

# Supplementary Material: Train in Germany, Test in The USA: Making 3D Object Detectors Generalize

Yan Wang<sup>\*1</sup>    Xiangyu Chen<sup>\*1</sup>    Yurong You<sup>1</sup>    Li Erran Li<sup>2,3</sup>

Bharath Hariharan<sup>1</sup>    Mark Campbell<sup>1</sup>    Kilian Q. Weinberger<sup>1</sup>    Wei-Lun Chao<sup>4</sup>

<sup>1</sup>Cornell University    <sup>2</sup>Scale AI    <sup>3</sup>Columbia University    <sup>4</sup>The Ohio State University

{yw763, xc429, yy785, bh497, mc288, kqw4}@cornell.edu    erranlli@gmail.com    chao.209@osu.edu

In this Supplementary Material, we provide details omitted in the main paper.

- **section S1:** data format conversion (section 3 of the main paper).
- **section S2:** evaluation metric (section 4.1 of the main paper).
- **section S3:** additional results on dataset discrepancy (section 4.4 and 4.5 of the main paper).
- **section S4:** object detection using PIXOR [9] (section 4.2 and 4.3 of the main paper).
- **section S5:** object detection using POINTRCNN with different adaptation methods (section 4.5 and 5 of the main paper).
- **section S6:** additional qualitative results (section 5 of the main paper).

## S1. Converting Datasets into KITTI Format

In this section we describe in detail how we convert Argoverse [3], nuScenes [2], Lyft [6], and Waymo [1] into KITTI [4, 5] format. As the formatting of images, point clouds and camera calibration information is trivial, and label fields such as  $\alpha$  and  $rotation_y$  have been well-defined, we only discuss the labeling process with non-deterministic definitions.

### S1.1. Object filtering

Due to the fact that KITTI focuses on objects that appear in the camera view, we follow its setting and discard all object annotations outside the frontal camera view. To allow truncated objects, we project the 8 corners of each object’s 3D bounding box onto the image plane. An object will be discarded if all its 8 corners fall out of the image boundary. To make other datasets consistent with KITTI, we do not consider labeled objects farther than 70 meters.

<sup>\*</sup>Equal contributions

Table S1: The original categories in each dataset that we include into the *car* and *truck* categories following the KITTI label formulation.

Dataset	Car	Truck
Argoverse	{VEHICLE}	{LARGE_VEHICLE, BUS, TRAILER, SCHOOL_BUS}
nuScenes	{car}	{bus, trailer, construction_vehicle, truck}
Lyft	{Car}	{other_vehicle, truck, bus, emergency_vehicle}
Waymo	{Car}	$\emptyset$

### S1.2. Matching categories with KITTI

Since the taxonomy of object categories among datasets are misaligned, it is necessary to re-label each dataset in the same way as KITTI does. As we focus on car detection, here we describe how we construct the new *car* and *truck* categories for each dataset except KITTI in Table S1. The *truck* category is also important since detected trucks are treated as false positives when we look at the *car* category. We would like to point out that Waymo labels all kinds of vehicles as *cars*. A model trained on Waymo thus will tend to predict trucks or other vehicles as cars. Therefore, directly applying a model trained on Waymo to other datasets will lead to higher false positive rates. For other datasets, the definition between categories can vary (e.g., Argoverse label Ford F-Series as cars; nuScenes labels some as trucks) and result in cross-domain accuracy drop even if the data are collected at similar locations with similar sensors.

### S1.3. Handling missing 2D bounding boxes

To annotate each object in the image with a 2D bounding box (the information is used by the original KITTI metric),

we first compute 8 corners of its 3D bounding box, and then calculate their pixel coordinates  $\{(cx_n, cy_n), 1 \leq n \leq 8\}$ . We then draw the smallest bounding box  $(x_1, y_1, x_2, y_2)$  that contains all corners whose projections fall in the image plane:

$$\begin{aligned} x_1 &= \max(\min_{0 \leq n < 8} cx_n, 0), \\ y_1 &= \max(\min_{0 \leq n < 8} cy_n, 0), \\ x_2 &= \min(\max_{0 \leq n < 8} cx_n, width), \\ y_2 &= \min(\max_{0 \leq n < 8} cy_n, height), \end{aligned} \quad (1)$$

where *width* and *height* denote the width and height of the 2D image, respectively.

### S1.4. Calculating truncation values

Following the KITTI formulation, the *truncation* value refers to how much of an object locates beyond image boundary. With Equation 1 we estimate it by calculating how much of the object’s 2D uncropped bounding box is outside the image boundary:

$$\begin{aligned} x'_1 &= \min_{0 \leq n < 8} cx_n, \\ y'_1 &= \min_{0 \leq n < 8} cy_n, \\ x'_2 &= \max_{0 \leq n < 8} cx_n, \\ y'_2 &= \max_{0 \leq n < 8} cy_n, \\ \text{truncation} &= 1 - \frac{(x_2 - x_1) \times (y_2 - y_1)}{(x'_2 - x'_1) \times (y'_2 - y'_1)}. \end{aligned} \quad (2)$$

### S1.5. Calculating occlusion values

We estimate the occlusion value of objects by approximating car shapes with corresponding 2D bounding boxes. The *occlusion* value is thus derived by computing the percentage of pixels occluded by bounding boxes from closer objects. We discretize the 0–1 occlusion value into KITTI’s  $\{0, 1, 2, 3\}$  labels by equally dividing the interval into 4 parts. We describe in algorithm 1 the detail of how we compute occlusion value for each object.

## S2. The New Difficulty Metric

In section 4.1 of the main paper, we develop a new difficulty metric to evaluate object detection (*i.e.*, how to define easy, moderate, and hard cases) so as to better align different datasets. Concretely, KITTI defines its *easy*, *moderate*, and *hard* cases according to truncation, occlusion, and 2D bounding box height (in pixels) of ground-truth annotations. The 2D box height (the threshold at 40 pixels) is meant to differentiate far-away and nearby objects: the *easy* cases

---

### Algorithm 1: Computing occlusion of objects from a single scene

---

**Input** : Image height  $H$ , image width  $W$ , object list  $objs$ .

```

1  $canvas \leftarrow Array([H, W]);$ 
2 for  $x \in \{0, 1, \dots, W - 1\}$  do
3   for  $y \in \{0, 1, \dots, H - 1\}$  do
4      $canvas[x, y] \leftarrow -1;$ 
5  $objs \leftarrow Sort(objs, key = depth, order =$ 
    $descending);$ 
6 for  $obj \in objs$  do
7    $[x_1, x_2, y_1, y_2] = obj.bounding\_bbox;$ 
8   for  $x \in \{x_1, x_1 + 1, \dots, x_2 - 1\}$  do
9     for  $y \in \{y_1, y_1 + 1, \dots, y_2 - 1\}$  do
10     $canvas[x, y] \leftarrow obj.id;$ 
11 for  $obj \in objs$  do
12    $[x_1, x_2, y_1, y_2] = obj.bounding\_bbox;$ 
13    $cnt \leftarrow 0;$ 
14   for  $x \in \{x_1, x_1 + 1, \dots, x_2 - 1\}$  do
15     for  $y \in \{y_1, y_1 + 1, \dots, y_2 - 1\}$  do
16       if  $canvas[x, y] = obj.id$  then
17          $cnt \leftarrow cnt + 1;$ 
18    $obj.occlusion \leftarrow 1 - \frac{cnt}{(x_2 - x_1) \times (y_2 - y_1)};$ 

```

---

only contain nearby objects. However, since the datasets we compare are collected using cameras of different focal lengths and contain images of different resolutions, directly applying the KITTI definition may not well align datasets. For example, a car at 50 meters is treated as a moderate case in KITTI but may be treated as a easy case in other datasets.

