# Semi-Supervised Domain Adaptation Using Target-Oriented Domain Augmentation for 3D Object Detection

Yecheol Kim[1*], Junho Lee[1*], Changsoo Park[2], Hyoung won Kim[2], Inho Lim[2], Christopher Chang[2], and Jun Won Choi[3]

*Abstract*—3D object detection is crucial for applications like autonomous driving and robotics. However, in real-world environments, variations in sensor data distribution due to sensor upgrades, weather changes, and geographic differences can adversely affect detection performance. Semi-Supervised Domain Adaptation (SSDA) aims to mitigate these challenges by transferring knowledge from a source domain, abundant in labeled data, to a target domain where labels are scarce. This paper presents a new SSDA method referred to as Target-Oriented Domain Augmentation (TODA) specifically tailored for LiDAR-based 3D object detection. TODA efficiently utilizes all available data, including labeled data in the source domain, and both labeled data and unlabeled data in the target domain to enhance domain adaptation performance. TODA consists of two stages: TargetMix and AdvMix. TargetMix employs mixing augmentation accounting for LiDAR sensor characteristics to facilitate feature alignment between the source-domain and target-domain. AdvMix applies point-wise adversarial augmentation with mixing augmentation, which perturbs the unlabeled data to align the features within both labeled and unlabeled data in the target domain. Our experiments conducted on the challenging domain adaptation tasks demonstrate that TODA outperforms existing domain adaptation techniques designed for 3D object detection by significant margins. The code is available at: https://github.com/rasd3/TODA.

*Index Terms*—Autonomous driving, 3D object detection, semi-supervised domain adaptation.

Fig. 1. **Performance evaluation in a domain adaptation task from Waymo dataset to nuScenes dataset:** 0.5%, 1%, and 5% labeled data in the target domain are used. A SSDA method using only 0.5% of the target label results in a remarkable performance gain over a UDA method (ST3D [8]). Our TODA also significantly outperforms SSDA3D [9] in all settings. Surprisingly, TODA even surpasses the *Oracle* performance with only 5% labels.

## I. INTRODUCTION

**3**D object detection is the task of detecting and localizing objects in 3D world coordinates using sensor measurements. 3D object detection has risen as a pivotal perception task in the field of autonomous vehicles and robotics. Recently, 3D point cloud data acquired by LiDAR sensor has been successfully used to achieve promising performance in 3D object detection. The recent progress of deep learning has sparked the development of a plethora of architectures for detecting objects from LiDAR point cloud. Widely used 3D object detectors include VoxelNet [1], PointPillar [2], SECOND [3], CenterPoint [4], PV-RCNN [5], PillarNet [6], and Voxel R-CNN [7].

A shift in the distribution of data often leads to notable decreases in the performance of 3D object detection [10], [11]. In the context of autonomous driving, shifts in distribution arise from change in sensor suites, fluctuations in weather conditions, disparities in geographical locations, and more. For instance, upgrades in sensor specifications, including resolution, field of view (FOV), and intensity, introduce shifts in data distribution. In this case, it is inefficient to collect new training data and retrain the model from scratch with each sensor replacement. Therefore, addressing the domain shift problem in 3D object detection is crucial for the commercial deployment of autonomous driving and still remains a significant open challenge.

Domain adaptation offers a solution to this problem by allowing models trained in source domains to adapt effectively to different but related target domains, thereby reducing the need for extensive data labeling. Domain adaptation strategies can be categorized into two main types: *unsupervised domain adaptation* (UDA) and *semi-supervised domain adaptation*

[1]Y. Kim and J. Lee are with Department of Electrical Engineering, Hanyang University, 04753 Seoul, Republic of Korea. (e-mail: yckim@spa.hanyang.ac.kr, jhlee@spa.hanyang.ac.kr)

[2]C. Park, H. Kim, I. Lim and C. Chang are with Kakao Mobility Corp, Pangyo 13529, Republic of Korea. (e-mail: teddy.p@kakaomobility.com, gemini.k@kakaomobility.com, ed.lim@kakaomobility.com, cswchang@alumni.caltech.edu)

[3]J. W. Choi is with Department of Electrical and Computer Engineering, College of Liberal Studies, Seoul National University, Seoul, 08826, Korea. (e-mail: junwchoi@snu.ac.kr)*(Corresponding author: Jun Won Choi)*
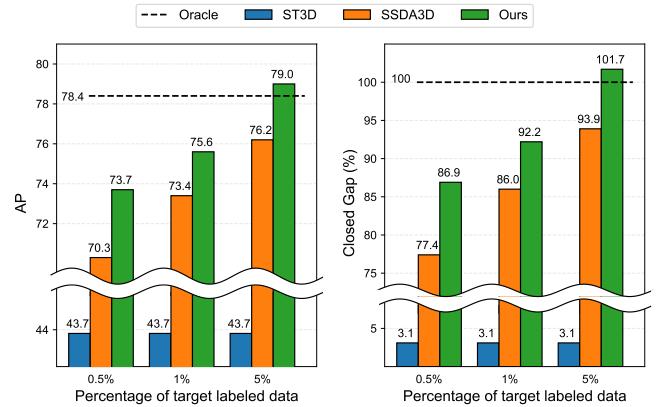
*denotes equal contribution.

(SSDA).

When labeled data is unavailable within the target domain, UDA transfers knowledge learned from labeled source domains to enhance performance in target domains. Recent UDA approaches for 3D object detection include CL3D [12], ST3D [8], and ST3D++ [13]. These methods have employed various domain-adaptive pseudo-labeling approaches to mitigate domain discrepancies.

While UDA can mitigate domain shift issues, addressing significant domain gaps between source and target domains remains challenging. While UDA is effective in narrowing the performance gap from the *Oracle*[1] by approximately 80% when applied to similar LiDAR specifications from the Waymo dataset [14] to the KITTI dataset [15], it achieves only a 3% reduction in the gap when dealing with markedly different LiDAR configurations, such as those from the Waymo dataset to the nuScenes dataset [16]. To address this limitation, semi-supervised domain adaptation (SSDA) has emerged as a cost-effective way to improve the effect of domain adaptation. Unlike UDA, SSDA uses a small amount of labeled target-domain data along with a substantial volume of unlabeled target-domain data to improve performance of domain adaptation. Fig. 1 illustrates that leveraging a small amount of labeled target-domain data, SSDA methods can yield substantial performance improvements over UDA.

SSDA3D [9], as of now, is the sole existing SSDA method specifically designed for 3D object detection. This method operates in a two-stage process. Initially, it incorporates an *Inter-domain Point-Cutmix* operation to reduce domain bias and thus learns domain-invariant representations. Following this, SSDA3D employs an *Intra-domain Point-MixUp* operation, which combines both labeled data and pseudo-labeled scenes on a global scale under a semi-supervised learning (SSL) framework. Despite its notable performance gains over UDA methods, we claim that SSDA3D has not fully exploited the distinct properties inherent in LiDAR point cloud data.

In this study, we present a novel SSDA framework referred to as *Target-Oriented Domain Augmentation* (TODA) for 3D object detection. TODA employs a two-stage data augmentation strategy for SSDA; *Target Sensor-Guided Mix Augmentation* (TargetMix) and *Adversarial-Guided Mix Augmentation* (AdvMix).

TargetMix reduces the disparity between the source and target domains by employing a cross-domain mixup strategy. Initially, TargetMix aligns the characteristics of LiDAR point clouds, such as Field of View (FOV) and beam configurations in the source-domain data with those in the target-domain data. Subsequently, a cross-domain mixup augmentation is applied in a polar coordinate system, incorporating an effective LiDAR distribution matching process that considers the scanning mechanism of LiDAR technology. By generating convex combinations of LiDAR point clouds from both source and target domains, TargetMix ensures smooth transitions between these domains.

TargetMix does not utilize the unlabeled data in the target domain, leaving room for further improvement. AdvMix employs a pseudo-labeling approach to leverage this unlabeled data. However, as pointed out in [17], the pseudo-labeling approach might suffer from intra-domain discrepancy, which arises when the teacher model, trained solely with labeled data, yields feature points divided into those attracted into the source domain and those that do not. Such inconsistency within the target domain data diminishes the effectiveness of pseudo labeling. To address this issue, we introduce AdvMix, a technique that integrates adversarial point augmentation with mixup augmentation. *Adversarial Point Augmentation*, derived from *Adversarial Augmentation* in [18], involves perturbing unlabeled data at the point level. These perturbations are informed by the negative gradient of the detection loss calculated in terms of the unlabeled points, effectively altering their distribution to enhance detection performance. This process generates perturbed pseudo-labeled samples that help minimize intra-domain discrepancies. Following this, we implement mixup augmentation, blending the labeled data with the pseudo-labeled data on a global scale.

