

CODA: A Real-World Road Corner Case Dataset for Object Detection in Autonomous Driving

Kaican Li^{1*}, Kai Chen^{3*}, Haoyu Wang^{1*}, Lanqing Hong^{1†}, Chaoqiang Ye¹, Jianhua Han¹, Yukuai Chen², Wei Zhang¹, Chunjing Xu¹, Dit-Yan Yeung³, Xiaodan Liang⁵, Zhenguo Li¹, and Hang Xu¹

¹ Huawei Noah’s Ark Lab

² Huawei Intelligent Automotive Solution BU

³ Hong Kong University of Science and Technology

⁴ Sun Yat-sen University

Appendix

The official website of CODA is at <https://coda-dataset.github.io>, where we have released 1000 of the 1500 annotated scenes of CODA. The remaining 500 scenes are reserved for the corner case challenge in the SSLAD workshop⁵ of ECCV 2022, and will be released after the challenge.

As a continued effort, CODA is further extended by 8711 additional scenes with more than 28k new corner cases during the reviewing and publishing process of this paper. The extension will also be made fully available after the ECCV challenge. Please stay tuned!

A Supplementary implementation details

A.1 Closed-world detectors

We re-implement the closed-world detectors in Tab. 2 of the main paper following the default configurations in MMDetection [2] if the corresponding official checkpoints⁶ are not publicly available. Detectors are by default trained for 12 epochs (1x) on BDD100K and Waymo, while for 24 epochs (2x) on SODA10M, except Deformable DETR and Sparse R-CNN, which are trained for 100 epochs on SODA10M training set, due to the limited labeled data size of SODA10M and data-hungry property of the Transformer layers.

* Equal contribution.

† Corresponding author at honglanqing@huawei.com.

⁵ <https://ssladcompetition.github.io/>

⁶ SODA10M: <https://soda-2d.github.io/> & BDD100K: <https://github.com/SysCV/bdd100k-models/>

A.2 ORE

The experiments are conducted using the released code⁷ of the original paper [7], without incremental learning. We use the training set of SODA10M [4] for training; the validation set of SODA10M and the whole CODA for fitting the energy distributions; the test set of SODA10M and the whole CODA for testing. The backbone of the network is ResNet-50 [6], and the whole network is trained for 24 epochs with a batch size of 8 under a learning rate of 0.02. More efficient training methods (*e.g.*, sparse training [18]) are potential solutions for better generalization, which will be explored in the future.

A.3 Anomaly detection

We utilize the segmentation model and GAN pre-trained on Cityscapes dataset [3] in method *synthesize then compare* during the test procedure. We crop the generated fake images and the original images according to the top-ranking bounding box proposals obtained from RPN, then we compare the cropped images by pixel-wise cosine similarity as illustrated in the paper and get the detection result according to the similarity rankings. As to *memory-based OOD detection*, we utilize the feature map of common ground truth and some conventional background (vegetation, sky, and pure black) extracted from ResNet-152 [6] as the memory bank. Likewise, We then compare the feature map extracted from off-the-shelf ResNet-152 of the top-ranking bounding boxes from RPN with the memory bank to filter out the anomaly detection results. More advanced OOD detection methods [16,17] will be explored in the future.

B Supplementary ablation studies on COPG

We conduct ablation studies on COPG by tuning its components one at a time. The experiment is conducted on CODA-ONCE whose final ground truths are utilized to compute the AP and AR (under COCO [9] protocol) of the output proposals of the modified COPGs. The results are shown in Tab. 1-3 below, where the **bold** hyperparameters are the ones in effect. The number of proposals (#Proposals) and the number of scenes (#Scenes) containing at least one proposal are also reported.

C Supplementary discussion on potential negative societal impact

CODA has no potential negative societal impact since it is constructed totally based on publicly available datasets with delicate privacy protection. For example, all objects containing personal information (*e.g.*, human faces, license plates) are blurred in ONCE [11].

⁷ <https://github.com/JosephKJ/OWOD>

Table 1. Ablation study on point-cloud clustering. The min. θ controls the minimum tolerance of the angle between two points being assigned to the same object (see Fig. 6 of the main paper). The cluster size controls the minimum tolerance to the number of points of a cluster being considered as a non-trivial object. The max. dist. controls the maximum tolerance to the distance from the nearest point of a cluster to the lidar sensor.

Min. θ	AP	AR	#Proposals	#Scenes
4°	1.3	3.1	1039	763
8°	2.0	5.1	1522	1040
12°	1.9	5.0	1620	1000
16°	1.7	4.7	1652	965

Cluster size	AP	AR	#Proposals	#Scenes
5	2.0	5.1	1522	1040
10	2.0	5.1	1522	1040
20	2.0	5.1	1513	1038

Max. dist.	AP	AR	#Proposals	#Scenes
25	1.8	4.5	1336	935
50	2.0	5.1	1522	1040
100	1.9	5.1	1548	1040

Table 2. Ablation study on background removal. The background (BG) ratio controls the maximum tolerance to the ratio of the proposal area overlapped with the background regions.

