

MixCycle: Mixup Assisted Semi-Supervised 3D Single Object Tracking with Cycle Consistency

Qiao Wu¹ Jiaqi Yang^{1*} Kun Sun² Chu'ai Zhang¹ Yanning Zhang¹ Mathieu Salzmann³

¹ Northwestern Polytechnical University ² China University of Geosciences, Wuhan

³ École Polytechnique Fédérale de Lausanne

qiaowu@mail.nwpu.edu.cn; jqyang@nwpu.edu.cn; sunkun@cug.edu.cn;

cazhang@mail.nwpu.edu.cn; ynzhang@nwpu.edu.cn; mathieu.salzmann@epfl.ch

Abstract

3D single object tracking (SOT) is an indispensable part of automated driving. Existing approaches rely heavily on large, densely labeled datasets. However, annotating point clouds is both costly and time-consuming. Inspired by the great success of cycle tracking in unsupervised 2D SOT, we introduce the first semi-supervised approach to 3D SOT. Specifically, we introduce two cycle-consistency strategies for supervision: 1) Self tracking cycles, which leverage labels to help the model converge better in the early stages of training; 2) forward-backward cycles, which strengthen the tracker's robustness to motion variations and the template noise caused by the template update strategy. Furthermore, we propose a data augmentation strategy named SOT-Mixup to improve the tracker's robustness to point cloud diversity. SOTMixup generates training samples by sampling points in two point clouds with a mixing rate and assigns a reasonable loss weight for training according to the mixing rate. The resulting MixCycle approach generalizes to appearance matching-based trackers. On the KITTI benchmark, based on the P2B tracker [16], MixCycle trained with 10% labels outperforms P2B trained with 100% labels, and achieves a 28.4% precision improvement when using 1% labels. Our code will be released at <https://github.com/Mumuqiao/MixCycle>.

1. Introduction

3D single object tracking (SOT) plays a critical role in the field of autonomous driving. For example, given object detection [15, 33] results as input, it can output the necessary information for trajectory prediction [8]. The goal of SOT is to regress the center position and 3D bounding-box (BBBox) of an object of interest in a search area, given the point cloud (PC) patch and BBBox of the object template.

*Corresponding author

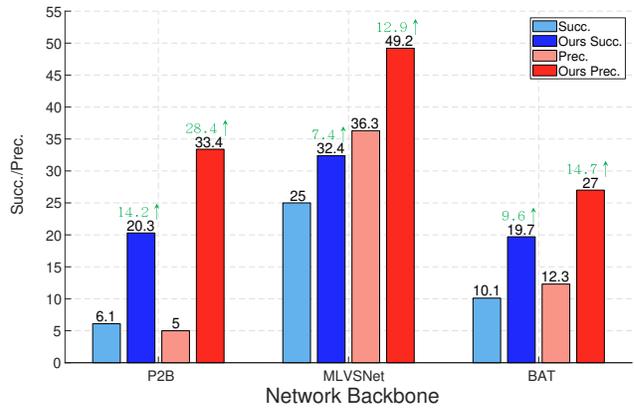


Figure 1. Comparison of MixCycle and fully-supervised methods [16, 25, 34], all trained with 1% labels on KITTI [5]. ‘Succ.’ and ‘Prec.’ represent Success and Precision, respectively.

This is a very challenging task because (i) point clouds obtained with, e.g., LiDAR sensors, suffer from occlusions and point sparsity, complicating the tracker’s task of finding the object of interest; (ii) the point distribution for an object may vary significantly, making it difficult for the model to learn discriminative object features.

To tackle the above challenges, existing 3D SOT models [6, 16, 4, 9, 10, 25, 34, 18, 35] rely on large scale annotated point cloud datasets for training. Unfortunately, obtaining annotations for this task, as for many 3D tasks, is extremely time-consuming. Furthermore, as shown in Fig. 1, the performance of these methods degrades dramatically as the number of labeled samples decreases. Nevertheless, no semi-supervised or unsupervised methods have been explored so far in 3D SOT.

As shown in Fig. 2, the matching-based tracker of [34] can still track the target at the very beginning of a sequence by predicting a motion offset relative to the reference coordinate, even though there are no points in the template for appearance matching. This indicates that the appearance

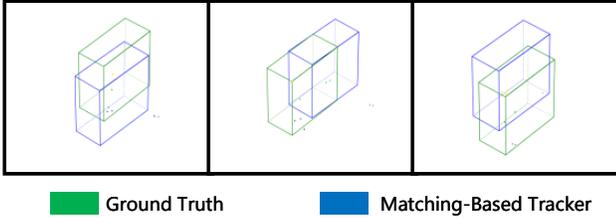


Figure 2. We observe that appearance matching-based trackers can learn the objects motion distribution and track them even in the absence of points for appearance matching. For instance, BAT [34] manages to track objects in extremely sparse point clouds.

matching-based trackers can learn motion information. By contrast, in 2D SOT, many trackers [23, 24, 36, 29] employ cycle consistency to leverage unsupervised data. Specifically, they encourage forward and backward tracking to produce consistent motions. In principle, we expect to apply this idea to 3D SOT and make trackers to learn the object’s motion distribution in unlabeled data. However, transferring these 2D methods directly to 3D is challenging. First, since the point cloud is sparse and the environment is cluttered with objects, it is hard to find meaningful patches to use as pseudo labels for training. Second, unsupervised 2D SOT methods rely on the assumption that the target appears in every frame of the sequence. Unfortunately, this assumption is not always satisfied in point cloud datasets such as KITTI [5], NuScenes [1], and Waymo [20]. This is because they are built for the multi-object tracking task, and cannot guarantee that the tracking object exists in the whole sequence.

In this paper, we introduce a label-efficient way to train 3D SOT trackers. We call it **MixCycle** - a 3D SOT approach based on a novel **SOTMixup** data augmentation strategy for semi-supervised **Cycle** tracking. Specifically, we first develop a tracking framework exploiting both self and forward-backward tracking cycles. Self tracking consistency is performed to cover the object point cloud appearance variation, and forward-backward consistency is built for learning the object’s motion distribution. Second, we present a data augmentation method for 3D SOT called **SOTMixup**, which is inspired by the success of *mixup* [31] and *Manifold mixup* [22]. Without changing the total number of points in the search area, **SOTMixup** samples points in a random point cloud and the search area point cloud according to the mixing rate and generates training samples. Specifically, the random point cloud is sampled from the labeled training set. **SOTMixup** thus increases the tracker’s robustness to point cloud variations. We evaluate **MixCycle** on KITTI, NuScenes, and Waymo. As shown in Fig. 1, our experiments clearly demonstrate the label efficiency, generalization and remarkable performance of our method on the 3D SOT task.

Contributions: (i) We propose the first semi-supervised

3D SOT framework. It exploits self and forward-backward consistency as supervision and generalizes to appearance matching-based trackers. (ii) We introduce a **SOTMixup** augmentation strategy that increases the tracker’s robustness to point distribution variations and allows it to learn motion information in extreme situations. (iii) Our framework demonstrates a remarkable performance in terms of label efficiency, achieving better results than existing supervised methods in our experiments on KITTI NuScenes, and Waymo when using fewer labels. In particular, we surpass P2B [16] trained on 100% labels while only using 10% labels.

2. Related Work

3D Single Object Tracking. Since LiDAR is insensitive to illumination, the appearance matching model has become the main choice in the field of 3D single object tracking. Giancola *et al.* [6] proposed SC3D which is the first method using a Siamese network to deal with this problem. However, it is very time-consuming and inaccurate due to heuristic matching. Zarzar *et al.* [30] built an end-to-end tracker by using 2D RPN in 2D bird’s eyes view (BEV). Unfortunately, the lack of information in one dimension leads to limited accuracy. The point-to-box (P2B) network [16] employs VoteNet[15] as object regression module to construct a point-based tracker. A number of works [4, 9, 10, 25, 34, 18] investigate different architectures of trackers based on P2B [16]. Zheng *et al.* [34] depicted an object using the point-to-box relation and proposed **BoxCloud**, which enables the model to better sense the size of objects. Hui *et al.* [9] discovered the priori information of object shapes in the dataset to obtain dense representations of objects from sparse point clouds. Zheng *et al.* [35] presented a motion centric method M^2 -Track, which is appearance matching-free and has made great progress in dealing with the sparse point cloud tracking problem. However, M^2 -track is limited by the LiDAR frequency of datasets as object motion in adjacent frames varies with the LiDAR sampling frequency.

All the above methods rely on large-scale labeled datasets. Unfortunately, 3D point cloud annotation is labor- and time-consuming. To overcome this, we propose **MixCycle**, a semi-supervised tracking method based on cycle consistency constraints with **SOTMixup** data augmentation. **Label-Efficient Visual Tracking.** Wang *et al.* [23] proposed unsupervised deep tracking (UDT) with cycle consistency, based on a Siamese correlation filter backbone network. UDT achieved remarkable performance, revealing the potential of unsupervised learning in visual tracking. Yuan *et al.* [29] improved the UDT approach to make the target features passed forward and backward as similar as possible. The self-supervised fully convolutional Siamese network [19] uses only spatially supervised learning of tar-

