

Self-Supervised Global-Local Structure Modeling for Point Cloud Domain Adaptation with Reliable Voted Pseudo Labels

Anonymous CVPR submission

Paper ID 308

Abstract

In this paper, we propose an unsupervised domain adaptation method for deep point cloud representation learning. To model the internal structures in target point clouds, we first propose to learn the global representations of unlabeled data by scaling up or down point clouds and then predicting the scales. Second, to capture the local structure in a self-supervised manner, we propose to project a 3D local area onto a 2D plane and then learn to reconstruct the squeezed area. Moreover, to effectively transfer the knowledge from source domains, we propose to vote pseudo labels for target samples based on their nearest source neighbors in the feature space. To avoid the noise of incorrect pseudo labels, we only select reliable target samples, whose voting consistencies are high enough, for enhancing adaptation. The voting method is able to adaptively select more and more target samples during training, which in return facilitates adaptation because the amount of labeled target data increases. Experiments on PointDA (ModelNet-10, ShapeNet-10 and ScanNet-10) and Sim-to-Real (ModelNet-11, ScanObjectNN-11, ShapeNet-9 and ScanObjectNN-9) demonstrate the effectiveness of our method.

1. Introduction

Large-scale learning methods based on deep neural networks [7, 8, 20, 29, 30] constitute the recent advances in 3D vision, and play an import role for visual perception in intelligent platforms such as robots, drones and self-driving cars. These intelligent platforms usually employ real-time depth sensors, such as LiDAR, to capture the accurate geometric information of scenes, which are represented by 3D point clouds. However, deep neural networks usually requires massive amounts of labeled point clouds for representation learning, which limits the scalability to the real world. To alleviate this problem, unsupervised point cloud domain adaption is recently attracting increasing attention from the community. Domain adaptation aims to transfer

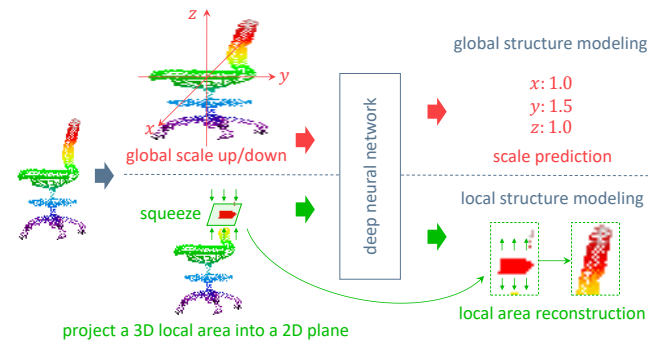


Figure 1. Illustration of the self-supervised global-local structure modeling for point cloud domain adaptation. The global structure is modeled by scaling up/down the point cloud and then predicting the scale. To capture the local structure, a random 3D local area is squeezed onto a 2D plane and then reconstructed by networks.

the knowledge from a labeled source domain to a related but unlabeled target domain, in which the source and target domains share the same feature space. However, due to different point scales, object sizes, densities, styles, sensor perspectives, etc., point cloud representations in the target domain usually inevitably deviate from the corresponding representations in the source domain, resulting in the domain shift or distributional shift problem.

To reduce domain shift, one solution is to directly learn from the target domain via self-supervised learning, *i.e.*, exploiting the relations or correlations between different input signals. However, most of the existing methods focus on leveraging only one of the global and local structure of unlabeled data, such as predicting global vertical rotation [19] and reconstructing point clouds from randomly rearranged object parts [24]. Moreover, the vertical rotation method ignore the horizontal structure (horizontal rotation is used for data augmentation) and the rearrangement method may fail when an object consists of similar parts. Besides, other methods [1, 34] try to indirectly capture structure by leveraging the relationships of two point clouds. However, because two objects are merged into

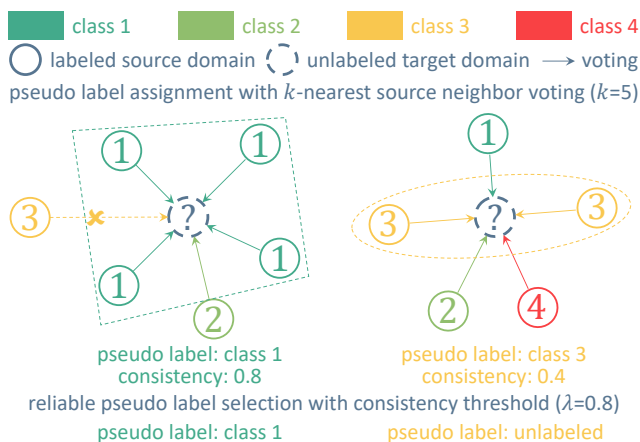


Figure 2. Illustration of reliable voted pseudo label generation. First, target point clouds’ pseudo labels are voted by a few of nearest source neighbors in the feature space. Then, the target samples whose nearest source neighbor labels are consistent enough are selected as reliable training data.

a single point cloud, they may neglect their own structure of each individual object. These methods also try to capture the local structure via complicated deformation-reconstruction/localization methods, *i.e.*, reconstructing a collapsed region via Chamfer Distance loss [1] or localizing a curvature-changed area [34]. Another domain adaptation solution is to transfer the knowledge from the source domain to the target domain via adversarial training [22] and pseudo labels [34]. However, the existing pseudo label method, which directly exploits the prediction on the target data, may be not reliable and add noise into training.

This paper is devoted to exploiting both the self-supervised and transfer learning for point cloud domain adaptation. First, as shown in Fig. 1, we propose to learn point cloud representations by scaling up or down point coordinates of one dimension and then predicting the scale based on the other two unchanged dimensions. In this way, the network is able to capture the global structure in a self-supervised manner. To model the local structure, we propose to project a 3D local area onto a 2D plane by simply setting point coordinates of a random selected dimension to the same value and recover the squeezed area via mean squared error loss. Because our self-supervised methods only need to change one dimension’s coordinates, they are simple compared to existing works [19, 22, 24, 34]. Second, as shown in Fig. 2, to enhance the knowledge transfer from the source domain, we propose a voting method to assign reliable pseudo labels to target samples. Specifically, pseudo labels are voted based on a few nearest source neighbors in the shared feature space. Then, only the target point clouds whose nearest source neighbor labels are consistent enough are selected as reliable training data. More-

over, with networks becoming stronger during training, our reliable voting method is able to adaptively select more target data, which in return facilitates learning because the labeled target data increases.

