

Self-Supervised Learning with Multi-View Rendering for 3D Point Cloud Analysis

Bach Tran¹, Binh-Son Hua¹, Anh Tuan Tran¹, and Minh Hoai^{1,2}

¹ VinAI Research, Hanoi, Vietnam

² Stony Brook University, New York, 11794, USA

{v.bachtx12,v.sonhb,v.anhtt152,v.hoainm}@vinai.io

Abstract. Recently, great progress has been made in 3D deep learning with the emergence of deep neural networks specifically designed for 3D point clouds. These networks are often trained from scratch or from pre-trained models learned purely from point cloud data. Inspired by the success of deep learning in the image domain, we devise a novel pre-training technique for better model initialization by utilizing the multi-view rendering of the 3D data. Our pre-training is self-supervised by a local pixel/point level correspondence loss computed from perspective projection and a global image/point cloud level loss based on knowledge distillation, thus effectively improving upon popular point cloud networks, including PointNet, DGCNN and SR-UNet. These improved models outperform existing state-of-the-art methods on various datasets and downstream tasks. We also analyze the benefits of synthetic and real data for pre-training, and observe that pre-training on synthetic data is also useful for high-level downstream tasks. Code and pre-trained models are available at https://github.com/VinAIRResearch/selfsup_pcd.git.

Keywords: Self-supervised learning · point cloud analysis · multiple-view rendering · 3D deep learning.

1 Introduction

Pixels and points are basic elements in computer vision for visual recognition. In the past decade, image collections have been successfully used to train neural networks for common visual recognition tasks, including object classification and semantic segmentation. Concurrently, advances in depth-sensing technologies, including RGB-D and LiDAR sensors, have enabled the acquisition of large-scale 3D data, facilitating the rapid development of visual recognition methods in 3D, notably neural networks for point cloud analysis in the last few years. Unlike images, annotation for point clouds are generally more expensive to acquire due to the laborious process of scene scanning, reconstruction, and annotation, and thus point cloud neural networks are often trained with datasets that are much smaller than image datasets.

A potential direction to improve the robustness for point cloud neural networks is self-supervised learning. By letting the point cloud network perform some pre-text tasks before supervised learning, a process commonly known as pre-training,

the network can perform more effectively than that trained from scratch. With self-supervised learning, the pre-text tasks are designed so that the pre-training does not use additional labels. In 3D deep learning, some initial effort has been spent on investigating this direction [46,54,58]. However, most previous works solely considered self-supervised learning in the 3D domain; only a few works exploited images to support the learning of point cloud neural networks. In an early work, Pham et al. [40] trained autoencoders on both images and point clouds and applied constraints on the latent space of both domains, allowing feature transfers between 2D and 3D. Inspired by the recently growing literature on network pre-training, we explore how to use images to more effectively (self-)supervise point cloud neural networks.

Particularly, in this paper, we propose a method that utilizes multi-view rendering to generate pixel/point and image/point cloud pairs for self-supervising a point cloud neural network. We train two neural networks, one for image and one for point cloud, respectively, and direct both networks to agree upon their latent features in the 2D and 3D domains. To achieve this, we use the pixel and point correspondences to formulate a local loss function that encourages features of the correspondences to be similar. This is well-motivated by projective geometry in 3D computer vision that defines the coordinate mapping between the image and 3D space. To further regularize the self-supervision, we utilize knowledge distillation to formulate a global loss that encourages the feature distribution between images and point clouds to be similar as well. Our method works even when there is big domain gap between the pre-train data and test data, e.g., between synthetic and real data.

In summary, we make three technical contributions in this paper: (1) a pre-training technique built upon multi-view rendering that advocates the use of multi-view image features to self-supervise the point cloud neural network; (2) a local loss function that exploits pixel-point correspondence in the pre-training; (3) a global loss function that performs knowledge distillation from the images to the 3D point clouds.

2 Related work

3D deep learning: Building a neural network to analyze 3D data is a non-trivial task. Early attempts involve extending neural networks for images to work with volumes [37], and projecting 3D data to 2D views that can be used with traditional neural networks [51]. Recent methods in deep learning with point clouds take a different approach by letting a network input point clouds directly. Two major directions can be taken to implement such idea: learning per-point features by pointwise MLP [43,44], and learning point features by examining points in a local neighborhood by custom convolutions tailored for point sets [32,26,30] and by graph neural networks [55,48,6]. Notable methods in such directions include PointNet [43] and DGCNN [55]. An inherent limitation of PointNet-based approaches is that they can only process a few thousands of points, limiting the ability to handle large-scale point clouds, where a sliding

window is often used as a workaround [43]. More recent developments include the use of sparse tensor and sparse convolution [9,13,14] on large-scale point clouds for semantic segmentation and 3D object detection. We refer readers to [16] for a survey of deep learning methods for point clouds.

Self-supervised learning: Unsupervised pre-training is a useful technique in deep learning with proven success in natural language processing [11] and representation learning [7,15,21,38,61,47,52,20]. For pre-training, one can use generative modeling techniques based on GANs [56,2] and autoencoders [20,54,19,60], or other self-supervised learning techniques [58,46,8,41,24,59,4,1]. Pre-training is relevant to knowledge distillation [23], a class of techniques for transferring features learned from a large network (teacher network) to a small network (student network). Here we assume that the pre-text task is rather general and can be very different to the downstream tasks, and so a subset of the layers in the pre-trained can be transferred depending on the downstream task.

Self-supervised learning techniques for pre-training 3D point cloud networks have been recently explored from multiple perspectives. Early works use a pre-text task for self-supervised learning. The pre-text task can be solving a jigsaw puzzle [46], where a point cloud is subdivided into randomly arranged voxels, and the task is to predict for each point the voxel ID the point belongs to. Another pre-text task is point cloud completion [54] (OcCo), where a mechanism similar to mask-based pre-training in natural language processing is utilized. As an extension of autoencoder on 3D point clouds, Eckart et al. [12] apply soft segmentation on point clouds and enforces these partitions to comply a latent parametric model in an encoder-decoder network paradigm. Recent contrastive learning is also shown to be effective for pre-training 3D point clouds [58,62]. PointContrast [58] create positive pairs and negative pairs for contrastive learning by the correspondences between two camera views of a point cloud. DepthContrast [62] learn the representation with multiple 3D data formats including points and voxels.

Self-supervised learning with other 3D data modalities [17,34,18,25,33,3,40] has also been explored. Jing et al. [28,29] (CM) use 3D data with multi-modality and build cross-modal and cross-view invariant constraints, maximizing cross-modal agreement of the features of point cloud, mesh, and images, and maximizing cross-view agreement with the image features. Hou et al. [25] use contrastive learning on multi-view images constraints and image-geometry constraint. However, they only focus on 2D downstream tasks. Huang et al. [27] (STRL) proposed self-supervised learning for a sequence of point clouds which utilizes spatio-temporal cues. Pham et al. [40] (LCD) leverages a 2D-3D correspondence dataset and a triplet loss to transfer features from 2D to 3D only on *small cropped regions* of images and 3D point clouds. Compared with LCD [40], our method is largely different as we self-supervise 3D point features via multi-view projection in the *entire* image space. LCD [40] is suitable for image matching and point cloud registration tasks, while our method is suitable for point cloud analysis tasks such as classification and segmentation.