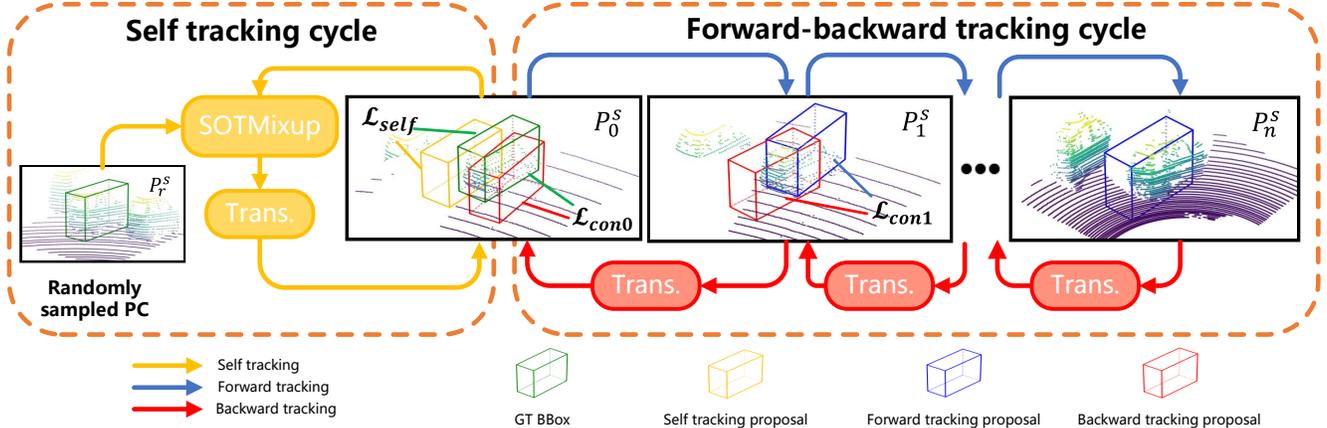


Figure 3. **MixCycle framework.** The label is only contained in P_0^s of point cloud (PC) sequence $\{P_0^s, P_1^s, \dots, P_n^s\}$. **1) Self tracking cycle:** we first sample a PC P_r^s from the labeled training set. Then, we generate pseudo labels by applying SOTMixup and random rigid transformation (Trans.) to P_0^s and P_r^s . SOTMixup directly mixes P_r^s and P_0^s based on the number of points with a mixing rate, assigning a reasonable loss weight corresponding to the mixing rate. We employ the consistency between self tracking proposals and pseudo labels to formulate the loss \mathcal{L}_{self} . **2) Forward-backward tracking cycle:** we leverage forward tracking proposals as pseudo labels and apply a random rigid transformation to them. Then, we employ the consistency between ground truth (GT)/pseudo labels and backward tracking proposals to formulate the losses $\mathcal{L}_{con0}/\{\mathcal{L}_{con1}, \dots, \mathcal{L}_{con(n-1)}\}$.

get correspondences in still video frames. Wu *et al.* [26] proposed a progressive unsupervised learning (PUL) network, which distinguishes the background by contrastive learning and models the regression result noise. PUL thus makes the tracker robust in long-time tracking. Unsupervised single object tracker [36] consists of an online-updating tracker with a novel memory learning scheme.

In essence, the above unsupervised trackers all make the implicit assumption that the tracked target exists in every frame of the sequence. Unfortunately, this is not necessarily true in KITTI [5], NuScenes [1], and Waymo [20]. Therefore, the above methods are not directly applicable to 3D SOT.

Mixup Data Augmentation. Data augmentation has become a crucial pre-processing step for many deep learning models. Zhang *et al.* [31] introduced a data augmentation method called *mixup*, which linearly interpolates two image samples, and *Manifold mixup* [22] transfers this idea to high-dimensional feature spaces. By interpolating a new sample, PointMixup [2] extends *mixup* to point clouds. Mix3D [14] introduces a scene-aware *mixup* by taking the union of two 3D scenes and their labels after random transformations. Lu *et al.* [13] developed a directed *mixup* based on the pixel values. Additionally, a variety of region *mixup* techniques have been proposed [32, 12, 21, 11]. In the case of outdoor scenes, Xiao *et al.* [28] combined two images using random rotation. CosMix[17] and structure aware fusion [7] combine point clouds using semantic structures. Fang *et al.* [3] turned a CAD into a point cloud to combat object occlusion.

However, the above-mentioned methods are made for the

multi-class classification scenario and are not suitable for SOT, which only contains a positive and negative sample.

3. Method

3.1. Overview

The purpose of 3D SOT is to continually locate the target in the search area point cloud sequence $\mathbf{P}^s = \{P_0^s, \dots, P_k^s, \dots, P_n^s | P_k^s \in \mathbb{R}^{N_s \times 3}\}$ given the tracking object template point cloud $P_0^o \in \mathbb{R}^{N_t \times 3}$ and the 3D BBox $B_0 \in \mathbb{R}^7$ in the initial frame. This can be described as

$$(\tilde{P}_{k+1}^o, \tilde{B}_{k+1}) = \mathcal{F}(P_{k+1}^s, \tilde{P}_k^o), \quad (1)$$

where \tilde{P}_k^o , \tilde{P}_{k+1}^o and \tilde{B}_{k+1} are the predicted target point cloud and 3D BBox in frame k and $k+1$, respectively. By referring to P2B [16] and its follow-ups [4, 9, 10, 18, 25, 34], we summarize the typically 3D SOT loss as

$$\mathcal{L} = \rho_1 \cdot \mathcal{L}_{cla} + \rho_2 \cdot \mathcal{L}_{prop} + \rho_3 \cdot \mathcal{L}_{reg} + \rho_4 \cdot \mathcal{L}_{box}, \quad (2)$$

where ρ is the manually-tuned hyperparameter, \mathcal{L}_{cla} , \mathcal{L}_{prop} , \mathcal{L}_{reg} and \mathcal{L}_{box} are the losses for foreground-background classification, confidences for the BBox proposals, voting offsets of the seed points, and offsets of the BBox proposals, respectively.

To address this task, we propose MixCycle, a novel semi-supervised framework for 3D SOT. Illustrated in Fig. 3, MixCycle relies on a SOTMixup data augmentation strategy to tackle data sparsity and diversity (Sec. 3.2). Further, it utilizes self and forward-backward cycle consistencies as sources of supervision to cover the object appear-

ance and motion variation (Sec. 3.3). Additionally, we apply rigid transformations to ground truth (GT) labels to generate search areas in unlabeled data (Sec. 3.4).

3.2. SOTMixup

Inspired by the great success of *mixup* [31], we develop SOTMixup to supply diverse training samples and deal with the point cloud diversity problem, providing a solution for the *mixup* application in binary classification in SOT task. With two image samples (I_A, I_B) , *mixup* can be simply describe as creating an image

$$I_A^m = \lambda I_A + (1 - \lambda) I_B, \quad (3)$$

$$y_A^m = \lambda y_A + (1 - \lambda) y_B, \quad (4)$$

where $\lambda \in [0, 1]$ is the mixing rate, and (y_A, y_B) are the image labels. Typically, λ follows a Beta distribution $\beta(\eta, \eta)$. A multi-class loss is then calculated as

$$\mathcal{L}_{muti.cal} = \lambda \cdot \mathcal{C}(\tilde{y}, y_A) + (1 - \lambda) \cdot \mathcal{C}(\tilde{y}, y_B), \quad (5)$$

where \tilde{y} is the predicted label, and \mathcal{C} is the criterion (usually being the cross-entropy loss). Vanilla *mixup* applies linear interpolation in aligned pixel space. However, this operation is not suitable to unordered point clouds. Furthermore, one of the key challenges for the SOT task is to determine whether the proposal is positive or negative. Multi-class label interpolation approach in *mixup* cannot be directly applied. Specifically, given the search area PC P_A and a random PC P_B sampled in the training set, we employ *mixup* to generate a foreground and background label pair $(\lambda \cdot y_A, (1 - \lambda) \cdot y_B)$. In practice, this label pair should be set to $(y_A, 0)$, as the points in P_B mismatch the template and should be considered as background.

We therefore develop a point cloud *mixup* strategy for SOT based on the number of points, called SOTMixup. As shown in Fig. 4, SOTMixup generates new samples and minimizes the gap between the generated samples and the real sample distribution. Specifically, SOTMixup mixes a point cloud randomly sampled from the training set and the search area point cloud by sampling points using a mixing rate, without changing the total number of points in the search area. First, given point cloud pair (P_A, P_B) , corresponding binary classification labels (y_A, y_B) , and a mixing rate λ , we separate the backgrounds and object points in P_A and P_B and obtain $(P_A^b, P_A^o, P_B^b, P_B^o)$, where $P_A^o \in \mathbb{R}^{N_A^o \times 3}$ and $P_B^o \in \mathbb{R}^{N_B^o \times 3}$. Second, we generate \hat{P}_A^o and \hat{P}_B^o by randomly sampling $\lambda \times N_A^o$ and $(1 - \lambda) \times N_A^o$ points from P_A^o and P_B^o , respectively. We then perform SOTMixup as

$$P_A^m = P_A^b + \hat{P}_A^o + \hat{P}_B^o, \quad (6)$$

where ‘+’ represents the concatenation operation.

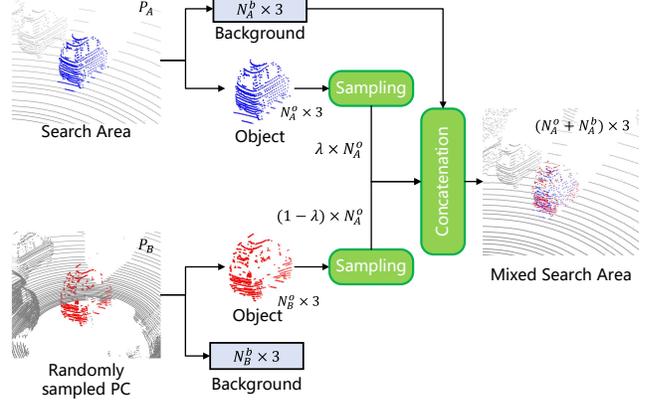


Figure 4. **SOTMixup**. First, the search area point cloud (PC) and a point cloud randomly sampled from the labeled training set are segmented into foreground and background, respectively. Second, a mixing rate λ is applied to sample two object PCs. Finally, we concatenate the sampled object PCs and search area background to generate the mixed search area.

