



Figure 1: **More visualizations on OpenLane val set.** Rows (a), (b), (c) illustrate the **ground truth** lanes, predictions generated by **LATR** and **Persformer** [2], respectively. The lanes are projected onto 2D images and different colors in the images represent specific categories in OpenLane. Missed lanes are indicated by arrows with dashed lines, while existing lanes are indicated by arrows with solid lines. In row (d), we compare the ground truths (**red**) and the predictions of our LATR (**green**) in 3D space. Best viewed in color (zoom in for details.)

in Tab. 2. While employing nearly half the parameters, LATR-Lite not only achieves a twofold increase in FPS but also exhibits a notable enhancement in F1 score, with an impressive improvement of +8.5.

# Layers	F1 / C.Acc.	X error (m)		Z error (m)	
		near	far	near	far
2	68.7 / 90.9	0.260	0.324	0.097	0.130
4	69.9 / 92.3	0.257	0.325	0.097	0.133
6	70.4 / 92.9	0.241	0.321	0.097	0.132
8	70.9 / 93.0	0.257	0.329	0.098	0.134

Table 1: **Ablation study on number of decoder layers.**

Input Sizes. Image resolution is a key factor that influences performance. To study the impact of different input shapes, we compared four resolutions, as detailed in Tab 3. Notably, the F1 score demonstrates improvement with increasing input size. This observation is consistent with our intuitive expectations, as larger images containing finer details can enhance the accuracy of lane location detection.

C.2. Model Complexity.

To comprehensively evaluate the performance of our proposed LATR, we compared its model parameters and FPS with those of the previous state-of-the-art model [2], as shown in Tab 2. Experimental results reveal that our model, LATR, achieves a superior F1 score of 61.9 and a frame rate of 11.34 FPS, outperforming Persformer, which exhibits a lower F1 score of 53.0 and operates at 6.92 FPS. Notably, our LATR-Lite significantly improves efficiency and effectiveness compared to Persformer, despite using almost half the number of parameters. Specifically, we achieve a more

than twofold increase in FPS, while obtaining a remarkable +8.5 improvement in F1 score. The F1 results are compared on OpenLane val set.

Model	Backbone	# Params	FPS	GPU Cost	F1
360 × 480					
Persformer	Efficient-B7	54.94 M	11.48	2.16GB	50.5
Persformer	Res50	62.54 M	9.96	2.24GB	52.6
LATR	Res50	44.35 M	12.63	2.28GB	58.6
720 × 960					
Persformer	Res50	63.19 M	6.92	3.00GB	53.0
LATR-Lite	Res50	38.78 M	17.75	2.26GB	61.5
LATR	Res50	44.35 M	11.34	2.55GB	61.9

Table 2: **Model complexity.** All models are tested on single A100 GPU and AMD EPYC 7351@2.60GHz CPUs. The reported F1 scores are based on OpenLane val set, aligning with the main results in the main paper.

Input Size	F1 / C.Acc.	X error (m)		Z error (m)	
		near	far	near	far
360 × 480	63.8 / 90.4	0.310	0.384	0.113	0.161
480 × 640	67.4 / 92.1	0.262	0.346	0.100	0.140
720 × 960	70.4 / 92.9	0.241	0.321	0.097	0.132
960 × 1280	71.2 / 92.9	0.253	0.315	0.096	0.129

Table 3: **Ablation study on input size.**

C.3. Visualizations

We present additional qualitative analysis in Fig. 1, highlighting differences using arrows of different colors. This analysis demonstrates that LATR can produce more accurate and robust 3D lane results.

References

- [1] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 213–229. Springer, 2020. [1](#)
- [2] Li Chen, Chonghao Sima, Yang Li, Zehan Zheng, Jiajie Xu, Xiangwei Geng, Hongyang Li, Conghui He, Jianping Shi, Yu Qiao, and Junchi Yan. Persformer: 3d lane detection via perspective transformer and the openlane benchmark. In *European Conference on Computer Vision (ECCV)*, 2022. [2](#)
- [3] Tianheng Cheng, Xinggang Wang, Shaoyu Chen, Wenqiang Zhang, Qian Zhang, Chang Huang, Zhaoxiang Zhang, and Wenyu Liu. Sparse instance activation for real-time instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4433–4442, 2022. [1](#)
- [4] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, pages 1–18. Springer, 2022. [1](#)
- [5] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. [1](#)
- [6] Fausto Milletari, Nassir Navab, and Seyed-Ahmad Ahmadi. V-net: Fully convolutional neural networks for volumetric medical image segmentation. In *2016 fourth international conference on 3D vision (3DV)*, pages 565–571. Ieee, 2016. [1](#)
- [7] Russell Stewart, Mykhaylo Andriluka, and Andrew Y Ng. End-to-end people detection in crowded scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2325–2333, 2016. [1](#)