# Unsupervised Domain Adaptation for Monocular 3D Object Detection via Self-Training

Zhenyu Li<sup>1</sup>, Zehui Chen<sup>2</sup>, Ang Li<sup>3</sup>, Liangji Fang<sup>3</sup>, Qinhong Jiang<sup>3</sup>, Xianming Liu<sup>1</sup>, and Junjun Jiang<sup>1 $\boxtimes$ </sup>

<sup>1</sup> Harbin Institute of Technology <sup>2</sup> University of Science and Technology <sup>3</sup> SenseTime Research {zhenyuli17, csxm, jiangjunjun}@hit.edu.cn lovesnow@mail.ustc.edu.cn {liang1, fangliangji, jiangqinhong}@senseauto.com

# 1 Dataset Comparisons

To provide more intuitive comparisons among different datasets (*e.g.*, KITTI [3], nuScense [1] and Lyft [4]), we present images with projected ground-truth labels in Fig. 1. One can easily observe cameras utilized in these datasets have different parameters, which are reflected in the image resolutions, FOV, *etc.* This work focuses on designing a general Mono3D UDA framework and solving the severe depth-shift caused by misaligned camera intrinsic parameters, which is the most crucial problem in Mono3D UDA. However, there are still numerous unsolved issues such as different image color styles, various distributions of object dimensions, different distributions of object depth, *etc.* Our proposed STMono3D can be a well-developed baseline for future research.

# 2 Visualizations of Pseudo Labels

Here, we present more visualizations of pseudo labels generated by the teacher model during the training stage. The images show the depth-shift issue caused by the misalignment of camera parameters can be well-solved. The reasonable pseudo labels provide regular supervision on the target domain and achieve Mono3D UDA in a teacher-student paradigm. In addition, we can find there is still a slight error of prediction locations or dimensions that can be improved by further development of Mono3D methods and enhancement of the UDA algorithms. There is still tremedous room for improvement of the Mono3D UDA.

#### **3** Detailed Training Settings

In this section, we introduce more detailed training settings. As for the model, we follow the basic config provided in MMDetection3D [2]. The only modification lies in the scaling of predicted object depth based on the pixel size (GAMS introduced in our paper). We then summary all the runtime settings in Tab. 1, including the number of interations, training schedule, threshold changing, *etc.* 



KITTI

 $\operatorname{nuScense}$ 

Lyft

Fig. 1. Dataset visualizations with ground-truth labels.



Fig. 2. Visualizations of pseudo labels generated by the teacher model.

4 Z. Li et al.

#### References

- Caesar, H., Bankiti, V., Lang, A.H., Vora, S., Liong, V.E., Xu, Q., Krishnan, A., Pan, Y., Baldan, G., Beijbom, O.: nuscenes: A multimodal dataset for autonomous driving. In: Computer Vision and Pattern Recognition (CVPR). pp. 11621–11631 (2020)
- 2. Contributors, M.: MMDetection3D: OpenMMLab next-generation platform for general 3D object detection. https://github.com/open-mmlab/mmdetection3d (2020)
- 3. Geiger, A., Lenz, P., Urtasun, R.: Are we ready for autonomous driving? the kitti vision benchmark suite. In: Computer Vision and Pattern Recognition (CVPR). pp. 3354–3361 (2012)
- Kesten, R., Usman, M., Houston, J., Pandya, T., Nadhamuni, K., Ferreira, A., Yuan, M., Low, B., Jain, A., Ondruska, P., Omari, S., Shah, S., Kulkarni, A., Kazakova, A., Tao, C., Platinsky, L., Jiang, W., Shet, V.: Level 5 perception dataset 2020. https://level5.global/level5/data/ (2019)

## STMono3D 5

Sche	edule			Optimizer		
number of iters	$880 \times 13$		0	optimizer type		
learning rate policy	5	step		learning rate	0.002	
warmup type	linear			gradient clip		
warmup iters	500		bat	batchsize per GPU		
warmup ratio	1/3		nu	number of GPUs		
$\operatorname{step}$	$[880 \times 8, 880 \times 11]$		source:de	source: domain samples per bs		
Mean Teac	her Par	adigm	Infer	Inference Settings (KITTI)		
interval		1		nms pre	100	
$\begin{array}{c c} & 1 \\ warmup \\ increasing thr per step \\ start iter \\ stop iter \\ \end{array} \begin{array}{c} 0 \\ 0.005 \\ 880 \times 8 \\ 880 \times 10 \end{array}$		0 0.005	student	nms thr score thr	$0.05 \\ 0.001 \\ 20$	
		teacher	score thr (only diff.)	0.35		

#### Table 1. Detailed training settings.

$\mathbf{Str}$	ong	Dat	a Au	gmentatior	n (type/prob/details)
n		7714	0.75		

RandomFlip3D	0.5	horizontal
Mono3DResize	1	$\begin{array}{c} (1600,840)(1600,900)\\ (1600,960)(1600,1020)\\ (1600,1080)(1600,1140)\\ (1600,1200)(1600,1260)\\ (1540,840)(1480,780)\\ (1420,720)(1380,680)\\ (1660,960)(1720,1020)\\ (1800,1080)(1880,1140) \end{array}$
OneOf	1	Identity AutoContrast RandEqualize RandSolarize RandColor RandContrast RandBrightness RandSharpness RandPosterize
RandErase	1	size= $[0, 0.2]$ n blocks= $(1, 5)$ squared=True

## Weak Data Augmentation (type/prob/details)

		, , _ ,
RandomFlip3D	0.5	horizontal
Mono3DResize	1	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$