To resolve this issue, we re-define detection difficulty based on object truncation, occlusion, and **depth range (in meters)**, which completely removes the influences of cameras. In developing this new metric we hope to achieve similar case partitions to the original metric of KITTI. To this end, we estimate the distance thresholds with

$$D = \frac{f_v \times H}{h}, \quad (3)$$

where  $D$  denotes depth,  $f_v$  denotes vertical camera focal length, and  $H$  and  $h$  are object height in the 3D camera space and the 2D image space, respectively. For a car of average height (1.53 meters) in KITTI, the corresponding depth for 40 pixels is 27.03 meters. We therefore select 30 meters as the new threshold to differentiate *easy* from *moderate* and *hard* cases. For *moderate* and *hard* cases, we disregard cars with depths larger than 70 meters since most of the annotated cars in KITTI are within this range. Table S3

Table S2: 3D object detection results across multiple datasets using the original KITTI evaluation metric (pixel thresholds). We apply POINTRCNN [8]. We report average precision (AP) of the *Car* category in bird’s-eye view and 3D ( $AP_{BEV}$  /  $AP_{3D}$ ,  $IoU = 0.7$ ) and compare object detection accuracy of different difficulties. The results are less comparable due to misaligned difficulty partitions among datasets. Red color: best generalization (per column and per setting); blue color: worst generalization; bold font: within-domain results.

Setting	Source \ Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>88.0 / 82.3</b>	44.2 / 21.4	27.5 / 7.1	72.3 / 45.5	42.1 / 10.6
	Argoverse	68.6 / 31.5	<b>69.9 / 43.6</b>	28.3 / 11.4	76.8 / 56.4	73.5 / 34.2
	nuScenes	49.4 / 13.2	57.0 / 16.5	<b>43.4 / 21.3</b>	83.0 / 31.8	71.7 / 28.2
	Lyft	72.6 / 38.9	66.9 / 33.2	35.5 / 13.1	<b>86.4 / 77.1</b>	78.0 / 54.6
	Waymo	52.0 / 13.1	64.9 / 29.4	31.5 / 14.3	82.5 / 68.8	<b>85.3 / 71.7</b>
Moderate	KITTI	<b>86.0 / 74.7</b>	44.9 / 22.3	26.2 / 8.3	63.2 / 36.3	43.9 / 12.3
	Argoverse	65.2 / 36.6	<b>69.8 / 44.2</b>	27.6 / 11.8	68.5 / 43.6	72.1 / 35.1
	nuScenes	45.4 / 12.1	56.5 / 17.1	<b>40.7 / 21.2</b>	73.4 / 26.3	68.1 / 30.7
	Lyft	67.3 / 38.3	62.4 / 35.3	33.6 / 12.3	<b>79.6 / 66.8</b>	77.3 / 53.1
	Waymo	51.5 / 14.9	64.4 / 29.8	28.9 / 13.7	75.5 / 58.2	<b>85.6 / 67.9</b>
Hard	KITTI	<b>85.7 / 74.8</b>	42.5 / 22.2	24.9 / 8.8	62.0 / 34.9	41.4 / 12.6
	Argoverse	63.5 / 37.8	<b>69.8 / 42.8</b>	26.8 / 14.5	65.9 / 44.4	68.5 / 36.7
	nuScenes	42.2 / 11.1	53.2 / 16.7	<b>40.2 / 20.5</b>	73.0 / 27.8	66.8 / 29.0
	Lyft	65.0 / 37.0	62.8 / 35.8	30.6 / 11.7	<b>79.7 / 67.3</b>	76.6 / 53.8
	Waymo	48.9 / 14.4	61.6 / 29.0	28.4 / 14.1	75.5 / 55.8	<b>80.2 / 67.6</b>

Table S3: Percentage (%) of data (total annotated cars with depths  $\in [0, 70]$  meters) in each difficult partition with old / new difficulty metric. The new *easy* threshold selects much fewer data than the old metric on all datasets except KITTI.

	Dataset	Easy	Moderate	Hard
Training Set	KITTI	21.7 / 21.6	55.5 / 67.8	76.1 / 91.0
	Argoverse	27.7 / 14.9	40.5 / 40.5	59.6 / 59.6
	nuScenes	31.9 / 13.9	47.2 / 47.2	64.8 / 64.8
	Lyft	25.0 / 15.5	50.4 / 54.9	64.9 / 70.5
	Waymo	29.0 / 10.7	40.1 / 40.1	58.7 / 58.7
	Validation Set	KITTI	20.4 / 20.4	55.4 / 65.5
Argoverse		29.2 / 14.3	41.7 / 41.7	60.6 / 60.6
nuScenes		38.3 / 18.4	53.6 / 53.6	68.5 / 68.5
Lyft		25.3 / 15.5	52.5 / 57.7	66.9 / 73.3
Waymo		30.3 / 10.3	42.3 / 42.3	60.9 / 60.9

shows the comparison between old and new difficulty partitions. The new metric contains fewer easy cases than the old metric for all but the KITTI dataset. This is because that the other datasets use either larger focal lengths or resolutions: the objects in images are therefore larger than in KITTI. We note that, the moderate cases contain all the easy cases, and the hard cases contain all the easy and moderate cases.

We also report in Table S2 the detection results within and across datasets using the old metric, in comparison to Table 2 of the main paper which uses the new metric. One notable difference is that for the easy cases in the old metric, both the within and across domain performances drop for

Table S4: The average size (meters) of 3D ground truth bounding boxes of the five datasets.

Dataset	Width	Height	Length
KITTI	1.62	1.53	3.89
Argoverse	1.96	1.69	4.51
nuScenes	1.96	1.73	4.64
Lyft	1.91	1.71	4.73
Waymo	2.11	1.79	4.80

all but KITTI datasets, since many far-away cars (which are hard to detect) in the other datasets are treated as easy cases in the old metric.

### S3. Dataset discrepancy

We have shown the box size distributions of each dataset in Figure 3 of the main paper. We also calculate the mean of the bounding box sizes in Table S4. There is a huge gap of size between KITTI and the other four datasets. In addition, we train an SVM classifier with the RBF kernel to predict which dataset a bounding box belongs to and present the confusion matrix result in Figure S1 (row: ground truth; column: prediction). The model has a very high confidence to distinguish KITTI from the other datasets.

