

Rail Detection: An Efficient Row-based Network and A New Benchmark

Xinpeng Li

Shenzhen Technology University

li.xin.peng@outlook.com

Xiaojiang Peng*

Shenzhen Technology University

pengxiaojiang@sztu.edu.cn

ABSTRACT

Rail detection, essential for railroad anomaly detection, aims to identify the railroad region in video frames. Although various studies on rail detection exist, neither an open benchmark nor a high-speed network is available in the community, making algorithm comparison and development difficult. Inspired by the growth of lane detection, we propose a rail database and a row-based rail detection method. In detail, we make several contributions: (i) We present a real-world railway dataset, Rail-DB, with 7432 pairs of images and annotations. The images are collected from different situations in lighting, road structures, and views. The rails are labeled with polylines, and the images are categorized into nine scenes. The Rail-DB is expected to facilitate the improvement of rail detection algorithms. (ii) We present an efficient row-based rail detection method, Rail-Net, containing a lightweight convolutional backbone and an anchor classifier. Specifically, we formulate the process of rail detection as a row-based selecting problem. This strategy reduces the computational cost compared to alternative segmentation methods. (iii) We evaluate the Rail-Net on Rail-DB with extensive experiments, including cross-scene settings and network backbones ranging from ResNet to Vision Transformers. Our method achieves promising performance in terms of both speed and accuracy. Notably, a lightweight version could achieve 92.77% accuracy and 312 frames per second. The Rail-Net outperforms the traditional method by 50.65% and the segmentation one by 5.86%. The database and code are available at: <https://github.com/Sampson-Lee/Rail-Detection>.

CCS CONCEPTS

- Computing methodologies → Computer vision.

KEYWORDS

datasets, neural networks, rail detection, efficient classification

ACM Reference Format:

Xinpeng Li and Xiaojiang Peng. 2022. Rail Detection: An Efficient Row-based Network and A New Benchmark. In *Proceedings of the 30th ACM International Conference on Multimedia (MM '22), October 10–14, 2022, Lisboa, Portugal*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3503161.3548050>

**Corresponding author.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

MM '22, October 10–14, 2022, Lisboa, Portugal

© 2022 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-9203-7/22/10.

<https://doi.org/10.1145/3503161.3548050>

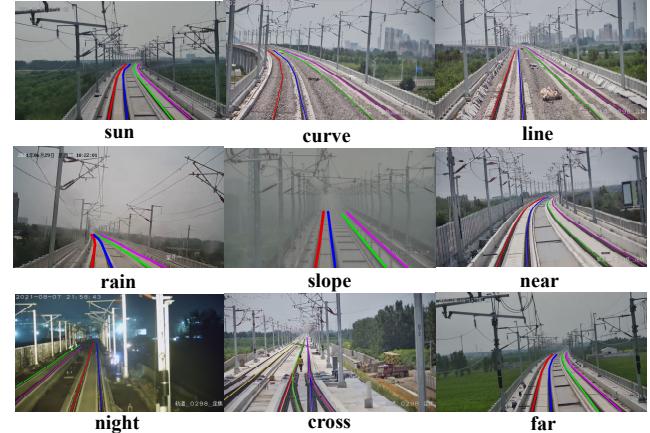


Figure 1: Samples on Rail-DB. The rails are labeled with polylines, where colors encode the rails' order numbers. The situations include different lighting, structures, and views: sun, rain, night, curve, slope, cross, line, near, and far.

1 INTRODUCTION

Rail detection is to identify the railroad region within a video frame. Given an input image, an algorithm outputs a segmentation or anchors to denote a railroad region. Recently, train driving safety has attracted increasing attention in the signal processing and multimedia community, including obstacle detection and damage detection [18, 19, 33]. Identifying the railroad region is vital for the above tasks. The main challenges of rail detection include lighting variety, rail changes, and high-speed requirements [32].

In the past decade, various algorithms have been proposed in rail detection, which can be mainly categorized into 1) traditional hand-crafted rail detection methods [2, 13, 15, 20, 24, 25, 27, 29, 36], and 2) deep learning based rail segmentation methods [7, 32]. The traditional hand-crafted rail detection methods usually first extract linear features of an image and then detect rails. These approaches achieve favorable results in simple conditions. Differently, the deep learning based rail segmentation ones adopt sufficient data and convolution neural networks for detection and segmentation. These methods show robustness in real-world environments.

Although various studies on rail detection exist, neither an open benchmark nor a high-speed network is available in the community. It leads to a significant gap in algorithm comparison and development. Recently, the algorithm of lane detection has been booming partly due to the rich and large databases [1, 3, 16, 34, 35] and high-speed frameworks [26, 28, 30, 37]. This paper proposes an open rail database, Rail-DB, inspired by the lane detection. The Rail-DB consists of 7432 pairs of images in diverse backgrounds

and high-quality annotations. Additionally, we present a row-based network, Rail-Net, to meet high-speed and accuracy requirements.

First, considering existing rail database [32] is unavailable for further research, we construct the Rail-DB and make it publicly available for the community. To build the Rail-DB, we first obtain videos from real-world trains in different conditions and collect 7432 images. Figure 1 illustrates some examples on the Rail-DB. The lighting covers sun, rain, and night. The structures range in curve, slope, cross, and line. The views include near and far. Experiments show that images with diverse environments contribute to the algorithm's robustness. Then, we label rails with polylines through coarse and refinement stages to gather high-quality annotations. At the coarse stage, eight annotators use *LabelMe*¹ to label every rail of an image. However, non-primary rails and railway tails are occasionally missed in the coarse annotations. Therefore, we need a refinement stage to correct the incomplete and inaccurate labels. An expert checks the above annotations and refines them. Besides, we annotate the images with nine scene properties based on backgrounds. Finally, we get the Rail-DB, a middle size, background-rich, and polyline annotated dataset, which will be public available.

Second, we propose the Rail-Net as a rail detection benchmark. Existing works mainly focus on segmentation algorithm [7, 32], which is computationally inefficient. Recently, row-based methods are widespread in efficient lane detection [26, 28, 30, 37], which inspires us to extend them to rail detection. The segmentation methods translate the labeled polylines to a multivalued mask. In contrast, we adopt anchors and gridings to further turn the mask into a series of horizontal locations. The locations are intersections between anchors and polylines, which is shrunk into class units by further gridding. Therefore, the process of rail detection can be treated as a row-based selecting problem. We produce a lower prediction number than outputting a segmentation map. Experiments suggest that the row-based selecting algorithm runs faster than existing segmentation methods while keeping higher accuracy.

