



# ULD-Net: 3D unsupervised learning by dense similarity learning with equivariant-crop

YU TIAN,<sup>1,2</sup> DA SONG,<sup>1,2</sup> MENGNA YANG,<sup>1,2</sup> JIE LIU,<sup>1,2,3</sup> GUOHUA GENG,<sup>1,2</sup> MINGQUAN ZHOU,<sup>1,2</sup> KANG LI,<sup>1,2</sup> AND XIN CAO<sup>1,2,4</sup> 

<sup>1</sup>School of Information Science and Technology, Northwest University, Xi'an, Shaanxi, China

<sup>2</sup>National and Local Joint Engineering Research Center for Cultural Heritage Digitization, Xi'an, Shaanxi 710127, China

<sup>3</sup>e-mail: jieliu2017@126.com

<sup>4</sup>e-mail: xin\_cao@163.com

Received 18 August 2022; revised 25 October 2022; accepted 27 October 2022; posted 1 November 2022; published 29 November 2022

Although many recent deep learning methods have achieved good performance in point cloud analysis, most of them are built upon the heavy cost of manual labeling. Unsupervised representation learning methods have attracted increasing attention due to their high label efficiency. How to learn more useful representations from unlabeled 3D point clouds is still a challenging problem. Addressing this problem, we propose a novel unsupervised learning approach for point cloud analysis, named ULD-Net, consisting of an equivariant-crop (equiv-crop) module to achieve dense similarity learning. We propose dense similarity learning that maximizes consistency across two randomly transformed global–local views at both the instance level and the point level. To build feature correspondence between global and local views, an equiv-crop is proposed to transform features from the global scope to the local scope. Unlike previous methods that require complicated designs, such as negative pairs and momentum encoders, our ULD-Net benefits from the simple Siamese network that relies solely on stop-gradient operation preventing the network from collapsing. We also utilize the feature separability constraint for more representative embeddings. Experimental results show that our ULD-Net achieves the best results of context-based unsupervised methods and comparable performances to supervised models in shape classification and segmentation tasks. On the linear support vector machine classification benchmark, our ULD-Net surpasses the best context-based method spatiotemporal self-supervised representation learning (STRL) by 1.1% overall accuracy. On tasks with fine-tuning, our ULD-Net outperforms STRL under fully supervised and semisupervised settings, in particular, 0.1% accuracy gain on the ModelNet40 classification benchmark, and 0.6% medium intersection of union gain on the ShapeNet part segmentation benchmark. © 2022 Optica Publishing Group

<https://doi.org/10.1364/JOSAA.473657>

## 1. INTRODUCTION

As a common 3D representation, the significant advantage of point cloud data over other representations (e.g., volumetric grids, meshes, and depth images) lies in its easy availability. With the advancement of 3D acquisition technologies, various types of 3D scanners, LiDARs, and red, green, and blue (RGB)-D cameras (e.g., cameras in Kinect and Apple devices) are becoming ever more accessible; thus, point cloud data can be quickly acquired without triangulating data into grids or voxel form. Hence, point cloud data is ideal for wide-ranging applications, such as autonomous driving [1], building information modeling [2], and digital preservation of ancient artifacts [3].

Recently, deep learning approaches became a dominant source in point cloud analysis in resolving various problems, including 3D shape classification, segmentation, object detection and tracking, registration, and so on. The remarkable advances in point cloud shape understanding rely on the large

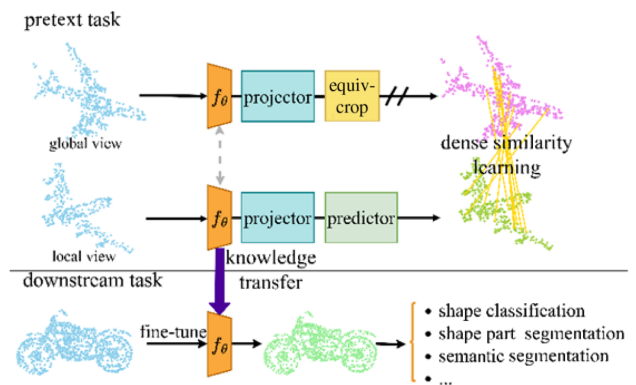
scale of labeled training data, and the performance improves logarithmically based on the size of annotated training data. The tedious and resource-consuming annotation process became a bottleneck for sustainable success due to the following reasons: (1) because of the sparsity, annotations for low-resolution point clouds are always ambiguous; (2) with the huge amount of points in dense objects, which can reach hundreds of millions, point-by-point annotation comes with significant costs; (3) the annotation for 3D objects are inherently more error prone than for 2D instances as its high complexity; (4) few works have focused on building automatic annotation tools for 3D point clouds, existing tools are still in the early stage manifested in their low accuracy and inconvenience.

In order to resolve the above practical difficulties, researchers explored unsupervised representation learning (URL) in the 3D point cloud analysis field based on the easy availability of unlabeled data. In common URL settings, the network learns

knowledge from pretext tasks without supervision in the pre-training stage, then, transfers the learned knowledge to other downstream tasks. Most existing works are based on generation tasks [4–8] that rely heavily on the specific architecture designation, such as folding-based decoders for completion and reconstruction tasks, the performance in downstream tasks degenerates when using a general multilayer perceptron (MLP-) based decoder. Meanwhile, generation-based tasks concentrate on geometric structures, which results in poor transferability on scene-level datasets. To eliminate the dependence on such specific components and improve model transferability to real-world scenes, several recent works consider context similarities [9–16] between samples to explore generic methods for URL.

Inspired by the huge success of self-supervised learning in the 2D computer vision domain, several efforts have been devoted to exploring context similarities in 3D point clouds based on Siamese networks. Most works built on top of contrastive learning rely on negative samples, Info3D [15] proposes to learn representations by maximizing mutual information between 3D objects and their local parts. Towards more discriminative features from local patches, Du *et al.* [16] introduced a hard negative sampling strategy into architecture. PointContrast [9] extracts dense correspondences across two views of scene point clouds for point-level contrastive learning. Without the requirements of negative samples, spatiotemporal self-supervised representation learning (STRL) [10] extends bootstrap your own latent (BYOL) [17] from 2D image processing to 3D point cloud analysis by learning features between original objects and their augmented views. However, previous works conduct unsupervised pretraining with complicated designations, such as negative pairs sampling [9], memory banks [15], and momentum encoders [10]. Additionally, most methods individually considered context similarities between transformed views at the instance level [10,15,16] or point level [9], separately maintaining consistency at both two levels was not taken into account.

To this end, we present a dense representation learning approach named 3D unsupervised learning by dense similarity learning with equivariant-crop (equiv-crop) (ULD-Net) based on three common-sense intuitions. First, purely considering instance-level similarity dismisses local spatial information; whereas, learning point-level similarity cannot extract representative abstract semantic information for the entire object. Thus, we jointly optimize the model at both levels, which helps to learn sufficient knowledge for downstream tasks. Second, the two branches of the network output point-level features within different scopes; whereas, point-level similarity learning aims to maximize corresponding features across views, the features should share the same scopes with one-to-one correspondence. Therefore, we propose an equiv-crop module equivariant with the cropping transformation to map the global features to the local scope. Third, it is proved that without redundant components which raise the computational cost, a simple stop-gradient (sg) design can get the network rid of collapse [18]. Using these inductive biases alone, we can train a Siamese network with a sg operation on top of SimSiam [18] to output point embeddings with objectives maximizing similarities between embeddings across local–global views, aiming at pretraining



**Fig. 1.** Illustration of the proposed method.

dense representations with strong transferability in downstream tasks.