We evaluate the performance of TODA on the challenging domain adapation tasks: from the Waymo dataset [14] to the nuScenes dataset [16] and from the nuScenes dataset to the KITTI dataset [15]. Our results demonstrate that TODA yields substantial performance improvements compared to the baseline method. Moreover, TODA outperforms existing domain adaptation techniques, including SSDA3D, by considerable margins.

The key contributions of TODA are summarized as follows:

- We propose a novel two-stage SSDA framework for 3D object detection.
- We present two novel data augmentation techniques, TargetMix and AdvMix, designed to improve the effectiveness of domain adaptation. TargetMix aims to reduce the domain disparity by leveraging labeled data from both the source and target domains. In contrast, AdvMix uses a combination of labeled and unlabeled data within the target domain to minimize the intra-domain gap.
- TargetMix method is specifically designed to take into account the unique characteristics of LiDAR point cloud data, thereby maximizing the impact of data augmentation.
- AdvMix generates adversarial examples to align the representation of unlabeled data effectively within the target domain. In our study, we are the first to pioneer the use of adversarial augmentation for point cloud-based SSDA.
- The proposed TODA achieves state-of-the-art performance on popular domain adaptation benchmarks. It attains performances on par with the *Oracle* performance even when utilizing merely 5% of labeled data in the target domain.

## II. RELATED WORK

### A. LiDAR-based 3D Object Detection

The advancement of deep learning has led to the development of various LiDAR-based 3D object detectors. These detectors
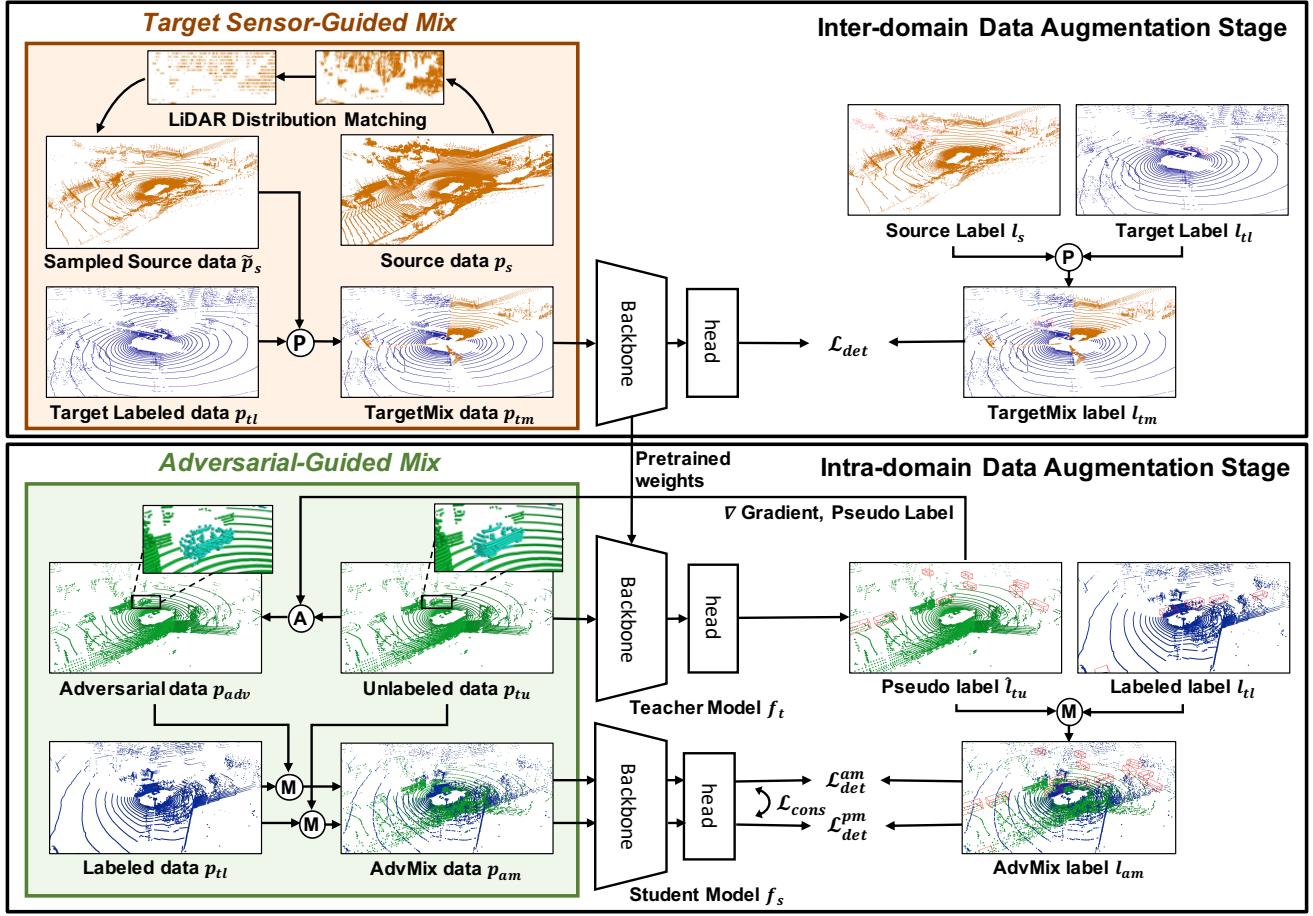
---

[1]The *Oracle* model denotes the fully-supervised model trained on the target domain.

Fig. 2. **Overall architecture of the proposed TODA:** First, TargetMix aligns the source-domain data with target-domain data by applying LiDAR Distribution Matching, followed by mixup augmentation in polar coordinates. Then, AdvMix utilizes Adversarial Point Augmentation to perturb the unlabeled data in the target domain, aiming to produce consistent representation of both labeled and unlabeled data. 'P', 'A', and 'M' denote *Polar Coordinate-based Mix*, *Adversarial Point Augmentation*, and *Point-Mixup* respectively.

utilize point encoding to extract high-level semantic features. Point encoding techniques can generally be categorized into two types: clustering-based and grid-based methods. Clustering-based methods group point clouds into clusters and encode points within each cluster in a hierarchical manner. Examples of clustering-based methods include PointRCNN [19], 3DDS [20], and PV-RCNN [5]. In contrast, grid-based encoding methods divide 3D point clouds into voxel or pillar grids, thereby producing features in a structured grid pattern. 3D object detectors that use grid-based encoding include VoxelNet [1], SECOND [3], PointPillar [2], and CenterPoint [4].

### B. Domain Adaptation for 3D Object Detection

Several UDA methods have proposed for 3D object detection. PointDAN [21] and SRDAN [22] transformed object features taking factors like the range and distance of point clouds into account. ST3D [8] and ST3D++ [13] tackled the challenge of learning instability induced by domain shift through the utilization of a memory bank and denoising techniques. Currently, only one SSDA method available in the literature is SSDA3D [9]. SSDA3D employed mixing augmentation strategies to reduce the distribution gap between the source and target domains.

### C. Mixing Augmentation for Domain Adaptation

Mixing augmentation has emerged as an effective strategy for addressing data limitations and enhancing robustness in semi-supervised learning (SSL) and domain adaptation. Mixup [23] generated new training data by creating convex combinations of input pairs, subsequently training the model on these blended inputs and their corresponding labels. CutMix [24], a variant of Mixup for image recognition, constructed new images by randomly cutting a rectangular patch from one image and pasting it onto another at a random location. These methods encouraged the model to learn interpolated decision boundaries between samples and hence improved generalization performance significantly.

In domain adaptation, mix augmentation techniques are employed to bridge inter-domain gaps by integrating data from both the source and target domains, thereby facilitating knowledge transfer from the source to the target domain. PolarMix [25] introduces a specialized mix augmentation

strategy for point cloud data that combines scene-level and object-level features in cylindrical coordinates.

### D. Adversarial Augmentation

Adversarial augmentation is a method for generating *adversarial examples* that are intended to challenge the accuracy of a model. Recently, AT [26] and VAT [27] have utilized a gradient-based perturbation generation method to enhance the robustness of models in both supervised and semi-supervised tasks. Furthermore, several studies [28]–[31] have explored using adversarial augmentation strategies to generate and translate point clouds.

### E. Semi-Supervised Learning for 3D Object Detection

SSL has also been actively studied for 3D object detection. SSL involves using a combination of a small amount of labeled data and a large amount of unlabeled data for training models. SESS [32] introduced a *Consistency Loss* as a means to align predicted 3D object proposals between the teacher and student networks. 3DIoUMatch [33] filtered low-quality pseudo-labels using the combination of confidence thresholds and 3D IoU predictions. Proficient Teachers [34] enhanced the precision of pseudo-labels through an augmented prediction approach that incorporates box voting-based ensembling.