We further train a point cloud classifier to tell which dataset a point cloud of car belongs to, using PointNet++ [7] as the backbone. For each dataset, we sample 8,000 object point cloud instances as training examples and 1,000 as test-

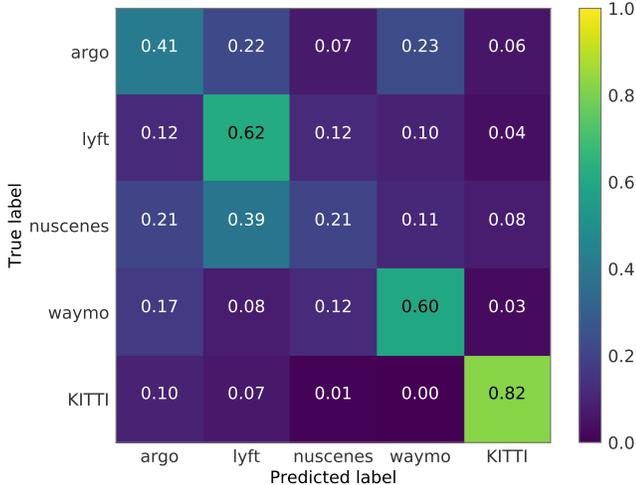


Figure S1: The confusion matrix of predicting which dataset an object belongs to using SVM with the RBF kernel. We take the (height, width, length) of car objects as inputs and the corresponding labels are the datasets the objects belong to.

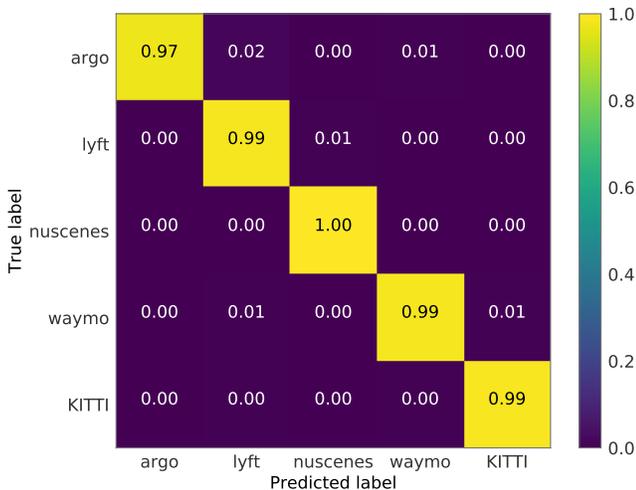


Figure S2: The confusion matrix of predicting which dataset a **Car** point cloud belongs to using PointNet++. We extract the points inside the ground truth **Car** bounding box as inputs, and the corresponding labels are the datasets the bounding boxes belong to.

ing examples. We show the confusion matrix in Figure S2. The classifier can almost perfectly classify the point clouds. Compared to Figure S1, we argue that not only the bounding box sizes, but also the point cloud styles (e.g., density, number of laser beams, etc) of cars contribute to dataset discrepancy. Interestingly, while the second factor seems to be more informative in differentiating datasets, the first factor is indeed the main cause of poor transfer performance

among datasets<sup>1</sup>.

## S4. PIXOR Results

We report object detection results using PIXOR [9], which takes voxelized tensors instead of point clouds as input. We implement the algorithm ourselves, achieving comparable results as the ones in [9]. We report the results in Table S5. PIXOR performs pretty well if the model is trained and tested within the same dataset, suggesting that its model design does not over-fit to KITTI.

We also see a clear performance drop when we train a model on one dataset and test it on the other datasets. The drop is more severe than applying the POINTRCNN detector in many cases. We surmise that the re-sampling operation used in POINTRCNN might make the difference. We therefore apply the same re-sampling operation on the input point cloud before inputting it to PIXOR. Table S6 shows the results: re-sampling does improve the performance in applying the Waymo detector to other datasets. This is likely because Waymo has the most LiDAR points on average and re-sampling reduces the number, making it more similar to that of the other datasets. We expect that tuning the number of points in re-sampling can further boost the performance.

## S5. Additional Results Using POINTRCNN

### S5.1. Complete tables

We show the complete tables across five datasets by replacing the predicted box sizes with the ground truth sizes (cf. section 4.5 in the main paper) in Table S7, by few-shot fine-tuning in Table S8, by statistical normalization in Table S9, and by output transformation in Table S10. For statistical normalization and output transformation, we see smaller improvements (or even some degradation) among datasets collected in the USA than between datasets collected in Germany and the USA.

### S5.2. Online sales data

In the main paper, for statistical normalization we leverage the average car size of each dataset. Here we collect car sales data from Germany and the USA in the past four years. The average car size  $(h, w, l)$  is  $(1.75, 1.93, 5.15)$  in the USA and  $(1.49, 1.79, 4.40)$  in Germany. The difference is  $(0.26, 0.14, 0.75)$ , not far from  $(0.20, 0.37, 0.78)$  between KITTI and the other datasets. The gap can be reduced by further considering locations (e.g., Argoverse from Miami and Pittsburgh, USA) and earlier data (KITTI was collected in 2011).

<sup>1</sup>As mentioned in the main paper, POINTRCNN applies point re-sampling so that every scene (in RPN) and object proposal (in RCNN) will have the same numbers of input points. Such an operation could reduce the point cloud differences across domains.

Table S5: 3D object detection across multiple datasets (evaluated on the validation sets). We report average precision (AP) of the *Car* category in the bird’s-eye view ( $AP_{BEV}$ ) at  $IoU = 0.7$ , using PIXOR [9]. We report results at different difficulties (following the KITTI benchmark, but we replace the 40, 25, 25 pixel thresholds on 2D bounding boxes with 30, 70, 70 meters on object depths, for *Easy*, *Moderate*, and *Hard* cases, respectively) and different depth ranges (using the same truncation and occlusion thresholds as KITTI *Hard* case). The results show a significant performance drop in cross-dataset inference. We indicate the best generalization results per column and per setting by red fonts and the worst by blue fonts. We indicate in-domain results by bold fonts.