To evaluate our method, we conduct experiments on the widely-used 3D domain adaptation benchmark PointDA [22], which consists of 10 shared classes from ModelNet40 [31], ShapeNet [2] and ScanNet [3]. Moreover, we also conduct experiments on a Sim-to-Real dataset [13], which consists of 11 shared classes from ModelNet40 and ScanObjectNN [28], and 9 shared classes from ShapeNet and ScanObjectNN, respectively. The Sim-to-Real dataset is created for point-cloud-based meta-learning methods. In this paper, we also employ it to evaluate point cloud domain adaption methods. The contributions of this paper are threefold:

- To model the internal structure of unlabeled target point clouds, we propose a global scaling-up-down prediction method and a local 3D-2D-3D projection-reconstruction method for point cloud domain adaptation.
- To transfer the knowledge from source domains, we propose a voting method to assign reliable pseudo labels to target samples. Without any supervision or control, our method is able to autonomously select more and more reliable target data during training, which in return facilitates learning.
- Extensive experiments on two datasets show that the proposed method effectively improves the accuracy of unsupervised domain adaptation on point clouds.

2. Related Work

Point Cloud Classification. Point-cloud-based object classification is one of the fundamental tasks for point cloud processing. Recently, a number of deep neural networks have been proposed to address this problem [20, 21, 26, 29, 30]. Most of these works aim to directly manipulate point clouds without converting irregular points into regular voxel grids, to avoid the quantization errors and high computational cost of voxelization. Since a point cloud is essentially a set of unordered points and invariant to permutations of its points, the key to point cloud processing is to design effective point-based spatial modeling operations that do not rely on point orderings. Our method is independent of these works and employ them as feature generators to encode point clouds.

Self-supervised Learning. To learn from internal structures in images, self-supervised learning tries to find or exploit the relations or correlations between different input signals [5, 6, 11, 16–18, 32], *e.g.*, modifying the input and (1) predict what changed or (2) ensure that the output representation does not change. These methods also inspire

the self-supervised learning on point clouds. For example, Sauder and Sievers [24] proposed to split an input point cloud into several parts with a random permutation and reconstruct the point cloud. Poursaeed *et al.* [19] proposed to rotate a point cloud and then predict the rotation angle. Achituve *et al.* [1] and Zou *et al.* [34] proposed to first mix up two point clouds and then predict the mixed labels and their angles, respectively. Besides, they also proposed to reconstruct and localize deformed local areas, respectively. When modeling local structure, our method is similar to the two deformation-based methods but simpler than them.

Domain Adaptation. Due to different illuminations, backgrounds, object styles, *etc.*, domain adaptation has been well developed on images [9, 15, 23, 25, 27, 33]. These methods can be divided into three categories. 1) Adversarial training [15, 23, 27], which aims to directly learn unbiased representations via a discriminator to judge whether the learned features are from the target domain or the source domain, and a feature generator to confuse the discriminator. 2) Style transfer [33], which employs Generative Adversarial Networks [12] to transfer source images as the target style for training. 3) Pseudo labels [4, 9]. Networks trained on source domains have a certain ability to recognize target images. Therefore, target data’s pseudo labels can be generated for training, in which self-paced learning [14] is usually employed to reduce noisy pseudo labels. The image-based methods can also be used for point-cloud-based domain adaptation. For example, Qin *et al.* [22] employed adversarial training to learn unbiased point cloud representations. Zou *et al.* [34] employed a pseudo-label-based method, equipped with self-paced learning, for point cloud domain adaptation. Our work is also based on adversarial training and pseudo label but we propose a new reliable-vote-based method for pseudo label generation.

3. Proposed Method

In this section, we first introduce and formulate the point cloud domain adaption problem in Sec. 3.1. Second, we introduce the scaling-up-down method for self-supervised global structure modeling in Sec. 3.2. Third, the 3D-2D-3D project-reconstruction method for local structure learning is described in Sec. 3.3. Fourth, we briefly introduce the adversarial training used in our paper. Fifth, the reliable voted pseudo label method is described in details in Sec. 3.5. Finally, the overall training procedure of our method is shown in Sec. 3.6.

3.1. Problem Formulation

The goal of unsupervised domain adaptation on point clouds is to transfer the knowledge for a labeled source domain $\mathcal{S} = \{(\mathbf{P}_i^s, y_i^s)\}_{i=1}^{n^s}$ to an unlabeled target domain $\mathcal{T} = \{(\mathbf{P}_i^t)\}_{i=1}^{n^t}$, where $\mathbf{P} \in \mathbb{R}^{m \times 3}$ and $y_i^s \in \mathcal{Y} =$

$\{1, \dots, c\}$ and m is the number of points and c is the number of shared classes. The n^s and n^t denote the number of source and target point clouds, respectively. The key to domain adaptation is to learn a mapping function or point cloud feature generator Φ that projects point clouds from different domains into a shared feature space. The feature generator Φ can be implemented by existing deep neural network, *e.g.*, PointNet [20] and DGCNN [29], which encode a point cloud to a vector, *i.e.*, $\mathbf{f} = \Phi(\mathbf{P})$.

In this paper, we suppose that point features, *e.g.*, color, norm or other information, are not available, which means the input is only about geometry. In this case, domain shifts can be caused by different point scales, point densities, object sizes, sensor perspectives, object styles, *etc.* Some of the shifts can be reduced by low-level data preprocessing and augmentation. For example, the point scale and object size problems can be addressed by normalizing object coordinates to a fixed range, *e.g.*, $[-1, 1]$. The density problem can be addressed by sampling, *e.g.*, Farthest Point Sampling (FPS), point clouds to the same number of points. The perspective shift problem can be addressed by rotation-based data augmentation. However, the other shifts, *e.g.*, object styles, have to be reduced via high-level representations. This paper focuses on adapting high-level representations.