## III. METHOD

In this section, we present the details of the proposed TODA method.

### A. Overview

In the SSDA framework, we train the model using labeled data from the source domain, as well as both labeled and unlabeled data from the target domain. First, the $N_S$ labeled point cloud samples in the source domain are denoted as $D_S = \{p_s^i, l_s^i\}_{i=1}^{N_S}$, where $p_s^i$ and $l_s^i$ are the $i$th point clouds and the corresponding 3D object detection labels. Similarly, the $N_{TL}$ labeled samples in the target domain are denoted as $D_{TL} = \{p_{tl}^i, l_{tl}^i\}_{i=1}^{N_{TL}}$ and the $N_{TU}$ unlabeled samples in the target domain are denoted as $D_{TU} = \{p_{tu}^i\}_{i=1}^{N_{TU}}$. The $i$th point clouds $p^i \in \mathbb{R}^{N_p \times 4}$ contain LiDAR points where each point measurement includes 3D coordinates $(x, y, z)$, and intensity $I$. The corresponding set of labels, $l^i$ contains descriptions of 3D object boxes, consisting of category, location, size, and heading angle. In typical SSDA setup, we assume that both $N_S$ and $N_{TU}$ are significantly larger than $N_{TL}$, i.e. $N_S >> N_{TL}$ and $N_{TU} >> N_{TL}$. The objective of SSDA is to maximize the object detection performance in the target domain through an effective use of $D_S$, $D_{TL}$, and $D_{TU}$.

Fig. 2 depicts the two-stage structure of the proposed TODA framework. In the first stage, TargetMix initially performs LiDAR Distribution Matching, which transforms the source-domain LiDAR data $D_S$ into the LiDAR data $D_S'$ such that the transformed LiDAR data $D_S'$ follows the configurations of the target domain LiDAR. Subsequently, TargetMix combines the transformed data $D_S'$ with $D_{TL}$ in polar coordinates using a weighted mixup operation as proposed in [24]. This operation
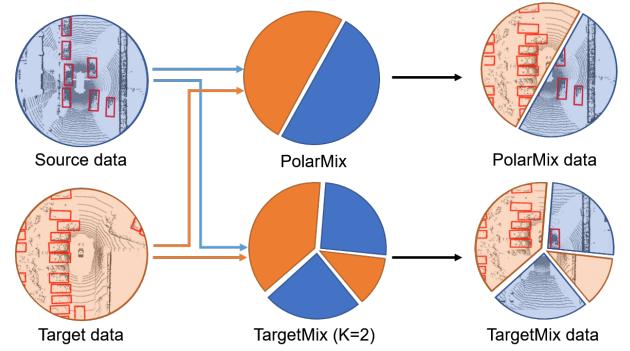


Fig. 3. **Comparison of TargetMix with PolarMix:** TargetMix divides the entire azimuth angle into $2K$ separate sectors while PolarMix divides it into two sectors.

results in a mixed dataset $D_{TM} = \{p_{tm}^i, l_{tm}^i\}_{i=1}^{N_{TM}}$. The mixed data $D_{TM}$ is then utilized to train a 3D object detection model, aimed at bridging the gap between the two domains.

In the second stage, we utilize two 3D object detection pipelines with identical structures. As illustrated in Fig. 2, one model $f_t$ serves as a teacher, while the other $f_s$ functions as a student model. The weights of the model trained in the first stage are initially copied to the teacher model. The teacher model is then employed to generate pseudo-labels $\{\hat{l}_{tu}^i\}_{i=1}^{N_{TU}}$ from the unlabeled data $D_{TU}$. At the same time, the corresponding unlabeled data $D_{TU}$ are perturbed by the Adversarial Point Augmentation. These perturbed samples are then mixed with the labeled data $D_{TL}$, resulting in the mixed data $D_{AM} = \{p_{am}^i, \hat{l}_{tu}^i\}_{i=1}^{N_{TU}}$. Finally, the student model is trained using the mixed data $D_{AM} = \{p_{am}^i, \hat{l}_{tu}^i\}_{i=1}^{N_{TU}}$. By incorporating adversarial examples, TODA ensures a consistent representation of the unlabeled samples to train the student model, facilitating the effective utilization of unlabeled target-domain samples.

### B. Target Sensor-Guided Mix

TargetMix combines LiDAR points from source and target domains. TargetMix performs two steps. The first step converts the source-domain LiDAR data $D_S$ into the data $D_S' = \{\tilde{p}_s^i, l_s^i\}_{i=1}^{N_S}$. This LiDAR Distribution Matching step adjusts the number of beam channels, the number of points per channel, and Vertical Field of View (VFOV) of $D_S$ to match the corresponding parameters of the LiDAR data in the target domain. Specifically, TargetMix converts the 3D points in $D_S$ from Cartesian coordinates $(x, y, z)$ to spherical coordinates:

$$\theta = \arctan \frac{z}{\sqrt{x^2 + y^2}}, \quad \phi = \arcsin \frac{y}{\sqrt{x^2 + y^2}}, \quad (1)$$

where $\theta$ and $\phi$ denote the azimuth and zenith angles, respectively. Using the point data in spherical coordinates, we then generate the range image $I_r \in \mathbb{R}^{H \times W}$, where each pixel in $I_r$ corresponds to specific $\theta$ and $\phi$ values. The image $I_r$ is downsampled based on the ratio of the number of beam channels, the number of points per channel, and the VFOV. The downsampled image $I_r'$ is transformed back into the Cartesian coordinate system to generate the transformed source data $D_S'$. For instance, consider the Waymo dataset as source-domain

data and the nuScenes dataset as target-domain data. LiDAR sensor used in Waymo dataset has 64 channels, about 2200 points per channel, and a VFOV of $[-17.6°, 2.4°]$, covering a range of $20°$. On the other hand, LiDAR sensor in nuScenes dataset has 32 channels, 1100 points per channel, and a VFOV of $[-30.0°, 10.0°]$, covering a range of $40°$. The image $I_r$ is downsized vertically by a factor of 4, considering both the VFOV range ratio $\frac{40°}{20°}$ and the channel ratio $\frac{64}{32}$. Additionally, the image is downsized horizontally by a factor of 2 to match the points per channel between the two datasets.

The second step of TargetMix applies mix augmentation between $D'_S$ and $D_{TL}$ in the polar coordinate system, where each LiDAR point is characterized by $(\theta, r, \phi)$: $\theta$ represents the azimuth angle, $r$ denotes the distance, and $\phi$ indicates the inclination angle between the z-axis and the point vector $(x, y, z)$. Inspired by PolarMix [25], TargetMix conducts a mix operation by partitioning the azimuth angle into two sections and filling one with data points from $D'_S$ and the other with those from $D_{TL}$. This mix operation is also applied to 3D box labels accordingly. Unlike PolarMix, which utilizes a single contiguous sector for $D_{TL}$, TargetMix assigns $K$ separate contiguous sectors to $D_{TL}$, accounting for the widespread distribution of objects. The remaining regions are assigned to $D'_S$. However, such partitioning can lead to ambiguity, as LiDAR points within objects may be divided by azimuth boundaries. To mitigate this issue, TargetMix not only eliminates the object boxes affected by these boundaries but also removes their corresponding point cloud data. This operation is referred to as *Enhanced Mix Strategy*. We denote the data generated by Enhanced Mix Strategy as $D_{TM} = \{p_{tm}^i, l_{tm}^i\}_{i=1}^{N_{TM}}$. Note that

$$p_{tm}^i = (M_p \odot \tilde{p}_s^i) \oplus ((1 - M_p) \odot p_{tl}^j) \quad (2)$$
$$l_{tm}^i = (M_l \odot l_s^i) \oplus ((1 - M_l) \odot l_{tl}^j) \quad (3)$$

where $M_p, M_l$ denotes a binary mask list indicating the $K$ azimuth ranges $[[\alpha^1, \beta^1], ... [\alpha^K, \beta^K]]$, $\odot$ is the element-wise multiplication operation, and $\oplus$ denotes the concatenation operation. TargetMix randomly selects either a mixed sample from $D_{TM}$ or a sample from $D'_S$ with a probability $P_{tm}$, and the chosen sample is fed into the 3D object detection model for training.

*C. Adversarial-Guided Mix*

The teacher model $f_t$ is initialized with the model trained by TargetMix in the first stage. Then, the teacher model $f_t$ generates pseudo labels $\{\hat{l}_{tu}^i\}_{i=1}^{N_{TU}}$ from the unlabeled target data $D_{TU}$. However, because the teacher model is trained on a limited subset of labeled data in both the source and target domains, it may not fully capture the distribution of the unlabeled target data. This results in intra-domain discrepancy, where the feature points from labeled data exhibit a different distribution from those of the unlabeled data. Since the teacher model is trained with labeled data, some unlabeled data may also be influenced by the labeled data while others are not. This discrepancy is evident in Fig. 4 (a), where a noticeable distribution gap exists between the representation of the labeled data and that of the unlabeled data. Such intra-domain inconsistency presents challenges when training the

student model using the inconsistently represented labeled and unlabeled data.