Usually, we consider the distance between the predicted object center and the ground truth to be positive if it is less than 0.3 meters, and negative if it is greater than 0.6 meters. The binary cross entropy loss for regression and foreground classification in SOTMixup can be written as

$$\mathcal{L}_{prop.mix} = -(\lambda \cdot y_A \cdot \log(s_i^p) + (1 - y_A) \cdot \log(1 - s_i^p)), \quad (7)$$

$$\mathcal{L}_{cla.mix} = -(\lambda \cdot y_A \cdot \log(b_j^p) + (1 - y_A) \cdot \log(1 - b_j^p)), \quad (8)$$

where $\mathcal{L}_{prop.mix}$ and $\mathcal{L}_{cla.mix}$ are the proposal confidence loss and foreground-background classification loss, respectively. s_i^p is the confidence score of proposal i , and b_j^p is the predicted foreground probability of point j in search area P^s . We replace losses \mathcal{L}_{prop} and \mathcal{L}_{cla} with $\mathcal{L}_{prop.mix} = \lambda \cdot \mathcal{L}_{prop}$ and $\mathcal{L}_{cla.mix} = \lambda \cdot \mathcal{L}_{cla}$. SOTMixup applies a loss weight λ to the positive proposals and foreground points, but does not change the loss weight of the negative proposals and background points. We reduce the loss penalty on the positive sample prediction scores to lessen the influence on the appearance matching ability of the tracker. We leave the loss weight unchanged for the negative samples. Because we intend the trackers to predict the motion offset of the object even if the object point cloud in the search area has dramatically changed.

3.3. Cycle Tracking

Self Tracking Cycle. In contrast with existing 2D cycle trackers [23, 24, 36, 29], which only consider forward-backward cycle consistency, we propose to create the self tracking cycle. The motivation is to leverage virtually infinite supervision information contained in the initial frame

itself. To this end, we first randomly sample a search area point cloud P_r^s and object 3D BBox B_r from the labeled training set. SOTMixup is then applied to generate the mixed point cloud

$$P_0^m = \text{SOTMixup}(P_0^s, B_0, P_r^s, B_r, \lambda), \quad (9)$$

where $\lambda \in [0, 1]$ is the mixing rate. Inspired by SC3D [6], we apply a random rigid transformation \mathcal{T} to $\{P_0^m, B_0\}$ and generate pseudo labels

$$(P_0^{mt}, B_0^t) = \mathcal{T}(P_0^m, B_0, \alpha), \quad (10)$$

where P_0^{mt} is the search area PC generated by applying SOTMixup and a rigid transformation on P_0^s and $\alpha = (\Delta x, \Delta y, \Delta z, \Delta \theta)$ is a transformation parameter with a coordinate offset $(\Delta x, \Delta y, \Delta z)$ and a rotation degree $\Delta \theta$ around the up-axis. This creates a self tracking cycle

$$(\tilde{P}_0^{ot}, \tilde{B}_0^t) = \mathcal{F}(P_0^{mt}, P_0^o), \quad (11)$$

where \tilde{B}_0^t is the predicted result on P_0^{mt} and P_0^o is the object template PC cropped from P_0^s . We then calculate the self consistency loss \mathcal{L}_{self} between \tilde{B}_0^t and B_0^t . \mathcal{L}_{self} has the same setting with \mathcal{L} of Eq. (2) while corresponding loss with $\mathcal{L}_{cal.mix}$ and $\mathcal{L}_{prop.mix}$.

In the self tracking cycle, the loss weight can be automatically quantified by the mixing rate in SOTMixup. This provides the tracker with simple training samples to make it converge faster in the early stage of training with a high mixing rate, and also allows us to improve the tracker’s robustness to point cloud variations using a low mixing rate.

Forward-Backward Tracking Cycle. In addition to self tracking cycles, we also use forward-backward consistency. Hence, we forward track the object in the given search area sequence \mathbf{P}^s , which can be written as

$$(\tilde{P}_1^o, \tilde{B}_1) = \mathcal{F}(P_1^s, P_0^o), \quad (12)$$

$$(\tilde{P}_2^o, \tilde{B}_2) = \mathcal{F}(P_2^s, \tilde{P}_1^o), \quad (13)$$

where $\{\tilde{P}_1^o, \tilde{P}_2^o\}$ and $\{\tilde{B}_1, \tilde{B}_2\}$ are the predicted forward tracking object point clouds and 3D BBoxes of P_1^s and P_2^s , respectively. Following this strategy lets us further predict $\{\tilde{P}_3^o, \dots, \tilde{P}_n^o\}$ and $\{\tilde{B}_3, \dots, \tilde{B}_n\}$ in $\{P_3^s, \dots, P_n^s\}$.

Then, we reverse the tracking sequence and perform backward tracking while applying random rigid transformations. This can be expressed as

$$(\tilde{P}_1^{o'}, \tilde{B}_1') = \mathcal{F}(\mathcal{T}(P_1^s, \tilde{B}_1, \alpha), \tilde{P}_2^{o'}), \quad (14)$$

$$(\tilde{P}_0^{o'}, \tilde{B}_0') = \mathcal{F}(\mathcal{T}(P_0^s, B_0, \alpha), \tilde{P}_1^{o'}), \quad (15)$$

where $\{\tilde{P}_0^{o'}, \tilde{P}_1^{o'}, \tilde{P}_2^{o'}\}$ and $\{\tilde{B}_1', \tilde{B}_0'\}$ are the predicted backward tracking object point clouds and 3D BBoxes.

We then measure the consistency losses \mathcal{L}_{con1} and \mathcal{L}_{con0} between \tilde{B}_1' and \tilde{B}_1 , as well as between \tilde{B}_0' and B_0 . Similarly, we can measure the consistency losses $\{\mathcal{L}_{con2}, \dots, \mathcal{L}_{con(n-1)}\}$. \mathcal{L}_{con} has the same setting with \mathcal{L} of Eq. (2).

The forward-backward tracking cycle provides real and diverse motion consistency, leading trackers to learn the object’s motion distribution between two neighboring frames. Furthermore, the tracker’s robustness is increased by training with a disturbed template generated by the template update strategy (Sec. 3.4).

3.4. Implementation Details

Training & Testing. We train MixCycle using the SGD optimizer with a batch size of 48 and an initial learning rate of 0.01 with a decay rate of $5e - 5$ at each epoch. All experiments are conducted using NVIDIA RTX-3090 GPUs. We set $n = 2$ and only measure the self and forward-backward cycle consistency losses \mathcal{L}_{self} and \mathcal{L}_{con0} for P_0^s due to the GPU memory limit. At test time, we track the object frame by frame in a point cloud sequence with labels in the first frame. For both training and testing, our default setting for the template update strategy is to merge the target in the first frame with the predicted previous result.

Input & Data Augmentation. Our MixCycle takes three frames $f, f + 1$ and $f + 2$ as input. For the initial frame f , we transform the BBox and point clouds to the object coordinate system. For the other frames in the tracking cycle, we transform the BBox and point clouds to the predicted object coordinate system in the last frame. We assume that the motion of the objects across neighboring frames is not significant. We apply random rigid transformations to BBoxes in labeled frames and use them to crop out search areas in the neighboring frames. Only the area within 2 meters around the reference object BBox is considered as input (search area) since we are only interested in the area where the target is expected to appear. The random rigid transformation parameter α is set to $(0.3, 0.3, 0.0, 5.0^\circ)$, and the β distribution parameter is set to $\eta = 0.5$.

Loss Function. The loss function of MixCycle is defined as $\mathcal{L}_{MixCycle} = \gamma_1 \mathcal{L}_{self} + \gamma_2 \mathcal{L}_{con0}$ containing self tracking cycle losses and forward-backward tracking cycle losses. Each of our cycle losses is set according to the original loss setting of the tracker to which we apply MixCycle. Additionally, the corresponding losses are replaced with $\mathcal{L}_{prop.mix}$ and $\mathcal{L}_{cla.mix}$. We empirically set $\gamma_1 = 1.0$ and $\gamma_2 = 2.0$, as we expect the tracker to focus more on learning the motion distribution of the object.

Table 1. Overall performance comparison between our MixCycle and the fully-supervised methods on the KITTI (left) and NuScenes (right) datasets, where the percentage of labels used for training is shown under the dataset names. Improvements based on the same tracker are shown in **green**. **Bold** and underline denote the best and the second-best performance, respectively.

	Dataset	KITTI			NuScenes								
		1%	5%	10%	0.1%	0.5%	1%						
Success	P2B [16]	6.1	25.5	34.3	15.2	23.0	24.3						
	MLVSNet [25]	<u>25.0</u>	35.5	36.6	21.5	29.9	34.0						
	BAT [34]	10.1	21.6	34.9	17.5	26.4	30.6						
	Ours(P2B)	20.3	14.2↑	36.7	11.2↑	<u>43.8</u>	9.5↑	23.2	7.9↑	<u>34.3</u>	11.3↑	34.3	10.0↑
	Ours(MLVSNet)	32.4	7.4↑	38.8	3.3↑	42.6	6↑	31.4	9.9↑	34.5	4.7↑	41.9	7.9↑
	Ours(BAT)	19.7	9.6↑	42.2	20.6↑	46.2	11.3↑	<u>24.4</u>	6.9↑	32.8	6.4↑	<u>34.4</u>	3.8↑
Precision	P2B [16]	5.0	39.5	52.7	13.0	21.2	22.6						
	MLVSNet [25]	36.3	53.2	54.7	19.5	30.4	<u>35.3</u>						
	BAT [34]	12.3	35.3	52.7	15.2	25.7	30.6						
	Ours(P2B)	33.4	28.4↑	55.3	15.8↑	<u>64.2</u>	11.5↑	21.9	8.9↑	<u>34.2</u>	13↑	34.0	11.4↑
	Ours(MLVSNet)	49.2	12.9↑	<u>56.6</u>	3.4↑	61.4	6.7↑	31.1	11.6↑	35.2	4.8↑	43.6	8.3↑
	Ours(BAT)	27.0	14.7↑	62.3	27.0↑	67.8	15.1↑	<u>22.7</u>	7.5↑	31.9	6.2↑	34.1	3.5↑

Table 2. Comparison of MixCycle against fully-supervised methods on each category. We train the models with 1%/0.1% sampling rate on KITTI/NuScenes. Improvements and decreases based on the same tracker are shown in **green** and **red**, respectively.

