

PointACL: Adversarial Contrastive Learning for Robust Point Clouds Representation under Adversarial Attack

Junxuan Huang
University at Buffalo
Junxuanh@buffalo.edu

Yatong An
Xmotors.ai
yatong.an@gmail.com

Cheng Lu
Xmotors.ai
luc@xiaopeng.com

Bai Chen
Xmotors.ai
chenbai@xiaopeng.com

Junsong Yuan
University at Buffalo
jsyuan@buffalo.edu

Chunming Qiao
University at Buffalo
qiao@buffalo.edu

Abstract

Despite recent success of self-supervised based contrastive learning model for 3D point clouds representation, the adversarial robustness of such pre-trained models raised concerns. Adversarial contrastive learning (ACL) is considered an effective way to improve the robustness of pre-trained models. In contrastive learning, the projector is considered an effective component for removing unnecessary feature information during contrastive pretraining and most ACL works also use contrastive loss with “projected” feature representations to generate adversarial examples in pretraining, while “unprojected” feature representations are used in generating adversarial inputs during inference. Because of the distribution gap between “projected” and “unprojected” features, their models are constrained of obtaining robust feature representations for downstream tasks. We introduce a new method to generate high-quality 3D adversarial examples for adversarial training by utilizing virtual adversarial loss with “unprojected” feature representations in contrastive learning framework. We present our robust aware loss function to train self-supervised contrastive learning framework adversarially. Furthermore, we find selecting high difference points with the Difference of Normal (DoN) operator as additional input for adversarial self-supervised contrastive learning can significantly improve the adversarial robustness of the pre-trained model. We validate our method, POINTACL on downstream tasks, including 3D classification and 3D segmentation with multiple datasets. It obtains comparable robust accuracy over state-of-the-art contrastive adversarial learning methods.

1. Introduction

Among various 3D representation methods, point clouds are popular for scene understanding and visual analysis. Tasks include 3D object classification, detection, and segmentation. Despite its popularity, adversarial robustness of its learned 3D perception models, namely, robustness against adversarial samples, is a major security concern in real-world application. Perturbed adversarial samples like point adding[34], point dropping[38] or point shifting[21] can easily mislead 3D perception models.

Adversarial training (AT) [22] and its variants[8, 4, 26] are considered effective defense strategies against adversarial attacks. However, they rely on class labels to generate adversarial samples that are used to supervise model training for robustness and it is difficult to conduct unsupervised training methods like self-supervised learning. Sun *et al.* [27] proposed a pretext tasks based self-supervised (SSL) adversarial training method, which forces the model to learn robust representation from solving a pre-designed pretext task without class labels. Some previous works like [32] build an additional networks connected to classifier as 3D point clouds purifier. Researchers later introduced contrastive learning framework as an extension to AT and SimCLR[9]. It projects feature representations into a different dimension space for conducting contrastive loss and the contrastive loss is also used to obtain adversarial examples of 2D images without using class labels [19, 17, 11]. But their contrastive loss guided adversarial examples rely on projected feature in contrastive learning framework, while unprojected feature are used in generating adversarial inputs in downstream task. Because of the distribution gap between unprojected feature and projected feature, contrastive loss guided adversarial examples can only provide limited robustness during adversarial pretraining stage.

In this work, we designed an adversarial contrastive

learning model specially for 3D point clouds, POINTACL. Specifically, in order to mitigate the distribution gap between training examples and testing adversarial inputs, we propose a novel method for generating adversarial examples with unprojected features. We introduce virtual adversarial loss[23] on the unprojected features to calculate gradient direction for perturbation, which is better than prior adversarial contrastive learning methods[19, 17, 11]. Leveraging the virtual adversarial loss, our method lowers the divergence between adversarial samples and perturbed adversarial inputs in testing, and experiments show that our model has an advantage in downstream tasks under adversarial robustness testing. Meanwhile, to enhance the robustness of the feature representation, we choose to add normal of point clouds surface information in adversarial pretraining. We first select point clouds from Difference of Normal(DoN) operator with information on surface gradient and treat those high difference points as additional input in pretraining. During adversarial contrastive learning, we extracted projected representations of high difference points and incorporated them into multi-view contrastive loss.

To verify the advantages of our proposed network, we applied our pre-trained network with standard linear finetuning on two downstream tasks : 3D classification and 3D segmentation. Specifically, for the classification network, we trained on ModelNet and tested on ModelNet and shapeNet. The robust accuracy on ModelNet achieved **27.51%** compared to **4.03%** without adversarial training and in shapeNet we achieved **13.34%** compared to **2.13%** without adversarial training. For 3D segmentation task, robust accuracy achieved was **39.08%** compared to **13.82%** without adversarial training.

Our major contributions can be summarized as follows:

① We proposed POINTACL, an adversarial contrastive learning framework for point clouds data, and we found that using the virtual adversarial loss to generate high-quality adversarial samples during pre-training stage can bring more robustness representation in downstream tasks.

② We verified that high-difference point clouds selected from the difference of normal (DoN) operator can contribute to the robustness of 3D representation learning. Experiments on downstream tasks verified the criticalness of high-difference points in improving the network’s adversarial robustness.

③ We extensively benchmarked our pre-trained model with other adversarial contrastive learning models. Our pre-trained models tested on two downstream tasks : *3D object classification* (on ModelNet40) and *3D segmentation* (on S3DIS) under standard linear finetuning POINTACL led to new improved state-of-the-art robust accuracy. For example, in the 3D object classification task with I-FGM attack under budget with $\epsilon=0.01m$ in 7 steps we achieve **18.72%** and **17.24%** robustness improvement over existing adver-

sarial contrastive learning methods[19] and [17].