Third, we perform extensive experiments on Rail-DB and Rail-Net. The results show that our method achieves promising performance in terms of both speed and accuracy. Notably, a lightweight version could even achieve 92.77% accuracy and 312 frames per second. The Rail-Net outperforms the traditional method by 50.65% and the segmentation one by 5.86%. Besides, ablation studies, including cross-scene setting, gridding number, metric threshold, rail and lane comparison, and various backbones, justify the effectiveness of every setting. Further, we find the cross structure and the cross-scene adaptation are confusing cases.

The contributions can be briefly summarized as follows.:

- 1) We propose the first open rail dataset, Rail-DB. The Rail-DB consists of 7432 pairs of multi-scenario images and high-quality annotations, expected to promote the rail detection.
- 2) We present a novel high-speed rail detection network, Rail-Net. We are the first to extend row-selecting for rail detection where only segmentation methods are available. The Rail-Net is demonstrated to be efficient and significantly better than the existing segmentation method.
- 3) We perform extensive experiments on the Rail-DB and the Rail-Net, including the comparison between rail and lane,

the comparison between segmentation and classification, and cross-scene adaptation. The proposed method achieves 92.3% accuracy and 312 frames per second.

2 RELATED WORK

2.1 Railroad Detection

The various algorithms in rail detection can be mainly categorized into 1) traditional hand-crafted feature based rail detection methods [2, 13, 15, 20, 24, 25, 27, 29, 36], and 2) deep learning based rail segmentation methods [7, 32].

Specifically, [7] performs visual rail inspection using material classification and semantic segmentation with deep convolutional neural networks. [32] builds a private railroad segmentation dataset and uses a pyramid structure to extract multi-scale features.

2.2 Lane Detection

Recently, the algorithm of lane detection has been booming partly due to the rich and large databases [1, 3, 16, 34, 35] and various frameworks [4, 6, 9, 10, 16, 23, 26, 28, 31, 34, 35, 37].

The Tusimple ² and CULane [34] are the two most famous datasets of lane datasets, promoting many algorithms in lane detection. Some efficient row-based methods take advantage of the slender characteristics of the rails to reduce computation costs. Specifically, [26] implements a small-sized anchor-based lane marking detection algorithm on ResNet. [28] proposes an anchor-based feature pooling and a novel attention mechanism.

3 RAIL-DB

Although many methods have been proposed for rail detection in the community, a public railroad dataset is unavailable. Recently, lane detection has been promoted partly due to various datasets, e.g., Tusimple and CULane [34]. A high-quality dataset is expected to contribute to algorithm performance and comparison in the deep learning era. To this end, we construct a high-quality rail dataset for rail detection. We will introduce the establishment of the Rail-DB, and the comparison between the Rail-DB and existing datasets.

3.1 Dataset collection and annotation

Figure 2 exhibits four steps to build the Rail-DB: the image collection stage, the coarse annotation stage, the refinement annotation stage, and the scene annotation stage.

First, we obtain videos from real-world trains in different conditions to collect rich rail data. Instead of downloading internet images, we capture videos of working trains on various railroads. As the train speed varies, we use adaptive frame intervals to guarantee sufficient difference between adjacent frames. We finally get 15 long-term videos and extract 7432 different frames. Figure 2 (a) shows some raw frames, which contain varied environments: lighting, road structures, and views. Second, we label rails with polylines at the coarse stage to gather annotations. At the coarse stage, eight annotators use *LabelMe* to label rails' polylines. A polyline represents a rail to fit its varying slender structure. Besides, we annotate different rails in an image with different numbers to encode the

¹<https://github.com/CSAILVision/LabelMeAnnotationTool>

²<https://github.com/TuSimple/tusimple-benchmark>

trails' priority order. The central pair of rails in an image get numbers 1 and 2 as they are the primary rails of the view. The other rails obtain pairs of increasing numbers according to their growing distance from the main rails. Figure 2 (b) illustrates the drawing and colored polylines. However, most of the coarse annotations are more or less incomplete and inaccurate in the non-primary trails, the tails, and the priority order. Therefore, we need a refinement stage to improve the quality of the coarse annotations. An expert checks all the rough annotations and refines them carefully. Figure 2 (c) displays that an incomplete annotation of the tails is made complete at the refinement stage. In addition to the rail annotation, we annotate all images with nine scenes properties based on contexts to explore algorithms' generalization. According to various lighting, structures, and views, the contexts are categorized into sun, rain, night, curve, slope, cross, line, near, and far, which account for 2977, 3733, 842, 2497, 1175, 3581, 288, 5534, and 1898 images, respectively. Some images include more than one scene property. Figure 2 (d) shows the statistics of our Rail-DB.

3.2 Dataset comparison

We compare the Rail-DB with existing lane and rail datasets in this section. Table 1 and Figure 3 show the property and visualization differences between existing datasets and the Rail-DB. The comparison shows the Rail-DB is a middle size, background-rich, polyline annotated, and open dataset.

The Tusimple³ and CULane [34] are two famous datasets of lane datasets [1, 3, 16, 34, 35]. The quantity, lighting, road structures, annotation and availability are shown in Table 1. In comparison, rail detection is more challenging than lane detection in three aspects: road structure, complex backgrounds, and dataset availability. From Figure 3, rails are slenderer and more curved than the lane; the rail background contains forests, fields, bridges, city buildings, etc. Worse, there is no public rail dataset. Therefore, a publicly available and multi-scenario rail dataset is urging for rail detection.

The RSDS [32] is a private rail dataset with 3000 pairs of images and segmentation labels. The lighting contains the sun, and the road structures include lines and curves. The RSDS is private and simple, although this dataset achieves practical results. In comparison, the Rail-DB is more extensive and informative than the RSDS. Specifically, the Rail-DB is an open rail dataset consisting of 7432 pairs of images and polyline labels. The lighting includes sun, night, and rain, and the road types vary in line, curve, cross, and slope. As is shown in Figure 3, the RSDS and the Rail-DB are made in different environments. The images of the RSDS are mainly captured in the city under the sun. In contrast, the pictures of the Rail-DB are recorded in the wild in various environments.