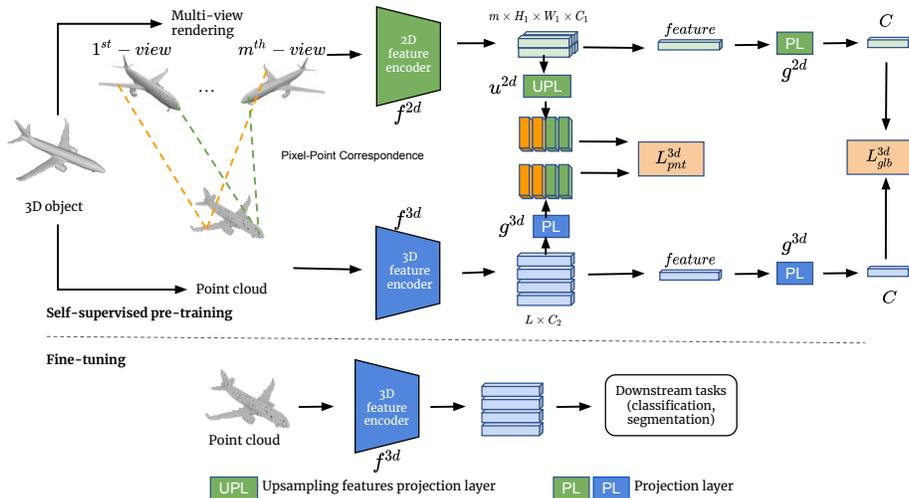


Fig. 1: Overview of our proposed method. The main proposal is pre-training steps that exploit multi-modal data, including multi-view images and point clouds, to learn a 3D feature encoder for effective point cloud representation. This model is then fine-tuned for different downstream tasks.

There are a few concurrent works [3,31]. In [3], the authors considered RGB rendering of the object surfaces but required the mesh textures in addition to the geometry for rendering. Our rendering is more practical in that we consider different rendering techniques and only require colorless point clouds. In [31], the authors focus on data from autonomous driving datasets including KITTI and nuScenes. Our method focuses on object datasets.

3 Self-supervised learning for 3D point clouds

In this section, we describe the proposed self-supervised learning with multi-view rendering for point clouds. Our goal is to leverage multi-modal data of 3D objects, in which each object is associated with a 3D point cloud and multiple 2D images from various view points to pre-train the point cloud network. We propose to use multi-view rendering to generate images for input 3D objects that pair with the point clouds for pre-training. Using rendered images to pre-train point cloud networks implies an advantage that different 3D data representations, including triangle mesh and point cloud, can all be converted into a unified 2D representation to pre-train the networks. To ease this pre-training process, we do not require annotation for the 3D objects and rely on self-supervised learning techniques for the pre-training.

Our method consists of two steps. First, we learn feature representation for 2D images with self-supervised learning by ensuring the similarity between the representation features of two transformed images of the same original image.

Second, we use the feature representation of 2D images to learn the 3D feature representation for 3D point clouds. We illustrate the overview of our method in Fig. 1, and we will describe the two steps in details in the next two subsections.

Let \mathcal{D} denote the pre-training data, $\mathcal{D} = \{P_i, \{\mathbf{y}_{ij}, M_{ij}\}_{j=1}^m\}_{i=1}^n$, where n is the number of objects in the dataset and m the number of 2D views for each object. P_i is the 3D point cloud of the i^{th} object, \mathbf{y}_{ij} is the projected image of the i^{th} object to the j^{th} view using the projection matrix M_{ij} .

3.1 Learning feature representation for 2D images

Let us start with learning the discriminative feature representation for multi-view images. In this step, we simply treat all object views $\{\{\mathbf{y}_{ij}\}_{j=1}^m\}_{i=1}^n$ as items of a set. Following SimCLR [7], we randomly sample a batch of k images from this image set in each training iteration. For each image in the batch, we randomly sample two augmentation operators (crop, color distortion, and Gaussian blur) to generate a pair of correlated images based on the original image. Given an image \mathbf{x}_i in the batch, let \mathbf{x}'_i and \mathbf{x}''_i be its augmented images, respectively. Our objective is to learn a feature extraction network f^{2d} so that the resulting feature vectors for different augmentations of the same image \mathbf{x}'_i and \mathbf{x}''_i are similar, while both \mathbf{x}'_i and \mathbf{x}''_i should be different from the feature vectors of other image augmentations \mathbf{x}'_j and \mathbf{x}''_j for $j \neq i$. We therefore define the loss for image \mathbf{x}_i as:

$$\begin{aligned} \mathcal{L}^{2d}(i) = & -\log \left(\frac{\psi(\mathbf{x}'_i, \mathbf{x}''_i)}{\psi(\mathbf{x}'_i, \mathbf{x}''_i) + \sum_{j \neq i} (\psi(\mathbf{x}'_i, \mathbf{x}'_j) + \psi(\mathbf{x}'_i, \mathbf{x}''_j))} \right) \\ & -\log \left(\frac{\psi(\mathbf{x}''_i, \mathbf{x}'_i)}{\psi(\mathbf{x}''_i, \mathbf{x}'_i) + \sum_{j \neq i} (\psi(\mathbf{x}''_i, \mathbf{x}'_j) + \psi(\mathbf{x}''_i, \mathbf{x}''_j))} \right). \end{aligned} \quad (1)$$

Here, $\psi(\mathbf{x}_i, \mathbf{x}_j)$ is the function that measures the similarity between two images, and we use the exponential cosine similarity of the two feature vectors, i.e.,

$$\psi(\mathbf{x}_i, \mathbf{x}_j) = \exp \left(\cos (g^{2d}(f^{2d}(\mathbf{x}_i)), g^{2d}(f^{2d}(\mathbf{x}_j))) / \tau \right), \quad (2)$$

where τ is the temperature hyper-parameter, and g^{2d} is the projection layer (nonlinear projection layer).

The loss function for a batch of k images is: $\mathcal{L}^{2d} = \frac{1}{k} \sum_{i=1}^k \mathcal{L}^{2d}(i)$. In each optimization iteration, we calculate the gradient of this loss to optimize for the parameters of the feature extractor network f^{2d} , which is a fully convolutional neural network. The input to the network is an RGB image of dimensions $H \times W \times 3$ and the output is a 3D tensor of size $H_1 \times W_1 \times C_1$, where C_1 is the number of output channels and $H_1 = H/2^s$, $W_1 = W/2^s$, with s being the number of down-sampling layers in the network. The output tensor can be vectorized to form a global representation vector for the entire image. This output tensor can also be up-sampled to yield a feature map having the same width and height as those of the input image; in this case, there is a corresponding C_1 -dim feature vector for each pixel of the input image.

3.2 Knowledge transfer from 2D to 3D

Once the feature extraction function f^{2d} for 2D images has been learned, we will use it to learn a point-wise feature extraction function f^{3d} for 3D point clouds. The input to this feature extraction is a point cloud of L points, and the output is a 2D tensor of size $L \times C_2$, where C_2 is the number of feature dimensions. Each point of the point cloud has a corresponding C_2 -dimensional feature vector. To learn the feature extraction f^{3d} , we develop a novel scheme that uses 2D-to-3D knowledge transfer. We use both global and point-wise knowledge transfer.

Global knowledge transfer. For global knowledge transfer, we minimize the distance between the aggregated 2D feature vector and the aggregated 3D feature vector by

$$\mathcal{L}_{glb}^{3d} = \frac{1}{n} \sum_{i=1}^n \left\| g^{2d}(\max_j f^{2d}(\mathbf{y}_{ij})) - g^{3d}(\max f^{3d}(P_i)) \right\|^2,$$

where P_i is the point cloud of the i^{th} object and \mathbf{y}_{ij} is the j^{th} view of the i^{th} object. f^{2d} is the feature extractor for 2D images, which was explained in Sec. 3.1. Function $\max_j f^{2d}(\mathbf{y}_{ij})$ is the pixel-wise max-pooling across different 2D views. Function f^{3d} is the feature extractor for 3D point cloud, which we seek to learn here. $\max f^{3d}(P_i)$ is element-wise max-pooling among all feature vectors of all points of point cloud P_i . Both g^{2d} and g^{3d} are nonlinear projection layers that transform 2D feature and 3D feature vectors to the feature space, respectively.