	Dataset	KITTI(1%)					Nuscenes(0.1%)														
		Car 6424	Pedestrian 6088	Van 1248	Cyclist 308	Mean 14068	Car 64159	Truck 13587	Trailer 3352	Bus 2953	Mean 84051										
Success	P2B [16]	8.1	3.6	8.1	5.6	6.1	15.8	13.1	12.8	16.1	15.2										
	MLVSNet [25]	<u>35.3</u>	15.2	22.9	12.8	<u>25.0</u>	21.0	25.2	22.5	13.5	21.5										
	BAT [34]	16.7	3.8	7.2	6.8	10.1	17.5	17.8	20.4	14.4	17.5										
	Ours(P2B)	20.6	12.5↑	22.8	19.2↑	8.0	0.1↓	16.6	11.0↑	20.3	14.2↑	23.0	7.2↑	25.2	12.1↑	22.4	9.6↑	<u>17.7</u>	1.5↑	23.2	7.9↑
	Ours(MLVSNet)	43.8	8.5↑	<u>20.7</u>	5.5↑	28.2	5.3↑	43.7	31.0↑	32.4	7.4↑	29.7	8.7↑	42.4	17.3↑	31.3	8.9↑	19.2	5.7↑	31.4	9.9↑
	Ours(BAT)	32.6	15.9↑	6.1	2.3↑	16.3	9.2↑	<u>34.1</u>	27.4↑	19.7	9.6↑	<u>24.3</u>	6.9↑	<u>26.9</u>	9.1↑	<u>23.7</u>	3.2↑	16.9	2.5↑	<u>24.4</u>	6.9↑
Precision	P2B [16]	7.4	2.2	6.1	4.4	5.0	14.5	8.2	6.8	8.4	13.0										
	MLVSNet [25]	<u>46.5</u>	28.8	<u>25.4</u>	16.6	<u>36.3</u>	20.5	20.0	11.3	6.4	19.5										
	BAT [34]	22.7	2.9	5.9	9.5	12.3	16.3	12.2	9.2	12.2	15.2										
	Ours(P2B)	30.0	22.6↑	43.7	41.5↑	6.1	0.0	11.1	6.7↑	33.4	28.4↑	23.5	9.0↑	18.9	10.7↑	11.2	4.4↑	14.0	5.6↑	21.9	8.9↑
	Ours(MLVSNet)	59.2	12.7↑	<u>40.7</u>	11.9↑	31.1	5.7↑	79.0	62.4↑	49.2	12.8↑	31.1	10.6↑	38.6	18.6↑	19.5	8.1↑	11.5	5.2↑	31.1	11.6↑
	Ours(BAT)	43.9	21.2↑	9.3	6.4↑	19.2	13.2↑	<u>57.3</u>	47.8↑	27.0	14.7↑	<u>24.1</u>	7.8↑	<u>21.1</u>	8.9↑	13.8	4.6↑	9.7	2.6↓	<u>22.7</u>	7.5↑

4. Experiments

4.1. Datasets

We evaluate our MixCycle on the challenging 3D visual tracking benchmarks of KITTI [5], NuScenes [1] and Waymo [20] for semi-supervised 3D single object tracking. Semi-supervision labels are generated by applying random sampling to the training set.

The KITTI tracking dataset contains 21 training sequences and 29 test sequences with 8 types of objects. Following previous works [6, 16, 34, 25], we split the training set into training/validation/testing: Sequences 0-16 are used for training, 17-18 for validation, and 19-20 for testing. The NuScenes dataset contains 1000 scenes and annotations for 23 object classes with accurate 3D BBoxes. NuScenes is officially divided into 700/150/150 scenes for training/validation/testing. Following the setting in [34], we train our MixCycle on the subset “training_track” of the training set, and test it on the validation set. Waymo includes 1150 scenes, 798/202/150 scenes for training/validation/testing. Following the setting in [35], we test trackers in the validation set. Compared to KITTI, NuScenes, and Waymo include larger data volumes and

more complex scenarios.

Evaluation Metrics. In this paper, we use One Pass Evaluation (OPE) [27] to evaluate the Success (Succ.) and Precision (Prec.) of different methods. *Success* is calculated as the overlap (Intersection Over Union, IoU) between the proposal BBox and the ground truth (GT) BBox. *Precision* represents the AUC of distance error between the centers of two BBoxes from 0 to 2 meters.

4.2. Comparison with Fully-supervised Methods

To the best of our knowledge, no other 3D single object trackers work in a semi-supervision fashion. Therefore, we choose P2B [16], the multi-level voting Siamese network (MLVSNet) [25] and the box-aware tracker (BAT) [34] to validate our method by sharing the same network backbone. Note that BAT is the state-of-the-art (SOTA) method in appearance matching-based trackers, and we regard it as our upper bound. The fully-supervised methods will be trained with labeled data, as their original approaches. We train MixCycle in a semi-supervised way, which uses both labeled and unlabeled data. We do not evaluate the motion-based tracker [35] since it requires 2 consecutive labeled point clouds for training and is not suitable for our training

Table 3. Overall performance comparison on KITTI/NuScenes between our MixCycle with 10%/1% sampling rates and the fully-supervised methods with 100% sampling rate.

Dataset	Method	Success	Precision
KITTI	P2B(100%) [16] vs Ours(10%)	42.4 vs 43.8	1.4↑
	MLVSNet(100%) [25] vs Ours(10%)	45.7 vs 42.6	3.1↓
	BAT(100%) [34] vs Ours(10%)	51.2 vs 46.2	5.0↓
NuScenes	P2B(100%) [16] vs Ours(1%)	39.7 vs 34.3	5.4↓
	MLVSNet(100%) [25] vs Ours(1%)	45.7 vs 41.9	3.8↓
	BAT(100%) [34] vs Ours(1%)	41.8 vs 34.4	7.4↓

Table 4. Comparison of MixCycle against BAT on Car in KITTI for different sampling rates. ‘Improv.’ denotes Improvement.

Sampling Rate		1%	5%	10%	30%	50%	70%	100%
Succ.	BAT [34]	16.7	24.3	44.0	48.0	48.7	55.5	60.5
	Ours(BAT)	32.6	49.2	55.2	56.2	56.6	60.9	64.7
	Improv.	15.9↑	24.9↑	11.2↑	8.2↑	7.9↑	5.4↑	4.2↑
Prec.	BAT [34]	22.7	34.8	57.3	63.1	65.3	69.5	77.7
	Ours(BAT)	43.9	62.1	70.0	70.5	70.3	75.7	77.9
	Improv.	21.2↑	27.3↑	12.7↑	7.4↑	5.0↑	6.2↑	0.2↑

set generation strategy. We employ different sampling rates for the two datasets. The first reason is that we account for the different scales of the dataset. The second reason is that the case of very limited labels is more practical in real applications, and we attempt to set the sampling rate as low as possible within the trainable range.

Results on KITTI. 1) We evaluate our MixCycle in 4 categories (Car, Pedestrian, Van and Cyclist) and compare it using 3 sampling rates: 1%, 5% and 10%. As shown in Tab. 1, our method outperforms the fully-supervised approaches under all sampling rates by a large margin. This confirms the high label efficiency of our proposed semi-supervised framework.

The performance gap between our MixCycle and the fully-supervised P2B becomes larger as the proportion of labeled samples decreases. In particular, in the extreme case of 1% labels usage, we achieve **14.2%** and **28.4%** improvement in success and precision, respectively, demonstrating the impact of our approach on the baseline method. Interestingly, our MixCycle based on the SOTA fully-supervised method BAT achieves the best results ((42.2%,46.2%)/(62.3%,67.8%) in succ./prec.) with 5% and 10% sampling rates, but the MLVSNet based MixCycle takes the first place (32.4%/49.2% in succ./prec.) with 1% sampling rate. We believe this is because the multi-scale approach in MLVSNet effectively enhances the robustness of its feature representation. The BoxCloud proposed by BAT further strengthens the reliance on labels, leading to a degradation of BAT and MixCycle performance at 1% sampling rate. 2) We further present the test results on each category with 1% sampling rate in Tab. 2. We achieve better performance in all categories, except Van on P2B. We assume this to be caused by the less labeled, huge, and moving fast feature of Van, which leads the tracker hard

Table 5. Comparison of MixCycle against BAT on Waymo. MixCycle(BAT) is trained only on KITTI with a 10% sampling rate; BAT* represents BAT trained on Waymo using all the labels.

Category	Vehicle		Pedestrian		Mean	
	Frame Number	1057651	510533	1568184		
Metrics	Succ.	Prec.	Succ.	Prec.	Succ.	Prec.
BAT [34]	26.5	28.2	16.5	31.1	23.2	29.1
Ours(BAT)	<u>31.1</u>	<u>33.5</u>	25.5	46.4	<u>29.3</u>	<u>37.7</u>
BAT* [34]	35.6	44.2	<u>22.1</u>	<u>36.8</u>	31.2	41.8

to predict a precise motion. The improvements on Cyclist (**31.0%/62.4%** and **27.4%/47.8%** in succ./prec. for P2B and BAT, respectively) and Pedestrian (**19.2%/41.5%** in succ./prec. on P2B) reveal the robustness of MixCycle to point cloud variations. As pedestrians and cyclists are usually considered to have the largest point cloud variations due to their small object sizes and the diversity of body motion. 3) In Tab. 3, we compare the performance between fully-supervised methods trained with 100% labels and MixCycle trained with 10% labels. Although our performance decreases slightly on MLVSNet and BAT, MixCycle still shows a remarkable result (43.8/64.2 in succ./prec.) on P2B. With only 10% of the labels, MixCycle based on P2B outperforms the fully-supervised method (1.4%/4.2% improvement in succ./prec.) using 100% of the labels. This confirms the strong ability of our approach to leverage data information and highlights its promise for future developments. 4) We provide the comparison of MixCycle against BAT on Car from KITTI with more sampling rates in Tab. 4. Our Mixcycle not only achieves great improvements in low sampling rates, but also boosts the tracker’s performance in the fully-supervised training (4.2/0.2 in succ./prec.). 5) In Fig. 5, we compare MixCycle with BAT trained with 10% labels. MixCycle achieves a better performance in both extremely sparse and complex point clouds.