4 RAIL-NET

We propose the Rail-Net to detect the rails in video frames automatically. The existing rail detection works include the traditional hand-craft methods [2, 13, 15, 20, 24, 25, 27, 29, 36] and the deep learning segmentation approaches [7, 32]. They are either accuracy effective or computation efficient. The efficient row-based methods are widespread in lane detection [26, 28, 30, 37]. These methods take advantage of the slender characteristics of the rails to reduce

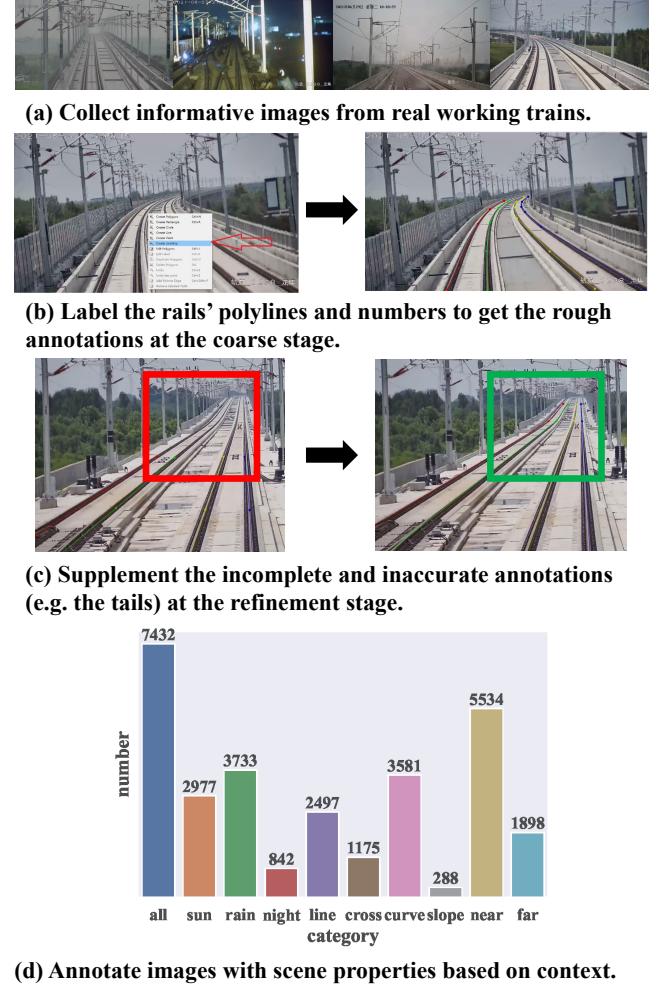


Figure 2: The establishment process of our Rail-DB: image collection, coarse annotation stage, refinement annotation stage, and scene annotation. Better to zoom in with PDF.

computation costs. To this end, we extend the row-based methods to rail detection. We will introduce the problem formulation and the network pipeline of the Rail-Net in this section.

4.1 Problem Formulation

We treat rail detection as a row-based selecting problem instead of a segmentation issue. The segmentation methods turn the polylines into multivalued maps. In contrast, we transfer the labeled polylines of rails into a series of horizontal locations. Figure 4 shows selecting areas of the rails in an image.

We adopt the anchors and uniform columns to grid the image into cells. The labeled polylines would intersect with the anchors in specific columns. As shown in Figure 4, we highlight the intersection areas with different colors for different rails, where the corresponding columns with the same colors represent the rails' locations. Therefore, we select these locations as the class labels for

³<https://github.com/TuSimple/tusimple-benchmark>

Database	Tusimple	CULane	RSDS	Rail-DB (ours)
quantity	3626	133,235	3000	7432
lighting	sun	sun, night, rain	sun	sun, night, rain
structures	line, curve	line, curve, cross	line, curve	line, curve, cross, slope
annotation	keypoints	keypoints	segmentation map	polylines
availability	public	public	private	public

Table 1: The comparison of properties between existing lane and rail datasets and the Rail-DB. The Rail-DB is a middle size, background-rich, polyline annotated, and open rail dataset.

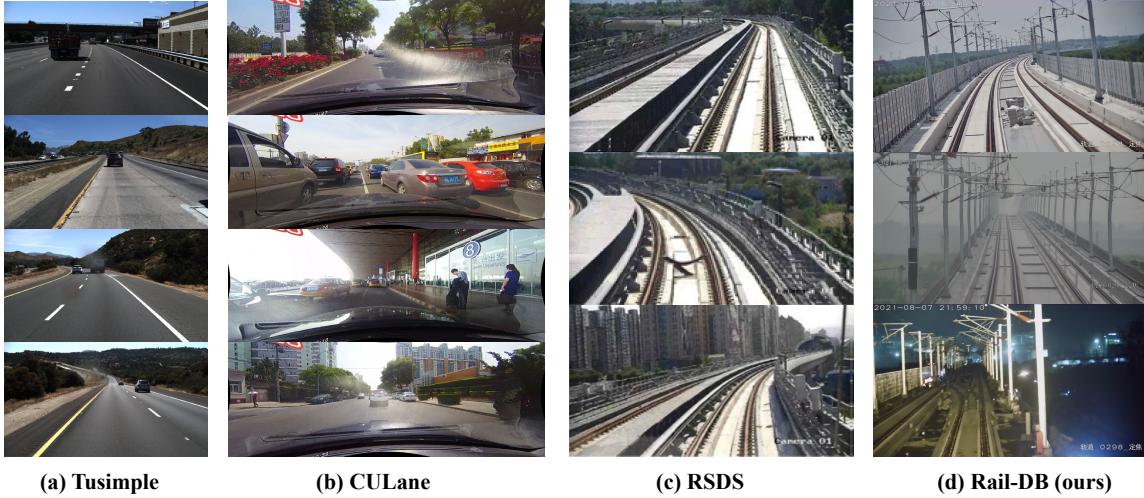


Figure 3: The comparison of visualization between existing database and the Rail-DB. The structures of rails are longer and thinner than lanes. Besides, the view of rails is different from lanes since trains are higher than cars. The images of the RSDS are mainly captured in the city under the sun. In contrast, the images of the Rail-DB are recorded in the wild in richer contexts.