Point-wise knowledge transfer: In addition to global knowledge transfer, we use contrastive learning that minimizes the distance between feature representation of a 3D point and its corresponding 2D pixel. To determine the point-to-pixel correspondences, we project each point of the point cloud P_i to each image view \mathbf{y}_{ij} using the camera projection matrix M_{ij} to have $\mathbf{y}_{ij}^{2d} = M_{ij}P_i$, where \mathbf{y}_{ij}^{2d} is a set of pixels from the rendered image \mathbf{y}_{ij} corresponding to P_i . For point-wise knowledge transfer, in each optimization iteration, we sample a batch of k corresponding pixel-point pairs, and let $\{(\mathbf{z}_i^{2d}, \mathbf{z}_i^{3d})\}_{i=1}^k$ be the corresponding set of feature vector pairs. For the i^{th} pixel-point pair, \mathbf{z}_i^{2d} is obtained by: (1) using the 2D feature extraction function f^{2d} on the image that contains the pixel; (2) passing the output to the upsampling feature projection module u^{2d} ; and (3) extracting the feature vector at the location of the pixel in the projected feature map. \mathbf{z}_i^{3d} is obtained by: (1) using the 3D feature extraction function f^{3d} on the point cloud containing that point; (2) passing the output through the projection function g^{3d} ; and (3) extracting the corresponding feature vector of the point in the point cloud.

For point-wise knowledge transfer, we maximize the similarity between the pixel representation vector and the point representation vector, using a loss

function inspired by SimCLR [7]. The loss term for the i^{th} pixel-point pair is:

$$\begin{aligned} \mathcal{L}_{pnt}^{3d}(i) = & -\log\left(\frac{\psi(\mathbf{z}_i^{2d}, \mathbf{z}_i^{3d})}{\psi(\mathbf{z}_i^{2d}, \mathbf{z}_i^{3d}) + \sum_{j \neq i} (\psi(\mathbf{z}_i^{2d}, \mathbf{z}_j^{2d}) + \psi(\mathbf{z}_i^{2d}, \mathbf{z}_j^{3d}))}\right) \\ & -\log\left(\frac{\psi(\mathbf{z}_i^{3d}, \mathbf{z}_i^{2d})}{\psi(\mathbf{z}_i^{3d}, \mathbf{z}_i^{2d}) + \sum_{j \neq i} (\psi(\mathbf{z}_i^{3d}, \mathbf{z}_j^{2d}) + \psi(\mathbf{z}_i^{3d}, \mathbf{z}_j^{3d}))}\right), \end{aligned} \quad (3)$$

where $\psi(\cdot, \cdot)$ is the exponential cosine function defined in Eq. (2). Intuitively, both 2D and 3D features can be regarded as augmentations of a common latent feature, so they form a positive pair of which similarity can be maximized with the contrastive loss. The total loss function for a batch of k pixel-point pairs is: $\mathcal{L}_{pnt}^{3d} = \frac{1}{k} \sum_{i=1}^k \mathcal{L}_{pnt}^{3d}(i)$.

Combined loss function. To pre-train the point cloud network, we minimize a loss function that is the weighted combination of the global knowledge transfer loss and the point-wise knowledge transfer loss:

$$\mathcal{L}^{3d} = \lambda_{glb} \mathcal{L}_{glb}^{3d} + \lambda_{pnt} \mathcal{L}_{pnt}^{3d}. \quad (4)$$

In our experiments, we simply use $\lambda_{glb} = \lambda_{pnt} = 1$. After pre-training, we can now use the pre-trained weights to initialize the training of downstream tasks.

4 Experiments

4.1 Implementation details

Dataset for pre-training. Unless otherwise mentioned, we use ModelNet40 [57] for pre-training. ModelNet40 is a synthetic dataset that includes 9,480 training samples and 2,468 test samples in 40 categories. ModelNet40 represents each object by a 3D mesh, making it suitable for our multi-view rendering purpose. For each object in the training set of ModelNet40, we generate its point cloud using farthest-point sampling. We render the object into multi-view images by moving a camera around the object. Unless otherwise mentioned, each point cloud has 1024 points rendered into 12 views with 32×32 resolution. We use 12 views as they tend to cover all object directions in general. We keep the views in low resolution of 32×32 to avoid out of memory usage at training.

Our multi-view rendering is implemented as follows. First, each object is normalized into a unit cube. To generate multi-view images from a mesh object, we used Blender [45] with fixed camera parameters (focal length 35, sensor width 32, and sensor height 32) and a light source. The camera positions are chosen to cover the surrounding views of the object, and the distances from each camera to its neighbor positions are equal.

2D feature representation learning. We use ResNet50 [22] as a 2D feature extractor f^{2d} in 2D self-supervised learning process. We use Adam optimizer

	DGCNN					
	Random	Jigsaw	OcCo	CM	STRL	Ours
MN40 [57]	92.7±0.1	92.9±0.1	92.9±0.0	93.0±0.1	93.1±0.1	93.2±0.1
SO [53]	82.8±0.5	82.1±0.2	83.2±0.2	83.0±0.2	83.2±0.2	84.3±0.6
SO BG [53]	81.4±0.5	82.0±0.4	82.9±0.4	82.2±0.2	83.2±0.2	84.5±0.6

Table 1: Comparison among random, Jigsaw [46], OcCo [54], CM [29], STRL [27], and our initialization to the object classification downstream task. We reported the mean and standard deviation at the best epoch over three runs.

with the initial learning rate 0.001 without learning decay. We then train the self-supervised model with 1000 epochs and batch size 512.

3D feature representation learning. We experiment with two common backbones PointNet [43] and DGCNN [55], which can be utilized for both classification and segmentation tasks. For PointNet [43], we use Adam optimizer with the initial learning rate 0.001, which decays 0.7 every 20 epochs. The momentum of batch normalization starts as 0.5, then divided by 2 every 20 epochs. For DGCNN [55], we use an SGD optimizer with the initial learning rate 0.1 and momentum 0.9. We use CosineAnnealingLR scheduler [35] for learning rate decay. For both backbones, we train the model with 250 epochs, 200 epochs, and 100 epochs for classification, part segmentation, and semantic segmentation task, respectively, with the same batch size as 32. After getting the pre-trained models, we test their effectiveness in training with different downstream tasks.

4.2 Object classification

We first experiment with object classification for 3D point cloud analysis. Two standard benchmarks are used, namely ModelNet40 [57] and ScanObjectNN [53] dataset. We follow the previous paper [54] to use ModelNet40 in both pre-training and downstream tasks. ScanObjectNN is a real-world dataset that has 15 categories with 2,321 and 581 samples for training and testing, respectively. We use the default variant (OBJ_ONLY, denoted by ScanObjectNN) and the variant with background (OBJ_BG, denoted by ScanObjectNN BG). We follow the experimental setting in the original PointNet [43].