3.2. Self-Supervised Global (G) Structure Modeling via Scaling-Up-Down Prediction

To enable Φ to capture the global structure without human-annotated class labels, we propose to scale up or down coordinates and then employ a regression Ω to predict the scale based on the point cloud feature \mathbf{f} . Specifically, suppose $\mathbf{s}_i = (s_i^x, s_i^y, s_i^z) \in \mathbb{R}^{+1 \times 3}$ is a random scale vector for the i -th target point cloud. Then, the coordinates \mathbf{P}_i^t of the point cloud is scaled by $\mathbf{s}_i \odot \mathbf{P}_i^t$, where \odot is element-wise multiplication and \mathbf{s}_i is broadcasted for the multiplication. Finally, regression Ω is used to predict the scale \mathbf{s}_i ,

$$\min_{\Phi, \Omega} \mathcal{L}_g, \quad \mathcal{L}_g = -\frac{1}{n^t} \sum_{i=1}^{n^t} \|\Omega(\Phi(\mathbf{s}_i \odot \mathbf{P}_i^t)) - \mathbf{s}_i\|_2^2. \quad (1)$$

Note that, to address the object size shift problem, we normalize object coordinates in data preprocessing. Therefore, when predicting the scale, regression Ω is actually based on the relative scales of the three dimensions, instead of the absolute change. For example, when we scale up two dimensions by 2 and leave the third dimension unchanged, regression Ω would misunderstand that the third dimension is scaled down by 0.5. To avoid this issue, we fix two of three dimensions and only scale up or down the other dimension. Therefore, \mathbf{s}_i is limited to $\{(s, 1, 1), (1, s, 1), (1, 1, s)\}$, where $s \in \mathbb{R}^+$. Moreover, because dramatically scaling up or down point clouds will

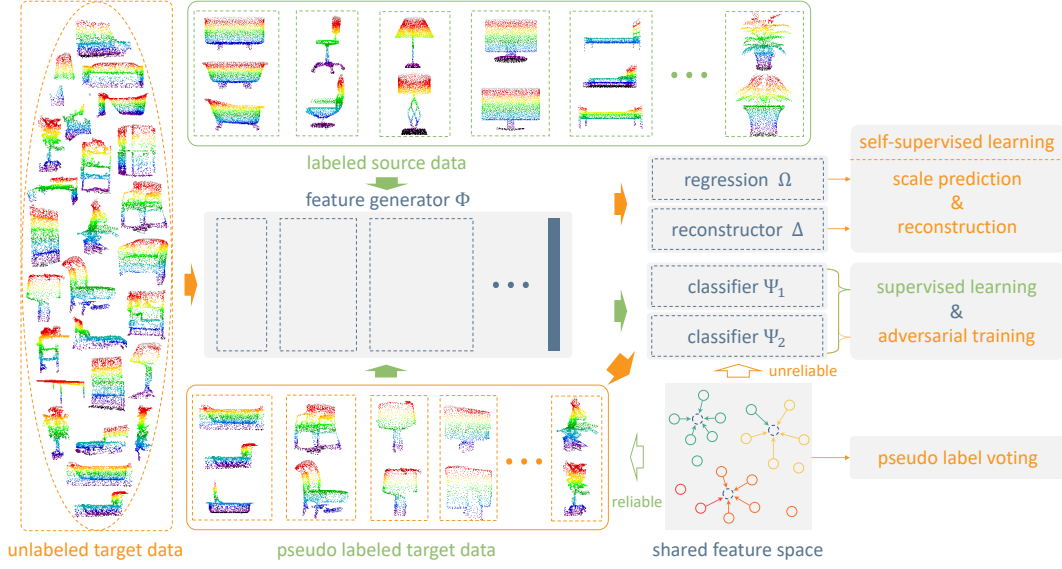


Figure 3. Illustration of the proposed self-supervised Global-Local structure modeling and reliable Voted pseudo label method (GLV) for point cloud domain adaption. The framework consists of a feature generator Φ to encode point clouds, a regression Ω to predict the scaling change for self-supervised global structure modeling, a reconstructor Δ to recover squeezed local area for self-supervised local structure modeling and two classifiers Ψ_1 and Ψ_2 for supervised learning and adversarial training. Besides, a reliable voting method is employed to obtain accurate target pseudo labels for enhancing source knowledge transfer.

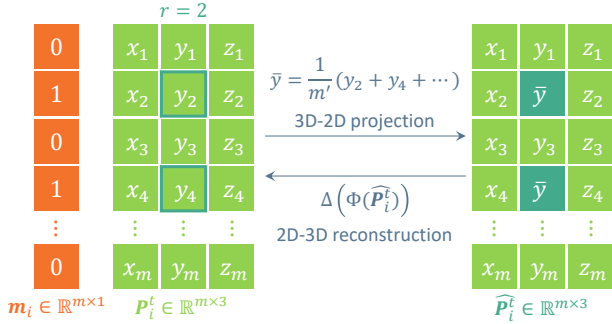


Figure 4. Illustration of the 3D-2D-3D projection-reconstruction process for the i -th target point cloud P_i^t . The points of a selected 3D local area is indicated by a mask vector m_i . Then, the local area is projected into 2D by squeezing the r -dimensional coordinates to their mean, resulting in the locally projected point cloud \hat{P}_i^t . Finally, the feature generator Φ and reconstructor Δ is asked to reconstruct the squeezed area based on \hat{P}_i^t .

change their structures, s is sampled from a small range, *i.e.*, $[0.5, 1.5]$, in this paper.