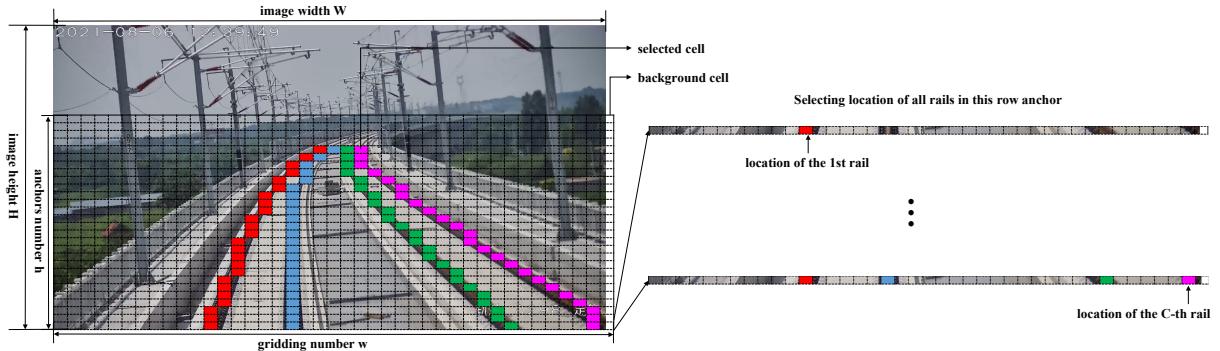


Figure 4: Illustration of selecting locations of the rails. Different colors refer to different rails.

each rail. The anchors and the column locations make rail detection a row-based selecting problem. Its algorithm outputs location classes, which is less than the one of segmentation pixels.

Suppose the number of rails is C , the number of row anchors is h , and the number of gridding columns is w . Suppose X is the feature extracted through the backbone and f^{ij} is the classifier to predict the rail location on the i -th rail, j -th row anchor. The formulation of rail prediction can be written as:

$$P_{i,j,:} = f^{ij}(X), \text{ s.t. } i \in [1, C], j \in [1, h], \quad (1)$$

where $P_{i,j,:}$ is the $(w + 1)$ -dimensional probability vector for selecting gridding columns of the i -th rail, j -th row anchor. Particularly, we extend w dimension to $(w + 1)$ dimension to indicate the background. For each rail and each anchor, the algorithm predicts the

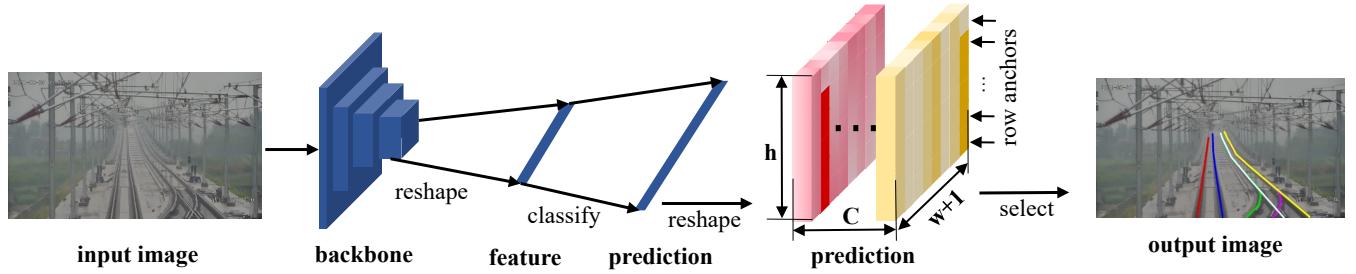


Figure 5: The pipeline of our Rail-Net. We first transfer an image into high-level feature maps through a backbone. Then, we reshape the high-level feature maps and classify them for the predictions of all the anchors and rails.

probability distribution of all locations. As a result, the right locations can be selected based on the probability distribution to represent the complete rails.

The row-based formulation significantly reduces computation costs compared to a segmentation map. The computation cost of our method is proportional to $C \times h \times (w + 1)$. Suppose the image size is $H \times W$, the computation cost of the segmentation method is proportional to $C \times H \times W$. In general, the numbers of gridding cells and row anchors are far smaller than the image size. That is to say, $h \ll H$ and $w \ll W$. In formulation, the computation reduction ratio (r) can be expressed as:

$$r = \frac{HW}{h \times (w + 1)}. \quad (2)$$

Take the setting of the Rail-DB as an example, the input size of the image is 288×800 , and the size of gridding cells is 52×200 . The computation reduction ratio achieves 22.

4.2 Network Pipeline

We propose the Rail-Net based on the row-based formulation. We discard up-sampling modules compared to the segmentation methods. Figure 5 shows the pipeline of the Rail-Net.

We first extract the high-level feature map for a given image through a backbone. Specifically, an 18-layers ResNet [8] serves as the backbone for its computation efficiency and performance effectiveness. The input size is $288 \times 800 \times 3$ in the Rail-DB, the output size is $9 \times 25 \times 512$ after a $1/32$ down-sampling in the backbone.

We classify the high-level feature for the predictions of all the anchors and rails. In detail, a 1×1 kernel convolutional layer reduces the high-level feature into $9 \times 25 \times 8$. Then a reshape operation flattens the high-level feature into a 1800 dimensions vector. These operations are more advanced than the global average pooling. The vector contains global spatial information and avoids computation redundancy, meeting the accuracy and high-speed requirements. Finally a classifier turns the vector into the $C \times h \times (w+1)$ predictions. The classifier consists of an $1800 - 2048$ linear layer, a ReLU layer, and a $2048 - C \times h \times (w+1)$ linear layer. In the setting of the Rail-DB, C is 4, h is 52, and w is 200.

We adopt a Cross-Entropy (CE) function for classification loss L_{cls} to regularize the parameters optimization in the training phase. Suppose $T_{i,j,:}$ is the one-hot label of correct locations. The formulation corresponds to:

$$L_{cls} = \sum_{i=1}^C \sum_{j=1}^h CE(P_{i,j,:}, T_{i,j,:}). \quad (3)$$

We use the expectation of predictions like [26] to obtain locations from the classification prediction in the testing phase. A softmax function gives the probability of different locations:

$$Prob_{i,j,:} = \text{softmax}(P_{i,j,1:w}), \quad (4)$$

where $P_{i,j,1:w}$ is a w -dimensional vector since the background gridding column is not included, and $Prob_{i,j,:}$ represents the normalized probability at each location. Then, the formulation of expected locations $Loc_{i,j}$ can be written as:

$$Loc_{i,j} = \sum_{k=1}^w k \cdot Prob_{i,j,k}, \quad (5)$$

in which $Prob_{i,j,k}$ is the probability of the i -th rail, the j -th row, and the k -th location.