2. Related works

In this section, we overview the progress on three related topics: adversarial attacking methods for point clouds, self-supervised Learning of point clouds and defensive methods in improving adversarial robustness.

2.1. Adversarial attack on Point Clouds

The robustness of deep learning model on 3D point clouds has attracted many researchers due to its applications in robotics and safe-driving cars. Existing 3D adversarial attack methods can be roughly divided into three classes: optimization-based, gradient-based, and generation-based. For gradient-based methods, Liu *et al.* [21] extended I-FGSM[20] into the 3D point clouds domain by perturbing the point coordinates. Zheng *et al.* [38] proposed an iterative point dropping attack by building a gradient-based saliency map. In optimization-based methods, xianget *al.* [34] first proposed to generate adversarial point clouds using C&W attack framework[5] by point perturbation and adding. However, those perturbed point clouds usually contain many outliers, which are not human-unnoticeable. To reduce those outliers, Wen *et al.* [30] focus on generating adversarial point cloud with much less outliers. LG-GAN [39] is a generation-based 3D attack method, which uses GANs[13] to generate adversarial point clouds that follows the input target labels.

2.2. Self-supervised Learning of Point Clouds

Many approaches [31, 1, 12, 36] have been proposed for unsupervised learning and generation of point clouds, but high-level downstream tasks like 3D object classification and segmentation is less discussed. Some recent work try to demonstrate the potentials for high-level tasks such as 3D object classification and 3D semantic segmentation. Those self-supervised learning methods can be divided into two classes: pretext task based approach and contrastive learning based approach. Poursaeed *et al.* [25] designed a pretext task that predicts rotations angle of point clouds object and achieved good performance in 3D object classification. Compare to pretext task based approach, contrastive learning methods have advantage in transferability and generalisability, they have achieved great success in 2D image domain. For 3D point clouds domain, Xie *et al.* [35] used contrastive loss between two different transformations results of point clouds. STRL[15] extended BYOL[14] structure to point clouds domain by utilizing the spatio-temporal contexts and structures of point clouds and it achieved state-of-the-art performance in various high-level downstream tasks.

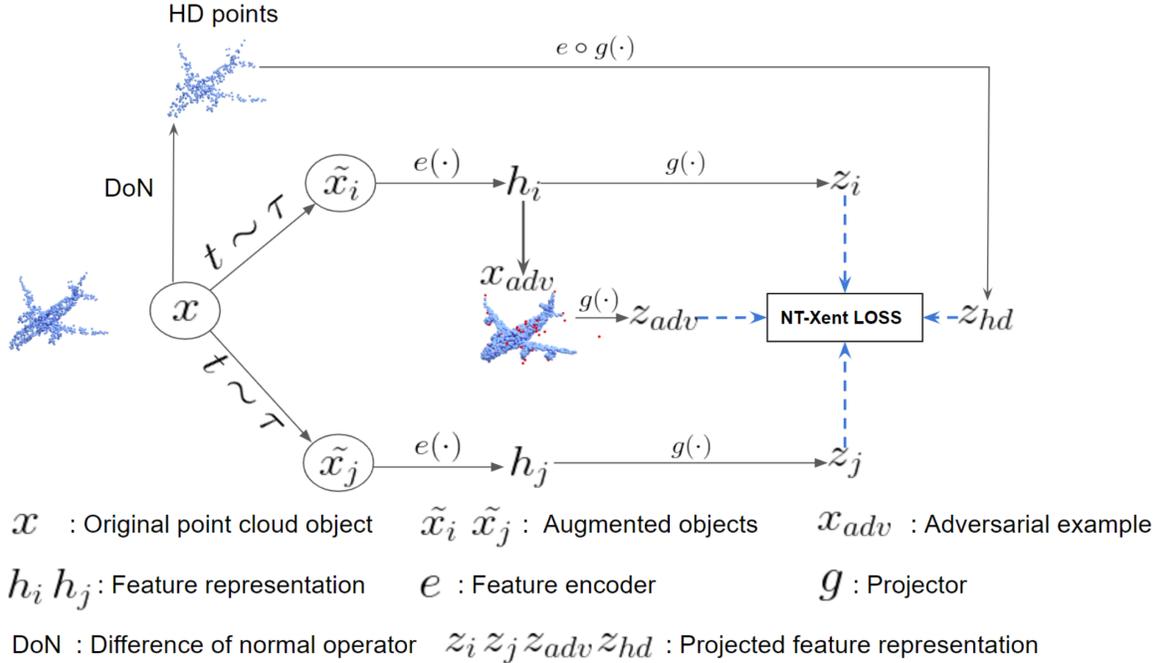


Figure 1: **Overview of PointACL.** x : Input point clouds object, At first, the input point clouds object x will generate three different versions of inputs: $\tilde{x}_i, \tilde{x}_j, x_{hd}$, corresponding to two augmented objects and high-difference point object. Adversarial examples will be generated based on one of the augmented object. Same as SimCLR[9], we have a mlp based nonlinear projection head g to improve representation quality and a multi-view normalized cross entropy loss with adjustable temperature (NT-Xent[9]) will be optimized based on projected feature representations z as the contrastive loss.