Results on NuScenes. Following the setting in BAT [34], we test our MixCycle in 4 categories (Car, Truck, Trailer and Bus). The results of P2B [16], MLVSNet [25] and BAT [34] on NuScenes are provided by M^2 -Track [35] and BAT [34]. 1) We use the published codes of the competitors to obtain results for each sampling rate. We compare them on 3 sampling rates: 0.1%, 0.5% and 1%, as NuScenes is larger than KITTI. As shown in Tab. 1, MixCycle still outperforms the fully-supervised approaches under all sampling rates. 2) Observing the individual categories in Tab. 2 evidences that MixCycle yields a remarkable improvement on Truck (**17.3%/18.6%** and **12.1%/10.7%** in succ./prec. on MLVSNet and P2B, respectively) and Car (8.7%/10.6% in succ./prec. on MLVSNet). However, MixCycle drops by 2.6% in precision on Bus, which has fewer labels, a greater size, and high velocity. 3) Moreover, in Tab. 3, we compare the performance of MixCycle trained with 1% labels and

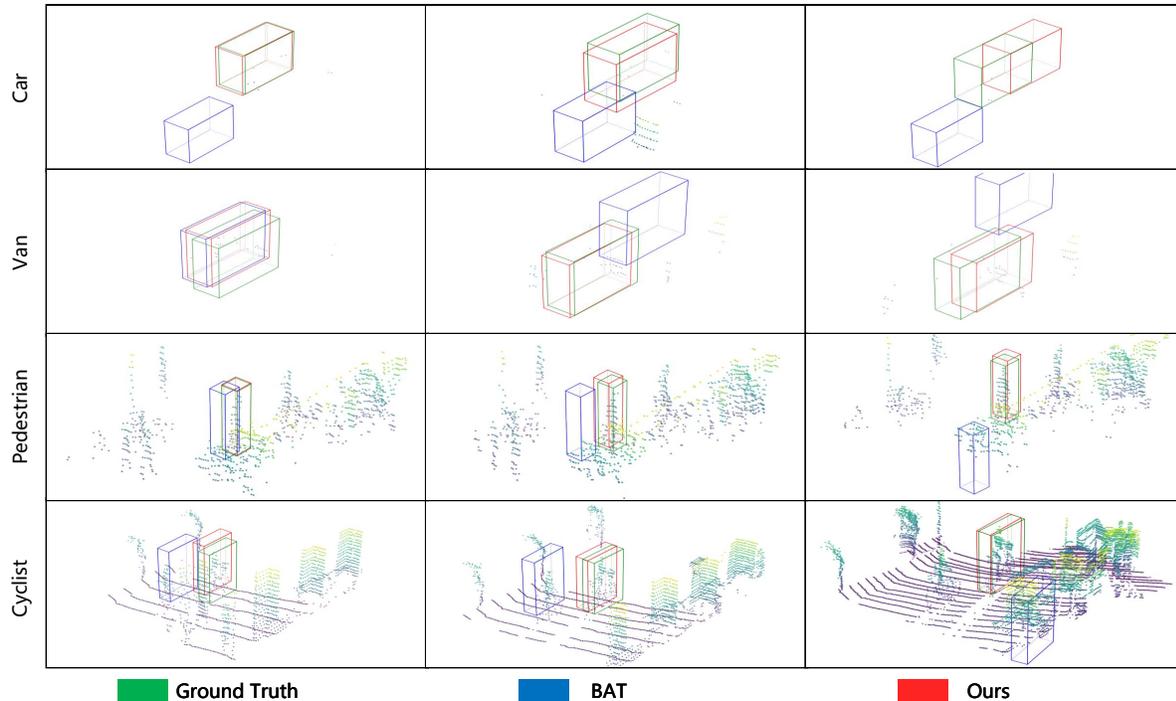


Figure 5. Visualization. **Car&Van**: Extremely sparse cases. **Pedestrian**: Medium density cases. **Cyclist**: complex environment cases.

fully-supervised methods trained with 100% labels. Our MixCycle based on MLVSNet surpasses the SOTA method BAT despite using significantly fewer labels. On such a challenging dataset with pervasive distractors and drastic appearance changes, our method exhibits even more competitive performance when using few labels.

Result on Waymo. Following the setting in [35], we test MixCycle in Vehicle and Pedestrian. We use the BAT backbone, and MixCycle(BAT) is only trained on KITTI with a 10% sampling rate. As shown in Tab. 5, MixCycle(BAT) outperforms BAT by 6.1% and 8.6% in terms of mean ‘succ.’ and ‘prec.’ values, respectively. *More impressively, our performance is close to the fully supervised BAT trained on Waymo (31.2%/41.75% in succ./prec. reported in [34]).* To summarize, our MixCycle still delivers excellent results on large-scale datasets.

4.3. Analysis Experiments

In this section, we extensively analyze MixCycle with a series of experiments. First, we study the effectiveness of each component in MixCycle. Second, we further analyze the influence of forward-backward cycle step sizes. Finally, we compare the various application ways of SOT-Mixup. All the experiments are conducted on KITTI with 10% sampling rate and with BAT as the backbone network, unless otherwise stated.

Ablation Study. We conduct experiments to analyze the effectiveness of different modules in MixCycle. First, we ver-

ify our assumption that appearance matching-based trackers can learn the object’s motion distribution. As shown in Tab. 6, the cycle tracking framework yields better performance when using only forward-backward cycle than when using only self cycle (3.5%/5.4% improvement in succ./prec.). Additionally, this supports the intuition that real and diverse motion information is helpful to appearance trackers. Furthermore, combining them boosts the results. Note that the random rigid transformation is necessary for self cycle, otherwise the GT BBox will always be fixed at the origin of the coordinate system. Second, we evaluate the effectiveness of random rigid transformation and SOTMixup in the framework. The performance grows significantly after applying SOTMixup (5.1%/6% improvement in succ./prec.), demonstrating the importance of SOT-Mixup in semi-supervised tasks. Random rigid transformation plays a negative role in the cycle tracking framework but is practical in MixCycle. We conjecture this to be due to the missing target in the search area caused by random rigid transformation. This phenomenon may occur when the target moves rapidly and the model makes wrong predictions. Applying SOTMixup to the cycle tracking framework can significantly improve the model’s tracking capabilities and address this issue.

Flexibility of MixCycle. We further explore the effect of forward-backward tacking cycle step size. As shown in Tab. 7, we conduct experiments with 2 and 3 step cycles, respectively. Compared to the 2 step cycle, 3 step cycle

Table 6. Results of MixCycle when different modules are ablated. ‘Self’, ‘F-B. Cycle’ and ‘Trans.’ stand for self cycle, forward-backward cycle and random rigid transformation in the forward-backward cycle, respectively.

Self	F-B. Cycle	Trans.	SOTMixup	Success	Precision		
✓				34.9	11.3↓	52.7	15.1↓
	✓			38.4	7.8↓	58.1	9.7↓
✓	✓			39.5	6.7↓	59.2	8.6↓
✓	✓	✓		38.8	7.4↓	59.3	8.5↓
✓	✓		✓	44.6	1.6↓	65.2	2.6↓
✓	✓	✓	✓	46.2		67.8	

Table 7. Analysis of the forward-backward cycle step size.

	Success	Precision
2 Steps	45.8	66.6
3 Steps	46.2	67.8
Improvement	0.4↑	1.2↑

Table 8. Results of SOTMixup with different settings. ‘Template’ and ‘Search Area’ indicate we apply SOTMixup with different inputs. ‘Self’ indicates we apply SOTMixup in the self cycle. ‘Backward’ means we apply SOTMixup in backward tracking P_1^s to P_0^s .

Self	Backward	Template	Search Area	Success	Precision		
	✓		✓	44.7	1.5↓	64.7	3.1↓
✓	✓		✓	45.4	0.8↓	67.6	0.2↓
✓		✓		42.3	3.9↓	63	4.8↓
✓		✓	✓	44.5	1.7↓	66.1	1.7↓
✓			✓	46.2		67.8	

achieves a better performance in both success and precision. This experiment demonstrates the potential of MixCycle for further growth in step size. We believe that larger step sizes can provide a template point cloud disturbed in long sequence tracking, leading to improved model robustness. Furthermore, according to [24], a larger step size more effectively penalizes inaccurate localization.

Influence of SOTMixup. In Tab. 8, we compare SOTMixup with different settings. First, we analyze the effect of applying SOTMixup to different inputs (Template and Search Area) while only using it in self cycle. The performance drops when we apply it to the template and to both the template and search area. We consider this to be due to the mismatch with real tracking. As the template is usually accurate while the search area point cloud varies significantly. Note that we share the same λ when taking SOTMixup in both the template and search area, and set $\lambda = 1$ in the SOTMixup loss. Second, we explore various ways of exploiting SOTMixup in MixCycle. SOTMixup leads to performance degradation when we apply it to backward tracking. This is caused by the disturbance of the template point cloud in backward tracking, making the loss weights misaligned. This further proves that the loss weights provided by SOTMixup are reliable.

5. Conclusion

In this paper, we have presented the first semi-supervised framework, MixCycle, for 3D SOT. Its three main components, self tracking cycle, forward-backward tracking cycle and SOTMixup, have been proposed to achieve robustness to point cloud variations and percept object’s motion distribution. Our experiments have demonstrated that MixCycle yields high label efficiency and outperforming fully-supervised approaches using scarce labels.

In the future, we plan to develop a more robust tracking network backbone for MixCycle, and thus further enhance its 3D SOT performance.

Acknowledgments. This work is supported in part by the National Natural Science Foundation of China (No. 62176242 and 62002295), and NWPU international cooperation and exchange promotion projects (No. 02100-23GH0501). Thank Haozhe Qi for the discussion and preliminary exploration.