The process of the proposed method is illustrated in Fig. 1. The pair of augmented point clouds (shown as blue dots) in global–local views are processed by the same encoder network and a projector network to extract features (shown as pentagons or crosses in other colors).

The equiv-crop module is applied to the global view side to project global features to the local scope. The predictor network is applied on one side, and the sg operation is on the other side. After taking dense similarity learning as a pretext task during pretraining, the trained encoder network transfers the learned knowledge to downstream tasks, such as shape classification, part segmentation, semantic segmentation, and so on. We theoretically prove the intuitions can improve the performance through serial experiments conducted, the method we proposed achieves competitive results to existing methods. Our contributions can be summarized as follows:

- (1) We propose a novel method for 3D point cloud unsupervised representation learning, which learns dense features by maximizing their local–global similarities at the point level and instance level, eliminating the need for negative samples or other complicated designs.
- (2) We introduce a novel point mapping strategy named equiv-crop for correspondence across views with local and global scopes, to provide the foundation for point-level feature learning. The local scope is produced by a cropping operation, and two augmented views are generated by integrating with an inv-aug strategy; whereas, the robustness is boosted.
- (3) We present a feature separability constraint that maximizes the separability of feature vectors from different dimensions; whereas, boosting the representability of features.

## 2. RELATED WORKS

### A. Deep Architectures for Point Cloud Processing

The advances in deep learning and learning-based point descriptors have been helpful to the impressive performance of recent point cloud processing for several 3D understanding tasks. Existing methods focus on alleviating the difficulty caused by the irregularity of 3D point clouds with most works extracting features directly from points.

PointNet [19] is the seminal work using deep learning that performs directly on raw point clouds, which achieved input-order invariance by symmetric functions. Since PointNet learns features independently through pointwise MLP for each point, later works paid attention to capturing local structural information using various methods. PointNet++ [20] learns local features by a hierarchical network, which is stacked by set abstraction layers. PointCNN [21] designed discrete convolutional kernel  $\chi$ -conv particular for point clouds. Considering each point-in-point clouds as a vertex, DGCNN [22] and RGCNN [23] construct graphs in spatial and spectral spaces. Kd-Net [24] learns features by constructing hierarchical data structures based on K-d trees. Recently, transformer-based methods [25,26] are proposed for long-range visual dependencies learning. In this paper, common architectures are suitable to be utilized as backbone networks because of our flexible designation.

**B. Deep Architectures for Point Cloud Processing**

URL is drawing increasing interest owing to its superiority in resolving the annotation bottleneck. Since annotations for 3D data take higher costs than 2D vision data, 3D tasks are supposed to benefit much more from URL. However, compared with natural language processing and 2D vision, the unsupervised pretext task defined for 3D point cloud data is much less mature.

Numerous pretext tasks have been proposed for strong presentation acquisition with specific objectives, which can be broadly divided into two categories: generation-based and context-based tasks. Generation-based tasks take point clouds themselves as supervised information, including reconstructing original input from low-dimensional vectors [4,5], generating new point clouds similar to training samples from random noise [6], up-sampling point clouds from sparse to dense [7], and completing missing parts [8]. Learning features through context-based methods is another rising research direction, including performing instance discrimination [9,10], solving 3D jigsaw puzzles [11], predicting rotation angles [12], predicting the next point in the sequence [13], and disentangling the mixed point clouds [14]. Considering multilevel similarity, a pretext task defined as optimizing the cosine similarities at both the instance level and the point level is proposed, accompanied by a feature separability constraint aiming at more representative features.

**C. Siamese Neural Networks**

The Siamese network consists of two identical artificial neural networks for comparing the projected representations of the two input vectors. The key challenge in Siamese methods is how to avoid collapsing solutions. SimCLR [27] and MoCo [28] are proposed based on the core idea of contrastive learning that drags positive sample pairs and pushes negative sample pairs away. Different from comparing samples in the current batch in SimCLR and MoCo builds a dynamic dictionary with a queue and a moving-averaged encoder to get rid of the dependence on large batch size and to improve the consistency of the queues. Clustering-based methods construct Siamese networks with

clustering intergraded, whereas, achieving competitive results without a memory bank. Specifically, SwAV [29] solves degenerate solutions through computing cluster assignments from one view playing as negative samples relying on the Sinkhorn–Knopp algorithm. Asymmetric methods prevent features from collapsing using asymmetric architecture. BYOL [17] uses a momentum encoder accompanied by sg and moving average; whereas, SimSiam [18] removes the momentum encoder and keeps minimum core architecture as an elegant realization. Inspired by SimSiam, we build our network based on the Siamese architecture with a sg operation as its computational advantage.

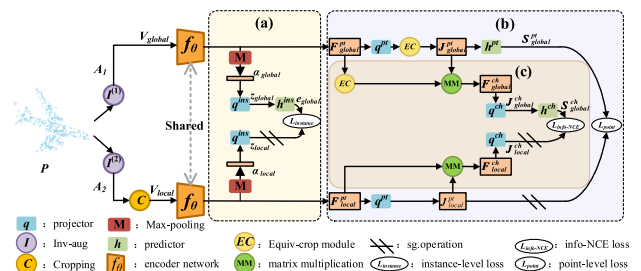
**3. METHOD**

The overall pipeline of our ULD-Net is depicted in Fig. 2. Taking unsupervised point cloud datasets as source data, our fundamental idea is to train an encoder network by modeling dense consistency between local–global features from transformed views to extract representations for better transferability on downstream tasks.

For each point cloud object  $P$ , we first transform the original input into two random augmented views in the global and local scope by inv-aug and cropping transformations. Then, we encode the point clouds to generate feature maps in high-dimensional space by a weight-shared encoder network  $f_\theta$ . Inspired by SimSiam, we promote the representational ability of the encoder by minimizing the dissimilarities between feature maps through dense similarity learning. We integrate the instance-level [Fig. 2(a)] and point-level [Fig. 2(b)] similarity learning into a unified framework, and we utilize feature separability constraint [Fig. 2(c)] for more discriminative features. Taking our approach to pretrain an encoder network from unlabeled data, the learned encoder can be transferred to various downstream tasks for feature extraction.

**A. Views and Features Generation**

Given each input point cloud  $P \in \mathbb{R}^{N \times 3}$  with  $N$  elements, we transform its geometric features ( $XYZ$  coordinates) by inv-aug and cropping operations with random factors. Inv-aug is a collection of data augmentations consisting of rotation, translation, scaling, and jittering. We start the transformation with two randomly inv-aug augmentations  $I^{(1)}$  and  $I^{(2)}$  that output two augmented point clouds  $A_1 = I^{(1)}(P)$  and  $A_2 = I^{(2)}(P)$ . Regarding one augmented point cloud  $A_1$  as the global view



**Fig. 2.** Overall overview of the proposed pretraining method. (a) Instance-level similarity learning. (b) Point-level similarity learning. (c) Feature separability constraint.