5 EXPERIMENTS

5.1 Implementations

Dataset splits. We randomly sample 80% images as the training set and 20% as the validation set. We resize images into 1280×720 in the evaluation following the Tusimple⁴. The predefined row anchors range from 200 to 720 with a step of 10.

Training Details. The optimizer is Adam. The learning rate strategy is cosine decay. The initialized learning rate is 4e-4. The batch size is 64. The total training epochs number is 50. We train and test all models with PyTorch [22] and NVIDIA Tesla V100 GPU.

Data Pre-process. In the optimizing process, images are resized to 800×288 following [21]. A row-based classification network is vulnerable to over-fit the training set because of the inherent structure of rails. Therefore, we augment the images with rotation and vertical and horizontal shifts.

Evaluation Metric. The evaluation metric is accuracy, which is

$$\text{accuracy} = \frac{\sum_{\text{clip}} C_{\text{clip}}}{\sum_{\text{clip}} S_{\text{clip}}}, \quad (6)$$

where C_{clip} is the number of rail points predicted correctly and S_{clip} is the total number of ground truth in each image. If the difference width between ground-truth and prediction is less than a

⁴<https://github.com/TuSimple/tusimple-benchmark>

Accuracy (%)	Hand-crafted	Segmentation	Rail-Net
sun	26.61	70.53	88.31
rain	56.55	95.62	95.42
night	31.09	76.61	96.22
line	43.57	87.97	95.82
cross	29.88	80.15	87.74
curve	45.58	88.85	92.76
slope	27.96	86.78	91.48
near	44.01	84.91	93.75
far	36.57	92.93	89.87
total	42.14	86.91	92.77
FLOPs (G)	-	13.02	8.42
memory (MB)	-	258.62	182.84
FPS	43	188	312

Table 2: Quantitative comparison of all methods on the Rail-DB. In accuracy, the Rail-Net outperforms the hand-crafted method and the segmentation approach by 50.65% and 5.86%. In FPS, the Rail-Net runs 7.16 times and 1.64 times faster than the hand-crafted method and the segmentation approach.

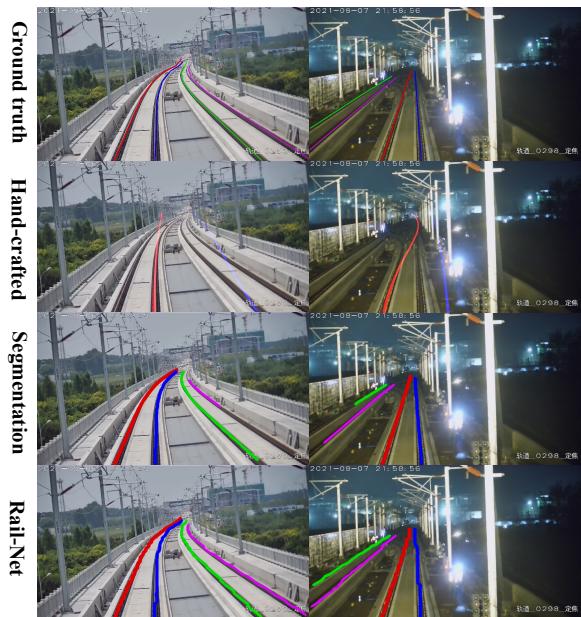


Figure 6: Qualitative comparison of all methods on the Rail-DB. Different colors refer to different rails. The Rail-Net is more accurate and complete in the sun and night situations than alternative methods.

threshold, the predicted point is a correct one. We set the threshold T_p by default to 6 pixels for the Rail-DB.

5.2 Method Comparison

We compare the Rail-Net with the traditional hand-crafted method and the deep learning segmentation approach on the Rail-DB in this section. The Rail-Net utilizes an 18-layers Resnet [8] for the

backbone and fully linear layers for the classifier. The hand-crafted method uses the Sobel operator [14] and polynomial fitting to match railroad rails. The segmentation method selects the same backbone as the Rail-Net, and convolutional blocks for the segmentation heads like FPN [17]. We show quantitative and qualitative comparison results in the following.

Quantitative comparison. We compare all methods in the evaluation accuracy and the frame per second (FPS) in these experiments. The validation accuracy of one run is reported as the evaluation accuracy. The FPS is recorded with the average time for 100 runs. We also calculate each situation’s accuracy and adopt torchinfo⁵ to get inference memory and floating point operations per second (FLOPs). The results are shown in Table 2.

From Table 2, we can see that the Rail-Net achieves dominant performance in accuracy and FPS. For the accuracy, the Rail-Net outperforms the hand-crafted method and the segmentation approach by 50.65% and 5.86%. For the FPS, the Rail-Net runs 7.16 times and 1.64 times faster than the hand-crafted method and the segmentation approach. The results show the Rail-Net exceeds alternative methods in most situations as the formulation of the Rail-Net is more suitable for the slender rails.

However, the segmentation approach outperforms the Rail-Net under the far scenes. Intuitively, the Rail-Net uses the slender property of rails and the contrast of an image row. In most cases, it is more certain and distinctive to output the rails’ positions. Nevertheless, the rail areas are more apparent in the far situation due to the enlarged perspective. It is better to segment the big area than locate its location. Additionally, the Rail-Net occupies less computational cost in FLOPs and memory than the segmentation approach, as the row-based formulation replaces segmenting each pixel with selecting the locations of each row. The reduced FLOPs and memory contribute to higher FPS.

Qualitative comparison. Figure 6 shows the visualization of all methods. We color the labels of the ground truth and the predictions under the sun and night conditions. We can see the Rail-Net performs better than other methods. The Rail-Net shows advanced completeness and accuracy in two situations. Specifically, the hand-crafted method fails to recognize the lane class or the distant trail position. The Sobel operator fails to process complicated contexts, leading to odd fitting curves. Besides, the segmentation methods are confused with the non-primary rails and the background because the secondary rails are much smaller than the main rails. In contrast, the row-based Rail-Net should select a specific location in one row, accounting for more complete results.