2.3. Adversarial robustness

Many defense methods have been proposed to improve model robustness against adversarial attacks. Adversarial training(AT)[22] provides one of the most effective defense methods by training the model over the adversarially perturbed training data. Zhou *et al.* [40] used Statistical Outlier Removal (SOR) method for removing points with a large kNN distance. They also proposed DUPNet, which is a combination of SOR and a point cloud upsampling network PU-Net[37]. In contrastive learning, RoCL[19] and ACL[17] AdvCL[11] extended AT by using contrastive loss to eliminate the need for class labels when generating adversarial samples. VAT[23] proposed a new regularization method based loss to generate adversarial samples on unlabeled data in Semi-Supervised Learning.

3. Our approach

3.1. Definition and Notations

Point Cloud. For classification task, a point cloud object is represented as $(\mathbf{x} \triangleq \{\mathbf{p}_i\}_{i=1\dots N}, y)$, where $\mathbf{p}_i \in \mathbb{R}^3$ is a 3D point and N is the number of points in the point cloud; $y \in \{1, 2, \dots, k\}$ is the ground-truth label, where

k is the number of classes. For semantic segmentation, a point cloud is $(\mathbf{x} \triangleq \{\mathbf{p}_i\}_{i=1\dots N}, \{y_i\}_{i=1\dots N})$, which $y_i \in \{1, 2, \dots, k\}$ represents ground-truth label for each point.

Contrastive learning. The main idea of contrastive learning is to learn representations without supervision by maximizing agreements of different augmented views of the same point cloud. Inspired by SimCLR[9] and STRL[15], we build our 3D contrastive learning framework. Specifically, consider augmentations \mathcal{T} (combination of rotation, translation, scaling, cropping, cutout, jittering and down-sampling) in Fig.1, each unlabeled point cloud x is augmented as \tilde{x}_i and \tilde{x}_j . The feature encoder generates output feature (h_i, h_j) from pair $(\tilde{x}_i, \tilde{x}_j)$. Then, MLP-based projector is applied for feature extraction. The projected feature representations (z_i, z_j) are optimized under a contrastive loss ℓ_{NT} (NT-Xent) to maximize their agreement. After training, we keep the the encoder part as the pre-trained model for downstream tasks.

Unsupervised Adversarial Training. Different from adversarial training (AT) which is supervised learning, Unsupervised Adversarial Training (UAT) does not use labels

to generate adversarial samples. Compared with other unsupervised adversarial training algorithms, ours improves robustness of 3D point cloud (self-supervised) pretraining model in an unsupervised manner.

3.2. Problem statement

Our method aims to improve the robustness of the 3D self-supervised contrastive learning pretraining model. Following [22], we add adversarial examples during the pretraining stage and aim to minimize contrastive loss between adversarial examples and normal inputs. So that we could achieve robust feature representations in the pretraining stage and, after linear finetuning, make more accurate predictions under adversarial attacks, improving the model’s adversarial robustness performance on downstream tasks such as 3D classification and 3D semantic segmentation. We can formulate our problem as:

$$\text{Pretraining: } \min_{\theta} \mathbb{E}_{x \in \mathcal{D}} \ell_{\text{CL}}(x + \delta, x; \theta) \quad (1)$$

$$\text{Linear finetuning: } \min_{\theta_c} \mathbb{E}_{(x,y) \in \mathcal{D}} \ell_{\text{CE}}(\phi_{\theta_c} \circ e_{\theta}(x), y) \quad (2)$$

where \mathcal{D} denotes the training set, θ represents the parameters of the model. x denotes the original point clouds object. ℓ_{CL} denotes the designed robustness aware contrastive loss with parameter θ , and the adversarial perturbation δ under budget ϵ . During linear finetuning phase (2), ℓ_{CE} is the cross-entropy loss that optimize parameters of linear prediction head ϕ_{θ_c} and e_{θ} is the fixed robustness aware feature encoder that we obtained after adversarial pretraining stage (1). $\phi_{\theta_c} \circ e_{\theta}$ denotes the classifier by equipping the linear prediction head ϕ_{θ_c} on top of the fixed feature encoder e_{θ} .

3.3. Robust Adversarial Contrastive Learning

Adversarial examples We now introduce our method on how to achieve adversarial robustness of representations in contrastive learning manner. Because the mechanism of adversarial training is to minimize loss of adversarial samples, the first step of our method is to generate label-free adversarial samples during contrastive learning. Fig.2 shows the visualization result of one adversarial example and we can see only a few perturbed points are human-noticeable.

A number of prior works [17, 11, 19] have proposed to use projected representation from self supervised model and contrastive loss to guide attacking algorithm’s generation of adversarial samples. Inspired by Virtual Adversarial Training (VAT)[23], we introduce a new method that uses Kullback-Leibler divergence (KLD) of unprojected representation between adversarial samples and augmented inputs to calculate gradient direction for updating adversarial perturbation. Following the idea of AT[22], the loss func-

Algorithm 1: Generate adversarial point clouds examples with untargeted I-FGM

Input: A set of point cloud objects x ;
 Augmentation family \mathcal{T} ; feature encoder e ;
 perturbation budget ϵ , number of steps t

Result: adversarial samples x_{adv}

Augment x to be $(\tilde{x}_i, \tilde{x}_j)$ with two augmentations sampled from \mathcal{T} .