We compare the performance of DGCNN [55] with and without pre-training. The results are shown in Table 13. We also provide comparisons with the PointNet backbone [43] in the supplementary material. As can be seen, models with pre-training outperform their randomly initialized counterparts. We further compare our method to previous point cloud pre-training methods, including Jigsaw [46], OcCo [54], CM [29], and STRL [27]. Particularly, Jigsaw [46] learns to solve jigsaw puzzles as a pretext task for pre-training. OcCo [54] is based on mask-based pre-training from natural language processing to propose a point cloud completion task for pre-training. CM [29] considered self-supervision from cross-modality and cross-view feature learning, which shares some similarity to ours. Our method differs in that we use a contrastive loss to learn 2D features and

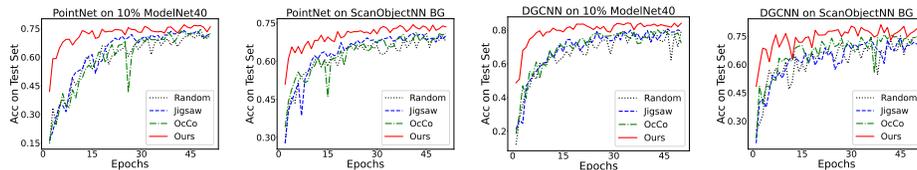


Fig. 2: Test-set accuracy over different training epochs in the object classification task. The plots show that previous pre-training methods are only marginally more effective than random initialization while our method shows a clear improvement.

	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
5%	73.2	73.8	73.9	77.9	82.0	82.1	82.3	84.9
10%	75.2	77.3	75.6	79.0	84.7	84.1	84.9	86.6
20%	81.3	82.9	81.6	84.6	89.4	89.2	89.1	90.2
50%	86.6	86.5	87.1	87.6	91.6	91.8	91.7	92.4
70 %	88.3	88.4	88.4	88.7	92.3	92.4	92.5	92.8
90 %	88.5	88.8	88.8	89.4	92.6	92.9	92.9	93.1

Table 2: Performance of the object classification task trained with fewer data. Our method has significant gains compared to other initialization methods. We reported the mean at the best epoch over three runs.

an L2 loss to match 2D-3D features while CM [29] uses a triplet loss for 2D features and a cross-entropy loss for matching 2D-3D features. STRL [27] explored self-supervision with spatial-temporal representation learning. In Table 13, it can be seen that our proposed self-supervision with contrastive loss and multi-view rendering outperforms other initialization methods on both ModelNet40 and ScanObjectNN dataset.

4.3 Network Analysis

Accuracy over epochs. Fig. 2 plots the accuracy on the test set over different training epochs. The proposed initialization helps both PointNet and DGCNN converge faster and obtain better accuracy than other initialization methods. For example, when we use ModelNet40 with 10% training dataset, the model with our initialization converges after around 15 epochs, while with other initialization methods, it takes around 40 epochs. For ScanObjectNN (OBJ_BG variant), the models converge after about 20 epochs with our initialization and about 45 epochs with other methods.

Training size. Our pre-training allows the point cloud network to be trained with less data compared to initialization with random weights. To demonstrate this, we reduce the number of samples in the training set of ModelNet40 to 5%, 10%, 20%, 50%, and 70%. We then use these datasets to train the model for the object classification task. Finally, we evaluate these learned models on the test set of ModelNet40. Table 16 shows the results with random initialization,

Jigsaw [46], OcCo [54], and our initialization, respectively. As can be seen, models using our proposed initialization outperform other models.

Number of views. We analyze the influence of multi-view rendering in our pre-training performance. We render the shapes with 4, 8, 12, and 24 views in object classification task. The results are shown in Table 3. For PointNet, the performance is best with 8 views, while for DGCNN it is generally enough to use 4 views, but DGCNN for ScanObjectNN performs best with 24 views.

	PointNet				DGCNN			
	4-views	8-views	12-views	24-views	4-views	8-views	12-views	24-views
MN40 [57]	88.9	89.2	88.9	88.9	92.8	92.3	92.5	92.3
SO [53]	79.0	80.4	79.3	79.1	82.7	82.6	82.8	84.9
SO BG [53]	74.2	77.1	75.7	76.6	82.8	81.9	82.6	81.4

Table 3: Performance of object classification tested with different number of views in multi-view rendering.

Classification with SVM. To evaluate the generalization ability of our pre-trained model, we follow a similar test scenario in [46]. First, we freeze the weights of the pre-trained model and pass the 3D object through this model to obtain their embeddings. Then, we train a Support Vector Machine (SVM) on the embeddings of the train set and evaluate it on the test set on three datasets ModelNet40, ScanObjectNN, and ScanObjectNN (OBJ_BG variant). For SVM, we used grid search to find the best parameter of SVM with a Radial Basis Function kernel. The results are shown in Table 4. The proposed initialization outperforms the other initialization methods, Jigsaw and OcCo, sometimes by a wide margin as in the ScanObjectNN (OBJ_BG variant). We also provide additional comparisons to previous self-supervised methods on ModelNet40 in Table 5. As can be seen, our proposed method outperforms almost other methods in both PointNet and DGCNN, except in PointNet, our method is ranked second while ParAE [12] performs best.

	PointNet			DGCNN		
	Jigsaw	OcCo	Ours	Jigsaw	OcCo	Ours
ModelNet40 [57]	82.5	87.2	89.7	83.1	89.5	91.7
ScanObjectNN [53]	49.7	62.1	70.2	57.8	69.0	76.3
ScanObjectNN BG [53]	48.9	61.7	69.5	51.1	67.5	74.2

Table 4: The result of SVM applied on the object embedding learned from different initializations. It shows that features learned by our method are more discriminative than other methods, resulting in more accurate classifications.

Ablation study of loss functions. Table 6 reports the performance of our method for the classification task when both global loss and pixel-point loss are

	PointNet	DGCNN
Rotation [41]	88.6	99.8
STRL [27]	88.3	90.9
ParAE [12]	90.3	91.6
CrossPoint [3]	89.1	91.2
Ours	89.7	91.7

Table 5: More comparisons of SVM classification on ModelNet40.

	PointNet			DGCNN		
	\mathcal{L}_{glb}^{3d}	\mathcal{L}_{pnt}^{3d}	\mathcal{L}^{3d}	\mathcal{L}_{glb}^{3d}	\mathcal{L}_{pnt}^{3d}	\mathcal{L}^{3d}
ModelNet40 [57]	88.5	88.5	88.9	92.4	92.1	92.5
ScanObjectNN [53]	77.6	78.8	79.3	81.8	81.1	82.8
ScanObjectNN BG [53]	74.5	74.2	75.7	81.6	81.6	82.6

Table 6: Effect of loss function choice to our pre-training.

	RGB	Silhouette	Shading
ModelNet40 [57]	88.3	88.9	88.9
ScanObjectNN [53]	79.7	78.8	79.3
ScanObjectNN BG [53]	75.1	75.6	75.7

Table 7: Effect of rendering techniques to the pre-training on PointNet [43].

used together or individually. The network achieves the best performance when trained with both losses. Using either global loss or pixel-point loss results in accuracy drop especially for the ScanObjectNN dataset [53]. This is because the global loss is only useful in distilling knowledge from an image view to a point cloud while the pixel-point loss encourages the model learn consistent features locally. Using both losses lets the model have the best of both worlds.

Multi-view rendering. Our pre-training requires multi-view image rendering, which can be implemented by a wide range of rendering techniques. We study the effect of images rendered from the object mesh, 3D position encoding, and object silhouette on the classification task (please refer to example rendering in Fig. 3). For the original object mesh, we use Blender [45] to render the object geometry with Phong shading, resulting in grayscale *shaded* images. For 3D position encoding, the images are rendered directly from point clouds. Specifically, we first assign a pseudo-color (RGB) to each point of a point cloud based on their 3D coordinates, then project the points to the image plane with preset camera projection matrices. For object silhouette, the process is generally similar except that we use black instead of the pseudo-color for each point in the point clouds. For pixel that has more than one corresponding point, we choose the point with minimum distance to the camera.

The performance of object classification with different rendering techniques is reported in Table 12, where there is no best technique overall. Using shaded images results in slightly higher accuracy than using position encoding and silhouette rendering in ModelNet40 and ScanObjectNN BG. This is because



Fig. 3: Object rendering with position encoding (RGB), silhouette, and shading.