Setting	Source\Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>87.2</b>	31.5	39.4	65.1	28.3
	Argoverse	57.4	<b>79.2</b>	49.0	89.4	69.3
	nuScenes	40.4	66.3	<b>56.8</b>	79.7	40.7
	Lyft	53.1	71.4	45.5	<b>90.7</b>	75.2
	Waymo	9.8	63.6	38.1	82.5	<b>87.2</b>
Moderate	KITTI	<b>72.8</b>	25.9	20.0	37.4	23.4
	Argoverse	42.5	<b>67.3</b>	24.4	57.9	55.0
	nuScenes	30.3	46.3	<b>30.1</b>	50.5	35.3
	Lyft	38.7	58.5	24.9	<b>78.1</b>	56.8
	Waymo	10.0	47.4	21.0	62.4	<b>76.4</b>
Hard	KITTI	<b>68.2</b>	28.4	18.8	34.5	24.1
	Argoverse	43.0	<b>64.7</b>	22.7	57.7	55.2
	nuScenes	27.1	46.0	<b>29.8</b>	50.6	35.8
	Lyft	35.8	54.3	24.8	<b>78.0</b>	53.7
	Waymo	11.1	46.9	21.3	63.4	<b>74.8</b>
0-30m	KITTI	<b>87.2</b>	39.5	38.9	62.2	32.1
	Argoverse	60.0	<b>82.4</b>	49.9	88.0	72.8
	nuScenes	38.7	63.0	<b>55.1</b>	79.4	43.9
	Lyft	50.7	73.5	48.4	<b>90.5</b>	76.4
	Waymo	12.9	65.7	42.7	83.1	<b>88.3</b>
30m-50m	KITTI	<b>50.3</b>	29.4	9.1	31.0	26.1
	Argoverse	23.7	<b>66.1</b>	0.8	54.5	56.5
	nuScenes	18.6	44.9	<b>12.3</b>	48.1	37.3
	Lyft	17.5	50.4	7.0	<b>77.0</b>	53.0
	Waymo	8.1	41.3	4.5	62.1	<b>78.0</b>
50m-70m	KITTI	<b>12.0</b>	3.0	3.0	9.0	10.1
	Argoverse	4.8	<b>31.7</b>	0.3	20.6	31.3
	nuScenes	9.1	13.4	<b>9.1</b>	21.6	20.9
	Lyft	6.5	19.1	9.1	<b>61.2</b>	29.9
	Waymo	1.7	20.6	9.1	39.7	<b>53.3</b>

In Table S11, we show the results of adapting a detector trained on KITTI to other datasets using statistical normalization with the car sales data:  $(\Delta h, \Delta w, \Delta l)$  is  $(0.26, 0.15, 0.75)$ . The performance is slightly worse than using the statistics of the datasets. Nevertheless, compared to directly applying the source domain detector, statistical normalization with the car sales data still shows notable improvements.

### S5.3. Pedestrian

We calculate the statistics of pedestrians, as in Table S12. There are smaller differences among datasets. We therefore

expect a smaller improvement by statistical normalization.

## S6. Qualitative Results

We further show qualitative results of statistical normalization refinement. We train a POINTRCNN detector on Waymo and test it on KITTI. We compare its car detection before and after statistical normalization refinement in Figure S3. Statistical normalization can not only improve the predicted bounding box sizes, but also reduce false positive rates.

Table S6: 3D object detection across multiple datasets (evaluated on the validation sets). The setting is exactly the same as Table S5, except that we perform POINTRCNN re-sampling on the input point cloud before applying the PIXOR detector.

Setting	Source\Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>85.9</b>	22.3	35.7	56.2	13.4
	Argoverse	59.4	<b>80.5</b>	47.1	89.3	66.5
	nuScenes	14.5	57.1	<b>66.2</b>	73.4	44.2
	Lyft	66.6	73.8	52.2	<b>90.7</b>	77.3
	Waymo	28.6	66.0	52.0	84.2	<b>86.7</b>
Moderate	KITTI	<b>70.3</b>	19.1	18.9	33.5	14.8
	Argoverse	43.0	<b>66.5</b>	24.1	57.9	52.6
	nuScenes	12.6	46.9	<b>36.5</b>	52.6	35.7
	Lyft	49.3	54.4	28.6	<b>79.4</b>	59.2
	Waymo	23.8	51.4	26.7	69.0	<b>77.1</b>
Hard	KITTI	<b>67.2</b>	20.0	17.4	33.1	15.0
	Argoverse	42.8	<b>63.8</b>	22.3	57.7	52.5
	nuScenes	14.4	44.6	<b>35.7</b>	53.1	36.0
	Lyft	45.5	54.5	27.6	<b>79.3</b>	58.5
	Waymo	24.0	54.2	26.4	70.3	<b>77.3</b>
0-30m	KITTI	<b>85.8</b>	28.6	33.2	56.6	14.7
	Argoverse	61.5	<b>82.7</b>	48.6	88.3	65.0
	nuScenes	20.2	61.5	<b>64.4</b>	75.4	48.4
	Lyft	62.9	71.9	54.3	<b>90.7</b>	78.3
	Waymo	31.0	65.9	55.4	85.9	<b>88.2</b>
30m-50m	KITTI	<b>48.8</b>	22.0	4.5	29.3	16.8
	Argoverse	21.1	<b>69.7</b>	2.3	55.1	54.6
	nuScenes	8.6	42.8	<b>15.9</b>	52.7	40.5
	Lyft	25.1	53.1	10.3	<b>78.6</b>	59.9
	Waymo	16.7	52.4	9.8	68.8	<b>78.8</b>
50m-70m	KITTI	<b>15.7</b>	3.4	0.4	7.5	12.0
	Argoverse	9.4	<b>29.5</b>	0.1	22.2	30.4
	nuScenes	0.7	12.8	<b>9.1</b>	23.1	16.4
	Lyft	7.3	18.8	3.0	<b>63.6</b>	30.8
	Waymo	2.3	23.5	9.1	45.1	<b>56.8</b>

Table S7: Cross-dataset performance by assigning ground-truth box sizes to detected cars while keeping their centers and rotations unchanged. We report  $AP_{BEV} / AP_{3D}$  of the *Car* category at  $IoU = 0.7$ , using POINTRCNN [8]. We indicate the best generalization results per column and per setting by red fonts and the worst by blue fonts. We indicate in-domain results by bold fonts.