Generate a initial small random perturbation δ and $x_{adv} = \tilde{x}_i + \delta$

for $i \in t$ **do**

Generate the corresponding adversarial point clouds with

$$\delta = \operatorname{argmax}_{\|\delta\|_{\infty} \leq \epsilon} KLD(e(x'), e(x_{adv}))$$

$$x_{adv} = x_{adv} + \delta$$

end

return: x_{adv}



Figure 2: Visualization of adversarial samples based on augmented input x_i , highly perturbed points are marked red

tion in leading attack algorithm can be written as

$$L_{adv}(x, \theta) = Div [q(y|x), p(y|x + \delta, \theta)] \quad (3)$$

$$\text{where } \delta = \operatorname{argmax}_{\|\delta\|_{\infty} \leq \epsilon} Div [q(y|x), p(y|x + r, \theta)], \quad (4)$$

where $Div[p, p']$ is a non-negative function that measures the divergence between two distributions p and p' . Because the true distribution of the output label, $q(y|x)$, is unknown, we use its *current* estimate $p(y|x, \theta)$. The goal of this loss function is to approximate the true distribution $q(y|x)$ by a parametric model $p(y|x, \theta)$ that is robust against adversarial attack to x .

Since we do not know the label y in self-supervised training, we rewrite eq:4 by approximating divergence of unprojected representation’ distributions between perturbed object $x + \delta$ and augmented original object x in contrastive learning framework.

$$L_{adv}(x, \theta) = Div [p(y|x, \theta), p(y|x + \delta, \theta)] \quad (5)$$

$$\delta = \operatorname{argmax}_{\|\delta\|_{\infty} \leq \epsilon} Div [e_{\theta}(x), e_{\theta}(x + \delta)], \quad (6)$$

High difference points cloud Because high-frequency information is a crucial contributing factor to improving the robustness of perception network in 2D image domain [28, 24, 11]. For 3D point clouds, we notice difference between normal of local surface and normal of global surface can reflect the frequency of surface gradient. We propose to use Difference of Normals(DoN)[16] to build a point saliency map according to multi-scale normal estimation and select higher difference points as an additional input during pretraining stage to enhance robustness of our model.

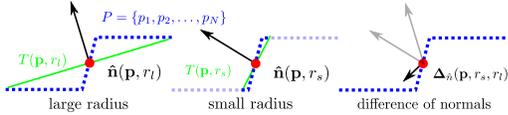


Figure 3: The normal support radius’ relation to scale.[16]

The Difference of Normals (DoN) operator $\Delta_{\hat{n}}$ for any point \mathbf{p} in a point cloud \mathbf{x} , is defined as:

$$\Delta_{\hat{n}}(p, r_1, r_2) = \frac{\hat{n}(p, r_1) - \hat{n}(p, r_2)}{2}, \quad (7)$$

where $r_1, r_2 \in \mathbb{R}$, $r_1 < r_2$, and $\hat{n}(\mathbf{x}, r)$ is the surface normal estimated at point \mathbf{p} , given the support radius r . For each point clouds object or scene x with number of points N , we remove low difference points based on $\Delta_{\hat{n}}(p, r_1, r_2)$ and keep high difference points with number $c \times N$, where $c \in (0, 1)$. Let \mathcal{D} denote DoN operator, an input point clouds object/scene X can then be decomposed into its high difference part X_{hd} and low-difference part X_{ld} :

$$[X_{hd}, X_{ld}] = \mathcal{D}(X). \quad (8)$$

The motivation is that surface normals estimated at any given radius reflect the underlying geometry of the surface at the scale of the support radius. By calculating the difference of multi-radius estimated surface normals, we can obtain the surface gradient. Thus, we can use DoN to select points which have high frequency information and use them as an additional view in robust contrastive learning eq:11.



Figure 4: Visualization result for selected points after removing 75% lower difference points in ModelNet40

Multi-view robust contrastive learning We first review the NT-Xent loss used in SimCLR[9]. The contrastive loss with a positive augmentation pair $(\tau_1(x), \tau_2(x))$ from each data input x is given by

$$\ell_{\text{CL}}(\tau_1(x), \tau_2(x)) = - \sum_{i=1}^2 \sum_{j \in \mathcal{P}(i)} \log \frac{\exp(\text{sim}(z_i, z_j)/t)}{\sum_{k \in \mathcal{N}(i)} \exp(\text{sim}(z_i, z_k)/t)}, \quad (9)$$

where $z_i = g \circ e(\tau_i(x))$ is the projected feature representation under the i th view, $\mathcal{N}(i)$ represents the set of augmented batch data not including the point $\tau_i(x)$. $\mathcal{P}(i)$ denote the set of positive views except i . $\text{sim}(z_{i1}, z_{i2})$ is the cosine similarity between projected representations z_{i1}, z_{i2} from two views of the same data x and t is a temperature parameter.

For involving adversarial examples and high-difference points as additional inputs, we follow the multi-view contrastive loss in [18]

$$\ell_{\text{CL}}(\tau_1(x), \tau_2(x), \dots, \tau_m(x)) = - \sum_{i=1}^m \sum_{j \in \mathcal{P}(i)} \log \frac{\exp(\text{sim}(z_i, z_j)/t)}{\sum_{k \in \mathcal{N}(i)} \exp(\text{sim}(z_i, z_k)/t)}, \quad (10)$$

where m denotes the number of views as input. By taking selected point clouds x_{hd} from DoN as third view input and adversarial samples $(x + \delta)$ as fourth view input, we propose our contrastive loss function as

$$\ell_{\text{CL}}(\tau_1(x), \tau_2(x), (x_{hd}), (x + \delta)) \quad (11)$$

To minimize the distance between representation from normal input x and other inputs, adversarial samples $x + \delta$ and high difference points x_{hd} , we add $\mathcal{KL}\mathcal{D}(h_1, h_{adv})$ and $\mathcal{KL}\mathcal{D}(h_{adv}, h_{hd})$ as the regularization term to our object loss function, where $h_1 = e(\tau_1(x))$ and $h_{adv} = e(x + \delta)$. Thus, we can force our backbone network to generate more stable and similar representations for both normal point clouds inputs and adversarial attacking point clouds inputs resulting in downstream task models that are more robust to adversarial samples. Eq:12 shows our loss function for updating parameters.