3.3. Self-Supervised Local (L) Structure Modeling via 3D-2D-3D Projection-Reconstruction

To enable Φ to learn the local structure without point-level human annotations, we propose a 3D-2D-3D projection-reconstruction method. Specifically, similar to [1], we first split the normalized 3D space into several

regions. In this way, a point cloud is divided into multiple parts. Then, we randomly select a part that contains enough points for projection-reconstruction. Suppose $m_i \in \{0, 1\}^{m \times 1}$ is a mask vector to indicate the selected points for the i -th target point cloud. Recall that m is the number of points in the cloud. When $m_i[j] = 1$, it indicates the j -point is selected. In this way, the selected local area can be denoted as $P_i^t[m_i] \in \mathbb{R}^{m'_i \times 3}$, where m'_i is the number of selected points, *i.e.*, $m'_i = \|m_i\|_1$. Third, we project the selected 3D area $P_i^t[m_i]$ onto a 2D panel. To do so, we randomly select a dimension r . Then, the r -dimensional coordinates of selected points $P_i^t[m_i, r] \in \mathbb{R}^{m'_i \times 1}$ are squeezed to their mean, *i.e.*, $P_i^t[m_i, r] \leftarrow \frac{1}{m'_i} \sum_{j=1}^{m'_i} P_i^t[m_i, r][j]$. In this way, we obtain the locally projected point cloud \hat{P}_i^t . Finally, a reconstructor Δ is employed to recover the projected area,

$$\min_{\Phi, \Delta} \mathcal{L}_l, \quad \mathcal{L}_l = -\frac{1}{n^t} \sum_{i=1}^{n^t} \|m_i \odot \Delta(\Phi(\hat{P}_i^t)) - m_i \odot P_i^t\|_F^2, \quad (2)$$

where m_i is broadcasted for the element-wise multiplication. We illustrate the 3D-2D-3D projection-reconstruction process in Fig. 4.

3.4. Adversarial Training for Unbiased Representation Learning

Like most domain adaptation works, we also employ adversarial training [12] to reduce domain shifts and learn unbiased representations. In this paper, we employ Maximum

Classifier Discrepancy (MCD) [23] for adversarial training. Specifically, MCD uses two classifiers Ψ_1 and Ψ_2 , which map the feature vector \mathbf{f} to two probability vectors of length c , i.e., $\Psi_1(\mathbf{f}) \in \mathbb{R}^c$ and $\Psi_2(\mathbf{f}) \in \mathbb{R}^c$, respectively. Recall that c is the number of classes. For labeled data, MCD performs supervised-learning-based classification,

$$\begin{aligned} & \min_{\Phi, \Psi_1, \Psi_2} \mathcal{L}_s, \\ \mathcal{L}_s = & -\frac{1}{n^s} \sum_{i=1}^{n^s} \sum_{j=1}^c \mathbb{1}_{[j=y_i^s]} \cdot \log(\Psi_1(\Phi(\mathbf{P}_i^s))[j]) \\ & -\frac{1}{n^s} \sum_{i=1}^{n^s} \sum_{j=1}^c \mathbb{1}_{[j=y_i^s]} \cdot \log(\Psi_2(\Phi(\mathbf{P}_i^s))[j]). \end{aligned} \quad (3)$$

For unlabeled data, MCD first tries to maximize the prediction discrepancy of the two classifiers with fixed Φ ,

$$\begin{aligned} & \min_{\Psi_1, \Psi_2} \mathcal{L}_s - \mathcal{L}_{adv}, \\ \mathcal{L}_{adv} = & \frac{1}{n^t c} \sum_{i=1}^{n^t} \|\Psi_1(\Phi(\mathbf{P}_i^t)) - \Psi_2(\Phi(\mathbf{P}_i^t))\|_1. \end{aligned} \quad (4)$$

Then, the generator Φ is trained to minimize the discrepancy with fixed classifiers.

$$\min_{\Phi} \mathcal{L}_{adv}. \quad (5)$$

In this way, Φ is enforced to learn unbiased representations.

3.5. Reliable Voted (V) Target Pseudo Label Generation for Enhancing Domain Adaptation

Although adversarial training provides a way to reduce domain shifts, its effectiveness is usually limited due to the complicated training process and indirect adaptation strategy. In this paper, we propose a pseudo-label-based method to directly transfer the knowledge from source to target domain. Specifically, our method employs a voting strategy to assign pseudo labels to target samples. Specifically, the pseudo labels of target point clouds are voted based on the labels of a few of their nearest source neighbors in the shared feature space. Suppose \mathbf{f}_i^t and \mathbf{f}_j^s are the features of the i -th target and the j -th source point clouds, respectively. Their similarity is calculated as

$$e_{ij}^{st} = \frac{\mathbf{f}_i^t \cdot \mathbf{f}_j^s}{\|\mathbf{f}_i^t\|_2 \times \|\mathbf{f}_j^s\|_2}. \quad (6)$$

Then, the k -nearest source neighbors are selected as follows,

$$\mathcal{N}(\mathbf{P}_i^t, k) = \{j \mid e_{ij}^{st} \in \text{top-}k(\{e_{i1}^{st}, \dots, e_{in^s}^{st}\})\}. \quad (7)$$

Third, the pseudo label of the i -th target point cloud is assigned with a voting mechanism,

$$\tilde{y}_i^t = \text{vote}(\{y_j^s \mid j \in \mathcal{N}(\mathbf{P}_i^t, k)\}), \quad (8)$$

where the vote function simply selects the majority as the output.

