
Point-GCC : Universal Self-supervised 3D Scene Pre-training via Geometry-Color Contrast

Guofan Fan¹

Zekun Qi¹

Wenkai Shi¹

Kaisheng Ma^{2†}

¹ Xi'an Jiaotong University

² Tsinghua University

Abstract

Geometry and color information provided by the point clouds are both crucial for 3D scene understanding. Two pieces of information characterize the different aspects of point clouds, but existing methods lack an elaborate design for the discrimination and relevance. Hence we explore a 3D self-supervised paradigm that can better utilize the relations of point cloud information. Specifically, we propose a universal 3D scene pre-training framework via **Geometry-Color Contrast** (Point-GCC), which aligns geometry and color information using a Siamese network. To take care of actual application tasks, we design (i) hierarchical supervision with point-level *contrast and reconstruct* and object-level *contrast* based on the novel deep clustering module to close the gap between pre-training and downstream tasks; (ii) architecture-agnostic backbone to adapt for various downstream models. Benefiting from the object-level representation associated with downstream tasks, Point-GCC can directly evaluate model performance and the result demonstrates the effectiveness of our methods. Transfer learning results on a wide range of tasks also show consistent improvements across all datasets. *e.g.*, new state-of-the-art object detection results on SUN RGB-D and S3DIS datasets. Codes will be released at <https://github.com/Asterisci/Point-GCC>.

1 Introduction

3D Self-supervised learning (SSL) has received abundant attention recently because of remarkable improvement on various downstream tasks. 3D scene datasets are tiny compared to the 2D field because 3D point cloud labeling is time-consuming and labor-intensive, which dramatically impedes the improvements of supervised methods. Hence many works [18, 31, 43, 49] explore pre-training models out of 3D labeled data to transfer knowledge for downstream tasks. The goal of self-supervised learning can be summarized as learning rich representations from unlabeled data and helping to improve performance on downstream tasks with labeled data. Most existing works follow the paradigm in the previous 2D field, such as contrastive learning [20, 21, 43, 51] and masked autoencoder (MAE) [28, 31, 47, 50]. After standing on the shoulders of giants in the 2D field, we could further see the particularity of 3D representation learning as follows:

- **Unique information.** 3D scene point cloud contains various information such as geometry and color, which makes 3D point cloud data different from 2D image data. Most existing methods [20, 43, 51] treat all information of each point as an entirety in model architecture design. We argue that directly concatenating all information can not adapt the model to discriminately learn different aspects of point clouds. Although some works [40, 46] propose the two-stream architecture that encodes point cloud by 3D network and images by 2D network, it needs extra 2D data, and 3D network can not clearly learn the discrimination between different information. Considering these additional differences may be beneficial for effective representation learning.

[†]Corresponding Author.

- **Mismatch between pre-training and downstream tasks.** Previous pre-training works [20, 28, 43, 50] design their self-supervised point-level tasks, such as contrast and reconstructing between specific points. However, 3D scene downstream tasks mostly focus the object representations such as object detection and instance segmentation. The gap in supervision level between pre-training and downstream tasks may hinder the improvements of 3D self-supervised learning.
- **Architecture diversity.** The 3D point cloud field has grown rapidly in recent years [13, 25, 30, 32, 36], and the popular architecture appears changeable and specific for downstream tasks. Hence a universal pre-training framework is important that can implement various existing methods for all kinds of tasks and is easy to adapt for future architecture.