Worth mentioning that the locations seem still precise when the rail extends far ahead since the line becomes thinner. On the one hand, we split the width 1280 into 200 columns with padding for rail detection, i.e., the horizontal resolution for our method is 6 pixels. It is enough in practical applications, and we do not observe serious degradation in far rail locations. On the other hand, the far rail locations can be accurate in the following frames.

⁵<https://github.com/TylerYep/torchinfo>

accuracy (%)	total	sun	rain	night	line	cross	curve	slope	near	far
total	92.77	56.31	60.75	36.14	64.13	55.22	79.53	36.56	78.11	48.32
sun	88.31	79.61	24.66	25.47	45.54	42.84	73.12	23.64	66.93	44.56
rain	95.42	41.35	95.84	36.86	73.16	60.81	92.56	49.61	83.39	54.58
night	96.22	37.77	25.96	75.8	91.14	75.66	35.02	20.91	92.25	33.78
line	95.82	54.54	63.09	42.62	95.52	57.61	71.24	35.61	76.96	54.01
cross	87.74	58.67	34.16	33.41	52.59	80.14	41.69	20.75	80.93	32.51
curve	92.76	56.11	68.94	33.21	48.22	46.69	90.79	41.65	78.62	49.77
slope	91.48	39.75	75.46	22.72	53.72	43.44	77.37	77.16	43.04	73.18
near	93.75	54.19	61.42	37.77	64.21	57.77	77.55	35.65	93.86	34.62
far	89.87	62.58	58.79	31.26	63.91	47.61	85.48	39.26	31.06	89.21

Table 3: The results of cross-scene experiments. The column keys and the row keys denote the training rail situations and testing rail situations. The training of the whole dataset exceeds the training of the partial one in most cases.

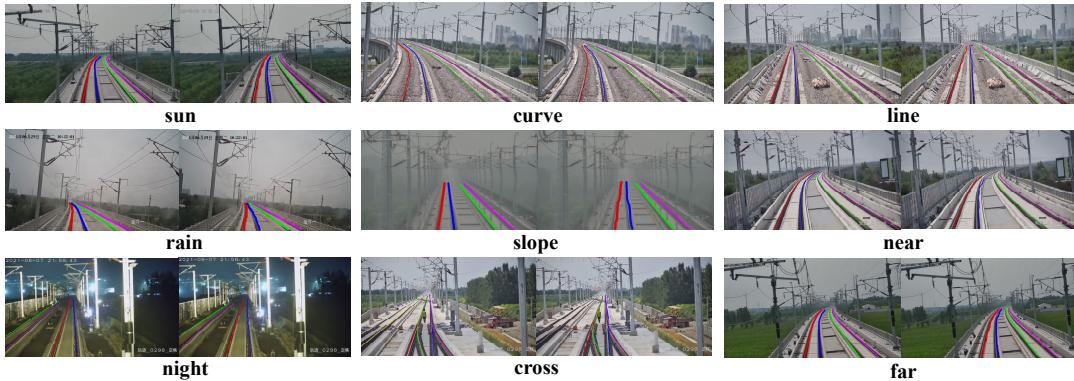


Figure 7: The visualization of the predictions under various scenes. The algorithm performs well under different conditions, even in the cross structure. This shows the situation diversity of the images enhances the method's stability.

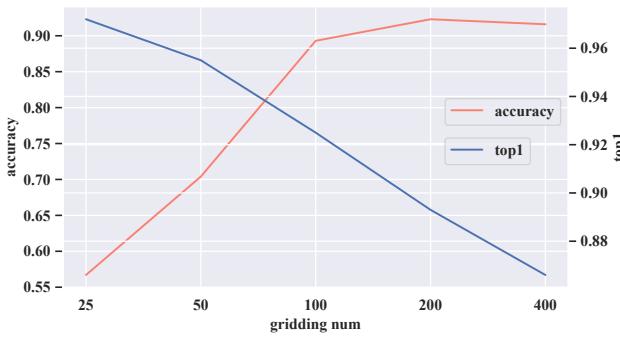


Figure 8: The results of different gridding numbers experiments. The evaluation accuracy arrives at the peak when the gridding number is 200.

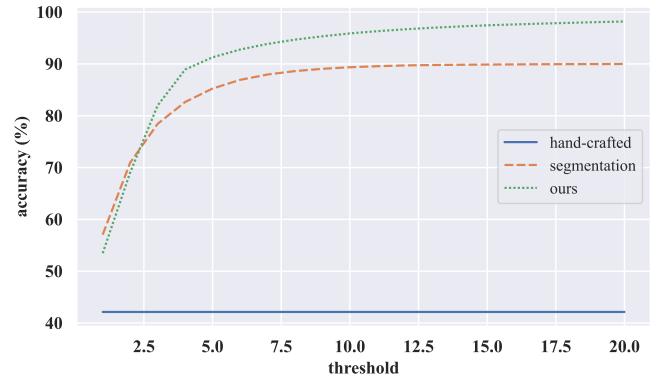


Figure 9: The results of different thresholds in metrics. The Rail-Net consistently outperforms the segmentation method with threshold larger than 3.

5.3 Ablation study

We conduct ablation studies for our Rail-Net on the Rail-DB to verify the effectiveness of every setting. We will show the effects of scene diversity, gridding number, backbone models, and thresholds.

The scene diversity of the Rail-DB. As described in Section 3.1, we collect railroad images from different conditions to improve

scene diversity. To survey the effects of scene variety, we carry out cross-scene experiments among nine situations. We adopt the partial dataset of one scene for training. The quantitative and qualitative results are shown in Table 3 and Figure 7.

backbone	accuracy (%)	FPS
ResNet (18)	92.77	312
ResNet (34)	93.09	169
ResNet (50)	92.71	126
SqueezeNet	90.99	370
MobileNet	92.24	147
ViT	82.41	109

Table 4: The results of different backbones for our Rail-Net. The ResNet (18) is sufficient for rail detection compared to deeper ResNets, efficient alternative networks and the ViT.