However, although we employ a k -NN based voting method, the pseudo labels can be still unreliable, which may add noise into training data and decrease performance. To address this problem, we propose to only exploit reliable target point clouds, of which nearest source neighbor labels are consistent enough, to train the model with their voted pseudo labels. In this paper, we use the consistency of source neighbors' labels to measure the reliability,

$$v_i^t = \begin{cases} 0, & \frac{\sum_{j \in \mathcal{N}(\mathbf{P}_i^t, k)} \mathbb{1}_{[y_j^s = \tilde{y}_i^t]}}{k} < \lambda, \\ 1, & \frac{\sum_{j \in \mathcal{N}(\mathbf{P}_i^t, k)} \mathbb{1}_{[y_j^s = \tilde{y}_i^t]}}{k} \geq \lambda, \end{cases} \quad (9)$$

where $\lambda \in (0, 1]$ is the consistency threshold and v_i^t indicates where the i -th target point cloud is selected as a reliable training sample. Finally, the selected reliable target data is used to train the feature generator in a supervised manner,

$$\begin{aligned} & \min_{\Phi, \Psi_1, \Psi_2} \mathcal{L}_t, \\ \mathcal{L}_t = & -\frac{1}{\|\mathbf{v}^t\|_1} \sum_{i=1}^{n^t} v_i^t \sum_{j=1}^c \mathbb{1}_{[j=\tilde{y}_i^t]} \cdot \log(\Psi_1(\Phi(\mathbf{P}_i^t))[j]) \\ & -\frac{1}{\|\mathbf{v}^t\|_1} \sum_{i=1}^{n^t} v_i^t \sum_{j=1}^c \mathbb{1}_{[j=\tilde{y}_i^t]} \cdot \log(\Psi_2(\Phi(\mathbf{P}_i^t))[j]), \end{aligned} \quad (10)$$

where $\mathbf{v}^t = (v_1^t, \dots, v_{n^t}^t)$ in the selection indicator vector.

Note that, k and λ are the only two hyper-parameters of our reliable voting method. Moreover, as shown in the following experiments, we can further simplify this method by fixing λ to 1. In this case, k is the only one hyper-parameter of this method and a bigger k indicates a higher reliability threshold.

Although with fixed k and λ , our reliable voting method has the ability to automatically and adaptively select more and more target data during training. At the early stage of training, because the feature generator Φ is weak and the domain shift is large, only a few target point clouds, which are similar to source samples and easy to be recognized, reach the consistency threshold and are selected as reliable training data. When Φ becomes stronger and domain shift reduces, target representations become more discriminative and the consistency of source neighbors' labels increases. Consequently, more target point clouds are added into the training set. In return, the increasing of labeled target data facilitates training. In this way, the feature generator Φ is improved progressively.

3.6. Overall Training

In summary, our approach includes two self-supervised learning methods, i.e., scaling-up-down prediction and 3D-

ALGORITHM 1: GLV Training Procedure

Input : labeled source dataset $\mathcal{S} = \{(\mathbf{P}_i^s, y_i^s)\}_{i=1}^{n^s}$,
unlabeled target dataset $\mathcal{T} = \{(\mathbf{P}_i^t)\}_{i=1}^{n^t}$,
number of source neighbors for voting k ,
reliability or consistency threshold λ ,
feature generator Φ , classifiers Ψ_1 and Ψ_2 ,
regression Ω , reconstructor Δ ,
number of training rounds R ,
number of epochs E for each round.
Output: Φ , Ψ_1 , Ψ_2 , Ω and Δ .
Initialization: randomly initialize Φ , Ψ_1 , Ψ_2 , Ω and Δ ;
randomly initialize pseudo labels $\{\tilde{y}_i^t\}_{i=1}^{n^t}$;
zero-initialize selection indicator $\mathbf{v}^t = \mathbf{0}$.

for 1 to R do
 for 1 to E do
 for $(\mathbf{P}_i^s, y_i^s), (\mathbf{P}_i^t, \tilde{y}_i^t)$ in $(\mathcal{S}, \mathcal{T})$ do
 $\min_{\Phi, \Psi_1, \Psi_2} \mathcal{L}_s$ with (\mathbf{P}_i^s, y_i^s) ;
 $\min_{\Phi, \Omega} \mathcal{L}_g$ with \mathbf{X}_i^t ;
 $\min_{\Phi, \Delta} \mathcal{L}_l$ with \mathbf{X}_i^t ;
 if $v_i = 1$ then
 $\min_{\Phi, \Psi_1, \Psi_2} \mathcal{L}_t$ with $(\mathbf{P}_i^t, \tilde{y}_i^t)$;
 else
 $\min_{\Psi_1, \Psi_2} \mathcal{L}_s - \mathcal{L}_{adv}$ with (\mathbf{P}_i^s, y_i^s) and \mathbf{P}_i^t ;
 $\min_{\Psi_1, \Psi_2} \mathcal{L}_{adv}$ with \mathbf{P}_i^t ;
 end
 end
 end
 update pseudo labels $\{\tilde{y}_i^t\}_{i=1}^{n^t}$ and selection indicator \mathbf{v}^t based on Φ , k and λ ;
end

2D-3D projection-reconstruction, and two transfer learning methods, *i.e.*, adversarial training and our reliable voted pseudo label method. The framework of our approach is illustrated in Fig. 3. The overall training process is shown in Alg. 1. The training contains multiple rounds. After each round, we perform reliable voted pseudo label generation. Each round contains several epochs. In each epoch, we perform supervised learning, scaling-up-down prediction, 3D-2D-3D projection-reconstruction and adversarial training.

4. Experiments

4.1. Datasets

PointDA. The PointDA [22] dataset is a widely-used benchmark for point cloud domain adaptation evaluation, which extracts the samples in 10 shared classes from ModelNet40 [31], ShapeNet [2] and ScanNet [3], respectively. In this way, PointDA consists of three subsets: ModelNet-10 (M10), ShapeNet-10 (S10) and ScanNet-10 (S*10). Given the three subsets, we can conduct six types of adaptation

Self-Supervised		Transfer			Accuracy
Scale	3D-2D-3D	Adversarial	Pseudo Label Vote	Reliable	
\times	\times	\times	\times	\times	64.2
\checkmark					69.8
	\checkmark				67.9
\checkmark	\checkmark				71.2
		\checkmark			67.7
			\checkmark		64.8
			\checkmark	\checkmark	74.8
		\checkmark	\checkmark	\checkmark	75.4
\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	77.2

Table 1. Ablation study on each component of our method. Experiments are conducted on PointDA with the S*10 \rightarrow S10 scenario. When none of the components is employed, the model is directly transferred to the target domain without adaptation.