To mitigate the aforementioned problems, we explore a 3D self-supervised paradigm that can better utilize the relations of point cloud information. Most 3D scene datasets [1, 16, 38] provide geometry and color information, representing different aspects of the point cloud. Geometry information describes the outline of objects and can easily distinguish between them, while color information refines the internal characteristics of objects and gives a more accurate view of each object. What’s more, different information has inherent relevance. For instance, we can roughly infer the geometric structure of the object from a color photo and vice versa. Motivated by the difference and relevance inherent in the information, we propose a self-supervised 3D scene pre-training framework via **Geometry-Color Contrast (Point-GCC)**, which uses a Siamese network to extract representations and implements elaborate hierarchical supervision. To bridge the gap between pre-training and downstream tasks, the hierarchical supervision contains point-level supervision that aims to align point-wise features and object-level supervision based on a novel deep clustering module to provide better object-level representations strongly associated with downstream tasks. Additionally, the universal Siamese network is designed as an architecture-agnostic backbone so that various downstream models can easily be adapted in a plug-and-play way.

In extensive experiments, we directly perform a fully unsupervised semantic segmentation task without fine-tuning to evaluate the quality of the pre-training model. The result outperforms the previous method with +7.8% mIoU on ScanNetV2, which proves that Point-GCC has learned rich object representations through our paradigm. Furthermore, we choose a broad downstream task to demonstrate our generality: object detection, semantic segmentation and instance segmentation on ScanNetV2 [16], SUN RGB-D [38] and S3DIS [1] datasets. Remarkably, our results indicate general improvements across all tasks and datasets. For example, we achieve new state-of-the-art results with 69.7% AP₂₅, 54.0% AP₅₀ on SUN RGB-D and 75.1% AP₂₅, 56.7% AP₅₀ on S3DIS datasets. Compared with previous pre-training methods, our method achieves higher AP₅₀ by +3.1% on ScanNetV2 and +1.1% on SUN RGB-D. Our contributions can be summarized as follows:

- We propose a new universal self-supervised 3D scene pre-training framework, called Point-GCC, which aligns geometry and color information via a Siamese network with hierarchical supervision. To the best of our knowledge, this is the first study to explore the alignment between geometry and color information of point cloud via the pre-training approach.
- We design a novel deep clustering module to generate object pseudo-labels based on the inherent feature consistency of the two pieces of information. The result demonstrates that Point-GCC has learned rich object representations by clustering.
- Extensive experiments show that Point-GCC is a general pre-training framework with an architecture-agnostic backbone, significantly improving performance on a wide range of downstream tasks and achieving new state-of-the-art on multiple datasets.

2 Related Work

2.1 3D Scene Understanding

Most 3D scene understanding works are still specially designed for downstream tasks, such as object detection [25, 27, 29, 36, 37], semantic segmentation [9, 30, 34, 41], and instance segmentation [11, 22, 23, 39]. The model architecture can be summarized as a backbone module extracting the features of point clouds, and a downstream head adapting for the special task. According to the processing method, these works can be roughly divided into two categories: point-based methods and voxel-based methods. Point-based methods [9, 25, 29] is widely used in point clouds thanks to the effectiveness of PointNet++ [30], which alternately use farthest point algorithm and multi-layer perceptron to sample and extract the features of point. Voxel-based methods [11, 22, 23, 36, 37, 39] is recently popular because of the better performance and efficiency on many downstream tasks than point-

based methods, which operate 3D sparse convolution on regular voxels transformed from irregular point clouds. We pre-train on both point-based PointNet++ and voxel-based 3D sparse convolution backbone and fine-tune on multiple downstream methods to give a comprehensive view of our work.

2.2 3D Self-supervised Learning

Compared to 2D vision or natural language, the 3D vision has a more serious problem of data scarcity [18] which limits the downstream performance of 3D tasks. To solve the raising problem, 3D self-supervised learning (SSL) has gotten more attention in recent years. The mainstream SSL methods can be roughly divided into two categories: contrastive learning and reconstructive learning. Contrastive learning is motivated to learn the invariant representation from different paired carriers such as view augmentation [12, 43] or different data formats [33, 45]. Reconstructive learning is designed to reconstruct the disturbed data to learn geometry knowledge between patches [4, 17]. Motivated by the success of masked autoencoder [19] in 2D, the MAE-style self-supervised method became popular in point cloud [28, 50]. Recently, some works find that the *pattern difference* between the two methods in attention area [44] and scaling performance [31]. Based on previous work, we consider the color and geometry of scene point clouds as two views for contrastive learning, and use a *swapped reconstruct* strategy for reconstructive learning. Therefore, Point-GCC achieves the integration of two methods and derives benefits from both of them.