	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
mAcc	93.3 \pm 0.2	93.0 \pm 0.0	93.3 \pm 0.1	93.4 \pm0.0	94.2 \pm 0.0	94.1 \pm 0.0	94.3 \pm0.0	94.2 \pm 0.1
mIoU	83.1 \pm 0.3	83.2 \pm 0.1	83.0 \pm 0.2	83.3 \pm0.1	84.7 \pm0.0	84.5 \pm 0.1	84.7 \pm0.1	84.7 \pm0.1

Table 8: The result of four initialization in the part segmentation task on the ShapeNetPart dataset [55]. We reported the mean and standard error of mAcc and mIoU over three runs.

	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
mAcc	83.9	82.5	83.6	85.0	86.8	86.8	87.0	87.0
mIoU	43.6	43.6	44.5	46.7	49.3	48.2	49.5	49.9

Table 9: Fold 1 of overall point accuracy (mAcc) and mean Intersection-over-Union (mIoU) on the S3DIS (Stanford Area 5 Test) [5].

shaded images often have more details than position encoding and silhouette images since shaded images are rendered from meshes. Exploring more robust rendering techniques for self-supervised learning is left for future work.

4.4 Part segmentation and scene segmentation

Beyond classification, we conduct experiments to validate our pre-training for semantic part segmentation and scene segmentation tasks. We first experiment with object part segmentation on the ShapeNetPart dataset [55] that includes 16,881 objects from 16 categories. Each object is represented by 2,048 points, and each point belongs to one of 50 part types. Most objects in the dataset are divided into two to five parts. As shown in Table 8, our initialization is slightly better than random initialization, Jigsaw, and OcCo in both overall point accuracy (mAcc) and mean Intersection-over-Union (mIoU) metric. We observed that the improvement is minor in the part segmentation task because the network architecture used for part segmentation is largely different from the pre-trained networks, and therefore only a few layers of the part segmentation networks can be initialized by the pre-trained model.

We also experiment with semantic scene segmentation on the Stanford Large-Scale 3D Indoor Spaces dataset (S3DIS) [5]. This dataset contains point clouds of 272 rooms from 6 areas and 13 categories. Each room is split into $1m \times 1m$ blocks. Each point is represented by a 9D vector including XYZ, RGB, and normalized location in the room. During training, each block is sampled with 4096 points, but during testing, all points are used. The results are shown in Table 15. We see that models initialized by our method outperform others in both PointNet [43] and DGCNN [55].

Dataset	Task (Metric)	Random	PC [58] ModelNet	Ours ModelNet	PC [58] ScanNet	Ours ScanNet
S3DIS (Area 5)	sem. seg. (Acc)	72.5	71.2 -1.3	73.2 +0.7	73.0 +0.5	73.0 +0.5
S3DIS (Area 5)	sem. seg. (IoU)	64.5	64.1 -0.4	66.0 +1.5	66.1 +1.6	66.5 +2.0
ScanNet	sem. seg. (Acc)	80.2	80.3 +0.1	81.1 +0.9	80.8 +0.6	81.0 +0.8
ScanNet	sem. seg. (IoU)	72.4	72.5 +0.1	73.3 +0.9	73.1 +0.7	73.6 +1.2
ScanNet	3D det. (AP ₅₀)	35.2	36.6 +1.4	38.2 +3.0	36.1 +0.9	39.2 +4.0
ScanNet	3D det. (AP ₂₅)	56.5	58.2 +1.7	58.4 +1.9	59.5 +3.0	60.3 +3.8
SUN RGB-D	3D det. (AP ₅₀)	32.3	34.8 +2.5	34.9 +2.6	34.8 +2.5	35.1 +2.8
SUN RGB-D	3D det. (AP ₂₅)	55.5	57.8 +2.3	58.1 +2.6	57.4 +1.9	57.8 +2.3

Table 10: Comparison to PointContrast (PC) [58] on the semantic segmentation and 3D object detection task on S3DIS dataset [5], ScanNet dataset [10], and SUN RGB-D dataset [50]. Our method outperforms PointContrast when pre-trained on both datasets. The subscript indicates the performance difference compared to the Random case.

4.5 Pre-training with synthetic vs. real-world data

Multi-view rendering can be easily used for self-supervised learning when working with synthetic data as we have shown with ModelNet40 [57]. Real-world 3D datasets, however, often do not provide such multi-view images, limiting our choices for pre-training. In this section, we investigate the role of synthetic and real-world data in pre-training by comparing to PointContrast [58] and DepthContrast [62] on the segmentation and detection task. We run different experiments using Sparse Residual U-Net (SR-UNet) [9] as the network backbone. Compared to PointNet and DGCNN backbone used in the previous sections, SR-UNet uses sparse convolutions to learn features on point clouds, discarding the need of a sliding window for processing large-scale point clouds. For pre-training data, we use ModelNet40 [57] as synthetic data and ScanNet [10] as real data.

Pre-training. As the original PointContrast [58] only supports ScanNet for pre-training, here we adapt ModelNet40 to PointContrast by using surface point cloud pairs, formed for every two continuous views, instead of the provided point cloud pairs from ScanNet. As for our model, we use two view images when their corresponding point cloud pairs have at least 30% overlapping. To define pixel-point pairs, we reconstruct a point cloud from the first depth image in an image pair, then project it to two color images to get pixel-point correspondences. During training, we follow original setting of PointContrast [58]. For our pre-trained model on ScanNet [10], we used the pre-trained ResNet50 [22] on ImageNet provided by Pytorch[39] as the 2D feature extractor. Besides, all images are resized to 240×320 and we only use the point-wise knowledge transfer loss for pre-training. We train the model with one GPU and four GPUs for ModelNet40 and ScanNet datasets, respectively.

Segmentation and detection results. We evaluate four pre-training configurations with the semantic segmentation task on two datasets S3DIS [5] and

ScanNet [10]. We show the results in Table 10 (comparisons to PointContrast [58]) and Table 11 (comparisons to DepthContrast [62]). In Table 10, on both datasets, the performance gap between our models pre-trained on synthetic and real data is small. When testing on S3DIS, our pre-trained network on ModelNet even provides a slightly better performance compared to the pre-trained model on ScanNet on Acc. metric, and it offers 2% increase when compared with the PointContrast counterpart on both Acc. and IoU metric. In Table 11, we also

Dataset	Task (Metric)	Random	DepthContrast [62]	Ours ModelNet	Ours ScanNet
S3DIS (Area 5)	sem. seg. (Acc)	70.9	72.1 +1.2	75.1 +4.2	74.5 +3.6
S3DIS (Area 5)	sem. seg. (IoU)	64.0	64.8 +0.8	66.8 +2.8	66.5 +2.5
ScanNet	sem. seg. (Acc)	77.2	77.6 +0.4	77.4 +0.2	78.3 +1.1
ScanNet	sem. seg. (IoU)	69.1	69.9 +0.8	69.2 +0.1	70.7 +1.6

Table 11: Comparison to DepthContrast [62] on the semantic segmentation task on S3DIS dataset [5] and ScanNet dataset [10]. The subscript indicates the performance difference compared to the Random case.

compare with DepthContrast on semantic segmentation task. For S3DIS, our pre-trained models on both synthetic and real data achieve better performance approximately 2% on IoU and 4% on Acc. For ScanNet, our pre-trained model on synthetic data outperforms the random setting but is slightly less effective than DepthContrast. However, our pre-trained model on real data outperforms both random and DepthContrast initialization.

We also perform comparison on the 3D object detection task on the ScanNet dataset [10] and SUN RGB-D dataset [50]. Following [58], we replace original PointNet++ [44] backbone of VoteNet [42] by SR-UNet [9]. The results are also shown in Table 10. As can be seen, our method outperforms PointContrast when pre-training on the same dataset. When using synthetic data, our model can obtain two points higher in mAP_{50} compared with the PointContrast counterpart. When using real data, the mAP scores increase slightly and achieve state-of-the-art performance.