$$\mathcal{L}_{\text{total}} := \ell_{\text{CL}}(\tau_1(x), \tau_2(x), (x_{hd}), (x + \delta)) + \alpha \mathcal{KL}\mathcal{D}(h_1, h_{adv}) + \beta \mathcal{KL}\mathcal{D}(h_{adv}, h_{hd}) \quad (12)$$

4. Experimental Results

Evaluation setting In this section, we validated our adversarial unsupervised learning method on public benchmark point cloud datasets. Specifically, we evaluated our method on two downstream tasks: classification and segmentation. We did not perform model full finetuning since

Algorithm 2: Algorithm of Pretraining

Input: A set of point clouds \mathbf{x} , a set of high difference point clouds \mathbf{x}_{hd} ; Augmentation family \mathcal{T} ; Network backbone and projection head e, g ;

Result: The parameters θ in e and g ;

for sampled mini-batch \mathbf{x} **do**

 Augment \mathbf{x} to be $(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j)$ with two augmentations sampled from \mathcal{T} .

 Generate the adversarial mini-batch $(\tilde{\mathbf{x}}_i + \delta)$ with algorithm 1

$\ell = \ell_{\text{CL}}(\tilde{\mathbf{x}}_i, \tilde{\mathbf{x}}_j, (\tilde{\mathbf{x}}_i + \delta)) + \alpha \mathcal{KL}\mathcal{D}(h_1, h_{adv}) + \beta \mathcal{KL}\mathcal{D}(h_{adv}, h_{hd})$

 Update parameters (θ_e, θ_g) to minimize ℓ .

end

tuning full network weights is not possible for the model to preserve robustness[10]. Instead, we froze the pretrained backbone weights from our self-supervised training to keep the robustness and performed standard finetuning (**SF**) of prediction heads in different downstream tasks. For attacking algorithm, we chose untarget ℓ_∞ I-FGSM[21] for point clouds with a fixed iteration number and budget ϵ . In 3D classification task, we ran 7 iteration steps and $\epsilon = 0.01m$ in attacking during robustness evaluation. In 3D segmentation, we ran 15 iteration steps and $\epsilon = 0.08m$ in attacking during robustness evaluation.

Pre-training For the backbone pretraining stage, we use the Adam optimizer with a cosine decay learning rate schedule, the exponential moving average parameter starts with $\tau_{start} = 0.996$ and is gradually increased to 1 during the training. For 3D object classification task, we train our PointNet backbone network in 50 epochs with learning rate 0.001 over 256 batch size and the number of input is 2048. For 3D object segmentation task, we train our DGCNN backbone network in 100 epochs with learning rate 0.0002 over 32 batch size and the number of input is 4096.

We used a combination of the following augmentation method to construct the augmentation family as in [15]: Random rotation, Random translation, Random scaling, Random cropping, Random cutout, Random jittering, Random drop-out, Down-sampling, Normalization.

- Random rotation: For each axis, we draw random angles within 15° and rotate around it.
- Random translation: We translate the point cloud globally within 10% of the point cloud dimension.
- Random scaling: We scale the point cloud with a factor $s \in [0.8, 1.25]$.

- Random cropping: A random 3D cuboid patch is cropped with a volume uniformly sampled between 60% and 100% of the original point cloud. The aspect ratio is controlled within $[0.75, 1.33]$.
- Random cutout: A random 3D cuboid is cut out. Each dimension of the 3D cuboid is within $[0.1, 0.4]$ of the original dimension.
- Random jittering: Each point’s 3D locations are shifted by a uniformly random offset within $[0, 0.05]$.
- Random drop-out: We randomly drop out 3D points by a drop-out ratio within $[0, 0.7]$.
- Down-sampling: We down-sample point clouds based on the encoder’s input dimension by randomly picking the necessary amount of 3D points.
- Normalization: We normalize the point cloud to fit a unit sphere while training on synthetic data.

4.1. 3D Object Classification

We first verified our proposed method on classification task. In this downstream task, we used PointNet [7] with SimCLR[9] as backbone. During evaluation, we apply two different settings: standard finetune and adversarial full finetune (AFF). For standard finetune, we add a linear layer to the representation obtained from the pretrained backbone in self-supervised training and only finetune the linear layer with training data. For adversarial full finetune (AFF), we first use attacking algorithm I-FGSM[21] to generate an adversarial example for each input of training data in a supervised manner. Then, we add a linear layer and we finetune the whole network including the pretrained backbone with original training data and adversarial examples.

We trained the backbone network on ModelNet40[33] dataset with $\alpha = 1, \beta = 1$, and evaluated downstream tasks on ModelNet40[33] dataset. For each point cloud object we selected 2048 points as input with only coordinates and removed 512 low difference points with DoN operator as HD input.

Experimental results are shown in Table.1. For standard finetune, our pretraining model PointACL achieved 27.51% robust accuracy under I-FGSM attack in testing set of ModelNet40, which improved **23.48%** over baseline pretraining model SimCLR. We also had more than **14%** robust accuracy improvement over other adversarial contrastive learning methods.

