimated Variance (RMV) in each bin to the Root Mean Squared empirical Error (RMSE) of the samples (lane points) in that bin. Equality between the two measures indicates well calibrated uncertainty. This evaluation method is described by Levi et al. [8] which also propose a single figure of merit, the Expected Normalized Calibration Error (ENCE) which averages the error between the RMV and the RMSE in each bin, normalized by the bins RMV. Fig. 8a shows the results for the above evaluation. In this analysis RMV is take as the maximal eigenvalue of our 3D covariance matrix which is reasonable since most of the error is in the lateral direction. We compare it to the empirical lateral RMSE taken between each detected point and the GT lane. We can see that our estimated uncertainty is very close to the ideal uncertainty line achieving an ENCE of 11%.

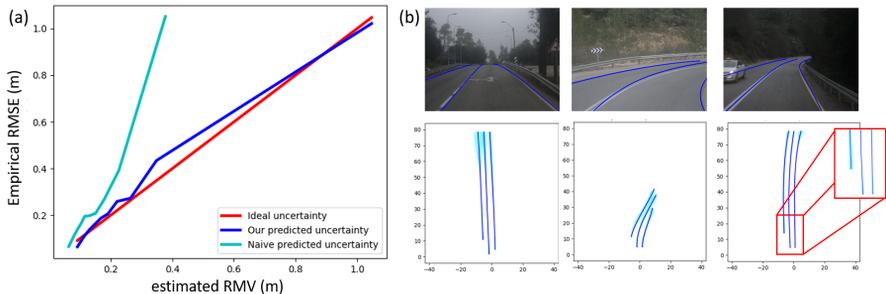


Fig. 8. uncertainty estimation results. (a) Shows comparison between our estimated uncertainty and a simple uncertainty baseline. (b) Examples of our predicted uncertainty. It is visible that the estimated uncertainty is larger on large curvatures (either lane curvature or surface curvature), and on occlusions (see bottom right inset)

In the absence of a previous baseline, we compare ourselves to an uncertainty model similar to ours only it is supervised by errors computed in the tile level, i.e. on the same tiles we train our segment parameters (orange tiles in Fig. 4a). As discussed in Sec. 3.3 we expect these errors to be bounded and generate a skewed distribution. Fig. 8a) demonstrated that indeed this is the case. The maximal estimated RMV reaches $\sim 40\%$ of the maximal empirical RMSE, and the resulting ENCE for this model is 60%. Fig. 8b shows qualitative results of the estimated uncertainty. We can see that the uncertainty captures detection errors occurring at large curvatures (either lane curvature or complex surface geometries), occlusions, or large distance. Note that other approaches based on offline error modeling and rule based look up tables wouldn't necessarily capture all the cases in which the uncertainty should be high. This is in contrast to our data driven approach.

6 Conclusions

We presented a novel 3D lane detection with uncertainty estimation framework. The method uses a semi-local representation that captures topology-invariant lane segments that are then clustered together using a learned global embedding into full lane curves. The efficacy of our approach was showcased in extensive experiments achieving SOTA results while demonstrating its ability to detect globally complex lanes having different topologies and curvatures and generalize well to unseen complex surface geometries and new cameras. In this work we also implemented the first learning based uncertainty estimation for lane detection. We show the importance of properly quantifying the detection errors, achieving almost ideal uncertainty results with respect to the real error statistics. Our work performs full 3D lane detection and uncertainty estimation thus closing the gap towards full lane detection requirements for autonomous driving.