$V_{\text{global}} = A_1$  of the input point cloud  $P$  [as in Eq. (3)], we conduct a cropping operation  $C$  on another augmented point cloud to generate the local view. Different from inv-aug, which keeps the object complete, the cropping operation transforms points to a random local scope.

### 1. Cropping Operation

The cropping operation  $C$  consists of two steps conducted at the augmented point set  $A_2$ . First, we compute the indices of, at least, 50% of points inside the coordinate range defined by a random 3D cuboid with computed indices  $m = \{j_i \in [1, \dots, N], i \in [1, \dots, L]\}$ , and the selected  $L$  points  $A_2[m]$  inside a cuboid local scope are kept as downsampled points. Since the point sequence is invariant to inv-aug operations  $I^{(1)}$  and  $I^{(2)}$ , the downsampled point set  $A_2[m]$  with  $L$  elements corresponds exactly to points with the same indices  $m$  in the global scope point set  $V_{\text{global}}$  (i.e.,  $A_1$ ). Second, we upsample the downsampled point set  $A_2[m]$  from size  $L$  to the predefined input size  $N$  of the encoder network. We choose the inverse distance weighted average based on  $k$ -nearest neighbors [as in Eq. (1), in default we use  $p = 2$ ,  $K = 3$ ] for interpolation, the outputs upsampled point set regarded as the local view [as in Eq. (4)],

$$V_{\text{local}}[i] = \frac{\sum_{r=1}^K w_r(A_2[i]) A_{\text{knn}}[i, r]}{\sum_{r=1}^K w_r(A_2[i])}, \quad i \in [1, \dots, N], \quad (1)$$

where  $A_{\text{knn}} \in \mathbb{R}^{N \times K \times 3}$  denotes the  $K$ -nearest neighbors of  $N$  points in the entire point set  $A_2$ , the neighbors are searched in the downsampled point set  $A_2[m]$ ,  $d(\cdot)$  denotes the Euclidean distance between two points, and  $w_r$  is computed for the weight of the  $r$ th neighbor,

$$w_r(A_2[i]) = \frac{1}{d(A_2[i], A_{\text{knn}}[i, r])^p}. \quad (2)$$

In conclusion, the correspondence between two scales is constructed by downsampled indices  $m$  and the  $K$ -neighbors  $A_{\text{knn}}$ . Thus, two different views in the local and global scopes with transformation in Eqs. (3) and (4) are produced

$$V_{\text{global}} = A_1 = I^{(1)}(P), \quad (3)$$

$$V_{\text{local}} = C(A_2) = C(I^{(2)}(P)). \quad (4)$$

### 2. Dense Feature Map Generation

The global–local views from the same point cloud are then processed by a backbone encoder network  $f_\theta$  with parameters  $\theta$ . The encoder shares the same weights between views. For local input  $V_{\text{local}}$ , the encoder  $f_\theta$  computes a point-level feature map  $F_{\text{local}}^{\text{pt}} = f_\theta(V_{\text{local}})$  including representations for each point in the local view  $V_{\text{local}}$ , the feature vector for the  $i$ th point is noted as  $F_{\text{local}}^{\text{pt}}[i]$ . Simultaneously,  $f_\theta$  yields a high-dimensional vector  $\alpha_{\text{local}}$  after max pooling describing the entire

local view  $V_{\text{local}}$  at the instance level. Following the same computation pipeline with local features, a point-level feature map  $F_{\text{global}}^{\text{pt}} = f_\theta(V_{\text{global}})$  and an instance-level representation  $\alpha_{\text{global}}$  from the global view are generated

$$\alpha_{\text{global}} = \text{maxpooling}(f_\theta(V_{\text{global}})), \quad (5)$$

$$\alpha_{\text{local}} = \text{maxpooling}(f_\theta(V_{\text{local}})). \quad (6)$$

### 3. Equiv-Crop

Denoting high-dimensional features from the global and local views without max pooling as  $F_{\text{global}}, F_{\text{local}} \in \mathbb{R}^{N \times D}$ , where  $D$  denotes the number of feature dimensions, the global view features  $F_{\text{global}}$  can be transformed to the local scope through a module EC equivalent to the cropping operation. Since the network is permutation invariant, the sequences of output features correspond to network input sequences. Namely, there is a one-to-one correlation between the input point and the output feature. For example, the  $i$ th point  $V_{\text{global}}[i]$  in the global view is represented by the feature vector  $F_{\text{global}}[i]$ . Based on such a principle, global features  $F_{\text{global}}$  can be directly mapped into the corresponding local scope using the same correspondence in the cropping operation performed in views transformation. Specifically, we gather the global features of points in the same local cuboid scope in cropping following the same downsampling and upsampling steps. First, we downsample the features by selecting indices  $m$  saved in cropping, the downsampled features  $F_{\text{global}}[m]$  represent features of points in downsampled points  $A_2[m]$  in cropping. Then, with the downsampled features  $F_{\text{global}}[m]$ , we upsample the features from the searched  $K$ -nearest neighbors  $A_{\text{knn}}$  the same as in cropping,

$$\text{EC}(F_{\text{global}}[i]) = \frac{\sum_{r=1}^K w_r(A_2[i]) F_{\text{knn}}[i, r]}{\sum_{r=1}^K w_r(A_2[i])}, \quad i \in [1, \dots, N], \quad (7)$$

where  $F_{\text{knn}} \in \mathbb{R}^{N \times K \times D}$  denotes features of neighbors  $A_{\text{knn}}$ ,  $d(\cdot)$  denotes the Euclidean distance between two vectors, and  $w_r$  denotes the weight of the  $r$ th neighbor. The equiv-crop module, then, transforms global features into the local scope defined in cropping.

### B. Dense Similarity Learning

To achieve a sophisticated similarity measurement, we learn local–global consistency through dense similarity learning. Solely learning consistency between local and global views at the instance level would cause most of the spatial information to be discarded during pooling. To tackle this question, we jointly learn instance-level and point-level similarities. Moreover, for point-level feature learning, we utilize the equiv-crop module in Section 3.A.3 towards mapping point embeddings from the global scope to the local one.

### 1. Instance-Level Similarity Learning

We learn instance-level similarity from representations  $\alpha_{\text{local}}$  and  $\alpha_{\text{global}}$ , the pipeline is shown in Fig. 2(a). We first transform features by the same projector network  $q^{\text{ins}}$ , which is a three-layer MLP head output with features  $z_{\text{global}} = q^{\text{ins}}(\alpha_{\text{global}})$  and  $z_{\text{local}} = q^{\text{ins}}(\alpha_{\text{local}})$ . Then, a predictor network  $h^{\text{ins}}$  transforms the projected feature from one view to predict another, outputs predictions  $e_{\text{global}} = h^{\text{ins}}(z_{\text{global}})$  and  $e_{\text{local}} = h^{\text{ins}}(z_{\text{local}})$ . Meanwhile, a sg operation is applied to the projected features from another view. We symmetrically minimize the distance of feature maps and predictions from another view,

$$L_{\text{instance}} = \frac{1}{2} D(e_{\text{local}}, \text{sg}(z_{\text{global}})) + \frac{1}{2} D(e_{\text{global}}, \text{sg}(z_{\text{local}})), \quad (8)$$

where sg avoids the outputs of the network collapsing to a constant and  $D(\cdot)$  in Eq. (9) is a distance function measuring negative cosine similarity in high-dimensional feature space,

$$D(e_{\text{local}}, z_{\text{global}}) = -\frac{e_{\text{local}}}{\|e_{\text{local}}\|_2} \cdot \frac{z_{\text{global}}}{\|z_{\text{global}}\|_2}, \quad (9)$$

where  $\|\cdot\|_2$  denotes  $l_2$  normalization.