2.3 Deep Clustering for Self-supervised Learning

Deep Clustering [3, 5, 7, 10, 26, 42, 48] aims to learn better features and discover data labels using deep neural networks, which has been broadly applied in self-supervised and semi-supervised learning. DeepCluster [6] uses the off-the-shelf K-means algorithm pseudo-labels as supervision which learns comparative representations for self-supervised learning. SeLa [2] proposes a simultaneous clustering and representation learning method using the Sinkhorn-Knopp algorithm to generate pseudo-labels with equal partitions quickly. SwAV [8] combines contrastive learning and deep clustering, which enforces consistency between cluster assignments from different views of the same image. In this work, we attempt to apply deep clustering in 3D self-supervised learning field, which generates pseudo-labels based on the inherent feature consistency of the geometry and color information of the point cloud.

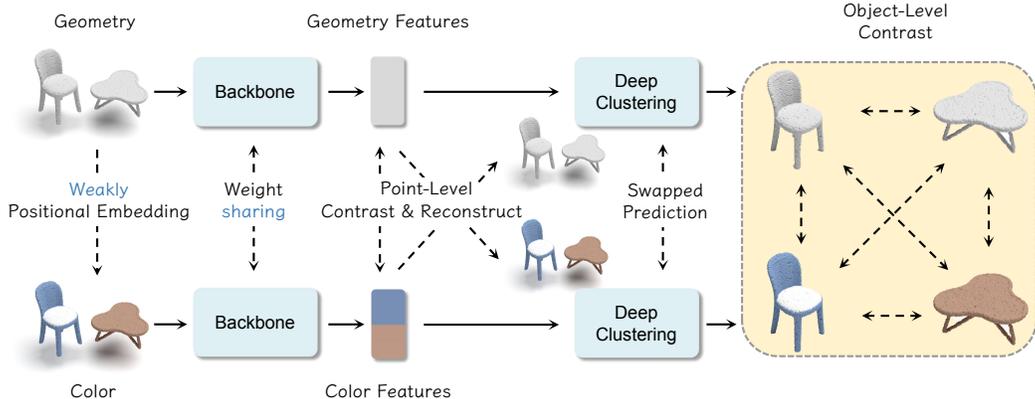


Figure 1: **Overview of our Point-GCC framework.** Point-GCC utilizes the Siamese network to extract the features of geometry and color with positional embedding respectively. Then we implement the hierarchical supervision on extracted features which contains point-level *contrast and reconstruct* and object-level *contrast* based on the deep clustering module.

3 Point-GCC: Pre-training via Geometry-Color Contrast

Existing methods mainly focus on geometric information, but our goal is to enhance the 3D representation capability by better utilizing all the information discriminately in scene point clouds. Therefore, a novel *Geometry-Color Contrast* method is proposed to address this motivation. Figure 1 illustrates the overall framework of Point-GCC. We first perform a Siamese backbone to extract the

features of the geometry and color information respectively in Section 3.1. To carefully align the features belonging to different information, we propose the point-level supervision via combining the contrastive and reconstructive learning in Section 3.2, then we design an unsupervised deep clustering module to generate object pseudo-labels and perform object-level contrastive learning between high-confidence object samples in Section 3.3. The final hierarchical supervision is described in Section 3.4. In Section 3.5, we propose a new method directly evaluating the pre-training model on unsupervised semantic segmentation to demonstrate the effectiveness of our method.