To address this issue, we employ adversarial augmentation as proposed in [18]. This technique aligns data distribution with the target distribution by perturbing input data in the direction of the negative gradient of the loss function derived from a model trained on the target distribution. AdvMix specifically implements Adversarial Point Augmentation on unlabeled LiDAR data within the target domain, altering the positions of LiDAR points in 3D space to optimize domain adaptation performance. More precisely, LiDAR points located within detection boxes provided by pseudo labels, are subject to perturbation. Within each bounding box, a subset of LiDAR points is randomly chosen to be perturbed, based on a probability $\rho$. Subsequently, for each selected point, one of three types of adversarial perturbations is applied

- Point translation: a perturbation $\delta$ is added to the $(x, y, z)$ coordinates of the selected point
- Point addition: a new point is generated by translating the coordinates of the selected point by $\delta$
- Point removal: the selected point is removed.

One of the three types is randomly chosen with equal probability, producing the adversarial samples $D_{adv} = \{p_{adv}^i, l_{adv}^i\}_{i=1}^{N_{TU}}$ The direction of the perturbation $\delta$ is determined by the negative gradient of the detection loss. For the point cloud input $p_{tu}$, the perturbation $\delta$ can be calculated as

$$\delta = \epsilon \frac{g}{||g||_2}, \quad (4)$$
$$g = -\nabla_{p_{tu}} \mathcal{L}_{det}(f_t(p_{tu}), \hat{l}_{tu}), \quad (5)$$

where $\mathcal{L}_{det}$ is the detection loss comprising the classification and regression losses, $\hat{l}_{tu}$ includes the pseudo labels obtained from $p_{tu}$, $g \in \mathbb{R}^3$ is the gradient of the detection loss with respect to the $(x, y, z)$ coordinates of the points in $p_{tu}$, and $\epsilon$ denotes the magnitude of the perturbation. This gradient indicates the direction of perturbation that reduces the detection loss. Fig. 4 (b) shows that the Adversarial Point Augmentation results in the reduced gap in distribution between the labeled and unlabeled data. Consequently, this can foster a more consistent and aligned learning process for the student model.

Finally, we mix the perturbed unlabeled target data $D_{adv}$ and the labeled target data $D_{TL}$ through *Point-MixUp* [9]. This step generates the adversarial mixed data $D_{AM} = \{\tilde{p}_{AM}^i, \hat{l}_{AM}^i\}_{i=1}^{N_S}$. In parallel, we also mix the unlabeled target data $D_{TU}$ with the labeled target data $D_{TL}$ without Adversarial Point Augmentation. This generates the point-mixed data $D_{PM} = \{\tilde{p}_{PM}^i, \hat{l}_{PM}^i\}_{i=1}^{N_S}$. Both mixed data $D_{AM}$ and $D_{PM}$ are used to train the student model, as described in the next section. Note that this mixup operation is conducted with the probability $P_{am}$. Without the mixup operation, we let $D_{AM} = D_{TU}$ and $D_{PM} = D_{adv}$.

*D. Training*

Our TODA is trained through two-stage training process. In the first stage, the teacher model $f_t$ is trained with the data $D_{TM}$ generated by TargetMix through the detection loss,

$$\mathcal{L}_{det}^{tm} = \mathcal{L}_{cls}^{tm} + \mathcal{L}_{reg}^{tm}, \quad (6)$$

---

**Algorithm 1** Target Oriented Data Augmentation (TODA)

---

**Input:** Source data $(p_s, l_s) \in D_S$, Target labeled data $(p_{tl}, l_{tl}) \in D_{TL}$, Target unlabeled data $p_{tu} \in D_{TU}$, Teacher model $f_t$, Student model $f_s$, the probability of TargetMix $P_{tm}$, the probability of AdvMix $P_{am}$, Regularization parameter $\lambda$, total epoch of TargetMix stage $T_{TM}$, total epoch of AdvMix stage $T_{AM}$

1: Apply DataTransformation: $D'_s$ ▷ TargetMix stage
2: **for** $T = 1$ to $T_{TM}$ **do**
3:    **for** $i = 1$ to $N_S + N_{TL}$ **do**
4:       **if** rand() $< P_{tm}$ **then**
5:          Sample $(p'^i_s, l^i_s)$ and $(p^i_{tl}, l^i_{tl})$
6:          Apply PolarEnhanceMix: $(p'^i_{tm}, l^i_{tm})$
7:       **else**
8:          Sample $(p'^i_s, l^i_s)$ if $i < N_S$ else $(p^i_{tl}, l^i_{tl})$
9:       **end if**
10:       Calculate $\mathcal{L}_{det}$ and update $f_t$
11:    **end for**
12: **end for**
13: Initialize $f_s$ with $f_t$ pre-trained weights ▷ AdvMix stage
14: Generate pseudo labels $\hat{l}_{tu}$ using $f_t$
15: **for** $T = 1$ to $T_{AM}$ **do**
16:    **for** $i = 1$ to $N_{TL} + N_{TU}$ **do**
17:       Sample $(p^i_{tl}, l^i_{tl})$ and $(p^i_{tu}, \hat{l}^i_{tu})$
18:       Generate adversarial example $p^i_{adv}$
19:       **if** rand() $< P_{am}$ **then**
20:          Apply PointMixUp: $(\tilde{p}^i_{AM}, \hat{l}^i_{AM}), (\tilde{p}^i_{PM}, \hat{l}^i_{PM})$
21:          Calculate $\mathcal{L}_{cons}(f_s(\tilde{p}^i_{AM}), f_s(\tilde{p}^i_{PM}))$
22:       **else**
23:          Calculate $\mathcal{L}_{cons}(f_s(p^i_{tu}), f_s(p^i_{adv}))$
24:       **end if**
25:       Calculate $\mathcal{L}_{det}$ and update $f_s$
26:    **end for**
27: **end for**

---

where $\mathcal{L}^{tm}_{reg}$ and $\mathcal{L}^{tm}_{cls}$ are the standard smoothed-L1 loss and focal loss [35], respectively. In the next stage, the student model was trained while freezing the teacher model. The loss function used to train the student model is given by

$$\mathcal{L}_{adv} = \mathcal{L}^{am}_{det} + \mathcal{L}^{pm}_{det} + \lambda \mathcal{L}_{cons}, \qquad (7)$$

where $\lambda$ is the regularization parameter, and $\mathcal{L}^{am}_{det}$ and $\mathcal{L}^{pm}_{det}$ are the detection loss terms associated with the adversarial mixed data $D_{AM}$ and the point mixed data $D_{PM}$, respectively.

We also consider the *Consistency Loss* $\mathcal{L}_{cons}$ to enforce the consistency between the detection results obtained from $D_{PM}$ and those obtained from $D_{AM}$. Specifically, $\mathcal{L}_{cons}$ is expressed as

$$\mathcal{L}^i_{cons} = \frac{\Sigma^{N^i_{AMB}}_{k=1}||b^{i,k}_{am} - \tilde{b}^{i,k}_{pm}||_2 + \Sigma^{N^i_{PMB}}_{k=1}||b^{i,k}_{pm} - \tilde{b}^{i.k}_{am}||_2}{N^i_{AMB} + N^i_{PMB}}, \qquad (8)$$

$$\mathcal{L}_{cons} = \frac{\Sigma^{N_{TU}}_{i=1}\mathcal{L}^i_{cons}}{N_{TU}}, \qquad (9)$$

where $b^{i,k}_{am}$ and $b^{i,k}_{pm}$ denote the $k$th predicted bounding boxes obtained from the $i$th sample, respectively. To quantify the
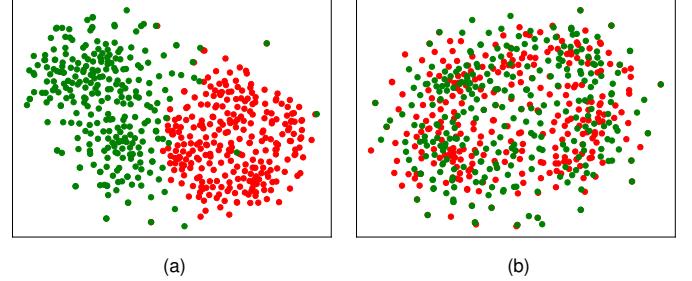


Fig. 4. **t-SNE visualization of features:** (a) unlabeled data (green) versus labeled data (red), (b) adversarial examples (green) versus labeled data (red) within the target domain. These features are extracted from the final layer of the teacher model trained with TargetMix.

distance of $b^{i,k}_{am}$ from $b^{i,k}_{pm}$, we find $\tilde{b}^{i,k}_{am}$ among the set $\{b^{i,k}_{am}\}^{N^i_{AMB}}_{k=1}$ which is closest to $b^{i,k}_{pm}$. Similarly, we also find $\tilde{b}^{i,k}_{pm}$ among the set $\{b^{i,k}_{pm}\}^{N^i_{PMB}}_{k=1}$ which is closest to $b^{i,k}_{am}$. Each bounding box is represented by center coordinates $(x, y, z)$ and size offset $(w, l, h)$ and $||.||_2$ denotes the $\ell_2$ norm, and $N^i_{AMB}$ and $N^i_{PMB}$ represent the number of detection boxes for the $i$th sample in $D_{AM}$ and that in $D_{PM}$, respectively. The overall training process is presented in Algorithm 1.