BG ratio	AP	AR	#Proposals	#Scenes
0.15	0.4	1.1	412	308
0.30	1.0	2.6	834	614
0.45	2.0	5.1	1522	1040
0.60	2.0	5.6	1834	1055
0.75	1.8	6.1	2315	1056

Table 3. Ablation study on common-class suppression. The IoU threshold controls the maximum tolerance to the IoU between every proposal and each detected common-class object.

IoU	AP	AR	#Proposals	#Scenes
0.00	1.2	2.7	834	643
0.25	2.0	5.1	1522	1040
0.50	1.8	5.2	1660	1041

D Supplementary benchmark results

Table 4. Detection results (%) on our CODA-ONCE dataset. The dramatic performance decrease still maintains, but compared with Tab. 2 of the main paper, we observe a decrease for the reported AR values of all detectors, suggesting that CODA-ONCE constructed with our automatic corner case proposal generation COPG in Sec. 4.2 of the main paper is the most challenging subset of CODA.

Table 5. More detection results(%) on our CODA-ONCE dataset. AR^s , AR^m and AR^l are the average recall for small, medium and large objects respectively, following COCO definition [9], while AR^1 represents the average recall when only one object prediction is allowed for each image. Here AR^s is also marked as “-”, since no small corner cases of common classes are collected in CODA-ONCE subset.

Table 6. More detection results(%) on our CODA dataset. AR^s, AR^m and AR^l are the average recall for small, medium and large objects respectively, following COCO definition [9], while AR¹ represents the average recall when only one object prediction is allowed for each image.

E COPG examples

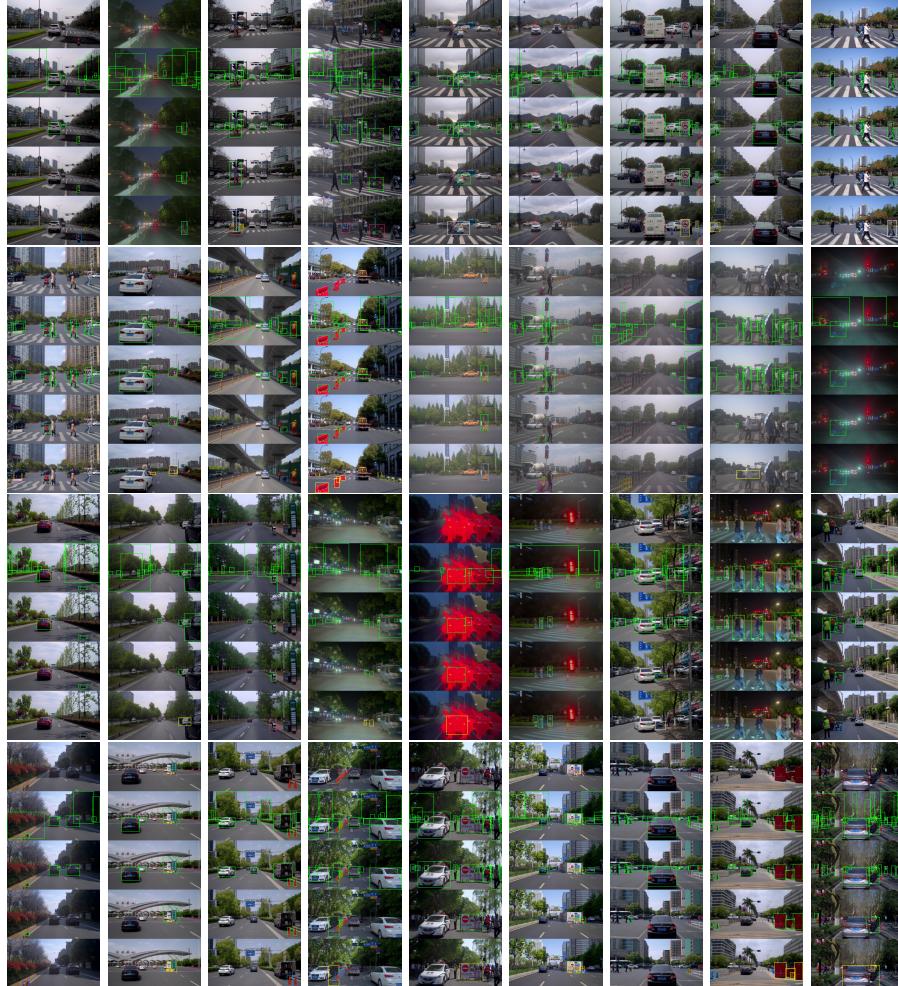


Fig. 1. Examples of COPG step-by-step proposals and final annotations. In each stack, the images from top to bottom correspond to the steps indicated in Fig. 5 of the main paper: **(b)** camera image, **(e)** initial proposal, **(f)** intermediate proposal, and **(g)** final proposal. The last row is the result after manual labeling. Best viewed with color and zoom in.