4.2. 3D Sementic Segmentation

In 3D Segmentic Segmentation, we used DGCNN [29] with SimCLR[9] structure as backbone. During evaluation, we added 2-layered MLP network to the representation obtained from the pretrained backbone in self-supervised training for segmentataion. We trained the bakcbone network on S3DIS[2] dataset on area 1-5 with $\alpha = 1, \beta = 1$,

Table 1: **Performance result of different methods in evaluation.** Standard Accuracy(SA) represent the accuracy in test dataset and Robust Accuracy (RA) means the accuracy evaluated under test dataset generated by untarget I-FGSM[21] for point clouds with $\epsilon=0.01m$ in 7 steps

| Training type | Method | Standard Accuracy(%) | Robust Accuracy(%) |
|--------------------------|-----------------------|----------------------|--------------------|
| Supervised | AT[22] | 82.54 | 44.49 |
| Self-supervised+finetune | SimCLR[9] | 86.33 | 4.03 |
| | RoCL[19] | 85.22 | 8.72 |
| | ACL[17] | 85.85 | 10.25 |
| | PointACL(Ours) w/o HD | 82.28 | 24.26 |
| | PointACL(Ours) | 80.71 | 27.51 |

and evaluated on area 6. We selected 4096 points as input with only coordinates and removed 1024 lower difference points with DoN eq:7 as HD input.

As Table.2 shows, our pretraining model PointACL increased **25.23%** robust accuracy and **13.52%** mIoU in robustness evaluation over baseline model SimCLR. Compared to other adversarial contrastive learning methods, our method had more than **15.13%** robust accuracy with RoCL[19] and more than **12.26%** robust accuracy with ACL[17].

Table 2: **Performance result of different methods in evaluation.** Standard Accuracy(SA) represent the accuracy and (S-mIoU) means the standard mean IoU in area 6 and Robust Accuracy (RA)/R-mIoU represents the accuracy/IoU evaluated under test dataset generated by untarget I-FGSM[21] for point clouds with $\epsilon=0.08m$ in 15 steps

| Training type | Method | SA(%) | S-mIoU(%) | RA(%) | R-mIoU(%) |
|--------------------------|-----------------------|-------|-----------|--------------|--------------|
| Supervised | AT[22] | 80.29 | 52.99 | 55.61 | 28.53 |
| Self-supervised+finetune | SimCLR[9] | 82.37 | 55.54 | 13.85 | 5.60 |
| | RoCL[19] | 79.53 | 52.06 | 23.95 | 11.37 |
| | ACL[17] | 78.10 | 49.79 | 26.82 | 12.41 |
| | PointACL(Ours) w/o HD | 79.06 | 50.27 | 36.16 | 19.01 |
| | PointACL(Ours) | 78.69 | 49.85 | 39.08 | 19.12 |

4.3. Tradeoff between Robust Accuracy (RA) and Standard Accuracy(SA)

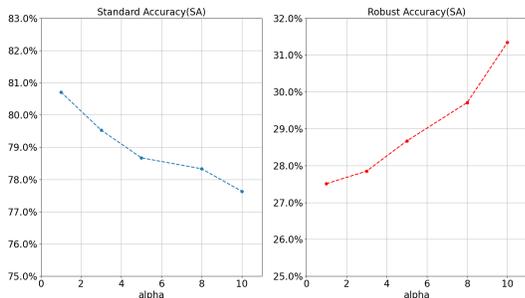


Figure 5: Standard Accuracy(SA) and Robust Accuracy(RA) with different α

We found that tuning the hyperparameter α yielded in-

creasing performance of robust accuracy (RA) at the cost of decreasing standard accuracy(SA). Fig.5 shows the RA and SA under different α in 3D object classification task. One possible reason is that increasing α forces the backbone network to generate similar unprojected representation h between adversarial sample and normal point clouds. This helps the model make more accurate predictions when facing attacking inputs but misleads it when the input data is clean.

4.4. Robustness evaluation vs. attack strength

The strength of attacking algorithm I-FGM[20] are affected by two things:(a) The number of iteration steps; increasing the number of iterations produced strong adversarial samples but increasing the number of steps did not generate stronger adversarial samples[3]. The left graph of Figure.6 shows that attacking strength slowed after iteration 5 and stopped increasing after iteration 25. Our method outperformed other methods in robust accuracy at different number of iterations; (b) The attacking budget(m) ϵ , which sets perturbation boundary of attacking, is also very important for attacking algorithm; increasing ϵ significantly enhances the attacking strength. We tested all the pretraining methods with budget from 0.001m to 0.02m with iteration 5. We found the performance of ROCL and ACL to be nearly 0% when $\epsilon > 0.02m$. Our method was better when $\epsilon > 0.002m$.

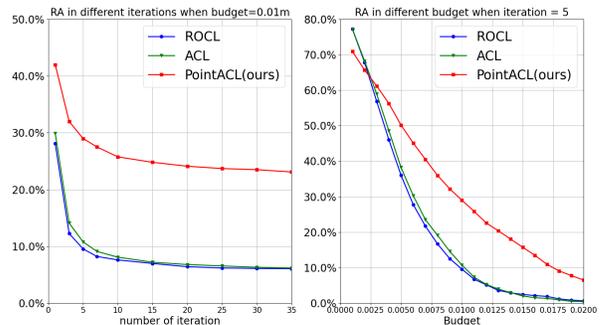


Figure 6: Robust accuracy for 3D classification in ModelNet with different steps and budget. **Left:** We fix attacking budget=0.01m and set different iteration.**Right:** We fix iteration=5 and set different attacking budget.