### 2. Point-Level Similarity Learning

We formulate point-level similarity learning as shown in Fig. 2(b) to maximize the similarity of point predictions. Input with point-level feature maps  $F_{\text{local}}^{\text{pt}}$  and  $F_{\text{global}}^{\text{pt}}$  following the same pipeline with instance-level similarity learning, a projector  $q^{\text{pt}}$  is used to transform the point-level features first. For each point, we predict its feature from another view. However, due to the input point represented by the  $i$ th feature mismatch between local and global scopes, it is incompatible with common sense to predict directly between two features from the same indices. To bridge the gap, an equiv-crop module EC maps the projected features from global to local scopes, and the projected features are noted as  $J_{\text{global}}^{\text{pt}} = \text{EC}(q^{\text{pt}}(F_{\text{global}}^{\text{pt}}))$ ,  $J_{\text{local}}^{\text{pt}} = q^{\text{pt}}(F_{\text{local}}^{\text{pt}})$ . After that, the predictions  $S_{\text{global}}^{\text{pt}} = h^{\text{pt}}(J_{\text{global}}^{\text{pt}})$  and  $S_{\text{local}}^{\text{pt}} = h^{\text{pt}}(J_{\text{local}}^{\text{pt}})$  for each point are outputted from the point-level predictor  $h^{\text{pt}}$ .

We symmetrically maximize the similarity between the projected feature for the  $i$ th point and its prediction,

$$L_{\text{point}} = \sum_{i=1}^N \frac{1}{2} D(S_{\text{local}}^{\text{pt}}[i], \text{sg}(J_{\text{global}}^{\text{pt}}[i])) + \frac{1}{2} D(S_{\text{global}}^{\text{pt}}[i], \text{sg}(J_{\text{local}}^{\text{pt}}[i])). \quad (10)$$

### C. Feature Separability Constraint

It is common that projected features and predictions (such as  $J_{\text{global}}^{\text{pt}}[i]$  and  $S_{\text{local}}^{\text{pt}}[i]$ ) contain different information after random augmentations, but similarity learning forces these embeddings to be close to each other, which leads to a risk of features from different dimensions degenerating to the same value. To address the degenerating issue, besides the sg operation, we further propose a feature separability constraint as illustrated in Fig. 2(c) to boost the expressiveness of features.

The channel embeddings are obtained by the sum of multiplication between the feature maps and the predictions,

$$F_{\text{global}}^{\text{ch}} = \sum_i^N J_{\text{global}}^{\text{pt}}[i] \cdot \text{EC}(F_{\text{global}}^{\text{pt}})[i], \quad (11)$$

$$F_{\text{local}}^{\text{ch}} = \sum_i^N J_{\text{local}}^{\text{pt}}[i] \cdot F_{\text{local}}^{\text{pt}}[i], \quad (12)$$

where  $F_{\text{global}}^{\text{ch}}, F_{\text{local}}^{\text{ch}} \in \mathbb{R}^{D \times D'}$ ,  $D$  and  $D'$  represent the number of output feature channels in the predictor and encoder.

Similar to similarity learning, we transformed the embeddings by a projector composed of MLP head  $q^{\text{ch}}$  and predictor  $h^{\text{ch}}$  output embeddings  $J_{\text{global}}^{\text{ch}} = q^{\text{ch}}(F_{\text{global}}^{\text{ch}})$ ,  $J_{\text{local}}^{\text{ch}} = q^{\text{ch}}(F_{\text{local}}^{\text{ch}})$  and predictions  $S_{\text{global}}^{\text{ch}} = h^{\text{ch}}(J_{\text{global}}^{\text{ch}})$ ,  $S_{\text{local}}^{\text{ch}} = h^{\text{ch}}(J_{\text{local}}^{\text{ch}})$ . By using the information-noise-contrastive estimation (info-NCE) loss [30], the similarities of features in different channels decreased, which leads to higher separability. Specifically, we optimize the feature separability by Eq. (13),

$$L_{\text{separability}} = \frac{1}{2} L_{\text{info-NCE}}(S_{\text{local}}^{\text{ch}}, \text{sg}(J_{\text{global}}^{\text{ch}})) + \frac{1}{2} L_{\text{info-NCE}}(S_{\text{global}}^{\text{ch}}, \text{sg}(J_{\text{local}}^{\text{ch}})), \quad (13)$$

where  $L_{\text{info-NCE}}$  is the info-NCE loss as

$$L_{\text{info-NCE}}(S, J) = -\sum_{r=0}^R \log \frac{\exp(S[r] \cdot J[r]/\tau)}{\sum_{r'=0}^R \exp(S[r] \cdot J[r']/\tau)}, \quad (14)$$

where  $\tau$  denotes the temperature coefficient of 0.1 in default and  $R$  denotes the number of dimensions of the prediction feature  $S$ .

## 4. EXPERIMENTS AND RESULTS

### A. Datasets

To validate the effectiveness and transferability of our method, three benchmarks [ModelNet40 [31], the ShapeNet part [32], and the Stanford 3D indoor spaces (S3DISs) [33]] are used in the experiments. In the pretraining stage, ModelNet40 is used for all experiments, and ShapeNet55 is additionally used for linear evaluation comparison. For downstream tasks, we use the ModelNet40 benchmark for shape classification, the ShapeNet part benchmark for the shape part segmentation, and the S3DIS benchmark for scene semantic segmentation.

**ModelNet40.** ModelNet40 includes 12,311 synthesized 3D objects (divided into 9843 training samples and 2468 testing samples) from 40 categories. We downsample each object to 2048 points whose  $XYZ$  coordinates normalized into a unit sphere following the preprocessing method from PointNet [19].

**ShapeNet Part.** ShapeNet55 [34] contains 57,748 synthetic 3D shapes from 55 categories. The ShapeNet part benchmark includes 16,881 shapes of 16 categories selected from ShapeNet55. Each sample is annotated with two to five parts, part labels for all categories amounted to 50. Intersection of union (IoU) is widely used for segmentation evaluation that

measures the ratio between pointwise ground truth and prediction. For the part segmentation task, we compute category mean IoU (mIoU) by averaging IoUs over parts of the same object category, and instance mIoU is obtained by averaging over all test shapes.

**S3DIS.** The S3DIS dataset contains 3D scans of six different places including 271 rooms, which cover over 6000 m<sup>2</sup>. Each point is represented by a nine-dimensional vector consisting of  $XYZ$  coordinates, RGB color values, and normalized location, and the individual point is labeled with 13 semantic categories. We use the same preprocessing procedures as the original work, each room is split into blocks with  $1m \times 1m$  areas, and each block contains 4096 points sampled. To evaluate semantic segmentation performance, mIoU is computed by averaging IoUs over all points.