A. Implementation Details

Framework Architecture. The overall pipeline with the grad flow of MixCycle is shown in Fig. 6. Due to the limitation of non maximum suppression (NMS) on gradient back-propagation, we only calculate the gradients of the directly supervised parts.

SOTMixup. Given the mix point cloud $P_A^m = P_A^b + \hat{P}_A^o + \hat{P}_B^o$ and Bounding-box B_A in label y_A at the SOTMixup, we only regard the points in B_A as the foreground points. Specifically, the points in \hat{P}_B^o are considered as background noise if they are outside the B_A . We believe that modifying the size of the tracking target is incompatible with real tracking.

B. More Analysis

Training & Inference Time. We compare MixCycle and fully-supervised methods [16, 25, 34] in training time shown in Tab. 9. They are trained on Car in KITTI with 10% labels using 2 NVIDIA RTX-3090 GPUs. Our MixCycle takes around 2.0 ~ 2.5 times as long as the fully-supervised methods. The experiments reveal that MixCycle requires a longer training time, but it is still in an acceptable range. Hence, we could expect a faster and more robust tracking network backbone for MixCycle.

Table 9. Training time comparison of MixCycle and fully-supervised methods on Car in KITTI with 10% labels using 2 NVIDIA RTX-3090 GPUs. Decreases based on the same tracker is shown in red.

Method	Time
P2B [16]	1h22m
MLVSNet [25]	1h47m
BAT [34]	1h22m
Ours(P2B)	3h35m 2h13m↓
Ours(MLVSNet)	3h32m 1h45m↓
Ours(BAT)	3h40m 2h18m↓

Frame Number of Cycle Tracking and Unlabeled Data Losses Balance. 1) Because of the limited memory of an NVIDIA RTX-3090 GPU, only a maximum of 2 cycle consistencies among 3 frames can be supervised. Therefore, we only present the losses for the self-supervised part. 2) For the one without labels part, we have made experiments to balance those losses. We try to supervise different consistencies in a two-stage training by supervising \mathcal{L}_{self} and \mathcal{L}_{con0} in stage 1, and \mathcal{L}_{self} and \mathcal{L}_{con1} in stage 2, based on BAT with a 10% sampling rate on KITTI. Without SOTMixup, the cycle framework achieves 38.8/59.3 and 41.0/60.6 in Succ./Prec. in stage 1 & 2, respectively. The performance drops in stage 2 if we use SOTMixup. We conjecture this to be due to conflicts between the delicate losses set by SOTMixup in Self Cycle and the ambiguous losses in

F.B. Cycle. We leave the design of a better training strategy for MixCycle as future work.

Fairness of Comparison with Fully-supervised Method.

Here we discuss fairness in the comparison experiments. The fully supervised method solely relies on labeled data, whereas our method utilizes both labeled and unlabeled data. 1) The intention of our work is to reduce the effort in data annotation. While reducing the cost of collecting data is also important, we constitute a different research task on its own. 2) We refer to a semi-supervised 3D object detection method SESS’s [33] experimental setting for comparison experiments. SESS directly reduces the usage of data of fully supervised methods for comparison experiments because no other method shares the same semi-supervised settings with it, which is very similar to our situation. 3) We present the performance comparison using the same amount of data but with different label usage in the Tab. 3 and Tab. 4 of the paper.

Further Details. We further demonstrate the test result on each category and sample rate on KITTI and Nuscenes shown in Tab. 10 and Tab. 11. We achieve great success on Cyclist. The maximum improvement on the Cyclist class is up to **44.77%/75.83%** in success/precision based on P2B [16] with 10% labels. For the most important class Car in KITTI and NuScenes, MixCycle also achieves a remarkable improvement in every sample rate.

C. Visualization

SOTMixup. Our MixCycle leverages SOTMixup to supply diverse training samples. As shown in Fig. 7, we present SOTMixup in a variety of categories. Our SOTMixup completes the point cloud of the occluded area in the Van in Fig. 7, making the training samples more diverse. In the Car in Fig. 7, SOTMixup almost removes the point cloud of the source object, allowing the trackers to regress the correct target center by learning the distribution of object motion in extreme cases.

KITTI Results. We present the visualization results of the comparison between Our MixCycle and BAT [34] with 10% sample rate in Fig. 8. The visualization results further validate the superiority of our approach in sparse and complex scenarios.

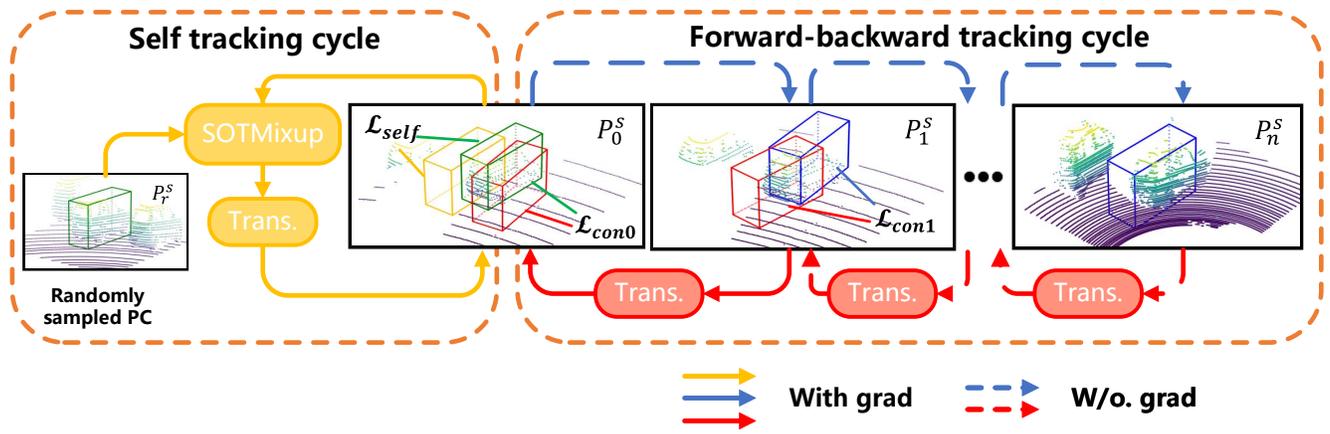


Figure 6. The framework of MixCycle. The gradient flow is represented by solid and dashed lines with arrows.

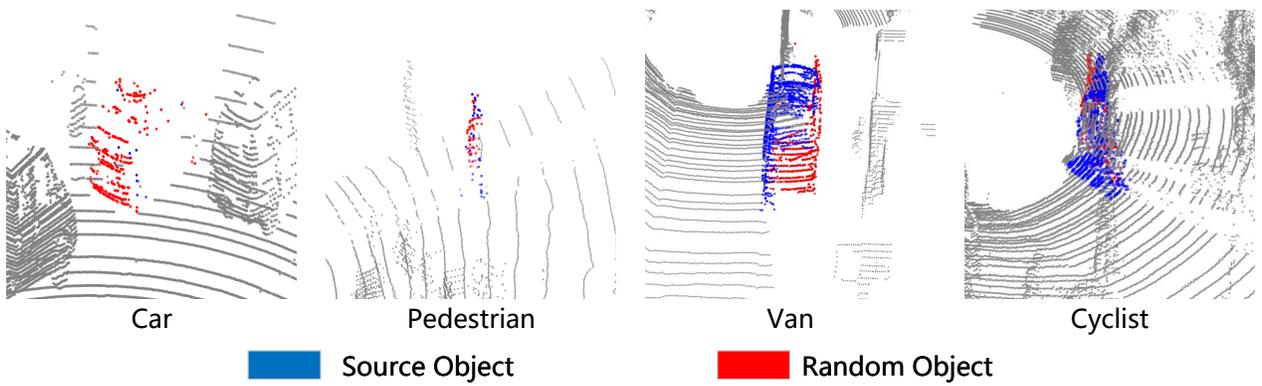


Figure 7. Visualization of SOTMixup.