References

1. Neven, D., Brabandere, B., Georgoulis, S., Proesmans, M., Van Gool, L.: Towards end-to-end lane detection: an instance segmentation approach. (06 2018) 286–291
2. Xingang Pan, Jianping Shi, P.L.X.W., Tang, X.: Spatial as deep: Spatial cnn for traffic scene understanding. In: AAAI Conference on Artificial Intelligence (AAAI). (February 2018)
3. Hou, Y., Ma, Z., Liu, C., Loy, C.C.: Learning lightweight lane detection cnns by self attention distillation. In: The IEEE International Conference on Computer Vision (ICCV). (October 2019)
4. Ghafoorian, M., Nugteren, C., Baka, N., Booij, O., Hofmann, M.: El-gan: Embedding loss driven generative adversarial networks for lane detection. In Leal-Taixé, L., Roth, S., eds.: Computer Vision – ECCV 2018 Workshops, Cham, Springer International Publishing (2019) 256–272
5. Lee, S., Kim, J., Shin Yoon, J., Shin, S., Bailo, O., Kim, N., Lee, T.H., Seok Hong, H., Han, S.H., So Kweon, I.: Vpnet: Vanishing point guided network for lane and road marking detection and recognition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 1947–1955
6. Choi, J., Chun, D., Kim, H., Lee, H.J.: Gaussian yolov3: An accurate and fast object detector using localization uncertainty for autonomous driving. In: The IEEE International Conference on Computer Vision (ICCV). (October 2019)
7. Bertoni, L., Kreiss, S., Alahi, A.: Monoloco: Monocular 3d pedestrian localization and uncertainty estimation. In: The IEEE International Conference on Computer Vision (ICCV). (October 2019)
8. Levi, D., Gispan, L., Giladi, N., Fetaya, E.: Evaluating and calibrating uncertainty prediction in regression tasks. arXiv preprint arXiv:1905.11659 (2019)
9. Garnett, N., Cohen, R., Pe’er, T., Lahav, R., Levi, D.: 3d-lanenet: end-to-end 3d multiple lane detection. In: Proceedings of the IEEE International Conference on Computer Vision. (2019) 2921–2930
10. Van Gansbeke, W., De Brabandere, B., Neven, D., Proesmans, M., Van Gool, L.: End-to-end lane detection through differentiable least-squares fitting. arXiv preprint arXiv:1902.00293 (2019)
11. Huval, B., Wang, T., Tandon, S., Kiske, J., Song, W., Pazhayampallil, J., Andriluka, M., Rajpurkar, P., Migimatsu, T., Cheng-Yue, R., Mujica, F.A., Coates, A., Ng, A.Y.: An empirical evaluation of deep learning on highway driving. CoRR **abs/1504.01716** (2015)
12. Bai, M., Mattyus, G., Homayounfar, N., Wang, S., Lakshminanth, S.K., Urtasun, R.: Deep multi-sensor lane detection. In: 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), IEEE (2018) 3102–3109
13. Gurghian, A., Koduri, T., Bailur, S.V., Carey, K.J., Murali, V.N.: Deeplanes: End-to-end lane position estimation using deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. (2016) 38–45
14. Gal, Y.: Uncertainty in deep learning. (2016)
15. Gal, Y., Ghahramani, Z.: Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In: Proceedings of the 33rd International Conference on Machine Learning (ICML-16). (2016)
16. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. In Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R., eds.: Advances in Neural Information Processing Systems 30. Curran Associates, Inc. (2017) 6402–6413

17. Nix, D.A., Weigend, A.S.: Estimating the mean and variance of the target probability distribution. In: Proceedings of 1994 IEEE International Conference on Neural Networks (ICNN'94). Volume 1. (June 1994) 55–60 vol.1
18. Phan, B., Salay, R., Czarnecki, K., Abdelzad, V., Denouden, T., Vernekar, S.: Calibrating uncertainties in object localization task. CoRR **abs/1811.11210** (2018)
19. Mahendran, S., Ali, H., Vidal, R.: A mixed classification-regression framework for 3d pose estimation from 2d images. In: BMVC. (2018)
20. Brabandere, B.D., Neven, D., Gool, L.V.: Semantic instance segmentation with a discriminative loss function. CoRR **abs/1708.02551** (2017)
21. Nix, D.A., Weigend, A.S.: Estimating the mean and variance of the target probability distribution. (1994)
22. Lakshminarayanan, B., Pritzel, A., Blundell, C.: Simple and scalable predictive uncertainty estimation using deep ensembles. In Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R., eds.: Advances in Neural Information Processing Systems 30. Curran Associates, Inc. (2017) 6402–6413
23. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In Precup, D., Teh, Y.W., eds.: Proceedings of the 34th International Conference on Machine Learning. Volume 70 of Proceedings of Machine Learning Research., International Convention Centre, Sydney, Australia, PMLR (06–11 Aug 2017) 1321–1330
24. : The tusimple lane challenge, <http://benchmark.tusimple.ai/>
25. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T., eds.: Computer Vision – ECCV 2014, Cham, Springer International Publishing (2014) 740–755
26. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: Conference on Computer Vision and Pattern Recognition (CVPR). (2012)
27. Homayounfar, N., Ma, W.C., Liang, J., Wu, X., Fan, J., Urtasun, R.: Dagmapper: Learning to map by discovering lane topology. In: The IEEE International Conference on Computer Vision (ICCV). (October 2019)
28. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2016) 770–778