## B. Implementation Details

**Architecture Parameters.** For a fair comparison with previous methods, the DGCNN backbone is used as the default encoder network, which outputs features with 1024 dimensions. All projectors and predictors are designed with the same architecture. Specifically, each projection MLP head consists of three fully connected layers with dimensions of [512,256,256], each prediction MLP head consists of two fully connected layers with dimensions of [512,256], each layer has batch normalization applied, and LeakyReLU activation with a negative slope of 0.2 is used except for the final output layer. For jointly learning instance-level similarity, point-level similarity, and feature separability, our ULD-Net optimizes the total loss  $L = \lambda_1 L_{\text{instance}} + \lambda_2 L_{\text{point}} + \lambda_3 L_{\text{separability}}$  to balance the significance of all tasks, we choose  $\lambda_1 = \lambda_2 = 100$  and  $\lambda_3 = 10$  based on the numbers of each loss to keep them in the same order of magnitude.

**Pretraining Setup.** We follow the settings of the STRL in unsupervised pretraining experiments. We implemented our paper with the deep learning library PyTorch using a single TITAN RTX GPU for all experiments. Specifically, the Adam optimizer is used in our model with an initial learning rate of 0.001, the learning rate is decayed by 0.7 every 20 epochs, and the batch size is 24 by default. We pretrain ULD-Net for 200 epochs on ModelNet40.

**Fine-tuning Setup.** As an end goal in URL, we verify the effectiveness of the pretrained features transferred to new tasks in a fully supervised fashion. For 3D shape classification on ModelNet40, we use a batch size of 24 for training and testing with 250 epochs, the stochastic gradient descent optimizer is used with an initial learning rate of 0.1, momentum 0.9, and weight decay 0.0001, and the learning rate is decayed with a cosine annealing scheduler. Slightly different from the above settings for the classification task, the batch size used for 3D part segmentation on the ShapeNet part is 16, and we train the network for 100 epochs with the Adam optimizer for 3D semantic segmentation on the S3DIS.

## C. Downstream Results

### 1. Evaluation for Shape Classification

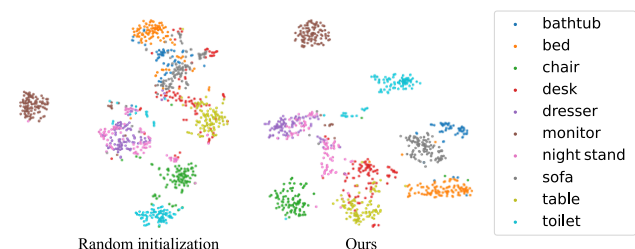
For 3D shapes classification, we train a linear support vector machine (SVM) on the target dataset ModelNet40 to evaluate the effectiveness of the learned instance-level features following the common protocol in prior URL works [8,10,11]. For the SVM classifier, the input features are obtained after the pretrained encoder network with the following pooling layer, and the weights of the feature extractor are frozen during evaluation. Following the settings in the DGCNN classification network, the pooling layer outputs concatenated features after max-pooling and average-pooling operations. The classification results compared with the state of the art are shown in Table 1, all methods tabulated are implemented with the DGCNN backbone as a feature extractor for a fair comparison. As shown in the table, the proposed method achieves 91.9% and 92.0% overall accuracy after pretrained on the ShapeNet55 and ModelNet40 datasets, which outperform the existing unsupervised method STRL [10] by 1.0% and OcCo [8] by 2.8%. These results suggest that the features attained by our pretraining method are discriminative that can easily achieve competitive performance even with little effort of training on the SVM.

Towards a better understanding of the capability of our method proposed, we visualize the learned features on the test dataset of the ModelNet10 as illustrated in Fig. 3, which is compared with features from a randomly initialized encoder network. Using  $t$ -distributed stochastic neighbor embedding [35] to project the instance-level high-dimensional features in 2D space, we observe that the learned features from instances of different categories are separable except dressers and night stands, which are difficult to distinguish even by a human.

**Table 1. Classification Accuracy Results (%) with Linear SVM in URL Methods on ModelNet40<sup>a</sup>**

Pretraining Dataset	Method	OA
ShapeNet	FoldingNet [4]	88.4
	Du <i>et al.</i> [16]	89.6
	Jigsaw3D [11]	90.6
	Rotation 3D [12]	90.8
	STRL [10]	90.9
	Ours	<b>91.9</b>
ModelNet40	FoldingNet [4]	84.4
	Jigsaw3D [11]	87.8
	MAP-VAE [5]	90.2
	OcCo [8]	89.2
	Ours	<b>92.0</b>

<sup>a</sup>“OA” denotes overall accuracy.



**Fig. 3.** Visualization of pretrained instance-level features.

Compared with the projected features from random initialization, our pretrained methods are dragged to further distances between features of distinct categories. Since the random initialized features can be regarded as the prior of the encoder network, the comparison proves our pretraining method can learn knowledge of 3D shapes without supervision.

## 2. Supervised Fine-Tuning for 3D Shape Classification

Further fine-tunes the encoder network on ModelNet40 without freezing, resulting in better classification accuracy. Following a common fine-tuning pipeline in URL methods, after an unsupervised pretraining stage aimed at maximizing dense similarities, we take the pretrained encoder network parameters as the initialization for the encoder network used in transfer learning, then optimize the network by the specific objective for the classification task in a supervised fashion. To produce predictions for the classification task, we train a classification MLP head during fine-tuning along with the encoder network, and the classification head takes instance-level features after pooling as input and output with classification scores for each object towards supervised validation on ModelNet40. Comparisons of fine-tuned classification results are illustrated in Table 2.

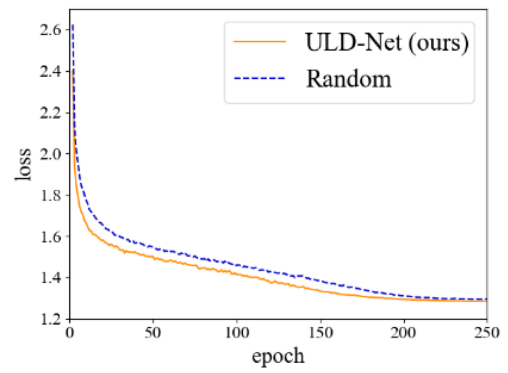
As shown in Table 2, after fine-tuning from our pretrained model, the proposed method achieves an additional 1.0% accuracy gain over the original DGCNN trained from randomly initialized parameters (93.2% versus 92.2%), which suggests our pretraining method can boost the ability of the feature extractor. Our method outperforms unsupervised methods OcCo and STRL by 0.2% and 0.1% in terms of overall accuracy and achieves the best fine-tuned performance on ModelNet40. The results indicate that our ULD-Net can attain a comparable performance with the state-of-the-art fully supervised methods.

Our method accelerates the convergence of the encoder framework during the fine-tuning stage. As shown in Fig. 4, compared with the random initialization, the loss number of our

**Table 2. Comparisons of Our Fine-Tuned Classification Accuracy (%) Result Against Other Methods on ModelNet40<sup>a</sup>**

Method	Sup.	OA
PointNet [19]	✓	89.2
RGCNN [23]	✓	90.5
PointNet++ [20]	✓	90.7
KD-Net [24]	✓	91.8
PointCNN [21]	✓	92.2
DGCNN [22]	✓	92.2
Point Cloud	✓	93.2
Transformer [26]		
PointTransformer [25]	✓	<b>93.7</b>
Jigsaw 3D [11]	✗	92.4
Info 3D [15]	✗	93.0
OcCo [8]	✗	93.0
FoldingNet [4]	✗	93.1
STRL [10]	✗	93.1
Ours	✗	<b>93.4</b>

<sup>a</sup>“Sup.” denotes supervised.