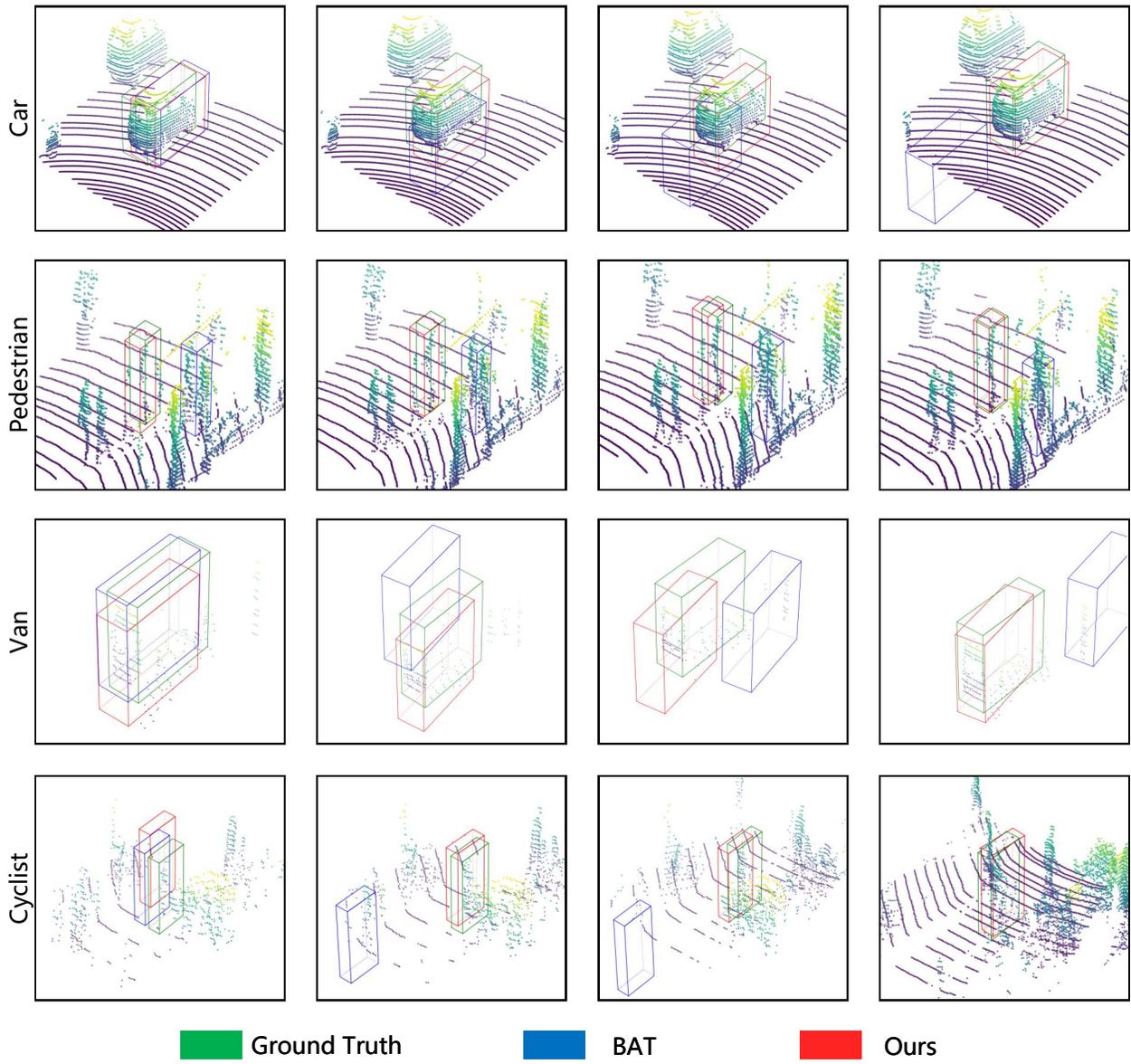


Figure 8. Visualization results. Our MixCycle and BAT are trained with 10% labels on KITTI.

Table 10. Comparison of MixCycle against fully-supervised methods on each category in KITTI. Improvements and decreases based on the same tracker are shown in **green** and **red**, respectively. **Bold** and underline denote the best and the second-best performance, respectively.

Category		Car	Pedestrian	Van	Cyclist	Mean						
Frame Number		6424	6088	1248	308	14068						
Success	1%	P2B [16]	8.11	3.61	8.10	5.60	6.11					
		MLVSNet [25]	<u>35.27</u>	15.15	<u>22.94</u>	12.76	<u>24.98</u>					
		BAT [34]	16.69	3.81	7.17	6.77	10.05					
		Ours(P2B)	20.56	12.45↑	22.76	19.15↑	7.97	0.13↓	16.62	11.02↑	20.31	14.20↑
		Ours(MLVSNet)	43.75	8.48↑	<u>20.68</u>	5.53↑	28.22	5.28↑	43.73	30.97↑	32.39	7.41↑
		Ours(BAT)	32.63	15.94↑	6.08	2.27↑	16.33	9.16↑	<u>34.12</u>	27.35↑	19.73	9.67↑
	5%	P2B [16]	33.99	20.31	12.10	5.73	25.51					
		MLVSNet [25]	43.50	28.09	<u>35.06</u>	19.77	35.55					
		BAT [34]	24.30	21.00	13.17	13.25	21.62					
		Ours(P2B)	44.13	10.14↑	<u>31.01</u>	10.70↑	26.15	14.05↑	36.77	31.04↑	36.70	11.19↑
		Ours(MLVSNet)	52.44	8.94↑	24.04	4.05↓	38.73	3.67↑	<u>46.54</u>	26.77↑	<u>38.80</u>	3.26↑
		Ours(BAT)	<u>49.24</u>	24.94↑	37.63	16.63↑	26.08	12.91↑	50.08	36.83↑	42.18	20.56↑
	10%	P2B [16]	41.94	30.63	19.61	7.37	34.31					
		MLVSNet [25]	48.21	24.76	37.90	24.89	36.64					
		BAT [34]	43.96	28.84	18.12	35.84	34.95					
		Ours(P2B)	45.82	3.88↑	41.59	10.96↑	42.59	22.98↑	<u>52.14</u>	44.77↑	<u>43.84</u>	9.53↑
		Ours(MLVSNet)	<u>54.08</u>	5.87↑	30.39	5.63↑	<u>41.29</u>	3.39↑	49.95	25.06↑	42.60	5.97↑
		Ours(BAT)	55.19	11.23↑	<u>38.62</u>	9.78↑	34.92	16.8↑	55.52	19.68↑	46.23	11.28↑
Precision	1%	P2B [16]	7.39	2.24	6.07	4.42	4.98					
		MLVSNet [25]	<u>46.54</u>	28.80	<u>25.41</u>	16.62	<u>36.33</u>					
		BAT [34]	22.66	2.92	5.94	9.54	12.35					
		Ours(P2B)	29.97	22.58↑	43.73	41.49↑	6.08	0.01↑	11.12	6.70↑	33.39	28.41↑
		Ours(MLVSNet)	59.24	12.7↑	<u>40.72</u>	11.92↑	31.08	5.67↑	79.03	62.41↑	49.16	12.83↑
		Ours(BAT)	43.87	21.21↑	9.32	6.40↑	19.18	13.24↑	<u>57.31</u>	47.77↑	27.02	14.67↑
	5%	P2B [16]	45.99	40.26	10.82	5.43	39.50					
		MLVSNet [25]	57.53	52.07	<u>42.30</u>	28.77	53.19					
		BAT [34]	34.81	40.35	15.55	25.52	35.30					
		Ours(P2B)	56.94	10.95↑	<u>58.04</u>	17.78↑	30.92	20.1↑	67.33	61.90↑	55.34	15.83↑
		Ours(MLVSNet)	66.61	9.08↑	47.15	4.92↓	45.26	2.96↑	<u>81.06</u>	52.29↑	<u>56.61</u>	3.42↑
		Ours(BAT)	<u>62.07</u>	27.26↑	68.05	27.70↑	30.81	15.26↑	82.63	57.11↑	62.33	27.04↑
	10%	P2B [16]	56.11	57.70	21.73	7.35	52.68					
		MLVSNet [25]	63.63	48.31	44.65	35.08	54.69					
		BAT [34]	57.25	56.08	21.48	19.69	52.75					
		Ours(P2B)	58.30	2.19↑	72.05	14.35↑	51.83	30.1↑	<u>83.18</u>	75.83↑	<u>64.22</u>	11.54↑
		Ours(MLVSNet)	<u>67.36</u>	3.73↑	56.28	7.97↑	<u>50.01</u>	5.36↑	82.52	47.44↑	61.36	6.67↑
		Ours(BAT)	70.02	12.77↑	<u>69.83</u>	13.75↑	42.28	20.8↑	85.37	65.68↑	67.81	15.06↑

Table 11. Comparison of MixCycle against fully-supervised methods on each category in NuScenes.

Category		Car	Truck	Trailer	Bus	Mean						
Frame Number		64159	13587	3352	2953	84051						
Success	0.1%	P2B [16]	15.77	13.09	12.81	16.12	15.23					
		MLVSNet [25]	20.99	25.16	22.46	13.53	21.46					
		BAT [34]	17.46	17.75	20.43	14.42	17.52					
		Ours(P2B)	23.01	7.24↑	25.22	12.13↑	22.37	9.56↑	17.65	1.53↑	23.15	7.92↑
		Ours(MLVSNet)	29.67	8.68↑	42.43	17.27↑	31.34	8.88↑	19.22	5.69↑	31.43	9.97↑
		Ours(BAT)	<u>24.32</u>	6.86↑	<u>26.88</u>	9.13↑	<u>23.66</u>	3.23↑	16.92	2.50↑	<u>24.45</u>	6.93↑
	0.5%	P2B [16]	24.42	19.21	20.30	12.38	22.99					
		MLVSNet [25]	29.82	32.25	27.40	<u>22.74</u>	29.87					
		BAT [34]	27.71	22.85	25.48	15.44	26.40					
		Ours(P2B)	36.85	12.43↑	28.23	9.02↑	21.75	1.45↑	21.14	8.76↑	<u>34.30</u>	11.31↑
		Ours(MLVSNet)	31.49	1.67↑	46.75	14.50↑	48.49	21.09↑	28.47	5.73↑	34.53	4.66↑
		Ours(BAT)	<u>32.20</u>	4.49↑	<u>38.22</u>	15.37↑	<u>31.04</u>	5.56↑	21.82	6.38↑	32.76	6.36↑
	1%	P2B [16]	23.95	27.83	25.84	14.57	24.32					
		MLVSNet [25]	33.23	<u>39.08</u>	39.62	22.23	34.04					
		BAT [34]	30.66	32.73	32.83	17.81	30.63					
		Ours(P2B)	<u>34.80</u>	10.85↑	35.24	7.41↑	30.40	4.56↑	22.61	8.04↑	33.43	9.10↑
		Ours(MLVSNet)	40.61	7.38↑	45.43	6.35↑	58.09	18.47↑	35.38	13.15↑	41.90	7.86↑
		Ours(BAT)	33.72	3.06↑	37.29	4.56↑	45.55	12.72↑	<u>24.26</u>	6.45↑	<u>34.44</u>	3.81↑
Precision	0.1%	P2B [16]	14.52	8.20	6.82	8.41	12.98					
		MLVSNet [25]	20.45	19.97	11.31	6.35	19.51					
		BAT [34]	16.31	12.16	9.19	<u>12.22</u>	15.21					
		Ours(P2B)	23.48	8.96↑	18.88	10.68↑	11.20	4.38↑	13.99	5.58↑	21.91	8.94↑
		Ours(MLVSNet)	31.05	10.60↑	38.57	18.60↑	19.45	8.14↑	11.53	5.18↑	31.12	11.60↑
		Ours(BAT)	<u>24.10</u>	7.79↑	<u>21.07</u>	8.91↑	<u>13.81</u>	4.62↑	9.67	2.55↓	<u>22.69</u>	7.48↑
	0.5%	P2B [16]	24.28	12.32	11.08	6.98	21.21					
		MLVSNet [25]	32.73	26.71	14.91	<u>15.35</u>	30.44					
		BAT [34]	28.69	18.06	15.09	8.89	25.73					
		Ours(P2B)	39.22	14.94↑	20.79	8.47↑	11.19	0.11↑	13.27	6.29↑	<u>34.21</u>	13.00↑
		Ours(MLVSNet)	<u>34.17</u>	1.44↑	42.78	16.07↑	38.71	23.8↑	19.59	4.24↑	35.23	4.80↑
		Ours(BAT)	33.35	4.66↑	<u>32.21</u>	14.15↑	<u>18.62</u>	3.53↑	14.49	5.60↑	31.92	6.18↑
	1%	P2B [16]	23.70	22.86	14.11	7.51	22.61					
		MLVSNet [25]	36.76	33.91	29.60	15.41	35.26					
		BAT [34]	32.47	28.36	20.42	11.19	30.58					
		Ours(P2B)	36.72	13.02↑	29.15	6.29↑	18.17	4.06↑	14.78	7.27↑	33.99	11.38↑
		Ours(MLVSNet)	45.07	8.31↑	40.17	6.26↑	46.28	16.68↑	25.01	9.60↑	43.62	8.36↑
		Ours(BAT)	35.29	2.82↑	32.30	3.94↑	<u>32.63</u>	12.21↑	<u>17.15</u>	5.96↑	34.06	3.49↑