Method	M \rightarrow S	M \rightarrow S*	S \rightarrow M	S \rightarrow S*	S* \rightarrow M	S* \rightarrow S
Rotation [19]	82.8	41.7	76.0	49.0	69.3	68.7
Scale (ours)	84.0	46.0	79.4	48.3	71.0	69.8

Table 2. Comparison of vertical rotation and our scaling-up-down method for self-supervised global structure modeling on PointDA.

Adaptation	Def-Rec [1]	Def-Loc [34]	3D-2D-3D (ours)
M10 \rightarrow S10	82.7	84.3	83.5
M10 \rightarrow S*10	43.9	46.2	48.1
S10 \rightarrow M10	79.8	69.8	75.7
S10 \rightarrow S*10	48.0	49.2	48.7
S*10 \rightarrow M10	66.0	66.6	68.2
S*10 \rightarrow S10	67.4	66.1	67.9

Table 3. Comparison of deformation-reconstruction (Def-Rec), deformation-localization (Def-Loc) and our 3D-2D-3D method for self-supervised local structure modeling on PointDA.

scenarios: M10 \rightarrow S10, M10 \rightarrow S10, S10 \rightarrow M10, S10 \rightarrow S10, S*10 \rightarrow M10 and S*10 \rightarrow S10.

Sim-to-Real. The Sim-to-Real [13] dataset is a fairly new benchmark, which consists of 11 shared classes from ModelNet40 and ScanObjectNN [28], and 9 shared classes from ShapeNet and ScanObjectNN, respectively. The dataset is built to evaluate meta-learning on point clouds. In this paper, we also employ the dataset for evaluating point cloud domain adaptation. The dataset consists of four subsets: ModelNet-11 (M11), ScanObjectNN-11 (S*O11), ShapeNet-9 (S9) and ScanObjectNN-9 (S*O9). Different from PointDA, Sim-to-Real asks models to transfer knowledge from simulated ModelNet or ShapeNet to real-world ScanObjectNN where the objects are even with back-

Method	M10 \rightarrow S10	M10 \rightarrow S*10	S10 \rightarrow M10	S10 \rightarrow S*10	S*10 \rightarrow M10	S*10 \rightarrow S10
w/o Adaptation	83.3 \pm 0.7	43.8 \pm 2.3	75.5 \pm 1.8	42.5 \pm 1.4	63.8 \pm 3.9	64.2 \pm 0.8
DANN [10]	74.8 \pm 2.8	42.1 \pm 0.6	57.5 \pm 0.4	50.9 \pm 1.0	43.7 \pm 2.9	71.6 \pm 1.0
PointDAN [22]	83.9 \pm 0.3	44.8 \pm 1.4	63.3 \pm 1.1	45.7 \pm 0.7	43.6 \pm 2.0	56.4 \pm 1.5
RS [24]	79.9 \pm 0.8	46.7 \pm 4.8	75.2 \pm 2.0	51.4 \pm 3.9	71.8 \pm 2.3	71.2 \pm 2.8
DefRec + PCM [1]	81.7 \pm 0.6	51.8 \pm 0.3	78.6 \pm 0.7	54.5 \pm 0.3	73.7 \pm 1.6	71.1 \pm 1.4
GAST [34]	84.8 \pm 0.1	59.8 \pm 0.2	80.8 \pm 0.6	56.7 \pm 0.2	81.1 \pm 0.8	74.9 \pm 0.5
GLV (ours)	85.4 \pm 0.4	61.0 \pm 0.4	78.8 \pm 0.6	57.7 \pm 0.4	76.0 \pm 1.3	77.2 \pm 0.6

Table 4. Accuracy on the PointDA dataset. Our GLV method achieves four best accuracies on the six adaptation scenarios.

Method	M11 \rightarrow S*O11		S9 \rightarrow S*O9	
	Object	Object & Background	Object	Object & Background
w/o Adaptation	61.68 \pm 1.26	57.61 \pm 0.44	57.42 \pm 1.01	54.42 \pm 0.80
PointDAN [22]	63.32 \pm 0.85	55.13 \pm 0.97	54.95 \pm 0.87	43.00 \pm 0.95
MetaSets [13]	72.42 \pm 0.21	65.66 \pm 1.06	60.92 \pm 0.76	59.08 \pm 1.01
GLV (ours)	75.16 \pm 0.34	66.74 \pm 0.59	62.46 \pm 0.55	59.76 \pm 0.84

Table 5. Accuracy on the Sim-to-Real dataset. Our GLV method achieves the best accuracies on the four adaptation scenarios.

grounds. Therefore, there are two types of adaptation scenarios Sim-to-Real: M11 \rightarrow S*O11 and S9 \rightarrow S*O9.

4.2. Implementation

Following existing works [1, 13, 34], we employ DGCNN [29] as feature generator. For PointDA, we use the setting of [1, 34]. For Sim-to-Real, we follow the setting of [13]. The training constrains 20 rounds, with 10 epochs in each round. Batch size is set to 32 and learning rate is set to 0.001. By default, the k and λ of reliable voting is set to 10 and 1.0, respectively.

4.3. Comparison to the State-of-the-art

For PointDA, we compare our method with the state-of-the-art point-based domain adaptation methods, including Domain Adversarial Neural Network (DANN) [10], Point Domain Adaptation Network (PointDAN) [22], Reconstruction Space Network (RS) [24], Deformation Reconstruction Network with Point Cloud Mixup (DefRec + PCM) [1] and Geometry-Aware Self-Training (GAST) [34]. We report the mean accuracy and standard error of the mean with three seeds in Table. 4. Our GLV method achieves four best accuracies on the six adaptation scenarios.

We also compare our scaling-up-down and 3D-2D-3D approaches with existing global (vertical rotation classification [19]) and local (deformation-reconstruction [1] and deformation-localization [34]) modeling methods, respectively. Results in Table. 2 and Table. 3 show the effectiveness of our method.