**Fig. 4.** Convergence curves during fine-tuning.

**Table 3. Fine-Tuned Results (%) under a Semisupervised Setting**

Method	1%	5%	10%	20%
DGCNN [22]	58.4	80.7	85.2	88.1
STRL [10]	60.5	<b>82.7</b>	86.5	89.7
Ours	<b>60.6</b>	82.5	<b>86.8</b>	<b>89.8</b>

method remains lower than random initialization at about 0.2 during training and convergence after fewer epochs.

## 3. Semisupervised Fine-Tuning for 3D Shape Classification

We further evaluate our pretrained model on the shape classification task under a semisupervised setting. We use the same setting as the STRL and report the overall accuracy on ModelNet40 as shown in Table 3. Specifically, we reduce the annotated input shapes to 1%, 5%, 10%, and 20% of the training data, and, at least, one shape is selected for each category. Then, we evaluate the model fine-tuned by the reduced training data on the full test dataset. The results show that our model surpasses the randomly initialized model by 2.2% and 1.7% when 1% and 20% of training shapes were sampled, and our ULD-Net slightly outperforms the STRL when the sampling ratios of 1%, 10%, and 20% indicate our pretraining method improves annotation efficiency.

## 4. Supervised Fine-Tuning for 3D Shape Part Segmentation

To validate the effectiveness of fine-grained point-level features gained from our method, we fine-tune the pretrained network for the part segmentation task. Different from classification, fine-tuning only transfers parameters from the encoder network, and segmentation fine-tuning uses parameters from the pretrained encoder and its attached point-level projector. We fine-tune them on the ShapeNet part dataset to verify the performance of our ULD-Net on the part segmentation task. The quantitative results compared with the state-of-the-art URL and supervised methods are shown in Table 4. It shows that our ULD-Net shows the best performance (85.7% instance mIoU) among other URL approaches and achieves top performance in six categories, such as aeroplane, car, and knife. Since the ShapeNet part is a long-tailed dataset, the instance of mIoU is mostly decided by shapes of large amounts (aeroplanes, chairs, lamps, tables, etc.), which leads to the unbalance performance

**Table 4.** Fine-Tuning Part Segmentation mIoU Results (%) on ShapeNet Part Dataset<sup>a</sup>

Shapes	Supervised Method				Unsupervised Method				
	DGCNN [22]	RSCNN [36]	PCT [26]	LGAN [6]	Method in [16]	Jigsaw 3D [11]	OcCo [8]	STRL [10]	Ours
Ins.	85.2	86.2	<b>86.4</b>	57.0	82.3	85.3	85.5	85.1	<b>85.7</b>
Aero	84.0	83.5	<b>85.0</b>	54.1	82.1	84.1	84.4	83.7	<b>84.7</b>
Bag	83.4	<b>84.8</b>	82.4	48.7	74.5	<b>84.0</b>	77.5	80.3	82.8
Cap	86.7	88.8	<b>89.0</b>	62.6	83.6	85.8	83.4	<b>87.6</b>	83.8
Car	77.8	79.6	<b>81.2</b>	43.2	74.9	77.0	77.9	77.7	<b>78.3</b>
Chair	90.6	91.2	<b>91.9</b>	68.4	87.9	90.9	<b>91.0</b>	90.9	90.9
Earphone	74.7	<b>81.1</b>	71.5	58.3	72.4	<b>80.0</b>	75.2	78.0	77.0
Guitar	91.2	<b>91.6</b>	91.3	74.3	89.9	91.5	<b>91.6</b>	91.4	91.3
Knife	87.5	<b>88.4</b>	88.1	68.4	85.4	87.0	<b>88.2</b>	87.7	<b>88.2</b>
Lamp	82.8	86.0	<b>86.3</b>	53.4	79.1	83.2	83.5	83.7	<b>83.8</b>
Laptop	95.7	<b>96.0</b>	95.8	82.6	95.2	95.8	<b>96.1</b>	<b>96.1</b>	95.6
Motor	66.3	<b>73.7</b>	64.6	18.6	67.3	<b>71.6</b>	65.5	66.7	68.6
Mug	94.9	94.1	<b>95.8</b>	75.1	93.3	94.0	94.4	<b>95.0</b>	94.3
Pistol	81.1	83.4	<b>83.6</b>	54.7	81.0	<b>82.6</b>	79.6	81.2	80.6
Rocket	<b>63.5</b>	60.5	62.2	37.2	58.2	60.0	58.0	58.2	<b>61.9</b>
Skateboard	74.5	<b>77.7</b>	77.6	46.7	74.0	<b>77.9</b>	76.2	75.3	75.1
Table	82.6	83.6	<b>83.7</b>	66.4	79.2	81.8	82.8	82.1	<b>83.4</b>

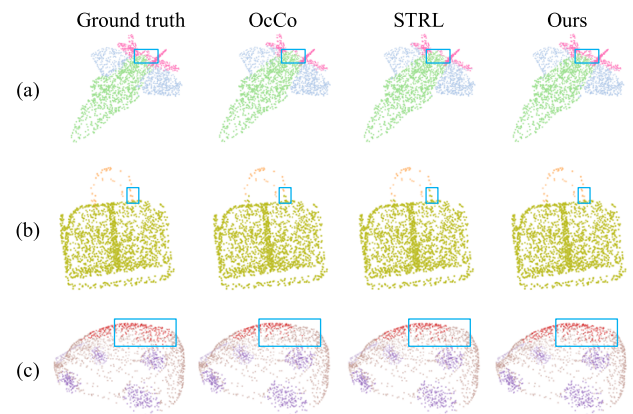
<sup>a</sup>“Ins.” denotes instance mIoU.

**Fig. 5.** Qualitative results on the ShapeNet part dataset.

on different categories of shapes. Compared with supervised methods, we also achieve comparable results.

The segmentation results of all shapes are qualitatively illustrated in Fig. 5. These visualization results show our method can segment one shape to clear parts close to the ground truths.

Furthermore, we compare our ULD-Net with the STRL and OcCo on shapes including aeroplanes, bags, and cars as illustrated in Fig. 6, which shows ULD-Net captures more local details than the STRL and OcCo. In confusing regions annotated with blue bounding boxes, containing points in the intersection of the main body and other parts of the different categories, such as the tail of aeroplanes demonstrated in Fig. 6(a), the handle of bags in Fig. 6(b), and the roof of cars in Fig. 6(c), show that our method distinguishes such regions better.

**Fig. 6.** Visual comparison of part segmentation on the ShapeNet part.

### 5. Supervised Fine-Tuning for 3D Semantic Segmentation

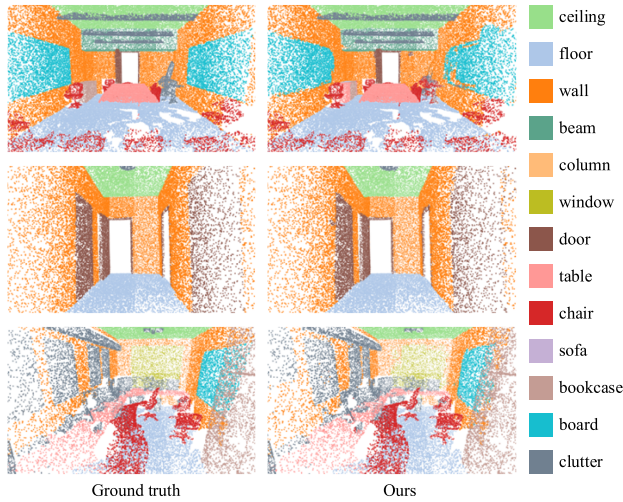
Transferring features pretrained on synthetic computer-aided design object models to real-world segmentation tasks is considered more challenging than tasks on synthetic shapes. To elucidate this problem, we also test our method for the indoor semantic segmentation task on the S3DIS dataset to validate the cross-domain generalizability of our pretrained features to a real-world dataset.