## IV. EXPERIMENTS

### A. Datasets

We consider the setup where domain adaptation is conducted for two scenarios:

- Transfer from Waymo dataset [14] to nuScenes dataset [16]
- Transfer from nuScenes dataset to KITTI dataset [15].

The Waymo dataset provides 160k labeled training samples collected with 64-beam LiDAR with a 20-degree VFOV of 20 degrees (from $-17.6°$ to $2.4°$) while nuScenes provides 28k frames of labeled training samples collected with 32-beam LiDAR and has a VFOV of 40 degrees (from $-30.0°$ to $10.0°$). KITTI provides 7,481 frames of labeled training samples collected with 64-beam LiDAR and has a VFOV of 30 degrees (from $-23.6°$ to $3.2°$). We use 0.5%, 1%, 5%, and 10% labels for nuScenes dataset and 1% labels for KITTI dataset, and the remaining data were utilized as unlabeled data.

We evaluated the performance of our method using both the nuScenes and KITTI 3D object detection metrics. For the nuScenes metric, we utilized the official nuScenes Detection Score (NDS) [16] and Average Precision (AP) for the car category. AP represents the average precision values obtained across different thresholds of $d = 0.5, 1, 2$, and 4 meters calculated based on the BEV center distance. NDS provides a comprehensive metric incorporating AP as well as errors in attributes, classification, localization, and velocity. For a fair comparison with existing methods, we also adopted the official KITTI 3D object detection metric. The KITTI metric calculates the AP for both BEV IoU and 3D IoU under an IoU threshold of 0.7 over 40 recall positions, specifically for the car category. Following the approach of ST3D [8], we measured the reduction in AP and NDS relative to the *Oracle* performance,

TABLE I
COMPARISON OF DOMAIN ADAPTATION PERFORMANCE ON WAYMO TO NUSCENES WITH DIFFERENT AMOUNTS OF TARGET LABELS FOR THE CAR CLASS.
WE REPORT AP, NDS, AND THE CLOSED GAP BASED ON NUSCENES METRIC. THE BEST ADAPTATION RESULT IS INDICATED IN **BOLD**.

| Method | 0.5% | | 1% | | 5% | | 10% | |
|---|---|---|---|---|---|---|---|---|
| | AP / NDS | Closed Gap (AP / NDS) | AP / NDS | Closed Gap (AP / NDS) | AP / NDS | Closed Gap (AP / NDS) | AP / NDS | Closed Gap (AP / NDS) |
| Source Only | 42.6 / 50.3 | +0% / +0% | 42.6 / 50.3 | +0% / +0% | 42.6 / 50.3 | +0% / +0% | 42.6 / 50.3 | +0% / +0% |
| ST3D [8] | 43.7 / 50.2 | +3.1% / -0.5% | 43.7 / 50.2 | +3.1% / -0.5% | 43.7 / 50.2 | +3.1% / -0.5% | 43.7 / 50.2 | +3.1% / -0.5% |
| Labeled Target | 36.0 / 37.7 | -18.4% / -64.2% | 37.2 / 38.1 | -15.1% / -62.2% | 61.0 / 53.2 | +51.4% / +14.8% | 65.6 / 58.2 | +64.2% / +40.3% |
| Co-training | 47.5 / 52.7 | +13.7% / +12.2% | 51.4 / 54.6 | +24.6% / +21.9% | 57.7 / 58.0 | +42.2% / +39.3% | 59.4 / 58.9 | +46.9% / +43.9% |
| SSDA3D [9] | 70.3 / 65.1 | +77.4% / +75.5% | 73.4 / 67.1 | +86.0% / +85.7% | 76.2 / 68.8 | +93.9% / +94.4% | 78.8 / 70.9 | +101.1% / +105.1% |
| Ours | **73.7 / 67.3** | **+86.9% / +86.7%** | **75.6 / 68.5** | **+92.2% / +92.9%** | **79.0 / 71.1** | **+101.7% / +106.1%** | **79.3 / 71.4** | **+102.5% / +107.7%** |
| Oracle | 78.4 / 69.9 | +100% / +100% | 78.4 / 69.9 | +100% / +100% | 78.4 / 69.9 | +100% / +100% | 78.4 / 69.9 | +100% / +100% |

TABLE II
COMPARISON OF DOMAIN ADAPTATION PERFORMANCE ON WAYMO TO
NUSCENES USING SECOND-IOU BASED ON KITTI METRIC. SSDA3D
AND TODA UTILIZED AN ADDITIONAL 1% OF TARGET-DOMAIN LABELED
DATA.

| Method | $AP_{BEV}$ / $AP_{3D}$ | Closed Gap |
|---|---|---|
| Source only | 32.9 / 17.2 | +0% / +0% |
| SN [11] | 33.2 / 18.6 | +1.7% / +7.5% |
| ST3D [8] | 35.9 / 20.2 | +15.9% / +16.7% |
| ST3D++ [13] | 35.7 / 20.9 | +14.7% / +20.9% |
| L.D [36] | 40.7 / 22.9 | +41.1% / +32.2% |
| DTS [37] | 41.2 / 23.0 | +43.7% / +32.8% |
| SSDA3D [9] | 46.6 / 29.6 | +72.1% / +70.1% |
| TODA (Ours) | **48.1 / 30.2** | **+80.0% / +73.4%** |
| Oracle | 51.9 / 34.9 | +100% / +100% |

TABLE III
COMPARISON OF DOMAIN ADAPTATION PERFORMANCE ON NUSCENES TO
KITTI USING SECOND-IOU BASED ON KITTI METRIC. SSDA3D AND
TODA UTILIZED AN ADDITIONAL 1% OF TARGET-DOMAIN LABELED DATA.

| Method | $AP_{BEV}$ / $AP_{3D}$ | Closed Gap |
|---|---|---|
| Source only | 51.8 / 17.9 | +0% / +0% |
| SN [11] | 59.7 / 37.6 | +25.1% / +35.4% |
| ST3D [8] | 75.9 / 54.1 | +76.6% / +59.5% |
| ST3D++ [13] | 80.5 / 62.4 | +91.1% / +80.0% |
| DTS [37] | 81.4 / 66.6 | +94.0% / +87.6% |
| SSDA3D [9] | 81.5 / 67.4 | +94.3% / +89.0% |
| TODA (Ours) | **82.7 / 68.6** | **+98.1% / +91.2%** |
| Oracle | 83.3 / 73.5 | +100% / +100% |

TABLE IV
ABLATION STUDY FOR EVALUATING THE CONTRIBUTION OF EACH
COMPONENT OF TODA ON THE NUSCENES VALIDATION DATASET.

| Method | TargetMix | AdvMix | AP / NDS |
|---|---|---|---|
| Source Only | | | 37.2 / 38.1 |
| Co-training | | | 51.4 / 54.6 |
| TODA | ✓ | | 71.2 / 65.9 |
| | ✓ | ✓ | 75.6 / 68.5 |

i.e., $\frac{AP_{model} - AP_{source\ only}}{AP_{oracle} - AP_{source\ only}}$ and $\frac{NDS_{model} - NDS_{source\ only}}{NDS_{oracle} - NDS_{source\ only}}$. We refer to this metric as the *closed gap*. The *Oracle* performance was obtained by conducting supervised learning using labels from the entire target-domain dataset.