For Sim-to-Real, we compared our method with a point cloud domain adaptation method, *i.e.*, PointDAN [22], and

a meta-learning method, *i.e.*, MetaSets [13]. We perform each adaptation scenarios three times and report the average and the standard deviation of the results in Table. 5. Our method outperforms both the domain adaptation and meta-learning methods.

4.4. Ablation Study

A) Influence of scaling-up-down prediction, 3D-2d-3D projection-reconstruction, adversarial training and reliable voted pseudo label.

To investigate the influence of each component in our method, we conduct an ablation study on PointDA with the S*10 \rightarrow S10 scenarios. The results are show in Table. 1. All the four components have important influences on domain adaption. Among them, the reliable voted pseudo label method (Vote + Reliable) is the most effective, which increases the baseline (64.2%) by 10.6%. Note that, without the consistency-based reliable target sample section, the pseudo method almost has no effect. This is because the single voting method fails to obtain accurate pseudo labels and adds a large amount of noise into training data.

B) Gradually increasing target data selection.

To verify the ability of the reliable voting method to autonomously and adaptively selecting more and more target data during training, we show the number of the selected reliable target samples, the accuracy of their pseudo labels and the accuracy on the test dataset in Fig. 5a \sim Fig. 5f. Under the premise of the fixed or slightly increasing accuracy of pseudo labels, more and more target data is selected, which indicates that the number of correctly labeled target point clouds increases. With the correct labeled target data

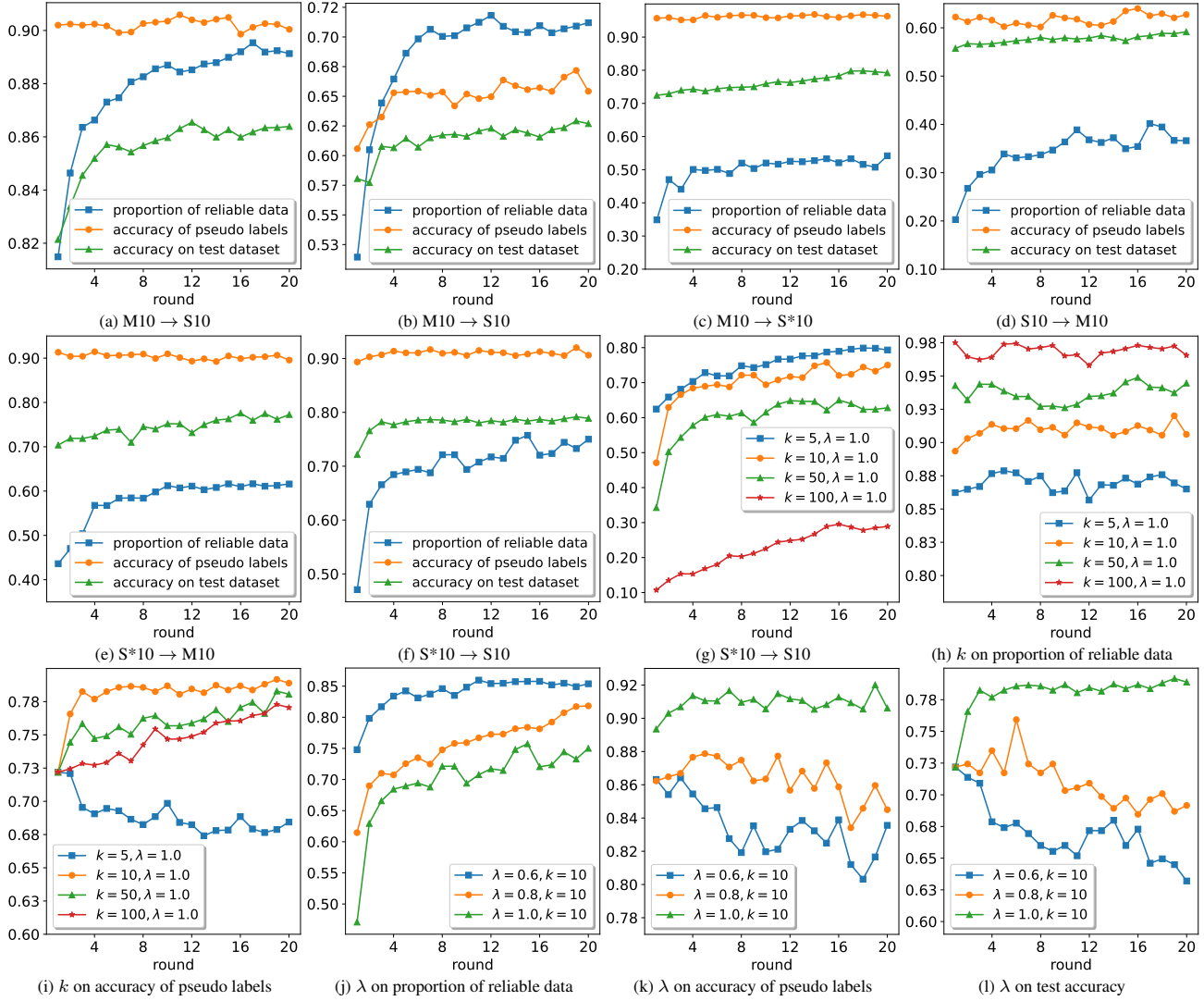


Figure 5. Influence of reliable voted pseudo labels. Experiments are conducted on PointDA. (a)~(f): During training, the proposed method is able to attentively select more and more target data and the test accuracy increases gradually. (g)~(l): Impact of different k and λ on voted pseudo label generation and adaptation performance (S*10 \rightarrow S10).

increasing, the feature generator is gradually improved on the target domain.

C) Impact of k and λ on voted pseudo label generation and adaptation performance.