4.5. Robustness transferability across datasets

In table.3, we evaluated the transferability of model robustness across different datasets for downstream 3D classification task. Following the same setting from Section. 4.1, where $A \rightarrow B$ denotes the transferability from pretraining on dataset A to finetuning on another dataset $B (\neq A)$. We chose ShapeNetCore[6] as our secondary dataset. It contained over 37000 3D point cloud objects with 55 categories

in the training set while ModelNet40 only had around 9840 objects. In table.3, *ShapeNetCore* \rightarrow *ModelNet40* scenario had very stable standard accuracy for all pretraining methods. Our method’s robust accuracy was 13.57% higher than baseline model and 8.18% higher than RoCL when $\alpha=10$. In *ModelNet40* \rightarrow *ShapeNetCore* scenario, our method significantly improved the robustness performance over baseline and other methods while suffering only a small drop to standard accuracy. This pattern of large RA gain and small SA drop was also observed in 4.3.

Table 3: **Performance result of different methods in cross dataset evaluation between ModelNet40 and ShapeNet-Core.** ModelNet40 \rightarrow ShapeNetCore means we pretrained the model on ModelNet40 and linear evaluate on ShapeNet-Core dataset. The robust accuracy evaluated under test dataset generated by untarget I-FGSM[21] for point clouds with $\epsilon=0.01m$ in 7 steps

| Pretraing method | ModelNet40 \rightarrow ShapeNetCore | | ShapeNetCore \rightarrow ModelNet40 | |
|------------------------------|---------------------------------------|--------------|---------------------------------------|--------------|
| | SA(%) | RA(%) | SA(%) | RA(%) |
| PointNet (SimCLR)[9] | 81.10 | 2.13 | 85.49 | 5.39 |
| RoCL[19] | 81.29 | 2.46 | 85.01 | 10.78 |
| ACL[17] | 81.93 | 2.99 | 85.29 | 16.09 |
| PointACL(Ours, $\alpha=1$) | 80.01 | 13.34 | 85.04 | 10.56 |
| PointACL(Ours, $\alpha=10$) | 77.16 | 26.12 | 84.56 | 18.96 |

5. Ablation Study

5.1. Projected Representation vs Unprojected Representation

In SimCLR[9], the authors showed that a nonlinear projection head improves the representation quality for contrastive learning. A number of prior works [19, 17, 11] used contrastive loss with projected feature representation z to generate adversarial samples during adversarial contrastive learning. In our pretraining method, we achieved better robustness performance (Table.4) using unprojected feature representation h . To further evaluate this approach, we also experimented with using unprojected feature representation h instead of z to calculate contrastive loss in generating adversarial samples for RoCL[19] and ACL[17]. Because we wanted to focus only on the importance of feature representation selection during adversarial contrastive learning, we didn’t include high difference(HD) point clouds input in this experiment.

The results in Table.4 showed that using unprojected feature representation h resulted in better robustness performance in our method but not in other adversarial contrastive learning methods. If we look at the regularization term $\mathcal{KLD}(h_1, h_{adv})$ in our pretraining loss function, we find an explanation for this. The regularization improves the model’s robustness by minimizing the distance between h_1 and h_{adv} (unprojected representation). Because the prediction result from robustness testing on downstream tasks is based on unprojected representation h and not projected

representation z , replacing h with z in our method will decrease the robustness performance of the model during adversarial contrastive training.

Table 4: **3D classification performance under ModelNet40 with different representation choice** The robust accuracy evaluated under test dataset generated by untarget I-FGSM[21] for point clouds with $\epsilon=0.01m$ in 7 steps

| Pretraining method | Standard Accuracy(%) | Robust Accuracy(%) |
|---------------------------------|----------------------|--------------------|
| SimCLR[9] | 86.33 | 4.03 |
| RoCL[19] | 85.22 | 8.72 |
| RoCL[19](use h) | 85.98 | 5.83 |
| ACL[17] | 85.85 | 10.25 |
| ACL[17](use h) | 86.46 | 6.24 |
| PointACL(Ours)(use z) w/o HD | 85.62 | 8.72 |
| PointACL(Ours)(use h) w/o HD | 82.28 | 24.26 |

5.2. Loss function analysis

To better understand the importance of $\mathcal{KLD}(h_1, h_{adv})$ and $\mathcal{KLD}(h_{adv}, h_{hd})$ in our pretraining method, we performed the following 3D classification experiment in ModelNet. From Table.5, we observed that the model has 13.83% robust accuracy advantage when we set ($\alpha = 1, \beta = 0$) compare with the baseline model ($\alpha = 0, \beta = 0$). From that result, we can find $\mathcal{KLD}(h_1, h_{adv})$ part plays a important role in improving robustness to our pretrained model. When we use setting ($\alpha = 1, \beta = 1$), we can see $\mathcal{KLD}(h_{adv}, h_{hd})$ part in our loss function further increased robustness, which proves the contribution of High-difference points.