## B. Implementation Details

We implemented the CenterPoint [4] and SECOND-IoU [3] using the OpenPCDet [38] codebase. We followed the respective training schedule used in CenterPoint [4] and SECOND-IoU. For the Waymo to nuScenes adaptation, the detection ranges for $x$, $y$, and $z$ axes were set to $[-54.0, 54.0]$, $[-54.0, 54.0]$, and $[-5.0, 4.8]$ meters, and the voxel size was set to $[0.075, 0.075, 0.2]$. In the case of nuScenes to KITTI adaptation, the detection ranges in $x$, $y$, and $z$ axes were set to $[-76.2, 76.2]$, $[-76.2, 76.2]$, and $[-3.0, 5.0]$ meters with the same voxel size of $[0.075, 0.075, 0.2]$. The intensity of the point cloud was normalized to have a value between $[0, 1]$. We applied data augmentation techniques, including random flip along X and Y axes, random rotation, random scaling, and GT sampling. GT sampling was only applied to the labeled data in the target domain for the second stage. The probability $P_{tm}$ used in TargetMix was set to 0.4. The probability $P_{am}$ in AdvMix was set to 0.6. The hyperparameters $\lambda$, $\rho$, and $\epsilon$ were empirically chosen to 1, 0.5, and 0.001, respectively.

All experiments were conducted on four 24GB RTX 3090TI GPUs.

## C. Main Results

We evaluated the performance of our SSDA model using CenterPoint. In the nuScenes dataset, we utilized labels for 0.5%, 1%, 5%, and 10% of the target domain data, corresponding to 141, 282, 1407, and 2813 frames, respectively. Table I presents the performance of TODA in the Waymo to nuScenes adaptation task, using the nuScenes metric. We compare TODA with the existing domain adaptation methods including *Co-training*, *Labeled Target*, ST3D [8], and SSDA3D [9]. *Co-training* utilized supervised learning using the labeled data in both source and target domains while *Labeled Target* employed supervised learning solely utilizing the labeled data in the target domain. Table I shows that TODA consistently outperforms all other methods by significant margins. When

TABLE V
ABLATION STUDY FOR THE TARGETMIX MODULE. POLAR, ENHANCE., AND MATCH., INDICATE POLAR COORDINATE-BASED MIX, ENHANCED MIX STRATEGY, AND LIDAR DISTRIBUTION MATCHING, RESPECTIVELY.

| Method | Polar | Enhance. | Match. | AP / NDS |
|--------|-------|----------|--------|----------|
| CutMix [9] | | | | 66.8 / 63.4 |
| | ✓ | | | 68.3 / 63.8 |
| TargetMix | ✓ | ✓ | | 69.6 / 64.5 |
| | ✓ | ✓ | ✓ | 71.2 / 65.9 |

TABLE VI
ABLATION STUDY FOR THE ADVMIX MODULE. ADV., CONS., MIXUP., INDICATE ADVERSARIAL POINT AUGMENTATION, CONSISTENCY LOSS, AND POINT-MIXUP, RESPECTIVELY.

| Method | Adv. | Cons. | MixUp | AP / NDS |
|--------|------|-------|-------|----------|
| TargetMix | | | | 71.2 / 65.9 |
| | ✓ | | | 73.2 / 67.1 |
| AdvMix | ✓ | ✓ | | 73.5 / 67.2 |
| | ✓ | ✓ | ✓ | 75.6 / 68.5 |

TABLE VII
ABLATION STUDY FOR THE PROBABILITY OF MIXING AUGMENTATION $P_{tm}$ AND $P_{am}$.

| $P_{tm}$ | AP |
|----------|------|
| 0.1 | 70.1 |
| 0.2 | **71.2** |
| 0.3 | 70.1 |
| 0.4 | 69.8 |

(a)

| $P_{am}$ | AP |
|----------|------|
| 0.4 | 78.1 |
| 0.5 | 78.5 |
| 0.6 | **79.0** |
| 0.7 | 78.8 |

(b)

0.5% of labeled data are used in the target domain, TODA demonstrates notable performance gains of 3.4% in AP and 2.2% in NDS over SSDA3D, the current state-of-the-art method. TODA also achieves more than 85% closed gap from the Oracle performance. Even with 1% labeled target-domain data, TODA achieves improvements of 2.2% in AP and 1.4% in NDS over SSDA3D. It also achieves above 90% closed gap in both metrics. Remarkably, when 5% and 10% labeled data are utilized, TODA even surpasses the Oracle performance. This would be possibly because TODA leverages both source-domain and target-domain data, while the Oracle performance is obtained using the target-domain dataset only. In summary, TODA achieves performance comparable to the Oracle performance while significantly reducing annotation costs.

We further validate the effectiveness of our method in comparison with various latest UDA techniques including SN [11], ST3D [8], ST3D++ [13], L.D [36], and DTS [37]. For fair comparison, we implemented TODA on the SECOND-IoU model and conducted evaluation using the KITTI metrics for the Waymo to nuScenes adaptation task. Table II shows that the proposed TODA method still maintains performance gains over other UDA methods. While UDA methods exhibit limitation in reducing the domain gap, TODA achieves a significantly higher closed gap of +80.0% using only 1% of the labeled
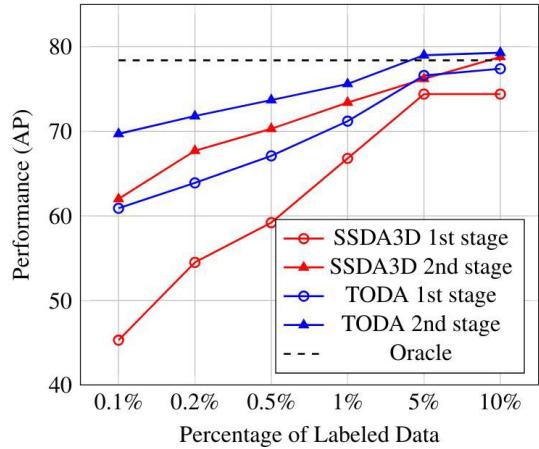


Fig. 5. **Comparison with SSDA3D in each stage for different percentages:** Performance comparison between TODA and SSDA3D across various sizes of labeled data (0.1%, 0.2%, 0.5%, 1%, 5%, and 10%)

TABLE VIII
THE PERFORMANCE OF TODA WITH 0.1% AND 0.2% OF TARGET-DOMAIN LABELED DATA.

| Method | 0.1% | 0.2% |
|--------|------|------|
| Labeled Target | fail | fail |
| Co-training | 50.3 | 52.2 |
| SSDA3D [9] | 62.0 | 67.6 |
| TODA | **69.7** | **71.8** |
| Oracle | 78.4 | 78.4 |

target data.

Table III shows the performance of TODA on the nuScenes to KITTI adaptation scenario as well. Although the UDA methods achieve performance close to the oracle, TODA generalizes well to the nuScenes to KITTI setup. Notably, TODA achieves a closed gap of 98.1/91.2% relative to the Oracle. These results highlight the effectiveness of our SSDA approach in adapting to different target domains, even with 1% of labeled data.

### D. Ablation Studies

*1) Contribution of Each Component :* We conducted ablation studies to assess the contribution of each component of TODA to the overall performance. The evaluation was conducted using the nuScenes validation set on the Waymo to nuScenes task. Table IV presents the results when 1% of the target-domain labeled data are used. We sequentially enabled *TargetMix* and *AdvMix* on top of the baseline method based on Co-training. The inclusion of *TargetMix* results in a notable performance gain of 19.8% in AP and 11.3% in NDS, highlighting its significant contribution. Incorporating *AdvMix* leads to a performance improvement of 4.4% in AP and 2.6% in NDS.