References

1. Cai, Z., Vasconcelos, N.: Cascade R-CNN: delving into high quality object detection. In: CVPR (2018) [4](#), [5](#), [6](#)
2. Chen, K., Wang, J., Pang, J., Cao, Y., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Xu, J., Zhang, Z., Cheng, D., Zhu, C., Cheng, T., Zhao, Q., Li, B., Lu, X., Zhu, R., Wu, Y., Dai, J., Wang, J., Shi, J., Ouyang, W., Loy, C.C., Lin, D.: MMDetection: Open mmlab detection toolbox and benchmark. arXiv preprint arXiv:1906.07155 (2019) [1](#)
3. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: CVPR (2016) [2](#)
4. Han, J., Liang, X., Xu, H., Chen, K., Hong, L., Mao, J., Ye, C., Zhang, W., Li, Z., Liang, X., Xu, C.: SODA10M: A large-scale 2d self/semi-supervised object detection dataset for autonomous driving. arXiv preprint arXiv:2106.11118 (2021) [2](#), [4](#), [6](#)
5. Han, J., Liang, X., Xu, H., Chen, K., Hong, L., Ye, C., Zhang, W., Li, Z., Liang, X., Xu, C.: Soda10m: Towards large-scale object detection benchmark for autonomous driving. arXiv:2106.11118 (2021) [5](#)
6. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016) [2](#)
7. Joseph, K., Khan, S., Khan, F.S., Balasubramanian, V.N.: Towards open world object detection. In: CVPR (2021) [2](#), [4](#), [5](#), [6](#)
8. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (2017) [4](#), [5](#), [6](#)
9. Lin, T., Maire, M., Belongie, S.J., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. In: ECCV (2014) [2](#), [5](#), [6](#)
10. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: ICCV (2021) [4](#), [5](#), [6](#)
11. Mao, J., Niu, M., Jiang, C., Liang, H., Chen, J., Liang, X., Li, Y., Ye, C., Zhang, W., Li, Z., et al.: One million scenes for autonomous driving: ONCE dataset. arXiv preprint arXiv:2106.11037 (2021) [2](#)
12. Ren, S., He, K., Girshick, R., Sun, J.: Faster R-CNN: Towards real-time object detection with region proposal networks. In: NeurIPS (2015) [4](#), [5](#), [6](#)
13. Sun, P., Kretzschmar, H., Dotiwalla, X., Chouard, A., Patnaik, V., Tsui, P., Guo, J., Zhou, Y., Chai, Y., Caine, B., Vasudevan, V., Han, W., Ngiam, J., Zhao, H., Timofeev, A., Ettinger, S., Krivokon, M., Gao, A., Joshi, A., Zhang, Y., Shlens, J., Chen, Z., Anguelov, D.: Scalability in perception for autonomous driving: Waymo open dataset. In: CVPR (2020) [4](#), [5](#), [6](#)
14. Sun, P., Zhang, R., Jiang, Y., Kong, T., Xu, C., Zhan, W., Tomizuka, M., Li, L., Yuan, Z., Wang, C., et al.: Sparse r-cnn: End-to-end object detection with learnable proposals. In: CVPR (2021) [4](#), [5](#), [6](#)
15. Yu, F., Chen, H., Wang, X., Xian, W., Chen, Y., Liu, F., Madhavan, V., Darrell, T.: Bdd100k: A diverse driving dataset for heterogeneous multitask learning. In: CVPR (2020) [4](#), [5](#), [6](#)
16. Zhou, X., Lin, Y., Pi, R., Zhang, W., Xu, R., Cui, P., Zhang, T.: Model agnostic sample reweighting for out-of-distribution learning. In: ICML (2022) [2](#)

17. Zhou, X., Lin, Y., Zhang, W., Zhang, T.: Sparse invariant risk minimization. In: ICML (2022) [2](#)
18. Zhou, X., Zhang, W., Xu, H., Zhang, T.: Effective sparsification of neural networks with global sparsity constraint. In: CVPR (2021) [2](#)
19. Zhu, X., Su, W., Lu, L., Li, B., Wang, X., Dai, J.: Deformable detr: Deformable transformers for end-to-end object detection. arXiv preprint arXiv:2010.04159 (2020) [4](#), [5](#), [6](#)