Using the pipeline similar to part segmentation, we transfer the parameters of the encoder and point-level projector to supervised fine-tuning for the semantic segmentation task. We test our model under sixfold cross validation over the six areas as in the original work [33]. As the quantitative results summarized in Table 5, our ULD-Net achieves the best segmentation result with 85.5% overall accuracy and 59.2% mIoU, which surpasses the state-of-the-art method OcCo by 0.4% overall accuracy and 0.7 mIoU. Compared with existing URL methods, these results demonstrate better transferability of our ULD-Net from



**Table 5. Semantic Segmentation Results (%) on S3DIS Dataset**

Method	Sup.	OA	mIoU
PointNet [19]	✓	78.6	47.6
PointNet++ [20]	✓	81.0	54.5
PointCNN [21]	✓	<b>88.1</b>	<b>65.4</b>
DGCNN [22]	✓	84.1	56.1
Jigsaw [11]	✗	84.4	56.6
OcCo [8]	✗	85.1	58.5
Ours	✗	<b>85.5</b>	<b>59.2</b>

**Fig. 7.** Visualization of semantic segmentation results on the S3DIS dataset.

synthetic shapes to real-world scene datasets. It is observed that our results even surpass the supervised PointNet, PointNet++, and DGCNN, and achieve competitive performance with other supervised models.

We show qualitative results of the S3DIS indoor semantic segmentation by visualizing selected rooms in Fig. 7. Empirically, we observe that our network is able to understand and classify semantic objects in a real-world scene, and our segmentation results are close to the ground truth.

#### D. Ablation Study

To investigate the effectiveness of our key components in ULD-Net, we study the impact of adopting different combinations of losses and transformations during the pretraining stage by validating the downstream SVM classification results using their pretrained features on ModelNet40.

##### 1. Transformations

We analyze the effectiveness of different transformations in inv-aug and cropping for view generation used in the pretraining stage. We remove certain transformations to produce augmented views when pretraining and validate the implication with SVM. As summarized in Table 6, our full model  $A_1$  uses all transformations and achieves the best result of 92.0%. Without any transformations (model  $B_1$ ), the network inputs of the two

**Table 6. Results (%) from Pretrained Features with Different Transformations<sup>a</sup>**

Model	C	Rot.	Trans.	Scal.	Jit.	OA
$A_1$	✓	✓	✓	✓	✓	<b>92.0</b>
$B_1$	✗	✗	✗	✗	✗	88.0
$C_1$	✗	✓	✓	✓	✓	89.6
$D_1$	✓	✗	✓	✓	✓	91.0
$E_1$	✓	✓	✗	✓	✓	91.0
$F_1$	✓	✓	✓	✗	✓	91.0
$G_1$	✓	✓	✓	✓	✗	91.0

<sup>a</sup>“Rot.” denotes rotation, “Trans.” denotes translation, “Scal.” denotes scaling, “Jit.” denotes jittering.

**Table 7. Ablation Study Results (%) of Different Pretraining Objectives**

Model	$L_{\text{instance}}$	$L_{\text{point}}$	$L_{\text{separability}}$	OA
$A_2$	✓	✗	✗	91.3
$B_2$	✓	✓	✗	91.7
$C_2$	✓	✗	✓	91.6
$D_2$	✓	✓	✓	<b>92.0</b>

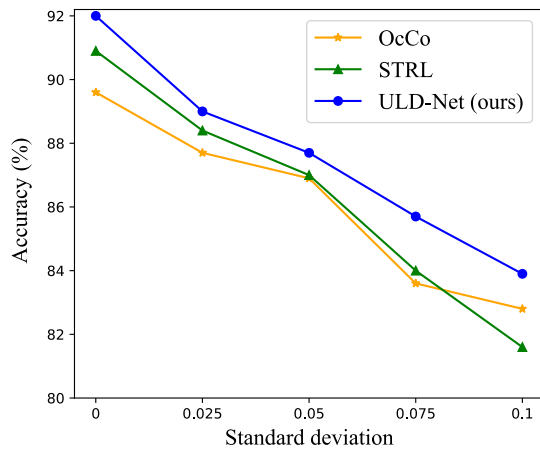
branches are exactly the same, which makes the network overfits pretraining samples due to too many task-irrelevant detailed features captured; hence, the classification result degenerates to 88.0%. The result reduces when one transformation is removed, proving that each adopted transformation schedule boosts the performance of pretrained features. Among transformations, removing the cropping transformation  $C$  (model  $C_1$ ) affects the performance the most by a 2.4% descent (92.0% vs 89.6%) compared with model  $A_1$ . Removing each transformation in inv-aug including rotation (model  $D_1$ ), translation (model  $E_1$ ), jittering (model  $F_1$ ), and scaling (model  $G_1$ ), the performance degenerates to 91.0%, 91.0%, 91.2%, and 91.5% respectively, which indicates the importance of each transformation is decreasing by the above order.

##### 2. Losses

We further study how the training objectives affect the performance of pretrained features. The results are shown in Table 7, the baseline model  $A_2$  is trained by the instance-level similarity loss, which closes the distance between the instance and its local parts in embedding space and gets a classification accuracy of 91.3%. Combined with one of the point-level similarity loss (model  $B_2$ ) or feature separability loss (model  $C_2$ ), we observed 0.4% and 0.3% improvements, respectively. Our full model joint learns with three objectives (model  $D_2$ ) and achieves a notable 92.0% on ModelNet40.

#### E. Robustness

To test the robustness of our method to random noise, we randomly jitter the  $XYZ$  coordinates of points with Gaussian noises in linear evaluation on ModelNet40 during test time. Each point cloud is jittered with randomly sampled Gaussian noises with zero mean and standard deviation  $\sigma \in \{0.025, 0.05, 0.075, 0.1\}$ . As shown in Fig. 8, we compare our ULD-Net with OcCo and the STRL under different



**Fig. 8.** Results with Gaussian noise.

noise levels. We can see that our ULD-Net remains robust with 83.9% accuracy even when noise is at a high level with a 0.1 standard deviation. It can also be observed that our ULD-Net gets competitive results with existing URL methods OcCo and the STRL.

## 5. DISCUSSION AND CONCLUSION

In this paper, we propose a novel URL method for point cloud analysis. Our method extracts features by dense similarity learning, which is composed of instance-level and point-level similarity learnings with the feature separability constraint. We also present the equiv-crop module to project point-level features from global to local scope to build correspondence across the transformed views. Without negative pairs, momentum encoder, or other complicated designs, ULD-Net pretrains the network that extracts representations with the best results on linear SVM validation. After fine-tuning the pretrained network on other downstream tasks including shape classification, shape part segmentation, and semantic segmentation, our ULD-Net also achieves competitive performances.