*2) Sub-components of TargetMix:* Table V illustrates the ablation study evaluating the contribution of three components within the TargetMix module: *Polar Coordinate-based Mix*, *Enhanced Mix Strategy*, and *LiDAR Distribution Matching*. The *Enhanced Mix Strategy* refers to the idea of segmenting the area into $K$ segments and allocating the LiDAR points accordingly. LiDAR Distribution Matching refers to the method
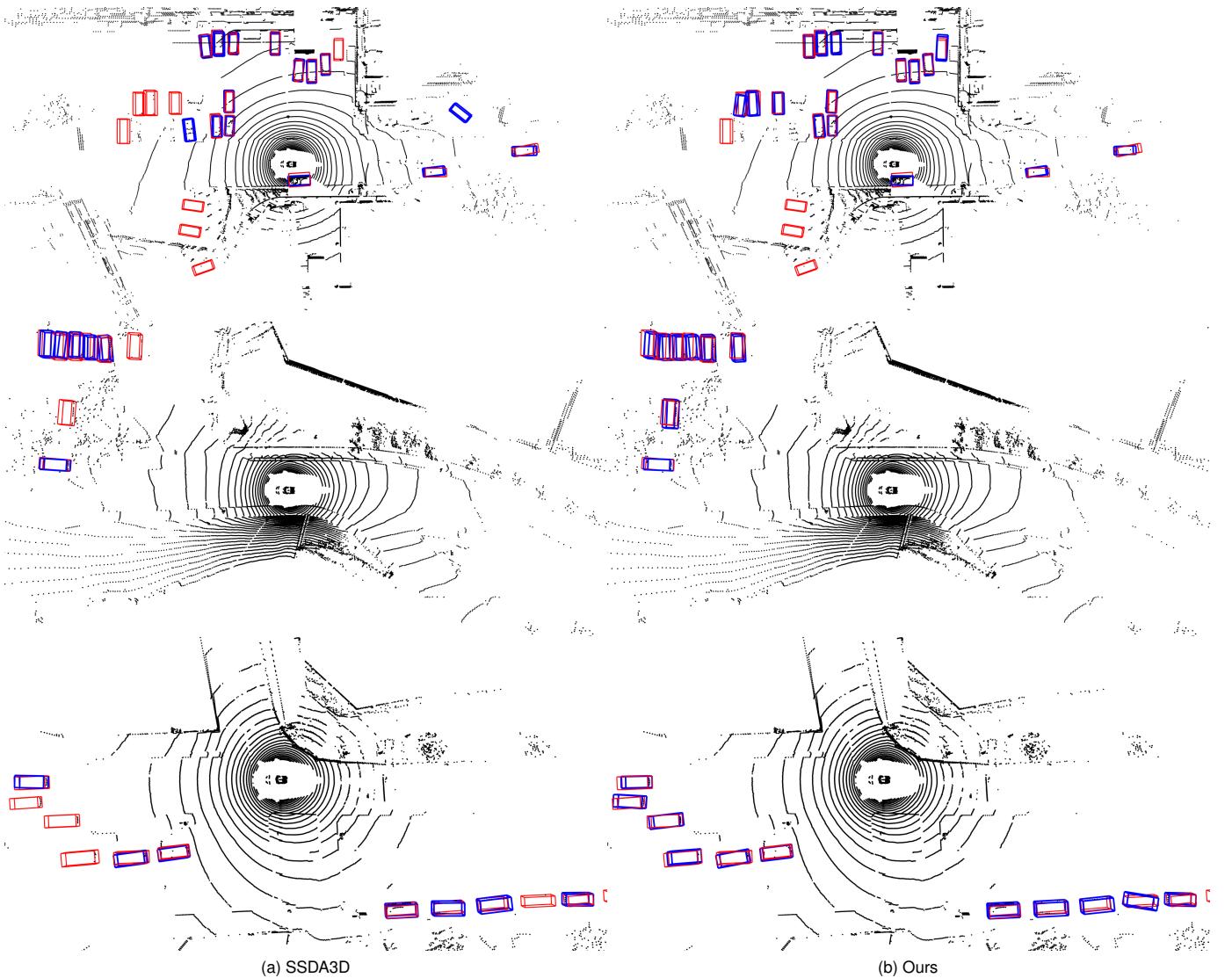
(a) SSDA3D

(b) Ours

Fig. 6. **Comparison of detection results:** (a) SSDA3D and (b) TODA. All samples are from nuScenes val split. The red box represents the ground truth, and the blue box indicates the predicted bounding box.

that adjusts LiDAR data to match the configuration of the target-domain LiDAR. While the naive *Polar Coordinate-based Mix* demonstrates only a slight improvement over the CutMix baseline, its impact becomes significant when combined with our *Enhanced Mix Strategy*. This *Enhanced Mix Strategy* achieves a performance gain of 2.8% in AP and 2.1% in NDA. Additionally, the *LiDAR Distribution Matching* method further enhances performance by 1.6% in AP and 1.4% in NDS.

*3) Sub-components of AdvMix:* Table VI presents the ablation study assessing the impact of each component within the AdvMix module. We utilized a model trained with TargetMix as our baseline. Three components were considered: *Adversarial Point Augmentation*, *Consistency Loss*, and *Point-Mixup*. *Adversarial Point Augmentation* results in a performance gain of 2.0% in AP and 1.2% in NDS. Additionally, *Consistency Loss* contributes to an increase in AP by 0.3% and NDS by 0.1%. Finally, *Point-MixUp* leads to further improvements of

2.1% in AP and 1.3% in NDS.

*4) Probabilities $P_{tm}$ and $P_{am}$ of TargetMix and AdvMix:* Table VII presents the performance of TODA as a function of the parameters $P_{tm}$ and $P_{am}$. Table VII (a) shows the performance of the first stage with varying $P_{tm}$ values, while Table VII (b) illustrates the final TODA performance with varying $P_{am}$ values using a model trained with $P_{tm} = 0.2$. The results demonstrate that TODA exhibits robust performance with respect to both $P_{tm}$ and $P_{am}$, with only slight degradation under different settings. The optimal performance is achieved with $P_{tm} = 0.2$ and $P_{am} = 0.6$.

*5) Scenarios Using Extremely Low Number of Labels:* In Table VIII, we evaluate the performance of TODA in scenarios where the percentage of labeled target-domain data is extremely low, e.g., 0.1% (28 frames) and 0.2% (56 frames). Training exclusively with such small number of labeled targets failed due to insufficient target-domain data. In contrast, TODA achieves AP values of 69.7% and 71.8% respectively, surpassing

SSDA3D by 6.3% and 4.2% respectively. Remarkably, TODA achieves 75% and 81% of closed gap even with 0.1% and 0.2% labeled data respectively. These results demonstrate the potential of TODA in highly data-constrained environments.

Figure 5 presents a performance comparison between TODA and SSDA3D across various sizes of labeled data (0.1%, 0.2%, 0.5%, 1%, 5%, and 10%). Performance was assessed at both the initial and secondary stages of each detector. We noted an increasing performance gap between TODA and SSDA3D at both stages as the percentage of labeled data decreased. This trend seems to be attributed to proposed LiDAR Distribution Matching that conducts the adaptation of LiDAR data distribution without using labeled data. This leads to notable performance improvements for TODA in the initial stage. Additionally, Adversarial Point Augmentation reshapes the distribution of unlabeled data in the target domain, enabling more effective feature alignment between labeled and unlabeled data. This boosts the performance of our pseudo-label based semi-supervised learning in the second stage.

### E. Qualitative Results

We present some qualitative results. We considered the setup where 1% of the labeled data are used in the target domain. The experiments were conducted on the nuScenes validation set.

In Fig. 4, we present visualizations of feature distributions utilizing t-SNE [39] under both pre- and post-application of Adversarial Point Augmentation. Fig. 4 (a) highlights the distribution gap between the labeled and unlabeled data, while Fig. 4 (b) shows the gap between the labeled data and the adversarially perturbed unlabeled data. Notably, Adversarial Point Augmentation demonstrates its effectiveness in mitigating the distribution shift.

Fig. 6 illustrates the detection results produced by both the existing method and the proposed TODA method. The figures in the left column depict the detection results attained by SSDA3D, while those in the right column show the results of the proposed TODA method. We observe that TODA demonstrates enhanced detection results, effectively identifying objects that were previously either missed or inaccurately detected by SSDA3D.

## V. CONCLUSIONS

In this paper, we introduced TODA, an SSDA framework for 3D object detection based on a target-oriented domain augmentation strategy. To mitigate the disparity in data distribution between the source and target domains, we introduced TargetMix. TargetMix utilizes an inter-domain mixup augmentation strategy within a polar coordinate system, considering the LiDAR scanning mechanism. Additionally, TargetMix incorporates LiDAR Distribution Matching to adapt the source domain data to align with the configurations of the target-domain LiDAR sensor. Additionally, we proposed AdvMix, which adds adversarial perturbation to the unlabeled data to mitigate intra-domain disparity. We optimized the perturbation direction to maximize detection performance, enabling AdvMix to generate consistent representations of both labeled and unlabeled data in the target domain. By integrating TargetMix and AdvMix,

TODA effectively utilizes both labeled and unlabeled data for domain adaptation. Our evaluation demonstrated that TODA achieved significant performance gains over existing domain adaptation methods and approached performance levels close to the Oracle performance.

Moving forward, we aim to enhance TODA's capabilities to cope with adverse weather conditions and low-resolution LiDAR environments, as well as explore its potential for other sparse data modalities like radar. These advancements will expand TODA's usefulness and offer promising avenues for future research in domain adaptation and 3D object detection.

## REFERENCES

[1] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4490–4499.

[2] A. H. Lang, S. Vora, H. Caesar, L. Zhou, J. Yang, and O. Beijbom, "Pointpillars: Fast encoders for object detection from point clouds," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 12 697–12 705.

[3] Y. Yan, Y. Mao, and B. Li, "Second: Sparsely embedded convolutional detection," *Sensors*, vol. 18, no. 10, p. 3337, Oct. 2018.

[4] T. Yin, X. Zhou, and P. Krahenbuhl, "Center-based 3d object detection and tracking," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 11 784–11 793.