5 Conclusion

We propose a self-supervised learning method based on multi-view rendering to pre-train 3D point cloud neural networks. Our self-supervision with multi-view rendering on global and local loss functions yield state-of-the-art performance on several downstream tasks including object classification, semantic segmentation and object detection. Our pre-training method works well on both synthetic and real-world data; it also proves the effectiveness of pre-training on synthetic data like ModelNet40 for downstream tasks with real data like semantic segmentation and 3D object detection.

References

1. Achituve, I., Maron, H., Chechik, G.: Self-supervised learning for domain adaptation on point clouds. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 123–133 (2021)
2. Achlioptas, P., Diamanti, O., Mitliagkas, I., Guibas, L.: Learning representations and generative models for 3d point clouds. In: International conference on machine learning. pp. 40–49. PMLR (2018)
3. Afham, M., Dissanayake, I., Dissanayake, D., Dharmasiri, A., Thilakarathna, K., Rodrigo, R.: Crosspoint: Self-supervised cross-modal contrastive learning for 3d point cloud understanding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9902–9912 (2022)
4. Alliegro, A., Boscaini, D., Tommasi, T.: Joint supervised and self-supervised learning for 3d real world challenges. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 6718–6725. IEEE (2021)
5. Armeni, I., Sener, O., Zamir, A.R., Jiang, H., Brilakis, I., Fischer, M., Savarese, S.: 3d semantic parsing of large-scale indoor spaces. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1534–1543 (2016)
6. Bruna, J., Zaremba, W., Szlam, A., LeCun, Y.: Spectral networks and locally connected networks on graphs. arXiv preprint arXiv:1312.6203 (2013)
7. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International conference on machine learning. pp. 1597–1607. PMLR (2020)
8. Chen, Y., Liu, J., Ni, B., Wang, H., Yang, J., Liu, N., Li, T., Tian, Q.: Shape self-correction for unsupervised point cloud understanding. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8382–8391 (2021)
9. Choy, C., Gwak, J., Savarese, S.: 4d spatio-temporal convnets: Minkowski convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3075–3084 (2019)
10. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: Proc. Computer Vision and Pattern Recognition (CVPR), IEEE (2017)
11. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
12. Eckart, B., Yuan, W., Liu, C., Kautz, J.: Self-supervised learning on 3d point clouds by learning discrete generative models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8248–8257 (2021)
13. Graham, B., Engelcke, M., Van Der Maaten, L.: 3d semantic segmentation with submanifold sparse convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 9224–9232 (2018)
14. Graham, B., van der Maaten, L.: Submanifold sparse convolutional networks. arXiv preprint arXiv:1706.01307 (2017)
15. Grill, J.B., Strub, F., Alché, F., Tallec, C., Richemond, P., Buchatskaya, E., Doersch, C., Avila Pires, B., Guo, Z., Gheshlaghi Azar, M., et al.: Bootstrap your own latent—a new approach to self-supervised learning. *Advances in neural information processing systems* **33**, 21271–21284 (2020)
16. Guo, Y., Wang, H., Hu, Q., Liu, H., Liu, L., Bennamoun, M.: Deep learning for 3d point clouds: A survey. *IEEE transactions on pattern analysis and machine intelligence* (2020)

17. Gupta, S., Hoffman, J., Malik, J.: Cross modal distillation for supervision transfer. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2827–2836 (2016)
18. Hafner, F., Bhuiyan, A., Kooij, J.F., Granger, E.: A cross-modal distillation network for person re-identification in rgb-depth. arXiv preprint arXiv:1810.11641 (2018)
19. Han, Z., Wang, X., Liu, Y.S., Zwicker, M.: 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: 2019 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 10441–10450. IEEE (2019)
20. Hassani, K., Haley, M.: Unsupervised multi-task feature learning on point clouds. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 8160–8171 (2019)
21. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9729–9738 (2020)
22. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
23. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. arXiv preprint arXiv:1503.02531 (2015)
24. Hou, J., Graham, B., Nießner, M., Xie, S.: Exploring data-efficient 3d scene understanding with contrastive scene contexts. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15587–15597 (2021)
25. Hou, J., Xie, S., Graham, B., Dai, A., Nießner, M.: Pri3d: Can 3d priors help 2d representation learning? In: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV). pp. 5693–5702 (October 2021)
26. Hua, B.S., Tran, M.K., Yeung, S.K.: Pointwise convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 984–993 (2018)
27. Huang, S., Xie, Y., Zhu, S.C., Zhu, Y.: Spatio-temporal self-supervised representation learning for 3d point clouds. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6535–6545 (2021)
28. Jing, L., Chen, Y., Zhang, L., He, M., Tian, Y.: Self-supervised modal and view invariant feature learning. arXiv preprint arXiv:2005.14169 (2020)
29. Jing, L., Zhang, L., Tian, Y.: Self-supervised feature learning by cross-modality and cross-view correspondences. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1581–1591 (2021)
30. Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B.: Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems* **31**, 820–830 (2018)
31. Li, Z., Chen, Z., Li, A., Fang, L., Jiang, Q., Liu, X., Jiang, J., Zhou, B., Zhao, H.: Simipu: Simple 2d image and 3d point cloud unsupervised pre-training for spatial-aware visual representations. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 1500–1508 (2022)
32. Liu, Y., Fan, B., Xiang, S., Pan, C.: Relation-shape convolutional neural network for point cloud analysis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8895–8904 (2019)
33. Liu, Y.C., Huang, Y.K., Chiang, H.Y., Su, H.T., Liu, Z.Y., Chen, C.T., Tseng, C.Y., Hsu, W.H.: Learning from 2d: Pixel-to-point knowledge transfer for 3d pretraining. arXiv preprint arXiv:2104.04687 (2021)

34. Liu, Z., Qi, X., Fu, C.W.: 3d-to-2d distillation for indoor scene parsing. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4464–4474 (2021)
35. Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983 (2016)
36. Manning, C.D., Raghavan, P., Schütze, H.: Introduction to Information Retrieval. Cambridge University Press (2008), <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>
37. Maturana, D., Scherer, S.: Voxnet: A 3d convolutional neural network for real-time object recognition. In: 2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 922–928. IEEE (2015)
38. Misra, I., Maaten, L.v.d.: Self-supervised learning of pretext-invariant representations. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6707–6717 (2020)
39. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style, high-performance deep learning library. In: Advances in Neural Information Processing Systems (2019)
40. Pham, Q.H., Uy, M.A., Hua, B.S., Nguyen, D.T., Roig, G., Yeung, S.K.: Lcd: Learned cross-domain descriptors for 2d-3d matching. In: Proceedings of the AAAI Conference on Artificial Intelligence (2020)
41. Poursaeed, O., Jiang, T., Qiao, H., Xu, N., Kim, V.G.: Self-supervised learning of point clouds via orientation estimation. In: 2020 International Conference on 3D Vision (3DV). pp. 1018–1028. IEEE (2020)
42. Qi, C.R., Litany, O., He, K., Guibas, L.J.: Deep hough voting for 3d object detection in point clouds. In: Proceedings of the IEEE International Conference on Computer Vision (2019)
43. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (July 2017)
44. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems* **30** (2017)
45. Roosendaal, T.: Blender - a 3D modelling and rendering package. Blender Foundation (2018), <http://www.blender.org>
46. Sauder, J., Sievers, B.: Self-supervised deep learning on point clouds by reconstructing space. *Advances in Neural Information Processing Systems* **32** (2019)
47. Sharma, C., Kaul, M.: Self-supervised few-shot learning on point clouds. *Advances in Neural Information Processing Systems* **33**, 7212–7221 (2020)
48. Simonovsky, M., Komodakis, N.: Dynamic edge-conditioned filters in convolutional neural networks on graphs. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3693–3702 (2017)
49. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556 (2014)
50. Song, S., Lichtenberg, S.P., Xiao, J.: Sun rgb-d: A rgb-d scene understanding benchmark suite. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 567–576 (2015)
51. Su, H., Maji, S., Kalogerakis, E., Learned-Miller, E.: Multi-view convolutional neural networks for 3d shape recognition. In: Proceedings of the IEEE international conference on computer vision. pp. 945–953 (2015)