3.1 Siamese Architecture

Information split and embedding. In 3D scene datasets, a point \mathbf{p} is usually associated geometry information represented by the coordinates \mathbf{p}_{geo} and color information represented by RGB value \mathbf{p}_{color} . Different from previous pre-training methods regarding a single point as an atom unit, we split the point cloud into two parts, the geometry and color respectively. Then we project them to universal embedding space e by Equation 1. Additionally, to distinguish similar colors in different coord, we add an extra weakly positional embedding e_{pos} to the color embedding with the Euclidean norm of coord. Note that we remove all embedding modules in fine-tuning stage to keep our framework plug-and-play in order that more existing methods can benefit from ours.

Siamese architecture-agnostic backbone. We use a symmetric Siamese network $\mathcal{F}(\cdot)$ to separately encode geometry features \mathbf{f}_{geo} and color features \mathbf{f}_{color} . Since we attempt to help more existing architectures learn better representations from the combination of geometry and color information, we do not modify any backbone architecture. So that we can directly reuse the core module for standard segmentation with any backbone architecture. In other words, the backbone encodes input $\mathbf{x} \in R^{N \times C_1}$ and extracts feature $\mathbf{y} \in R^{N \times C_2}$. To align the two information, Siamese backbone $\mathcal{F}(\cdot)$ encodes the geometry embedding e_{geo} and color embedding e_{color} with weakly positional embedding e_{pos} to geometry features \mathbf{f}_{geo} and color features \mathbf{f}_{color} respectively:

$$e_{geo} = \mathcal{E}_{geo}(\mathbf{p}_{geo}), \quad e_{color} = \mathcal{E}_{geo}(\mathbf{p}_{color}), \quad e_{pos} = \mathcal{E}_{pos}(\|\mathbf{p}_{geo}\|_2^2), \quad (1)$$

$$\mathbf{f}_{geo} = \mathcal{F}(e_{geo}), \quad \mathbf{f}_{color} = \mathcal{F}(e_{geo} + e_{pos}), \quad (2)$$

where \mathcal{E} is corresponding linear layer of each embedding, $\mathcal{F}(\cdot)$ is the Siamese network.

3.2 Point-level Supervision

Inspired by the success of associating contrastive learning and reconstructive learning in recent work [31], We propose our point-level supervision elaborately designed for our Siamese architecture, which first contrast and then *swapped reconstruct* the features to benefit from different paradigms.

Contrastive learning. The geometry features \mathbf{f}_{geo} and color feature \mathbf{f}_{color} are point-wise aligned because they are split from the same point cloud \mathbf{p} and extracted by the Siamese segmentation-style backbone network. We apply the InfoNCE loss aiming to pull positive pairs close, and push negative pairs away across the geometry features and color features:

$$\mathcal{L}_{pc} = - \sum_i^N \log \frac{\exp(\mathbf{z}_{geo}^{iT} \cdot \mathbf{z}_{color}^i / \tau)}{\sum_j^N \exp(\mathbf{z}_{geo}^{iT} \cdot \mathbf{z}_{color}^j / \tau)}, \quad (3)$$

where τ is the temperature hyper-parameter, we follow the previous works [43] to set it as 0.4. \mathbf{z}_{geo}^i and \mathbf{z}_{color}^i correspond to matched ℓ_2 -normalized feature \mathbf{f}_{geo}^i and \mathbf{f}_{color}^i from same point \mathbf{p}^i , which represent a pair of positive sample. And \mathbf{z}_{geo}^i with other \mathbf{z}_{color}^j except \mathbf{z}_{color}^i represent negative pairs.