[5] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, "Pv-rcnn: Point-voxel feature set abstraction for 3d object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 10 529–10 538.

[6] G. Shi, R. Li, and C. Ma, "Pillarnet: Real-time and high-performance pillar-based 3d object detection," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 35–52.

[7] J. Deng, S. Shi, P. Li, W. Zhou, Y. Zhang, and H. Li, "Voxel r-cnn: Towards high performance voxel-based 3d object detection," in *Proc. AAAI Conf. Artif. Intell.*, vol. 35, no. 2, 2021, pp. 1201–1209.

[8] J. Yang, S. Shi, Z. Wang, H. Li, and X. Qi, "St3d: Self-training for unsupervised domain adaptation on 3d object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 10 368–10 378.

[9] Y. Wang, J. Yin, W. Li, P. Frossard, R. Yang, and J. Shen, "Ssda3d: Semi-supervised domain adaptation for 3d object detection from point cloud," in *Proc. AAAI Conf. Artif. Intell.*, vol. 37, no. 3, 2023, pp. 2707–2715.

[10] B. Sun, J. Feng, and K. Saenko, "Return of frustratingly easy domain adaptation," in *Proc. AAAI Conf. Artif. Intell.*, vol. 30, no. 1, 2016, pp. 2058–2065.

[11] Y. Wang, X. Chen, Y. You, L. E. Li, B. Hariharan, M. Campbell, K. Q. Weinberger, and W.-L. Chao, "Train in germany, test in the usa: Making 3d object detectors generalize," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11 713–11 723.

[12] X. Peng, X. Zhu, and Y. Ma, "Cl3d: Unsupervised domain adaptation for cross-lidar 3d detection," in *Proc. AAAI Conf. Artif. Intell.*, vol. 37, no. 2, 2023, pp. 2047–2055.

[13] J. Yang, S. Shi, Z. Wang, H. Li, and X. Qi, "St3d++: Denoised self-training for unsupervised domain adaptation on 3d object detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 45, no. 5, pp. 6354–6371, 2022.

[14] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine *et al.*, "Scalability in perception for autonomous driving: Waymo open dataset," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 2446–2454.

[15] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit,*, 2012, pp. 3354–3361.

[16] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11 621–11 631.

[17] T. Kim and C. Kim, "Attract, perturb, and explore: Learning a feature alignment network for semi-supervised domain adaptation," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 591–607.

[18] P. Jiang, A. Wu, Y. Han, Y. Shao, M. Qi, and B. Li, "Bidirectional adversarial training for semi-supervised domain adaptation." in *Proc. 29th Int. Joint Conf. Artif. Intell.*, 2020, pp. 934–940.

[19] S. Shi, X. Wang, and H. Li, "Pointrcnn: 3d object proposal generation and detection from point cloud," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 770–779.

[20] Z. Yang, Y. Sun, S. Liu, and J. Jia, "3dssd: Point-based 3d single stage object detector," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11 040–11 048.

[21] C. Qin, H. You, L. Wang, C.-C. J. Kuo, and Y. Fu, "Pointdan: A multiscale 3d domain adaption network for point cloud representation," *Adv. Neural Inf. Process. Syst.*, vol. 32, pp. 7192–7203, 2019.

[22] W. Zhang, W. Li, and D. Xu, "Srdan: Scale-aware and range-aware domain adaptation network for cross-dataset 3d object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 6769–6779.

[23] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *Proc. Int. Conf. Learn. Representations*, 2018.

[24] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 6023–6032.

[25] A. Xiao, J. Huang, D. Guan, K. Cui, S. Lu, and L. Shao, "Polarmix: A general data augmentation technique for lidar point clouds," *Adv. Neural Inf. Process. Syst.*, vol. 35, pp. 11 035–11 048, 2022.

[26] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," in *Proc. Int. Conf. Learn. Representations*, 2015.

[27] T. Miyato, S.-i. Maeda, M. Koyama, and S. Ishii, "Virtual adversarial training: a regularization method for supervised and semi-supervised learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 8, pp. 1979–1993, Aug. 2019.

[28] D. Liu, R. Yu, and H. Su, "Extending adversarial attacks and defenses to deep 3d point cloud classifiers," in *Proc. IEEE Int. Conf. Image Process.*, 2019, pp. 2279–2283.

[29] ——, "Adversarial shape perturbations on 3d point clouds," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 88–104.

[30] C. Xiang, C. R. Qi, and B. Li, "Generating 3d adversarial point clouds," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 9136–9144.

[31] A. Hamdi, S. Rojas, A. Thabet, and B. Ghanem, "Advpc: Transferable adversarial perturbations on 3d point clouds," in *Proc. Eur. Conf. Comput. Vis.*, 2020, pp. 241–257.

[32] N. Zhao, T.-S. Chua, and G. H. Lee, "Sess: Self-ensembling semi-supervised 3d object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2020, pp. 11 079–11 087.

[33] H. Wang, Y. Cong, O. Litany, Y. Gao, and L. J. Guibas, "3dioumatch: Leveraging iou prediction for semi-supervised 3d object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 14 615–14 624.

[34] J. Yin, J. Fang, D. Zhou, L. Zhang, C.-Z. Xu, J. Shen, and W. Wang, "Semi-supervised 3d object detection with proficient teachers," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 727–743.

[35] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 2980–2988.

[36] Y. Wei, Z. Wei, Y. Rao, J. Li, J. Zhou, and J. Lu, "Lidar distillation: Bridging the beam-induced domain gap for 3d object detection," in *Proc. Eur. Conf. Comput. Vis.*, 2022, pp. 179–195.

[37] Q. Hu, D. Liu, and W. Hu, "Density-insensitive unsupervised domain adaption on 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17 556–17 566.

[38] O. D. Team, "Openpcdet: An open-source toolbox for 3d object detection from point clouds," https://github.com/open-mmlab/OpenPCDet, 2020.

[39] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *J. Mach. Learn. Res.*, vol. 9, no. 11, pp. 2579–2605, Nov. 2008.

## BIOGRAPHY SECTION

**Yecheol Kim** received a B.S. degree in Electrical Engineering from the Hanyang University, Seoul, South Korea, in 2018. He is currently pursuing a Ph.D. degree in the Hanyang University. His research interests include multi-modal object detection, domain adaptation, and efficient architecture for autonomous driving.

**Junho Lee** received a B.S. degree in Electrical Engineering from the Hanyang University, Seoul, South Korea, in 2018. He is currently pursuing a Ph.D. degree in the Hanyang University. His research interests reinforcement learning, domain adaptation, and domain generalization for autonomous driving.

**Changsoo Park** received B.S, M.S and Ph.D. degrees from the Electrical Engineering, KAIST. In 2015, he joined Samsung Electronics, where he participated in research on autonomous driving and robot. Since 2019, he has been researching autonomous driving at Kakao mobility.

**Hyung Won Kim** received the B.S., M.S. and Ph.D. degrees in Electronic Engineering from Hanyang University. From 2015 to 2021, he was with the Hyundai Mobis where his major research topics included target tracking and sensor fusion. He has been working at Kakao Mobility since 2021. His research area includes sensor fusion and machine learning.

**Inho Lim** received B.S degres from the Information and Technology Department, Ajou University. In 2012, he joined Samsung Electronics, Suwoun, Korea, where he participated in the software engineering research team. Since 2017, he has participated in autonomous driving SW research at Samsung Advanced Institute of Technology. Since 2019, he has been researching the Localization and Perception technologies for autonomous driving in Kakao Mobility Corp.

**Christopher Chang** Christopher Chang received a B.S. degree from Department of Electrical Engineering, Seoul National University, and M.S. and Ph.D. degrees in Electrical Engineering from California Institute of Technology, Pasadena CA. Since 2012, he had been working with Samsung Electronics and Hyundai Motor Company, where he developed corporate strategies and business roadmaps. In 2020, he joined Kakao Mobility Corp. where he has been leading strategy, business development, and R&D for the next generation mobilities including autonomous driving, robotics, digital twin, and urban air mobility.

**Jun Won Choi** earned his B.S. and M.S. degrees from Seoul National University and his Ph.D. from the University of Illinois at Urbana-Champaign. Following his studies, he joined Qualcomm in San Diego, USA, in 2010. From 2013 to 2024, he served as a faculty member in the Department of Electrical Engineering at Hanyang University. Since 2024, he has held a faculty position in the Department of Electrical and Computer Engineering at Seoul National University. He currently serves as an Associate Editor for both IEEE Transactions on Intelligent Transportation Systems, IEEE Transactions on Vehicular Technology, International Journal of Automotive Technology. His research spans diverse areas including signal processing, machine learning, robot perception, autonomous driving, and intelligent vehicles.