Setting	Source \ Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>95.6 / 84.6</b>	80.5 / 65.7	66.5 / 33.5	89.8 / 74.8	90.3 / 77.1
	Argoverse	80.0 / 59.2	<b>83.1 / 77.3</b>	54.2 / 26.7	87.5 / 75.6	89.1 / 74.7
	nuScenes	80.5 / 63.9	77.4 / 52.5	<b>74.8 / 46.4</b>	89.4 / 65.2	85.6 / 62.9
	Lyft	83.9 / 58.4	80.2 / 67.2	65.2 / 29.2	<b>90.3 / 87.3</b>	89.9 / 73.9
	Waymo	86.1 / 78.2	79.2 / 72.7	63.1 / 30.0	88.3 / 86.1	<b>90.2 / 86.2</b>
Moderate	KITTI	<b>81.4 / 72.6</b>	64.5 / 50.9	35.0 / 18.2	74.6 / 54.3	79.4 / 63.0
	Argoverse	66.9 / 51.0	<b>73.6 / 60.1</b>	28.2 / 17.6	67.6 / 52.3	77.3 / 61.5
	nuScenes	61.4 / 47.3	59.0 / 36.2	<b>41.7 / 25.4</b>	72.4 / 45.1	69.2 / 50.6
	Lyft	71.4 / 49.4	68.5 / 49.3	34.6 / 17.4	<b>84.2 / 66.9</b>	79.7 / 64.7
	Waymo	73.7 / 60.6	68.0 / 54.9	30.8 / 18.4	75.0 / 63.2	<b>86.4 / 74.4</b>
Hard	KITTI	<b>82.5 / 71.9</b>	64.0 / 49.3	31.4 / 17.7	73.1 / 53.0	77.2 / 59.1
	Argoverse	65.6 / 52.5	<b>73.6 / 59.2</b>	27.5 / 16.6	65.3 / 52.2	75.8 / 58.4
	nuScenes	61.3 / 45.7	55.3 / 33.5	<b>40.9 / 25.4</b>	72.6 / 43.6	68.3 / 46.2
	Lyft	72.0 / 52.0	65.4 / 49.8	31.2 / 16.5	<b>84.8 / 67.2</b>	78.2 / 63.6
	Waymo	75.3 / 60.7	67.8 / 51.9	30.2 / 17.0	75.2 / 61.9	<b>80.8 / 68.9</b>
0-30m	KITTI	<b>89.2 / 86.7</b>	82.3 / 70.2	62.8 / 35.1	89.8 / 76.2	90.4 / 78.7
	Argoverse	83.3 / 68.7	<b>86.1 / 80.2</b>	56.8 / 31.9	88.3 / 77.3	89.8 / 78.1
	nuScenes	76.5 / 62.5	80.9 / 55.2	<b>74.2 / 49.1</b>	89.4 / 67.7	87.5 / 62.6
	Lyft	86.7 / 62.7	84.2 / 69.3	63.1 / 31.7	<b>90.5 / 88.5</b>	90.2 / 77.2
	Waymo	88.0 / 75.8	82.8 / 76.3	62.0 / 32.9	88.7 / 86.8	<b>90.5 / 88.1</b>
30m-50m	KITTI	<b>71.6 / 56.4</b>	63.1 / 40.4	11.1 / 9.1	74.3 / 52.9	80.4 / 64.3
	Argoverse	43.2 / 27.0	<b>74.5 / 53.9</b>	9.5 / 9.1	67.4 / 49.3	78.8 / 62.9
	nuScenes	37.1 / 25.0	49.0 / 18.3	<b>17.4 / 10.4</b>	71.0 / 42.2	75.4 / 50.5
	Lyft	52.8 / 31.4	61.9 / 35.6	11.3 / 9.1	<b>85.1 / 65.9</b>	80.4 / 65.4
	Waymo	57.6 / 38.5	63.5 / 45.8	10.0 / 9.1	75.4 / 62.5	<b>87.7 / 74.9</b>
50m-70m	KITTI	<b>30.8 / 15.1</b>	23.6 / 9.7	1.6 / 1.0	50.0 / 24.0	53.3 / 31.2
	Argoverse	13.7 / 10.1	<b>34.2 / 11.7</b>	0.5 / 0.0	37.2 / 19.8	51.9 / 31.7
	nuScenes	9.2 / 5.7	15.8 / 4.3	<b>9.8 / 9.1</b>	46.9 / 20.1	43.4 / 23.2
	Lyft	17.2 / 8.1	28.2 / 11.4	1.1 / 0.1	<b>64.2 / 39.9</b>	57.8 / 36.3
	Waymo	13.1 / 4.9	29.2 / 11.4	0.9 / 0.0	53.0 / 29.4	<b>65.9 / 45.2</b>

Table S8: Cross-dataset performance by few-shot fine-tuning using 10 labeled target domain instances (average over five rounds of experiments). We report  $AP_{BEV}/AP_{3D}$  of the *Car* category at  $IoU = 0.7$ , using POINTRCNN [8]. We indicate the best generalization results per column and per setting by red fonts and the worst by blue fonts. We indicate in-domain results by bold fonts.