Reconstructive learning. Based on our Siamese architecture, we apply the reconstructive learning by *swapped reconstruct* strategy instead of mask strategy, which solves the raising problem about the distribution mismatch between training and testing data in masked autoencoding for point cloud [24]. Specifically, we simply project the geometry features \mathbf{f}_{geo} and color features \mathbf{f}_{color} to reconstruct color $\hat{\mathbf{p}}_{geo}$ and geometry $\hat{\mathbf{p}}_{color}$. The reconstructive loss is the mean squared error (MSE) between the reconstructed and original information of each point:

$$\mathcal{L}_{pr} = \frac{1}{N} \sum \|\mathbf{p}_{geo}^{i'} - \hat{\mathbf{p}}_{geo}^i\|_2^2 + \frac{1}{N} \sum \|\mathbf{p}_{color}^{i'} - \hat{\mathbf{p}}_{color}^i\|_2^2, \quad (4)$$

where N is the number of points, $\hat{\mathbf{p}}_{geo}^i$ and $\hat{\mathbf{p}}_{color}^i$ represent the reconstruct prediction, $\mathbf{p}_{geo}^{i'}$ and $\mathbf{p}_{color}^{i'}$ represent the reconstruct targets which both scale to between 0 and 1 for stability training loss.

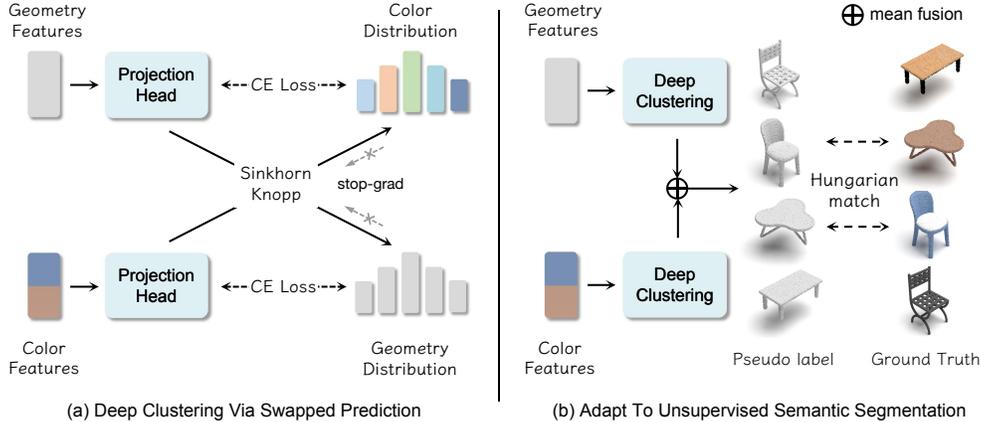


Figure 2: (a) The deep clustering module obtains pseudo prediction for different features and enforces consistent with the swapped partition distribution from the Sinkhorn-Knopp algorithm. (b) Point-GCC generates the pseudo-labels by utilizing cluster prediction from both branches and projects to ground-truth labels for unsupervised semantic segmentation using Hungarian matching alignment.

3.3 Object-level Supervision

Point-level supervision is widely applied in 3D self-supervised learning, which provides rich representations for downstream tasks. However, the object representation strongly associated with downstream tasks hasn't been noticed before. We propose our object-level supervision driven by the novel unsupervised deep clustering module. The clustering module generates pseudo-label predictions \mathcal{P}_{geo} and \mathcal{P}_{color} for the geometry features \mathbf{f}_{geo} and color features \mathbf{f}_{color} respectively, and enforces consistent prediction between geometry prediction \mathcal{P}_{geo} and color prediction \mathcal{P}_{color} of same point \mathbf{p} . We argue that the pseudo-labels represent more various object features, which are not restricted by human annotations with fixed object classes. To achieve robust supervision among these object-level pseudo labels, we sample the high-confidence object features based on the prediction confidence score and apply object-level contrastive learning according to pseudo labels.