52. Thabet, A., Alwassel, H., Ghanem, B.: Self-supervised learning of local features in 3d point clouds. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 938–939 (2020)
53. Uy, M.A., Pham, Q.H., Hua, B.S., Nguyen, D.T., Yeung, S.K.: Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In: International Conference on Computer Vision (ICCV) (2019)
54. Wang, H., Liu, Q., Yue, X., Lasenby, J., Kusner, M.J.: Unsupervised point cloud pre-training via occlusion completion. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9782–9792 (2021)
55. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)* **38**(5), 1–12 (2019)
56. Wu, J., Zhang, C., Xue, T., Freeman, W.T., Tenenbaum, J.B.: Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In: Proceedings of the 30th International Conference on Neural Information Processing Systems. pp. 82–90 (2016)
57. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2015)
58. Xie, S., Gu, J., Guo, D., Qi, C.R., Guibas, L., Litany, O.: Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In: European Conference on Computer Vision. pp. 574–591. Springer (2020)
59. Yamada, R., Kataoka, H., Chiba, N., Domae, Y., Ogata, T.: Point cloud pre-training with natural 3d structures. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 21283–21293 (2022)
60. Yu, X., Tang, L., Rao, Y., Huang, T., Zhou, J., Lu, J.: Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19313–19322 (2022)
61. Zhang, L., Zhu, Z.: Unsupervised feature learning for point cloud understanding by contrasting and clustering using graph convolutional neural networks. In: 2019 International Conference on 3D Vision (3DV). pp. 395–404. IEEE (2019)
62. Zhang, Z., Girdhar, R., Joulin, A., Misra, I.: Self-supervised pretraining of 3d features on any point-cloud. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10252–10263 (2021)

Supplementary Material

In this document, we provide more details for our proposed method. Firstly, we present the complete experiment of using different rendering styles, including position encoding, silhouette, and shading for both PointNet and DGCNN (Section A). We also include additional results for rest of object classification task with PointNet backbone(Section B.1), 6-fold cross validation results for the semantic segmentation task on S3DIS dataset (Section B.3). Secondly, we also report the performance of training with limited data on both ModelNet40 and ScanObjectNN (Section C). Finally, we report detail settings, runtime statistics and more insights into the proposed method by analyzing the t-SNE embedding and the critical

A Evaluation of multi-view rendering

We report the performance of PointNet [43] and DGCNN [55] on different rendering styles in Table 12. It can be seen that shaded images yield slightly higher performance than other renderings on both datasets. However, other rendering styles such as position encoding (RGB) and silhouette still produce competitive results. It implies that in cases where only point clouds are available for pre-training, RGB and silhouette rendering can be used while not causing a significant performance difference compared to mesh-based rendering.

	PointNet			DGCNN		
	RGB	Silhouette	Shading	RGB	Silhouette	Shading
ModelNet40 [57]	88.3 ± 0.2	88.9 ± 0.2	88.9 ± 0.1	92.5 ± 0.2	92.5 ± 0.2	92.5 ± 0.1
ScanObjectNN [53]	79.7 ± 0.5	78.8 ± 0.6	79.3 ± 0.3	82.8 ± 0.5	82.0 ± 0.2	82.8 ± 1.0
ScanObjectNN BG [53]	75.1 ± 0.3	75.6 ± 0.4	75.7 ± 0.5	81.0 ± 0.2	81.8 ± 0.9	82.6 ± 0.7

Table 12: Effect of different rendering techniques to our pre-training

B Details of downstream tasks

B.1 Object classification

Similar to the comparison with the DGCNN backbone in the main paper, we provide comparisons with the PointNet backbone. The results are shown in Table 13. As can be seen, our method outperforms random initialization as well as other pre-training methods, including Jigsaw [46], OcCo [54], and CM [29].

	PointNet				
	Random	Jigsaw	OcCo	CM	Ours
MN40 [57]	88.9±0.0	89.2±0.0	89.2±0.1	89.1±0.1	89.5±0.2
SO [53]	78.2±0.1	79.4±0.1	79.5±0.1	79.3±0.5	80.5±0.4
SO BG [53]	76.4±0.0	76.4±0.4	76.4±0.1	74.1±0.2	78.5±0.5

Table 13: Comparison among random, Jigsaw [46], OcCo [54], CM [29], and our initialization to the object classification downstream task. We reported the mean and standard deviation at the best epoch over three runs.

	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
mAcc	83.9	82.5	83.6	85.0	86.8	86.8	87.0	87.0
mIoU	43.6	43.6	44.5	46.7	49.3	48.2	49.5	49.9

Table 14: Fold 1 of overall point accuracy (mAcc) and mean Intersection-over-Union (mIoU) on the S3DIS (Stanford Area 5 Test) [5]

B.2 A note on the OcCo baseline

It can be seen that in our paper, some experiment results of OcCo are lower than the results reported by its original paper. We did our best to reproduce the results of OcCo but unfortunately, we were not able to match the results with the original paper. We confirmed this issue by using the docker image provided by the OcCo authors and rerun the experiments, but still could not reproduce the results exactly as in the OcCo paper. For fair comparison and reproducibility, we decided to report the results based on our own runs. Additionally, the pre-training time of OcCo is about 7x slower than our method.

B.3 Semantic segmentation

In addition to the Area-5 results reported in the main paper, we further report the results of 6-fold cross-validation over the 6 areas on the S3DIS dataset. For completeness, all results are shown in Table 14 (Area-5), and Table 15 (6 folds). In both cases, we can see that models initialized by our method outperform others in both PointNet [43] and DGCNN [55].

B.4 Details of PointContrast baseline

Semantic segmentation: We evaluate on two datasets S3DIS [5] and ScanNet [10]. We use SGD optimizer with the initial learning rate 0.1 and 0.8 for S3DIS and ScanNet respectively. We use PolynomialLR scheduler with a power factor 0.9. For ScanNet dataset, we train the model with 15000 iterations and batch size

	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
mAcc	82.8	82.8	82.7	83.2	86.9	86.6	87.1	87.5
mIoU	50.6	51.4	51.1	52.1	58.4	58.1	58.7	59.0

Table 15: Average of 6-fold cross validation of overall point accuracy (mAcc) and mean Intersection-over-Union (mIoU) on the S3DIS [5]

48 on 4 GPUs. For S3DIS dataset, we train the model with 20000 iterations and batch size 32 on 1 GPU.

Object detection: For object detection task, in the training we follow the configuration of PointContrast [58]. We use Adam optimizer with the initial learning rate 0.001 and train the model on 1 GPU with 180 epochs. Specifically, we train the model with batch size 32 and 64 for ScanNet and SUN RGB-D, respectively. Before voxelization, we sample 40000 and 20000 points from original point of ScanNet and SUN RGB-D and the voxel sizes are 2.5 cm and 5 cm respectively.