	Ultra-Fast	LaneATT	Rail-Net
Tusimple (Acc %)	95.82	95.57	95.64
CULane (F1 %)	68.40	75.09	68.1
Rail-DB (Acc %)	84.5	85.6	92.7

Table 5: The comparison of rail and lane detection. The results show it is not trivial to apply the row-selecting idea of lane detection to rail detection since those SOTA methods based on the idea are largely inferior to ours.

From Table 3, the training of the whole dataset exceeds the training of the partial one in most cases except the rain and near scenes, as the diversity of training images often contributes to higher accuracy for one situation and better generalization. It demonstrates that the Rail-DB with diversified situations is more advanced than the existing single-scene rail dataset. Figure 7 shows the algorithm performs well under different conditions, even in the cross structure, because the situation diversity enhances the method’s stability.

In addition, the accuracy suffers from a sharp degradation for cross-scene conditions. For cross-lighting, the performance falls dramatically, e.g., the accuracy suffers a 54.14% decline from the sun to the night scenes. The performance also has a rapid cut for cross-structure, e.g., the accuracy decreases 37.91% from the line to the cross. For cross-view, the performance degrades seriously, e.g., the accuracy gets a 58.15% reduction from the far to the near. The algorithm is trained to distinguish the difference between original rails and backgrounds. Intuitively, the performance drops when adopting other structures, lighting, and views. In the future, we will pay more attention to solving this challenge.

The gridding number of the Rail-Net. As described in Section 4.1, we adopt gridding to establish classification-based formulation. We make the gridding numbers 25, 50, 100, 200, and 400 to explore the effects of different gridding numbers. Figure 8 shows the results of the top1 accuracy and the evaluation accuracy. Specifically, the top1 accuracy refers to the classification accuracy of gridding cells, while the evaluation accuracy means the location accuracy of rails.

From Figure 8, the top-1 accuracy decreases as the gridding number increases. Because the more the gridding cells, the higher the requirement of the fine-grained classification. However, the evaluation accuracy gets a lift and a slight drop instead of suffering from a monotonic decrease. It is because the gridding number is a trade-off between classification and location. When the gridding number is too small, the gridding cell is too large to represent a precise location. When the gridding number is too big, the classification is too difficult, leading to a slight drop in the evaluation accuracy. The

evaluation accuracy arrives at the peak when the gridding number is 200. Therefore, we choose 200 as the default gridding number.

The thresholds of metric. As described in Section 5.1, the performance is closely relevant to the threshold T_p of the metric. Figure 9 illustrates the results in the experiment of Table 2 with varied T_p . The Rail-Net consistently outperforms the segmentation method with a threshold larger than 3. The hand-crafted method maintains the same performance since it only provides the curves. However, the segmentation method is superior to our Rail-Net with small thresholds. Intuitively, the gridding of the Rail-Net results in coarse grain outputs, while segmentation provides pixel-wise regions. Worth noting that a threshold like 6 makes sense in practice.

The comparison between rail and lane detection. Table 5 conduct experiments on lane and rail datasets. a) We adopt the Rail-Net on two public lane datasets, TuSimple and CULane, and respectively achieves 95.64 in accuracy and 68.1 in F1. The Rail-Net is slightly inferior to lane detection methods like Ultra-Fast and LaneATT on lane datasets. b) We adopt Ultra-Fast [26] and LaneATT [28] for rail detection on the Rail-DB. The Ultra-Fast achieves 84.5, and the LaneATT obtains 85.6, which is much lower than the one of our Rail-Net (i.e., 92.7). These results indicate that a straightforward transfer from lane detection method to rail detection leads to degradation due to the inherent difference between the two tasks. On the contrary, our elaborately-designed yet simple row-selecting Rail-Net achieves the best performance.

The backbones of Rail-Net. As described in Section 4.2, we adopt an 18-layers ResNet as the backbone. We evaluate various backbones for our Rail-Net to show their effects, including ResNet [8] with different layers, SqueezeNet [12], MobileNet [11], and Vision Transformer (ViT) [5]. Table 4 shows the accuracy and FPS results of different backbones. In accuracy, all the ResNet models perform similarly, indicating that a small network like ResNet18 is sufficient for the rail feature extraction. SqueezeNet runs faster than other networks but slightly degrades accuracy as a lighter computation leads to a weaker performance. ViT performs worse than the convolutional networks, partly indicating that the square patching of the picture is unsuitable for slender rails.

6 CONCLUSION

We have proposed an open benchmark, Rail-DB, and an efficient row-based framework, Rail-Net, for rail detection. Experiments show that a lightweight version of the Rail-Net achieves 312 frames per second and 92.77% accuracy on the Rail-DB.

In the future, we will pay more attention to solving these challenges: confusing situations (e.g., the night lighting and crossroad structure) and cross-scene generalization. Besides, considering the slender coherence of rows for smoother outputs and combining segmentation for more fine-grain results are potential directions.

Acknowledgement. This work is partially supported by the National Natural Science Foundation of China (62176165), and the Natural Science Foundation of Top Talent of SZTU (grant no.GDRC202131). Many thanks to all the annotators: Shuyi Mao, Junyao Huang, Xubin Huang, Shaoshuo Liu, Zhengchun Yang, Shao-long Zhao, Ping Huang and Bin Lin.