Deep clustering via swapped prediction. We apply the swapped prediction [8] in 2D contrastive learning to our model, which predicts the pseudo label of an image from the clustering result of another view. In our framework, we swap the cluster target of different information features, and predict the pseudo label from the other information feature based on the inherent consistency of the two types of information as shown in Figure 2(a). For pseudo label classes K , we use a learnable matrix $\mathcal{C} = [c_1, \dots, c_K]$ to represent the cluster centroids, and calculate the similarity \mathcal{S} between the ℓ_2 -normalized features \mathbf{f} and cluster centroids \mathbf{c} . To avoid the degeneration problem that all features collapse into the same prediction, the Sinkhorn-Knopp algorithm [15] is used to generate the equal partition cluster distribution \mathcal{Q} from the similarity \mathcal{S} by converting pseudo-label generation to an optimal transport problem. And the learnable prediction \mathcal{P} is computed by $\text{softmax}(\mathcal{S}/\tau)$, where τ is the temperature hyper-parameter. We set all hyper-parameter in swapped prediction same to the previous works [8] in 2D. Finally, The swapped prediction loss is the cross entropy losses between the learnable prediction \mathcal{P} and swapped equal partition distribution \mathcal{Q} :

$$\mathcal{L}_{clu} = \ell(\mathcal{Q}_{geo}, \mathcal{P}_{color}) + \ell(\mathcal{Q}_{color}, \mathcal{P}_{geo}), \quad (5)$$

where ℓ is the cross-entropy loss between the prediction and target.

Object-level contrastive learning. For the features \mathbf{f} with corresponding pseudo prediction \mathcal{P} and confidence score from deep clustering, we pick features with confidence scores higher than the picking threshold to alleviate the noise from unsupervised clustering. Then we compute the mean features of high-confidence samples from geometry and color branches, respectively. We take the two types of mean features with the same pseudo-label as positive pairs, oppositely with different pseudo-label as negative pairs, and apply the InfoNCE loss at object-level:

$$\mathcal{L}_{oc} = - \sum_i^N \log \frac{\exp(\mathbf{z}_{geo}^{iT} \cdot \mathbf{z}_{color}^i / \tau)}{\sum_j^N \exp(\mathbf{z}_{geo}^{iT} \cdot \mathbf{z}_{color}^j / \tau)}, \quad (6)$$

Method	ScanNetV2		SUN RGB-D		S3DIS	
	AP ₂₅	AP ₅₀	AP ₂₅	AP ₅₀	AP ₂₅	AP ₅₀
Supervised Only						
VoteNet [29]	58.6	33.5	57.7	-	-	-
GroupFree-3D [25]	66.3	47.8	-	-	-	-
FCAF3D [36]	71.5	57.3	64.2	48.9	66.7	45.9
TR3D [37]	72.9	59.3	67.1	50.4	74.5	51.7
Self-supervised Pre-training						
VoteNet [29]	58.6	33.5	57.7	-	-	-
+ PointContrast [43]	59.2	38.0	57.5	34.8	-	-
+ DepthContrast [51]	62.1	39.1	60.4	35.4	-	-
+ CSC [20]	-	39.3	-	36.4	-	-
+ Ponder [21]	63.6	41.0	61.0	36.6	-	-
+ Point-GCC*	65.3 (+3.0)	44.1 (+3.3)	61.3 (+1.5)	37.7 (+2.0)	-	-
VoteNet+FF [37]	-	-	64.5	39.2	-	-
+ Point-GCC	-	-	64.9 (+0.4)	41.3 (+2.1)	-	-
GroupFree-3D [25]	66.3	47.8	-	-	-	-
+ Point-GCC	68.1 (+1.8)	49.2 (+1.4)	-	-	-	-
TR3D [37]	72.9	59.3	67.1	50.4	74.5	51.7
+ Point-GCC	73.1 (+0.2)	59.6 (+0.3)	67.7 (+0.6)	51.0 (+0.6)	74.9 (+0.4)	53.2 (+1.5)
+ Point-GCC [†]	-	-	-	-	75.1 (+0.6)	56.7 (+5.0)
TR3D+FF [37]	-	-	69.4	53.4	-	-
+ Point-GCC	-	-	69.7 (+0.3)	54.0 (+0.6)	-	-