References

- [1] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscnets: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019. [2](#), [3](#), [6](#)
- [2] Yunlu Chen, Vincent Tao Hu, Efstratios Gavves, Thomas Mensink, Pascal Mettes, Pengwan Yang, and Cees GM Snoek. Pointmixup: Augmentation for point clouds. In *Proceedings of the European Conference on Computer Vision*, pages 330–345, 2020. [3](#)
- [3] Jin Fang, Xinxin Zuo, Dingfu Zhou, Shengze Jin, Sen Wang, and Liangjun Zhang. Lidar-aug: A general rendering-based augmentation framework for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4710–4720, 2021. [3](#)
- [4] Zheng Fang, Sifan Zhou, Yubo Cui, and Sebastian Scherer. 3d-siamrpn: An end-to-end learning method for real-time 3d single object tracking using raw point cloud. *IEEE Sensors Journal*, 21(4):4995–5011, 2020. [1](#), [2](#), [3](#)
- [5] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3354–3361, 2012. [1](#), [2](#), [3](#), [6](#)
- [6] Silvio Giancola, Jesus Zarzar, and Bernard Ghanem. Leveraging shape completion for 3d siamese tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1359–1368, 2019. [1](#), [2](#), [5](#), [6](#)
- [7] Frederik Hasecke, Martin Alsfasser, and Anton Kummert. What can be seen is what you get: Structure aware point cloud augmentation. In *2022 IEEE Intelligent Vehicles Symposium (IV)*, pages 594–599. IEEE, 2022. [3](#)
- [8] Christopher Hazard, Akshay Bhagat, Balarama Raju Budharaju, Zhongtao Liu, Yunming Shao, Lu Lu, Sammy Omari, and Henggang Cui. Importance is in your attention: agent importance prediction for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2532–2535, 2022. [1](#)
- [9] Le Hui, Lingpeng Wang, Mingmei Cheng, Jin Xie, and Jian Yang. 3d siamese voxel-to-bev tracker for sparse point clouds. *Advances in Neural Information Processing Systems*, 34:28714–28727, 2021. [1](#), [2](#), [3](#)
- [10] Le Hui, Lingpeng Wang, Linghua Tang, Kaihao Lan, Jin Xie, and Jian Yang. 3d siamese transformer network for single object tracking on point clouds. *arXiv preprint arXiv:2207.11995*, 2022. [1](#), [2](#), [3](#)
- [11] Dogyoon Lee, Jaeha Lee, Junhyeop Lee, Hyeongmin Lee, Minhyeok Lee, Sungmin Woo, and Sangyoun Lee. Regularization strategy for point cloud via rigidly mixed sample. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15900–15909, 2021. [3](#)
- [12] Sanghyeok Lee, Minkyu Jeon, Injae Kim, Yunyang Xiong, and Hyunwoo J Kim. Sagemix: Saliency-guided mixup for point clouds. *arXiv preprint arXiv:2210.06944*, 2022. [3](#)
- [13] Yuheng Lu, Fangping Chen, Ziwei Zhang, Fan Yang, and Xiaodong Xie. Directed mix contrast for lidar point cloud segmentation. In *Proceedings of the IEEE International Conference on Multimedia and Expo*, pages 1–6, 2022. [3](#)
- [14] Alexey Nekrasov, Jonas Schult, Or Litany, Bastian Leibe, and Francis Engelmann. Mix3d: Out-of-context data augmentation for 3d scenes. In *2021 International Conference on 3D Vision (3DV)*, pages 116–125. IEEE, 2021. [3](#)
- [15] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9277–9286, 2019. [1](#), [2](#)
- [16] Haozhe Qi, Chen Feng, Zhiguo Cao, Feng Zhao, and Yang Xiao. P2b: Point-to-box network for 3d object tracking in point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6329–6338, 2020. [1](#), [2](#), [3](#), [6](#), [7](#), [10](#), [13](#), [14](#)
- [17] Cristiano Saltori, Fabio Galasso, Giuseppe Fiameni, Nicu Sebe, Elisa Ricci, and Fabio Poiesi. Cosmix: Compositional semantic mix for domain adaptation in 3d lidar segmentation. In *Proceedings of the European Conference on Computer Vision*, pages 586–602, 2022. [3](#)
- [18] Jiayao Shan, Sifan Zhou, Zheng Fang, and Yubo Cui. Ptt: Point-track-transformer module for 3d single object tracking in point clouds. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1310–1316, 2021. [1](#), [2](#), [3](#)
- [19] Chon Hou Sio, Yu-Jen Ma, Hong-Han Shuai, Jun-Cheng Chen, and Wen-Huang Cheng. S2siamfc: Self-supervised fully convolutional siamese network for visual tracking. In *Proceedings of the ACM International Conference on Multimedia*, pages 1948–1957, 2020. [2](#)
- [20] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2446–2454, 2020. [2](#), [3](#), [6](#)
- [21] Ardian Umam, Cheng-Kun Yang, Yung-Yu Chuang, Jen-Hui Chuang, and Yen-Yu Lin. Point mixswap: Attentional point cloud mixing via swapping matched structural divisions. In *Proceedings of the European Conference on Computer Vision*, pages 596–611, 2022. [3](#)
- [22] Vikas Verma, Alex Lamb, Christopher Beckham, Amir Najafi, Ioannis Mitliagkas, David Lopez-Paz, and Yoshua Bengio. Manifold mixup: Better representations by interpolating hidden states. In *Proceedings of the International Conference on Machine Learning*, pages 6438–6447, 2019. [2](#), [3](#)
- [23] Ning Wang, Yibing Song, Chao Ma, Wengang Zhou, Wei Liu, and Houqiang Li. Unsupervised deep tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1308–1317, 2019. [2](#), [4](#)
- [24] Ning Wang, Wengang Zhou, Yibing Song, Chao Ma, Wei Liu, and Houqiang Li. Unsupervised deep representation learning for real-time tracking. *International Journal of Computer Vision*, 129(2):400–418, 2021. [2](#), [4](#), [9](#)
- [25] Zhoutao Wang, Qian Xie, Yu-Kun Lai, Jing Wu, Kun Long, and Jun Wang. Mlvsnet: Multi-level voting siamese net-

- work for 3d visual tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3101–3110, 2021. [1](#), [2](#), [3](#), [6](#), [7](#), [10](#), [13](#), [14](#)
- [26] Qiangqiang Wu, Jia Wan, and Antoni B Chan. Progressive unsupervised learning for visual object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2993–3002, 2021. [3](#)
- [27] Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. Online object tracking: A benchmark. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2411–2418, 2013. [6](#)
- [28] Aoran Xiao, Jiaying Huang, Dayan Guan, Kaiwen Cui, Shijian Lu, and Ling Shao. Polarmix: A general data augmentation technique for lidar point clouds. *arXiv preprint arXiv:2208.00223*, 2022. [3](#)
- [29] Weihao Yuan, Michael Yu Wang, and Qifeng Chen. Self-supervised object tracking with cycle-consistent siamese networks. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 10351–10358. IEEE, 2020. [2](#), [4](#)
- [30] Jesus Zarzar, Silvio Giancola, and Bernard Ghanem. Efficient bird eye view proposals for 3d siamese tracking. *arXiv preprint arXiv:1903.10168*, 2019. [2](#)
- [31] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. [2](#), [3](#), [4](#)
- [32] Jinlai Zhang, Lyujie Chen, Bo Ouyang, Binbin Liu, Jihong Zhu, Yujin Chen, Yanmei Meng, and Danfeng Wu. Pointcutmix: Regularization strategy for point cloud classification. *Neurocomputing*, 505:58–67, 2022. [3](#)
- [33] Na Zhao, Tat-Seng Chua, and Gim Hee Lee. Sess: Self-ensembling semi-supervised 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11079–11087, 2020. [1](#), [10](#)
- [34] Chaoda Zheng, Xu Yan, Jiantao Gao, Weibing Zhao, Wei Zhang, Zhen Li, and Shuguang Cui. Box-aware feature enhancement for single object tracking on point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13199–13208, 2021. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [10](#), [13](#), [14](#)
- [35] Chaoda Zheng, Xu Yan, Haiming Zhang, Baoyuan Wang, Shenghui Cheng, Shuguang Cui, and Zhen Li. Beyond 3d siamese tracking: A motion-centric paradigm for 3d single object tracking in point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8111–8120, 2022. [1](#), [2](#), [6](#), [7](#), [8](#)
- [36] Jilai Zheng, Chao Ma, Houwen Peng, and Xiaokang Yang. Learning to track objects from unlabeled videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13546–13555, 2021. [2](#), [3](#), [4](#)