Setting	Source \ Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>88.0 / 82.5</b>	75.8 / 49.2	54.7 / <b>21.7</b>	<b>89.0</b> / 78.1	87.4 / <b>70.9</b>
	Argoverse	<b>80.0</b> / <b>49.7</b>	<b>74.2</b> / <b>42.0</b>	54.0 / 19.2	<b>86.6</b> / <b>63.5</b>	86.6 / <b>56.3</b>
	nuScenes	83.8 / 58.7	<b>68.7</b> / <b>33.7</b>	<b>73.4</b> / <b>38.1</b>	88.4 / 67.7	<b>84.3</b> / 59.8
	Lyft	<b>85.3</b> / <b>72.5</b>	73.5 / 48.9	<b>56.5</b> / 17.7	<b>90.2</b> / <b>87.3</b>	<b>89.1</b> / 70.4
	Waymo	81.0 / 67.0	<b>76.9</b> / <b>55.2</b>	<b>51.0</b> / <b>16.7</b>	88.3 / <b>81.0</b>	<b>90.1</b> / <b>85.3</b>
Moderate	KITTI	<b>80.6</b> / <b>68.9</b>	60.7 / 37.3	28.7 / <b>12.5</b>	74.2 / 53.4	75.9 / 55.3
	Argoverse	68.8 / <b>42.8</b>	<b>66.5</b> / <b>34.4</b>	27.5 / 11.2	<b>65.4</b> / <b>40.2</b>	75.3 / <b>46.7</b>
	nuScenes	67.2 / 45.5	<b>54.5</b> / <b>24.2</b>	<b>40.7</b> / <b>21.2</b>	71.9 / 44.0	<b>72.8</b> / 47.0
	Lyft	<b>73.9</b> / <b>56.2</b>	61.0 / 35.3	<b>30.3</b> / <b>10.6</b>	<b>83.7</b> / <b>65.5</b>	<b>78.3</b> / <b>57.9</b>
	Waymo	<b>66.8</b> / 51.8	<b>65.7</b> / <b>41.8</b>	<b>26.7</b> / 11.0	<b>75.1</b> / <b>54.8</b>	<b>85.9</b> / <b>67.9</b>
Hard	KITTI	<b>81.9</b> / <b>66.7</b>	59.8 / 36.5	27.5 / <b>12.4</b>	71.8 / 52.9	70.1 / 54.4
	Argoverse	66.3 / <b>43.0</b>	<b>67.9</b> / <b>37.3</b>	26.9 / 11.8	<b>66.0</b> / <b>42.0</b>	70.3 / <b>43.9</b>
	nuScenes	<b>64.7</b> / 44.5	<b>52.0</b> / <b>23.4</b>	<b>40.2</b> / <b>20.5</b>	71.0 / 44.3	<b>68.7</b> / 44.3
	Lyft	<b>74.1</b> / <b>56.2</b>	61.9 / 37.0	<b>28.6</b> / <b>11.1</b>	<b>79.3</b> / <b>65.5</b>	<b>76.9</b> / <b>55.6</b>
	Waymo	68.1 / 52.9	<b>62.3</b> / <b>39.3</b>	<b>26.7</b> / 11.7	<b>74.7</b> / <b>55.2</b>	<b>80.4</b> / <b>67.7</b>
0-30m	KITTI	<b>88.8</b> / <b>84.9</b>	73.6 / 55.2	54.0 / <b>23.6</b>	<b>89.3</b> / 77.6	88.7 / 74.1
	Argoverse	84.0 / <b>56.9</b>	<b>81.2</b> / <b>52.2</b>	54.0 / 22.6	<b>87.7</b> / <b>68.7</b>	88.3 / <b>60.7</b>
	nuScenes	<b>81.2</b> / 59.8	<b>70.5</b> / <b>40.1</b>	<b>73.2</b> / <b>42.8</b>	88.8 / 69.6	<b>86.2</b> / 62.4
	Lyft	<b>87.5</b> / <b>73.9</b>	78.1 / 54.3	<b>56.9</b> / 21.2	<b>90.4</b> / <b>88.5</b>	<b>89.4</b> / <b>74.8</b>
	Waymo	84.8 / 71.0	<b>79.4</b> / <b>56.6</b>	<b>52.8</b> / <b>20.8</b>	88.8 / <b>79.1</b>	<b>90.4</b> / <b>87.2</b>
30m-50m	KITTI	<b>70.2</b> / <b>51.4</b>	59.0 / 29.9	<b>9.5</b> / 6.1	73.7 / 50.4	78.1 / 57.2
	Argoverse	47.9 / <b>23.8</b>	<b>70.8</b> / <b>34.0</b>	7.3 / <b>2.0</b>	<b>65.4</b> / <b>36.9</b>	78.1 / 48.5
	nuScenes	<b>45.0</b> / 25.1	<b>51.4</b> / <b>17.1</b>	<b>17.1</b> / <b>4.1</b>	71.5 / 41.5	<b>74.2</b> / <b>48.0</b>
	Lyft	<b>57.7</b> / <b>33.3</b>	<b>62.4</b> / 29.5	<b>6.5</b> / 3.3	<b>83.8</b> / <b>62.7</b>	<b>79.7</b> / <b>59.9</b>
	Waymo	49.2 / 29.2	60.6 / <b>34.7</b>	9.4 / <b>6.3</b>	<b>75.1</b> / <b>52.6</b>	<b>87.5</b> / <b>68.8</b>
50m-70m	KITTI	<b>28.8</b> / <b>12.0</b>	20.1 / 6.3	<b>3.3</b> / <b>1.2</b>	46.8 / 19.4	45.2 / 24.3
	Argoverse	<b>8.1</b> / <b>3.8</b>	<b>33.0</b> / <b>12.7</b>	<b>0.4</b> / <b>0.0</b>	<b>38.0</b> / <b>10.3</b>	51.1 / 23.4
	nuScenes	12.9 / 5.7	<b>15.5</b> / <b>2.6</b>	<b>9.1</b> / <b>9.1</b>	47.0 / 14.9	<b>44.3</b> / <b>19.3</b>
	Lyft	<b>17.5</b> / <b>8.0</b>	26.8 / <b>9.1</b>	2.5 / 0.0	<b>62.7</b> / <b>33.1</b>	<b>54.0</b> / <b>27.2</b>
	Waymo	10.5 / 4.8	<b>27.6</b> / 7.3	1.3 / 0.0	<b>51.2</b> / <b>19.9</b>	<b>63.5</b> / <b>41.1</b>

Table S9: Cross-dataset performance by fine-tuning with source data after statistical normalization. We report  $AP_{BEV}/AP_{3D}$  of the *Car* category at  $IoU = 0.7$ , using POINTRCNN [8]. We indicate the best generalization results per column and per setting by red fonts and the worst by blue fonts. We indicate in-domain results by bold fonts.

Setting	Source\Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>88.0 / 82.5</b>	74.7 / 48.2	60.8 / 23.9	88.3 / 73.3	84.6 / 53.3
	Argoverse	76.2 / 46.1	<b>79.2 / 57.8</b>	48.3 / 18.6	84.8 / 65.0	84.8 / 49.2
	nuScenes	83.2 / 35.6	72.0 / 25.3	<b>73.4 / 38.1</b>	88.7 / 38.1	76.6 / 43.3
	Lyft	83.5 / 72.1	74.4 / 44.0	57.8 / 21.1	<b>90.2 / 87.3</b>	86.3 / 66.4
	Waymo	82.1 / 48.7	75.0 / 44.4	54.9 / 20.7	85.7 / 80.0	<b>90.1 / 85.3</b>
Moderate	KITTI	<b>80.6 / 68.9</b>	61.5 / 38.2	32.9 / 16.4	73.7 / 53.1	74.9 / 49.4
	Argoverse	67.2 / 40.5	<b>69.9 / 44.2</b>	24.7 / 11.1	63.3 / 38.9	72.0 / 43.6
	nuScenes	67.4 / 31.0	55.6 / 17.9	<b>40.7 / 21.2</b>	71.1 / 24.5	66.6 / 32.2
	Lyft	73.6 / 57.9	59.7 / 33.3	30.4 / 10.9	<b>83.7 / 65.5</b>	75.5 / 51.3
	Waymo	71.3 / 47.1	62.3 / 31.7	28.8 / 11.5	71.5 / 52.6	<b>85.9 / 67.9</b>
Hard	KITTI	<b>81.9 / 66.7</b>	60.6 / 37.1	31.9 / 15.8	73.1 / 53.5	69.4 / 49.4
	Argoverse	68.5 / 41.9	<b>69.9 / 42.8</b>	24.3 / 10.9	61.6 / 40.2	68.2 / 42.7
	nuScenes	65.2 / 30.8	52.5 / 17.2	<b>40.2 / 20.5</b>	67.3 / 28.6	65.7 / 30.4
	Lyft	75.2 / 58.9	60.8 / 31.8	29.5 / 14.4	<b>79.3 / 65.5</b>	75.5 / 53.2
	Waymo	73.0 / 49.7	60.2 / 32.5	28.4 / 10.9	71.6 / 53.3	<b>80.4 / 67.7</b>
0-30m	KITTI	<b>88.8 / 84.9</b>	73.1 / 54.2	60.0 / 29.2	88.8 / 75.4	87.1 / 60.1
	Argoverse	83.3 / 53.9	<b>83.3 / 63.3</b>	51.5 / 23.0	86.3 / 68.4	87.3 / 59.7
	nuScenes	83.6 / 42.8	72.8 / 27.2	<b>73.2 / 42.8</b>	88.9 / 47.1	78.5 / 45.9
	Lyft	87.4 / 73.6	78.7 / 51.8	58.7 / 26.8	<b>90.4 / 88.5</b>	87.9 / 72.4
	Waymo	85.7 / 59.0	79.9 / 50.5	57.6 / 24.3	87.2 / 75.8	<b>90.4 / 87.2</b>
30m-50m	KITTI	<b>70.2 / 51.4</b>	61.5 / 31.5	11.0 / 2.3	73.8 / 52.2	78.1 / 54.9
	Argoverse	48.9 / 25.7	<b>72.2 / 39.5</b>	5.0 / 4.5	61.0 / 32.4	74.4 / 46.2
	nuScenes	44.9 / 18.6	45.6 / 7.3	<b>17.1 / 4.1</b>	70.1 / 18.1	67.9 / 31.6
	Lyft	58.3 / 38.0	57.2 / 18.5	6.5 / 4.5	<b>83.8 / 62.7</b>	77.2 / 52.4
	Waymo	57.3 / 36.3	54.9 / 20.1	9.1 / 1.5	71.3 / 48.4	<b>87.5 / 68.8</b>
50m-70m	KITTI	<b>28.8 / 12.0</b>	23.8 / 5.6	3.0 / 2.3	49.9 / 22.2	46.8 / 25.1
	Argoverse	9.1 / 2.6	<b>29.9 / 6.9</b>	0.2 / 0.1	28.9 / 8.8	46.2 / 21.2
	nuScenes	9.4 / 5.1	14.8 / 2.3	<b>9.1 / 9.1</b>	40.7 / 5.2	36.4 / 14.9
	Lyft	21.1 / 6.7	21.2 / 4.9	4.5 / 0.0	<b>62.7 / 33.1</b>	52.1 / 25.3
	Waymo	14.4 / 5.7	27.7 / 11.0	1.0 / 0.0	46.9 / 22.0	<b>63.5 / 41.1</b>