REFERENCES

- [1] Mohamed Aly. 2008. Real time detection of lane markers in urban streets. In *2008 IEEE Intelligent Vehicles Symposium*. IEEE, 7–12.
- [2] M. Arastouria. 2015. Automated Recognition of Railroad Infrastructure in Rural Areas from LIDAR Data. *Remote Sensing* (2015).
- [3] Karsten Behrendt and Ryan Soussan. 2019. Unsupervised labeled lane markers using maps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*. 0–0.
- [4] Shriyash Chougule, Nora Koznek, Asad Ismail, Ganesh Adam, Vikram Narayan, and Matthias Schulze. 2018. Reliable multilane detection and classification by utilizing CNN as a regression network. In *proceedings of the european conference on computer vision (ECCV) workshops*. 0–0.
- [5] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [6] Mohsen Ghafoorian, Cedric Nugteren, Nóra Baka, Olaf Booij, and Michael Hofmann. 2018. El-gan: Embedding loss driven generative adversarial networks for lane detection. In *proceedings of the european conference on computer vision (ECCV) Workshops*. 0–0.
- [7] X. Giben, V. M. Patel, and R. Chellappa. 2015. Material classification and semantic segmentation of railway track images with deep convolutional neural networks. In *IEEE International Conference on Image Processing*.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [9] Yuenan Hou, Zheng Ma, Chunxiao Liu, Tak-Wai Hui, and Chen Change Loy. 2020. Inter-region affinity distillation for road marking segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12486–12495.
- [10] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. 2019. Learning lightweight lane detection cnns by self attention distillation. In *Proceedings of the IEEE/CVF international conference on computer vision*. 1013–1021.
- [11] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. 2017. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861* (2017).
- [12] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. 2016. SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5 MB model size. *arXiv preprint arXiv:1602.07360* (2016).
- [13] F. Kaleli and Y. S. Akgul. 2009. Vision-based railroad track extraction using dynamic programming. In *International IEEE Conference on Intelligent Transportation Systems*.
- [14] Nick Kanopoulos, Nagesh Vasanthavada, and Robert L Baker. 1988. Design of an image edge detection filter using the Sobel operator. *IEEE Journal of solid-state circuits* 23, 2 (1988), 358–367.
- [15] B. Le Saux, A. Beaupère, A. Boulch, J. Brossard, A. Manier, and G. Villemain. 2018. Railway Detection: From Filtering to Segmentation Networks. In *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium*. 4819–4822. <https://doi.org/10.1109/IGARSS.2018.8517865>
- [16] Seokju Lee, Junsik Kim, Jae Shin Yoon, Seunghak Shin, Oleksandr Bailo, Namil Kim, Tae-Hee Lee, Hyun Seok Hong, Seung-Hoon Han, and In So Kweon. 2017. Vpgnet: Vanishing point guided network for lane and road marking detection and recognition. In *Proceedings of the IEEE international conference on computer vision*. 1947–1955.
- [17] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.
- [18] Jun Lu, Bo Liang, Qujiang Lei, Xiuhao Li, Junhao Liu, Ji Liu, Jie Xu, and Weijun Wang. 2020. SCueU-net: Efficient damage detection method for railway rail. *IEEE Access* 8 (2020), 125109–125120.
- [19] Hiroki Mukojima, Daisuke Deguchi, Yasutomo Kawanishi, Ichiro Ide, Hiroshi Murase, Masato Ukai, Nozomi Nagamine, and Ryuta Nakasone. 2016. Moving camera background-subtraction for obstacle detection on railway tracks. In *2016 IEEE international conference on image processing (ICIP)*. IEEE, 3967–3971.
- [20] B. T. Nassu and M. Ukai. 2011. Rail extraction for driver support in railways. In *IEEE*.
- [21] Xingang Pan, Jianping Shi, Ping Luo, Xiaogang Wang, and Xiaoou Tang. 2018. Spatial as deep: Spatial cnn for traffic scene understanding. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [22] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- [23] Jonah Philion. 2019. Fastdraw: Addressing the long tail of lane detection by adapting a sequential prediction network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 11582–11591.
- [24] A. I. Purica, B. Pesquet-Popescu, and F Dufaux. 2017. A Railroad Detection Algorithm for Infrastructure Surveillance using Enduring Airborne Systems. In *IEEE International Conference on Acoustics*.
- [25] Zhiqian Qi, Yingji Tian, and Yong Shi. 2013. Efficient railway tracks detection and turnouts recognition method using HOG features. *Neural Computing and Applications* 23, 1 (2013), 245–254.
- [26] Zequn Qin, Huanyu Wang, and Xi Li. 2020. Ultra fast structure-aware deep lane detection. In *European Conference on Computer Vision*. Springer, 276–291.
- [27] Laf Rodriguez, J. A. Uribe, and Jfv Bonilla. 2012. Obstacle detection over rails using hough transform. In *Image, Signal Processing, Artificial Vision*.
- [28] Lucas Tabelini, Rodrigo Berriel, Thiago M Paixao, Claudine Badue, Alberto F De Souza, and Thiago Oliveira-Santos. 2021. Keep your eyes on the lane: Real-time attention-guided lane detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 294–302.
- [29] Z. Teng, F. Liu, and B. Zhang. 2016. Visual railway detection by superpixel based intracellular decisions. *Multimedia Tools and Applications* 75, 5 (2016), 2473–2486.
- [30] Tian, Yan, Gelerner, Judith, Wang, Xun, Chen, Weigang, Gao, and Junxiang and. 2018. Lane marking detection via deep convolutional neural network. *Neurocomputing* (2018).
- [31] Wouter Van Gansbeke, Bert De Brabandere, Davy Neven, Marc Proesmans, and Luc Van Gool. 2019. End-to-end lane detection through differentiable least-squares fitting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*. 0–0.
- [32] Yin Wang, Lida Wang, Yu Hen Hu, and Ji Qiu. 2019. RailNet: A segmentation network for railroad detection. *IEEE Access* 7 (2019), 143772–143779.
- [33] Xiukun Wei, Dehua Wei, Da Suo, Limin Jia, and Yujie Li. 2020. Multi-target defect identification for railway track line based on image processing and improved YOLOv3 model. *IEEE Access* 8 (2020), 61973–61988.
- [34] Ping Luo Xiaogang Wang Xingang Pan, Jianping Shi and Xiaoou Tang. 2018. Spatial As Deep: Spatial CNN for Traffic Scene Understanding. In *AAAI Conference on Artificial Intelligence (AAAI)*.
- [35] Hang Xu, Shaoyu Wang, Xinyue Cai, Wei Zhang, Xiaodan Liang, and Zhenguo Li. 2020. Curvelane-nas: Unifying lane-sensitive architecture search and adaptive point blending. In *European Conference on Computer Vision*. Springer, 689–704.
- [36] B Yang and L Fang. 2015. Automated Extraction of 3-D Railway Tracks from Mobile Laser Scanning Point Clouds. *IEEE Journal of Selected Topics in Applied Earth Observations Remote Sensing* 7, 12 (2015), 4750–4761.
- [37] Seungwoo Yoo, Hee Seok Lee, Heesoo Myeong, Sungtrack Yun, Hyoungwoo Park, Janghoon Cho, and Duck Hoon Kim. 2020. End-to-end lane marker detection via row-wise classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 1006–1007.