ModelNet40 [57]	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
5%	73.2	73.8	73.9	77.9	82.0	82.1	82.3	84.9
10%	75.2	77.3	75.6	79.0	84.7	84.1	84.9	86.6
20%	81.3	82.9	81.6	84.6	89.4	89.2	89.1	90.2
50%	86.6	86.5	87.1	87.6	91.6	91.8	91.7	92.4
70 %	88.3	88.4	88.4	88.7	92.3	92.4	92.5	92.8
ScanObjectNN [53]	PointNet				DGCNN			
	Random	Jigsaw	OcCo	Ours	Random	Jigsaw	OcCo	Ours
5%	52.1	51.8	53.7	60.8	48.3	46.7	51.4	60.9
10%	63.0	62.3	62.5	69.0	58.7	58.0	61.5	69.5
20%	69.0	68.5	67.1	72.2	69.8	68.7	71.6	74.7
50%	73.7	75.1	72.6	77.0	76.3	77.1	78.0	81.6
80 %	76.1	77.9	76.7	78.4	79.9	78.1	80.8	82.1

Table 16: Performance of the object classification task trained with fewer data on ModelNet40 [57] and ScanObjectNN [53]. Our method has significant gains compared to other initialization methods. We reported the mean at the best epoch over three runs

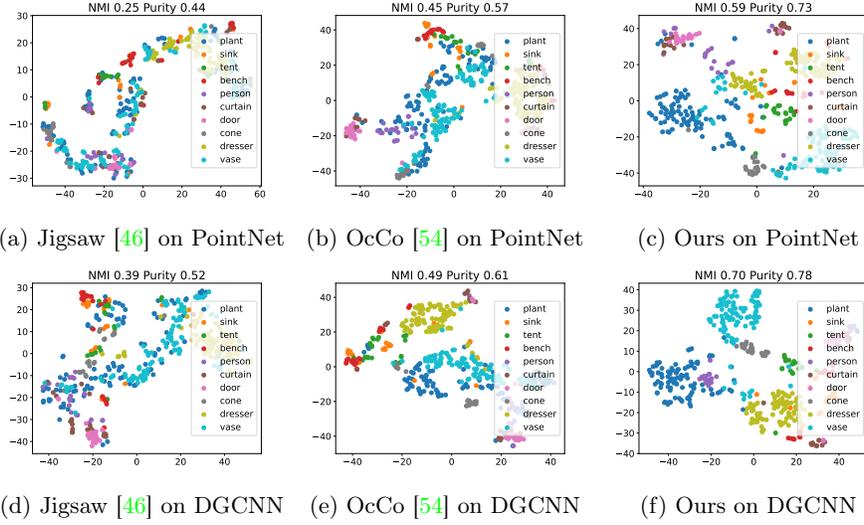


Fig. 4: t-SNE visualization of the object embedding of the test data of ModelNet40. Our method has better cluster quality measured by NMI and purity.

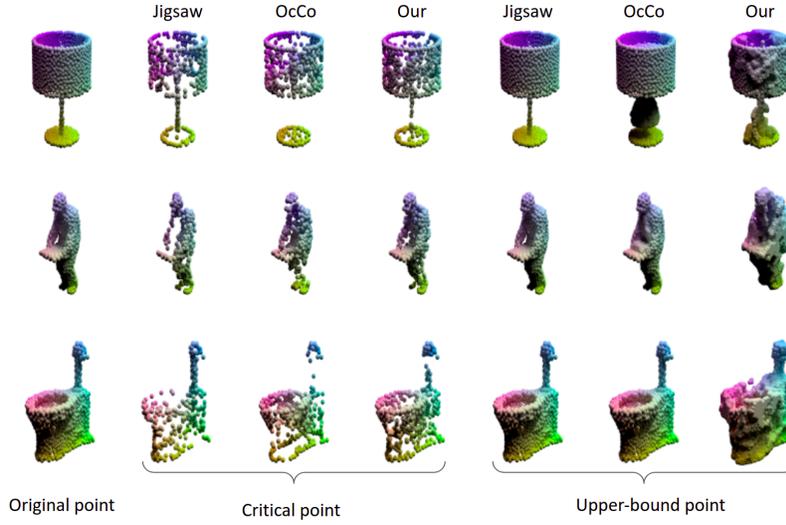


Fig. 5: Critical and upper-bound point visualizations for models pretrained with Jigsaw, OcCo, and our method.

C Training with limited data

To prove the effectiveness of our pre-training, we supervise the downstream task with fewer data when the network is pre-trained and compare to other initializations. We show both results on ModelNet40 (also reported in the paper) and ScanObjectNN. In this experiment, we decrease the number of training samples to 5%, 10%, 20%, 50%, and 80%, and evaluate the model on the original test set. The results are reported in Table 16, which shows that the performance of our method outperforms Random, Jigsaw [46], and OcCo [54] in most cases except DGCNN on 80% of ScanObjectNN.

D Visualization

D.1 t-SNE embedding

We further visualize learned object embeddings of the ModelNet40 test set by using t-SNE with perplexity 15 and 1000 iterations in Figure 4. We observe that the embeddings learned from using our initialization for different classes are well clustered than those acquired with OcCo and Jigsaw initialization indicated by normalized mutual information (NMI) and purity [36].

D.2 Critical point sets

We then visualize the critical point sets and upper-bound shapes by following PointNet [43] for selected samples in Figure 5. To recap, a critical point set is a set of points that contribute directly to the last max pooled feature, i.e., the global feature. Perturbing the critical point set can lead to changes in the global features and thus classification results. The upper-bound shape is the largest possible point set that does not directly affect the global feature of the original shape. From Figure 5, we found that in our method, the critical point sets can represent the object skeleton more faithfully (e.g., the toilet example) than other methods. Jigsaw sometimes causes sparse critical points, and OcCo tends to discard points along thin geometric features. We also found that the upper-bound shape of our initialization appears thicker than that of Jigsaw and OcCo, which we hypothesize that our model can be more robust to point perturbations than Jigsaw and OcCo.

E Running time

Following the request, we provide the pre-training time of three methods on **an NVIDIA Tesla V100 GPU** in Table 17. As can be seen, the pre-training time of our method is slightly longer than Jigsaw and significantly shorter than OcCo. Despite such, our method achieves better performance than the others.

Additionally, we provide more statistics of our training process. Specifically, it takes 2.3, 2.4, and 6.2 hours to render RGB, silhouette, and shaded images,

	Jigsaw	OcCo	Ours
PointNet	2.6	24.8	3.8
DGCNN	4.1	35.1	5.7

Table 17: Pre-training time of three methods (hours).

respectively. For the 2D self-supervision, we train the model for 80 hours on an NVIDIA Tesla V100 GPU. Knowledge distillation takes 3.8, 5.7, 26 and 62 hours of training for PointNet [43], DGCNN [55], SR-UNet on ModelNet40 [57] and SR-UNet on ScanNet [10], respectively. As for downstream-task training, the PointNet classification model takes 18.5 hours, and the DGCNN classification model takes 75.0 hours. The segmentation models require longer training time, with 32.0 hours and 90.0 hours for PointNet and DGCNN backbone, respectively. For SR-UNet backbone [9], in semantic segmentation task, it consumes 32 and 22 hours for S3DIS [5] and ScanNet [10], respectively. For object detection task on ScanNet [10], it takes 8.5 hours.

F Future Work

Our method is not without limitations. First, our image encoder is trained from scratch without leveraging existing popular feature extractors such as VGG [49] or ResNet [22]. Further utilizing such pre-trained networks on natural images could potentially improve the performance of the downstream tasks, which could be interesting for future work. Second, the multi-view rendering used in our method could potentially be further explored. While we attempted with position encoding, silhouette, and shaded rendering, there are many other rendering styles that could be experimented, e.g., rendering with colors and textures when applicable, rendering with depth completion, etc. Applying advanced techniques to enhance multi-view rendering is thus a good avenue for future research.