Table S10: Cross-dataset performance by output transformation: directly adjusting the predicted box size by adding the difference of mean sizes between domains. We report  $AP_{BEV} / AP_{3D}$  of the *Car* category at  $IoU = 0.7$ , using POINTRCNN [8]. We indicate the best generalization results per column and per setting by red fonts and the worst by blue fonts. We indicate in-domain results by bold fonts.

Setting	Source \ Target	KITTI	Argoverse	nuScenes	Lyft	Waymo
Easy	KITTI	<b>88.0 / 82.5</b>	72.7 / 9.0	55.0 / 10.4	88.2 / 23.5	86.1 / 16.2
	Argoverse	53.3 / 5.7	<b>79.2 / 57.8</b>	52.6 / 21.3	87.1 / 66.1	87.6 / 56.1
	nuScenes	75.4 / 31.5	73.3 / 27.9	<b>73.4 / 38.1</b>	89.2 / 44.3	78.4 / 35.5
	Lyft	71.9 / 4.7	77.1 / 48.0	63.1 / 24.5	<b>90.2 / 87.3</b>	89.2 / 73.9
	Waymo	64.0 / 3.9	74.3 / 54.8	58.8 / 25.2	88.3 / 85.3	<b>90.1 / 85.3</b>
Moderate	KITTI	<b>80.6 / 68.9</b>	59.9 / 7.9	30.8 / 6.8	70.1 / 17.8	69.1 / 13.1
	Argoverse	52.2 / 7.3	<b>69.9 / 44.2</b>	27.5 / 11.7	66.9 / 42.1	74.3 / 45.5
	nuScenes	58.5 / 27.3	56.8 / 20.4	<b>40.7 / 21.2</b>	71.3 / 27.3	67.8 / 26.2
	Lyft	60.8 / 5.6	62.7 / 37.6	33.5 / 12.5	<b>83.7 / 65.5</b>	78.4 / 60.8
	Waymo	54.9 / 3.7	62.9 / 40.4	30.1 / 14.5	74.3 / 59.8	<b>85.9 / 67.9</b>
Hard	KITTI	<b>81.9 / 66.7</b>	59.3 / 9.3	27.8 / 7.6	66.5 / 19.1	68.7 / 13.9
	Argoverse	53.5 / 8.6	<b>69.9 / 42.8</b>	26.7 / 14.5	64.6 / 43.0	70.0 / 44.2
	nuScenes	59.5 / 27.8	53.6 / 19.9	<b>40.2 / 20.5</b>	67.6 / 28.5	66.3 / 26.0
	Lyft	63.1 / 6.9	63.4 / 38.6	30.4 / 13.3	<b>79.3 / 65.5</b>	77.3 / 57.3
	Waymo	58.0 / 4.1	60.5 / 39.2	29.4 / 14.6	74.0 / 57.2	<b>80.4 / 67.7</b>
0-30m	KITTI	<b>88.8 / 84.9</b>	73.0 / 13.7	56.2 / 13.9	88.4 / 27.5	87.7 / 22.2
	Argoverse	64.9 / 10.1	<b>83.3 / 63.3</b>	55.2 / 27.0	87.8 / 69.9	87.9 / 62.6
	nuScenes	74.6 / 36.6	73.7 / 32.0	<b>73.2 / 42.8</b>	89.2 / 46.2	79.6 / 41.6
	Lyft	74.8 / 9.1	81.2 / 55.8	61.2 / 27.2	<b>90.4 / 88.5</b>	89.6 / 77.2
	Waymo	71.3 / 4.4	78.4 / 55.7	60.5 / 25.8	88.7 / 85.0	<b>90.4 / 87.2</b>
30m-50m	KITTI	<b>70.2 / 51.4</b>	56.1 / 5.4	10.8 / 9.1	67.4 / 10.7	73.6 / 10.4
	Argoverse	35.1 / 9.1	<b>72.2 / 39.5</b>	9.5 / 0.3	66.3 / 39.1	77.5 / 44.9
	nuScenes	35.5 / 15.5	47.4 / 7.8	<b>17.1 / 4.1</b>	69.9 / 22.5	68.7 / 21.1
	Lyft	43.3 / 3.9	60.8 / 25.4	11.2 / 9.1	<b>83.8 / 62.7</b>	79.5 / 61.4
	Waymo	39.8 / 4.5	58.1 / 34.9	9.9 / 9.1	74.5 / 57.5	<b>87.5 / 68.8</b>
50m-70m	KITTI	<b>28.8 / 12.0</b>	20.5 / 1.0	1.5 / 1.0	41.3 / 6.8	42.6 / 4.2
	Argoverse	8.0 / 0.8	<b>29.9 / 6.9</b>	0.5 / 0.0	35.6 / 14.2	49.2 / 20.3
	nuScenes	7.8 / 5.1	15.3 / 3.0	<b>9.1 / 9.1</b>	41.4 / 5.6	37.0 / 12.0
	Lyft	12.7 / 0.9	25.6 / 6.0	1.1 / 0.0	<b>62.7 / 33.1</b>	54.9 / 30.4
	Waymo	7.7 / 1.1	25.5 / 6.5	0.9 / 0.0	50.8 / 22.3	<b>63.5 / 41.1</b>