Table 5: **Ablation study of PointACL**

| Pretraining method | Standard Accuracy(%) | Robust Accuracy(%) |
|-------------------------------------|----------------------|--------------------|
| PointACL($\alpha = 0, \beta = 0$) | 86.79 | 11.94 |
| PointACL($\alpha = 1, \beta = 0$) | 82.02 | 25.77 |
| PointACL($\alpha = 1, \beta = 1$) | 80.71 | 27.51 |

6. Conclusion

In this paper, we have studied methods to make contrastive learning pretrained model more robust in the 3D point clouds domain. We have showed that using virtual adversarial loss to generate adversarial samples are beneficial towards robustness. We have further showed that using difference of normal (DoN) operator to select high difference points as additional input view can enhance the robustness. Our proposed approaches can achieve state-of-the-art robust accuracy using standard linear finetuning in two downstream tasks: 3D object classification and 3D segmentic segmentation. Extensive experiments involving cross-datasets and attacking strength have also been made to demonstrate universality of our method in improving robustness.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Representation learning and adversarial generation of 3d point clouds. 2018.
- [2] Iro Armeni, Ozan Sener, Amir Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. pages 1534–1543, 06 2016.
- [3] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *arXiv preprint arXiv:1802.00420*, 2018.
- [4] Akhilan Boopathy, Sijia Liu, Gaoyuan Zhang, Cynthia Liu, Pin-Yu Chen, Shiyu Chang, and Luca Daniel. Proper network interpretability helps adversarial robustness in classification. In *ICML*, 2020.
- [5] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks, 2016.
- [6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015.
- [7] R. Qi Charles, Hao Su, Mo Kaichun, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 77–85, 2017.
- [8] Jiefeng Chen, Xi Wu, Vaibhav Rastogi, Yingyu Liang, and Somesh Jha. Robust attribution regularization. In *Advances in Neural Information Processing Systems*, pages 14300–14310, 2019.
- [9] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations, 2020.
- [10] Tianlong Chen, Sijia Liu, Shiyu Chang, Yu Cheng, Lisa Amini, and Zhangyang Wang. Adversarial robustness: From self-supervised pre-training to fine-tuning, 2020.
- [11] Lijie Fan, Sijia Liu, Pin-Yu Chen, Gaoyuan Zhang, and Chuang Gan. When does contrastive learning preserve adversarial robustness from pretraining to finetuning?, 2021.
- [12] Matheus Gadelha, Rui Wang, and Subhransu Maji. Multiresolution tree networks for 3d point cloud processing. In *ECCV*, 2018.
- [13] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
- [14] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised learning, 2020.
- [15] Siyuan Huang, Yichen Xie, Song-Chun Zhu, and Yixin Zhu. Spatio-temporal self-supervised representation learning for 3d point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6535–6545, 2021.
- [16] Y. Ioannou, B. Taati, R. Harrap, and M. Greenspan. Difference of normals as a multi-scale operator in unorganized point clouds. *IEEE*, oct 2012.
- [17] Ziyu Jiang, Tianlong Chen, Ting Chen, and Zhangyang Wang. Robust pre-training by adversarial contrastive learning, 2020.
- [18] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschiot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *arXiv preprint arXiv:2004.11362*, 2020.
- [19] Minseon Kim, Jihoon Tack, and Sung Ju Hwang. Adversarial self-supervised contrastive learning. In *Advances in Neural Information Processing Systems*, 2020.
- [20] Alexey Kurakin, Ian J. Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *CoRR*, abs/1607.02533, 2016.
- [21] Daniel Liu, Ronald Yu, and Hao Su. Extending adversarial attacks and defenses to deep 3d point cloud classifiers, 2019.
- [22] Aleksander Madry, Aleksandar Makelev, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks, 2017.
- [23] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.
- [24] Hanieh Naderi, Arian Etemadi, Kimia Noorbakhsh, and Shohreh Kasaei. Lpf-defense: 3d adversarial defense based on frequency analysis, 2022.
- [25] Omid Poursaeed, Tianxing Jiang, Han Qiao, Nayun Xu, and Vladimir G. Kim. Self-supervised learning of point clouds via orientation estimation, 2020.
- [26] Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! In *Advances in Neural Information Processing Systems*, pages 3353–3364, 2019.
- [27] Jiachen Sun, Yulong Cao, Christopher B Choy, Zhiding Yu, Anima Anandkumar, Zhuoqing Morley Mao, and Chaowei Xiao. Adversarially robust 3d point cloud recognition using self-supervisions. *Advances in Neural Information Processing Systems*, 34, 2021.
- [28] Haohan Wang, Xindi Wu, Zeyi Huang, and Eric P Xing. High-frequency component helps explain the generalization of convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8684–8694, 2020.
- [29] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph cnn for learning on point clouds, 2018.
- [30] Yuxin Wen, Jiehong Lin, Ke Chen, C. L. Philip Chen, and Kui Jia. Geometry-aware generation of adversarial point clouds, 2019.

- [31] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *NeurIPS*, 2016.
- [32] Ziyi Wu, Yueqi Duan, He Wang, Qingnan Fan, and Leonidas J. Guibas. If-defense: 3d adversarial point cloud defense via implicit function based restoration, 2020.
- [33] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes, 2014.
- [34] Chong Xiang, Charles R Qi, and Bo Li. Generating 3d adversarial point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9136–9144, 2019.
- [35] Saining Xie, Jiatao Gu, Demi Guo, Charles R. Qi, Leonidas J. Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding, 2020.
- [36] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *CVPR*, 2018.
- [37] Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. Pu-net: Point cloud upsampling network, 2018.
- [38] Tianhang Zheng, Changyou Chen, Junsong Yuan, Bo Li, and Kui Ren. Pointcloud saliency maps, 2018.
- [39] Hang Zhou, Dongdong Chen, Jing Liao, Kejiang Chen, Xiaoyi Dong, Kunlin Liu, Weiming Zhang, Gang Hua, and Nenghai Yu. Lg-gan: Label guided adversarial network for flexible targeted attack of point cloud based deep networks. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [40] Hang Zhou, Kejiang Chen, Weiming Zhang, Han Fang, Wenbo Zhou, and Nenghai Yu. Dup-net: Denoiser and upsampler network for 3d adversarial point clouds defense, 2018.