Although our ULD-Net can generalize representations across domains and achieve competitive results on real-world scene understanding tasks, there still exists a domain gap for transferring from synthetic to scene-level data due to the large point numbers and complicated structures. In the future, we will further explore how to extend our method to domain adaptive analysis of point clouds with the domain gap bridged. We hope the dense similarity learning, the feature separability constraint, and the equiv-crop module proposed could provide insights into future context-based discriminative URL methods.

**Funding.** National Key Research and Development Program of China (2019YFC1521102, 2019YFC1521103); Key Research and Development Program of Shaanxi Province (2019GY215, 2021ZDLSF06-04, 2019ZDLGY10-01); Major Research and Development Project of Qinghai (2020-SF-143, 2020-SF-140); China Postdoctoral Science Foundation (2018M643719).

**Disclosures.** The authors declare no conflicts of interest.

**Data availability.** Data underlying the results presented in this paper are available in Refs. [31–33].

## REFERENCES

1. A. Oliver, S. Kang, B. C. Wünsche, and B. MacDonald, "Using the Kinect as a navigation sensor for mobile robotics," in *Proceedings of the 27th Conference on Image and Vision Computing New Zealand* (2012), pp. 509–514.
2. X. Lu, X. Mao, H. Liu, X. Meng, and L. Rai, "Event camera point cloud feature analysis and shadow removal for road traffic sensing," *IEEE Sens. J.* **22**, 3358–3369 (2022).
3. C. Rausch, M. Nahangi, C. Haas, and J. West, "Kinematics chain based dimensional variation analysis of construction assemblies using building information models and 3D point clouds," *Autom. Constr.* **75**, 33–44 (2017).
4. Y. Yang, C. Feng, Y. Shen, and D. Tian, "Foldingnet: point cloud auto-encoder via deep grid deformation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018), pp. 206–215.
5. Z. Han, X. Wang, Y.-S. Liu, and M. Zwicker, "Multi-angle point cloud-VAE: unsupervised feature learning for 3D point clouds from multiple angles by joint self-reconstruction and half-to-half prediction," in *IEEE/CVF International Conference on Computer Vision (ICCV)* (IEEE, 2019), pp. 10441–10450.
6. P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, "Learning representations and generative models for 3d point clouds," in *International Conference on Machine Learning* (PMLR, 2018), pp. 40–49.
7. R. Li, X. Li, C.-W. Fu, D. Cohen-Or, and P.-A. Heng, "PU-GAN: a point cloud upsampling adversarial network," in *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2019), pp. 7203–7212.
8. H. Wang, Q. Liu, X. Yue, J. Lasenby, and M. J. Kusner, "Unsupervised point cloud pre-training via occlusion completion," in *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 9782–9792.
9. S. Xie, J. Gu, D. Guo, C. R. Qi, L. J. Guibas, and O. Litany, "Pointcontrast: unsupervised pre-training for 3d point cloud understanding," in *European Conference on Computer Vision* (Springer, 2020), pp. 574–591.
10. S. Huang, Y. Xie, S.-C. Zhu, and Y. Zhu, "Spatio-temporal self-supervised representation learning for 3d point clouds," in *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 6535–6545.
11. J. Sauder and B. Sievers, "Self-supervised deep learning on point clouds by reconstructing space," in *Advances in Neural Information Processing Systems* (2019), Vol. **32**.
12. O. Poursaeed, T. Jiang, H. Qiao, N. Xu, and V. G. Kim, "Self-supervised learning of point clouds via orientation estimation," in *International Conference on 3D Vision (3DV)* (IEEE, 2020), pp. 1018–1028.
13. A. Thabet, H. Alwassel, and B. Ghanem, "Self-supervised learning of local features in 3d point clouds," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops* (2020), pp. 938–939.
14. C. Sun, Z. Zheng, X. Wang, M. Xu, and Y. Yang, "Point cloud pre-training by mixing and disentangling," arXiv:2109.00452 (2021).
15. A. Sanghi, "Info3d: representation learning on 3d objects using mutual information maximization and contrastive learning," in *European Conference on Computer Vision* (Springer, 2020), pp. 626–642.
16. B. A. Du, X. Gao, W. Hu, and X. Li, "Self-contrastive learning with hard negative sampling for self-supervised point cloud learning," in *Proceedings of the 29th ACM International Conference on Multimedia* (2021), pp. 3133–3142.
17. J.-B. Grill, F. Strub, F. Altché, C. Tallec, P. H. Richemond, E. Buchatskaya, C. Doersch, B. A. Pires, Z. D. Guo, M. G. Azar, B. Piot, K. Kavukcuoglu, R. Munos, and M. Valko, "Bootstrap your own latent—a new approach to self-supervised learning," in *Advances in Neural Information Processing Systems* (2020), Vol. **33**, pp. 21271–21284.
18. X. Chen and K. He, "Exploring simple siamese representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 15750–15758.

19. C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017), pp. 652–660.
20. C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: deep hierarchical feature learning on point sets in a metric space," in *Advances in Neural Information Processing Systems* (2017), Vol. 31.
21. Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: convolution on x-transformed points," in *Advances in Neural Information Processing Systems* (2018), Vol. 31.
22. Y. Wang, Y. Sun, Z. Liu, E. S. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph CNN for learning on point clouds," *ACM Trans. Graph.* **38**, 1–12 (2019).
23. G. Te, W. Hu, A. Zheng, and Z. Guo, "RGCNN: regularized graph CNN for point cloud segmentation," in *Proceedings of the 26th ACM International Conference on Multimedia* (2018), pp. 746–754.
24. R. Klokov and V. Lempitsky, "Escape from cells: deep kd-networks for the recognition of 3d point cloud models," in *Proceedings of the IEEE International Conference on Computer Vision* (2017), pp. 863–872.
25. H. Zhao, L. Jiang, J. Jia, P. Torr, and V. Koltun, "Point transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2021), pp. 16259–16268.
26. M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, "Pct: Point cloud transformer," *Comput. Vis. Media* **7**, 187–199 (2021).
27. T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International Conference on Machine Learning* (PMLR, 2020), pp. 1597–1607.
28. K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020), pp. 9729–9738.
29. M. Caron, I. Misra, J. Mairal, P. Goyal, P. Bojanowski, and A. Joulin, "Unsupervised learning of visual features by contrasting cluster assignments," in *Advances in Neural Information Processing Systems* (2020), Vol. 33, pp. 9912–9924.
30. A. V. D. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," arXiv:1807.03748 (2018).
31. Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3D ShapeNets: A Deep Representation for Volumetric Shapes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 1912–1920.
32. L. Yi, V. G. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas, "A scalable active framework for region annotation in 3D shape collections," *ACM Trans. Graph.* **35**, 1–12 (2016).
33. I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, "3D semantic parsing of large-scale indoor spaces," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 1534–1543.
34. A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "ShapeNet: an information-rich 3D model repository," arXiv:1512.03012 (2015).
35. L. van der Maaten and G. Hinton, "Visualizing data using t-SNE," *J. Mach. Learn. Res.* **9**, 2579–2605 (2008).
36. Y. Liu, B. Fan, S. Xiang, and C. Pan, "Relation-shape convolutional neural network for point cloud analysis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2019), pp. 8895–8904.