Our reliable voting method contains two hyper-parameters, *i.e.*, k and λ . To investigate the influence of the two hyper-parameters, we conduct the S*10 \rightarrow S10 adaptation with different k and λ . The results are shown in Fig. 5j ~ Fig. 5l. When k becomes bigger, only a few considerably discriminative target samples satisfy the reliability requirement and are selected. Although the pseudo label accuracy increases, less target data is selected. Consequently, the influence of pseudo labels diminishes and the corresponding improvement drops. When λ becomes smaller, the selected

target data becomes less reliable. Because noise is added, the accuracy decreases.

5. Conclusion

In this paper, we propose two self-supervised learning methods, *i.e.*, scaling-up-down prediction and 3D-2D-3D projection-reconstruction, and one reliable voted pseudo label method for point cloud domain adaptation. Experiments on two datasets demonstrate the effectiveness of our approach. However, when selecting target data, our reliable voting method does not take the class balance problem into consideration. A promising improvement is to integrate class diversity into selection, instead of only reliability.

References

- [1] Idan Achituve, Haggai Maron, and Gal Chechik. Self-supervised learning for domain adaptation on point clouds. In *WACV*, pages 123–133, 2021. 1, 2, 3, 4, 6, 7
- [2] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing 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. *arXiv*, 1512.03012, 2015. 2, 6
- [3] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas A. Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, pages 2432–2443, 2017. 2, 6
- [4] Yuhang Ding, Hehe Fan, Mingliang Xu, and Yi Yang. Adaptive exploration for unsupervised person re-identification. *ACM Trans. Multim. Comput. Commun. Appl.*, 16(1):3:1–3:19, 2020. 3
- [5] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, pages 1422–1430, 2015. 2
- [6] Alexey Dosovitskiy, Philipp Fischer, Jost Tobias Springenberg, Martin A. Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(9):1734–1747, 2016. 2
- [7] Hehe Fan, Yi Yang, and Mohan S. Kankanhalli. Point 4d transformer networks for spatio-temporal modeling in point cloud videos. In *CVPR*, pages 14204–14213, 2021. 1
- [8] Hehe Fan, Xin Yu, Yuhang Ding, Yi Yang, and Mohan S. Kankanhalli. Pstnet: Point spatio-temporal convolution on point cloud sequences. In *ICLR*, 2021. 1
- [9] Hehe Fan, Liang Zheng, Chenggang Yan, and Yi Yang. Unsupervised person re-identification: Clustering and fine-tuning. *ACM Trans. Multim. Comput. Commun. Appl.*, 14(4):83:1–83:18, 2018. 3
- [10] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor S. Lempitsky. Domain-adversarial training of neural networks. *J. Mach. Learn. Res.*, 17:59:1–59:35, 2016. 7
- [11] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018. 2
- [12] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pages 2672–2680, 2014. 3, 4
- [13] Chao Huang, Zhangjie Cao, Yunbo Wang, Jianmin Wang, and Mingsheng Long. Metasets: Meta-learning on point sets for generalizable representations. In *CVPR*, pages 8863–8872, 2021. 2, 6, 7
- [14] M. Pawan Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *NeurIPS*, pages 1189–1197, 2010. 3
- [15] Jingjing Li, Erpeng Chen, Zhengming Ding, Lei Zhu, Ke Lu, and Heng Tao Shen. Maximum density divergence for domain adaptation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 43(11):3918–3930, 2021. 3
- [16] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, pages 69–84, 2016. 2
- [17] Mehdi Noroozi, Hamed Pirsiavash, and Paolo Favaro. Representation learning by learning to count. In *ICCV*, pages 5899–5907, 2017. 2
- [18] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *CVPR*, pages 2536–2544, 2016. 2
- [19] Omid Poursaeed, Tianxing Jiang, Han Qiao, Nayun Xu, and Vladimir G. Kim. Self-supervised learning of point clouds via orientation estimation. In *3DV*, pages 1018–1028, 2020. 1, 2, 3, 6, 7
- [20] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 77–85, 2017. 1, 2, 3
- [21] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, pages 5099–5108, 2017. 2
- [22] Can Qin, Haoxuan You, Lichen Wang, C.-C. Jay Kuo, and Yun Fu. Pointdan: A multi-scale 3d domain adaption network for point cloud representation. In *NeurIPS*, pages 7190–7201, 2019. 2, 3, 6, 7
- [23] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *CVPR*, pages 3723–3732, 2018. 3, 5
- [24] Jonathan Sauder and Bjarne Sievers. Self-supervised deep learning on point clouds by reconstructing space. In *NeurIPS*, pages 12942–12952, 2019. 1, 2, 3, 7
- [25] Ozan Sener, Hyun Oh Song, Ashutosh Saxena, and Silvio Savarese. Learning transferrable representations for unsupervised domain adaptation. In *NeurIPS*, pages 2110–2118, 2016. 3
- [26] Hugues Thomas, Charles Ruizhongtai Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *ICCV*, pages 6410–6419, 2019. 2
- [27] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *CVPR*, pages 2962–2971, 2017. 3
- [28] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Duc Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *ICCV*, pages 1588–1597, 2019. 2, 6
- [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. *ACM Trans. Graph.*, 38(5):146:1–146:12, 2019. 1, 2, 3, 7

- [30] Wenxuan Wu, Zhongang Qi, and Fuxin Li. Pointconv: Deep convolutional networks on 3d point clouds. In *CVPR*, pages 9621–9630, 2019. 1, 2
- [31] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, pages 1912–1920, 2015. 2, 6
- [32] Richard Zhang, Phillip Isola, and Alexei A. Efros. Colorful image colorization. In *ECCV*, pages 649–666, 2016. 2
- [33] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by GAN improve the person re-identification baseline in vitro. In *ICCV*, pages 3774–3782, 2017. 3
- [34] Longkun Zou, Hui Tang, Ke Chen, and Kui Jia. Geometry-aware self-training for unsupervised domain adaptation object point clouds. In *ICCV*, pages 6403–6412, 2021. 1, 2, 3, 6, 7

1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079