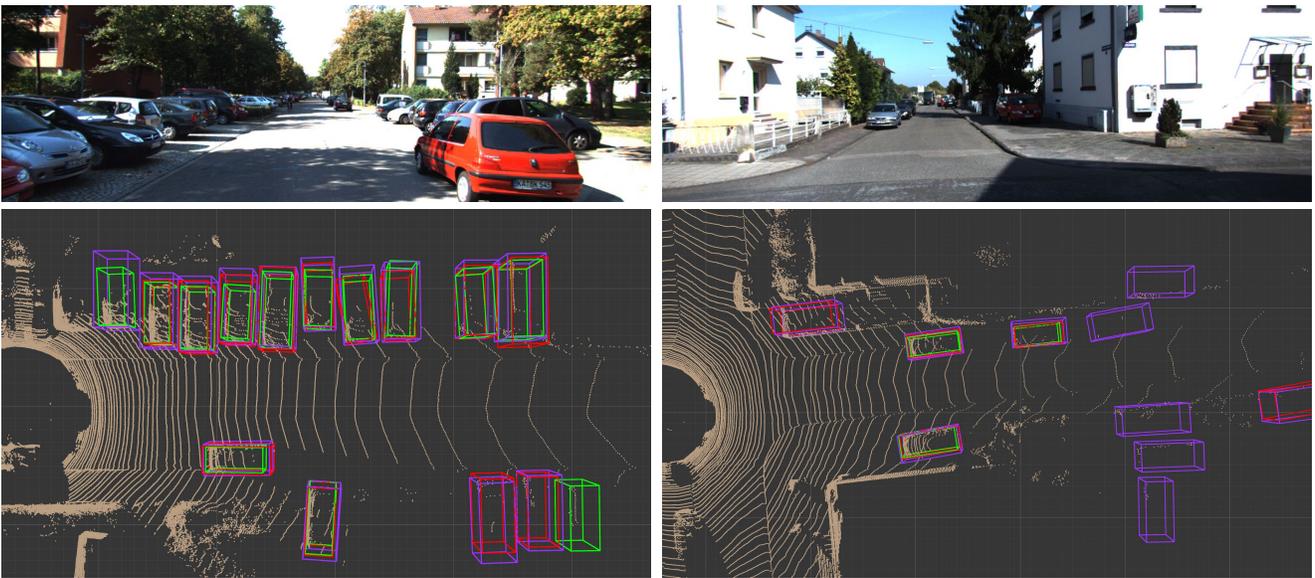


Figure S3: 3D prediction on KITTI using the model trained on Waymo before and after statistical normalization refinement. The green boxes are the ground truth *car* bounding boxes. Purple boxes and magenta boxes are predictions by the model before and after statistical normalization refinement, respectively. The left column demonstrates that statistical normalization is able to resize bounding box predictions to the correct sizes, while the right case shows that it also can reduce false positive rates.

Table S11: Statistical normalization using the mean sizes of datasets versus car sales data. Direct: directly applying the source domain detector.

Setting	Dataset	From KITTI (KITTI as the source)		
		Direct	Datasets	Car sales data
Easy	Argoverse	55.8 / 27.7	74.7 / 48.2	68.6 / 32.8
	nuScenes	47.4 / 13.3	60.8 / 23.9	62.0 / 24.4
	Lyft	81.7 / 51.8	88.3 / 73.3	88.9 / 69.9
	Waymo	45.2 / 11.9	84.6 / 53.3	66.7 / 22.8
Moderate	Argoverse	44.9 / 22.3	61.5 / 38.2	57.7 / 29.1
	nuScenes	26.2 / 8.3	32.9 / 16.4	32.6 / 13.0
	Lyft	61.8 / 33.7	73.7 / 53.1	72.6 / 47.6
	Waymo	43.9 / 12.3	74.9 / 49.4	61.8 / 22.9
Hard	Argoverse	42.5 / 22.2	60.6 / 37.1	54.0 / 30.0
	nuScenes	24.9 / 8.8	31.9 / 15.8	29.8 / 13.2
	Lyft	57.4 / 34.2	73.1 / 53.5	71.7 / 45.7
	Waymo	41.5 / 12.6	69.4 / 49.4	62.7 / 25.1
0-30m	Argoverse	58.4 / 34.7	73.1 / 54.2	71.0 / 44.0
	nuScenes	47.9 / 14.9	60.0 / 29.2	60.1 / 26.1
	Lyft	77.8 / 54.2	88.8 / 75.4	89.2 / 72.5
	Waymo	48.0 / 14.0	87.1 / 60.1	72.4 / 30.2
30m-50m	Argoverse	46.5 / 19.0	61.5 / 31.5	57.4 / 20.0
	nuScenes	9.8 / 4.5	11.0 / 2.3	5.7 / 3.0
	Lyft	60.1 / 34.5	73.8 / 52.2	72.2 / 42.7
	Waymo	50.5 / 21.4	78.1 / 54.9	66.8 / 35.5
50m-70m	Argoverse	9.2 / 3.0	23.8 / 5.6	16.8 / 4.5
	nuScenes	1.1 / 0.0	3.0 / 2.3	1.0 / 0.1
	Lyft	33.2 / 9.6	49.9 / 22.2	46.0 / 18.8
	Waymo	27.1 / 12.0	46.8 / 25.1	44.2 / 18.0

Table S12: Dataset statistics on pedestrians (meters)

	KITTI	Argoverse	nuScenes	Lyft	Waymo
H	1.76±0.11	1.84±0.15	1.78±0.18	1.76±0.18	1.75±0.20
W	0.66±0.14	0.78±0.14	0.67±0.14	0.76±0.14	0.85±0.15
L	0.84±0.23	0.78±0.14	0.73±0.19	0.78±0.17	0.90±0.19

## References

- [1] Waymo open dataset: An autonomous driving dataset, 2019. [1](#)
- [2] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019. [1](#)
- [3] Ming-Fang Chang, John W Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, and James Hays. Argoverse: 3d tracking and forecasting with rich maps. In *CVPR*, 2019. [1](#)
- [4] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. [1](#)
- [5] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. [1](#)
- [6] R. Kesten, M. Usman, J. Houston, T. Pandya, K. Nadhamuni, A. Ferreira, M. Yuan, B. Low, A. Jain, P. Ondruska, S. Omari, S. Shah, A. Kulkarni, A. Kazakova, C. Tao, L. Platinsky, W. Jiang, and V. Shet. Lyft level 5 av dataset 2019. [url=https://level5.lyft.com/dataset/](https://level5.lyft.com/dataset/), 2019. [1](#)
- [7] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NIPS*, 2017. [3](#)
- [8] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointnet: 3d object proposal generation and detection from point cloud. In *CVPR*, 2019. [3](#), [7](#), [8](#), [9](#), [10](#)
- [9] Bin Yang, Wenjie Luo, and Raquel Urtasun. Pixor: Real-time 3d object detection from point clouds. In *CVPR*, 2018. [1](#), [4](#), [5](#)