Table 1: 3D Object detection results on ScanNetV2, SUN RGB-D, S3DIS validation set. The overall best results are **bold**, and the best results with the same baseline model are underlined. + means fine-tuning with pre-training on the corresponding dataset. * means that we evaluate the performance on VoteNet with the stronger MMDetection3D implementation for a fair comparison. [†] means with extra training dataset ScanNetV2.

where τ is the temperature hyper-parameter, we set it to 0.4 following the above-mentioned setting. z^i is the ℓ_2 -normalized mean feature with pseudo-label i . z_{geo}^i and z_{color}^i represent a pair of positive sample with same pseudo-label i . And z_{geo}^i with z_{color}^j corresponding different pseudo-label j represent negative samples.

3.4 Overall Hierarchical Loss

Our framework contains hierarchical supervision in point-level and object-level, and the final loss is a combination of the four losses above-mentioned:

$$\mathcal{L}_{over} = \mathcal{L}_{pc} + \alpha\mathcal{L}_{pr} + \beta\mathcal{L}_{clu} + \gamma\mathcal{L}_{oc}, \quad (7)$$

where α , β and γ are the loss weight hyper-parameters, we set them to 100, 100 and 1 respectively to balance the magnitude of losses.

3.5 Adapt to unsupervised semantic segmentation

Due to the pseudo-label from object-level supervision, Point-GCC can adapt to unsupervised downstream tasks without fine-tuning. Meanwhile, previous pre-training methods evaluate the performance by transfer learning on downstream tasks. The results can be greatly affected by the fine-tuning setting and are not intuitive between different baselines. As shown in Figure 2(b), we generate the final pseudo-labels by utilizing cluster prediction from geometry and color branch. During the evaluation stage, we use the Hungarian matching alignment [48] to project the pseudo-labels to ground-truth labels because we are agnostic to the ground truth in pre-training. Although our method is not specifically designed for unsupervised downstream tasks, we find that the process is more intuitive and fair for evaluating the performance of pre-training methods.

Method	Supervision	Backbone	Pseudo Classes	mIoU
Unsupervised Method				
SL3D [9]	unsupervised	PointNet++	400	8.5
SL3D [9]	unsupervised	Point Transformer	800	10.5
Point-GCC	unsupervised	PointNet++	20	18.3
Weakly-supervised Method				
WyPR [34]	scene-level	PointNet++	20	29.6
MPRM [41]	subcloud-level	KPConv	20	41.0
Supervised Method				
PointNet++(SSG) [30]	supervised	PointNet++	20	54.4
+ Point-GCC	supervised	PointNet++	20	59.8 (+5.4)

Table 2: 3D semantic segmentation results on ScanNetV2 dataset by different level of supervision. The overall best results are **bold**. + means fine-tuning with pre-training on the corresponding dataset.

4 Experiments

To analyze the 3D representation learned by Point-GCC, we conduct extensive experiments on multiple datasets and tasks described in Section 4.1. First we evaluate on fully unsupervised semantic segmentation tasks to validate the effectiveness of object representation in Section 4.2. Then we expand experiments by transfer learning on multiple downstream tasks and datasets in Section 4.3.

4.1 Experiment setting

Dataset. We use three popular indoor scene datasets: ScanNetV2 [16], SUN RGB-D [38], S3DIS [1] in our experiments. **ScanNetV2** is a 3D reconstruction dataset, which provides 1513 indoor scans with a total of 20 classes. **SUN RGB-D** is a monocular RGB-D image dataset, which provides 10335 RGB-D images from four different sensors with a total of 37 classes. **S3DIS** is another 3D indoor scene dataset, which provides 271 point cloud scenes across 6 areas with 13 classes.

Implementation details. We implement Point-GCC built upon the MMDetection3D [14] framework. We use the AdamW optimizer with an initial learning rate of 0.001 and weight decay of 0.0001. Other implementation details are followed the default scheme. To ensure fair comparability of results, we refer to selecting downstream models implemented by MMDetection3D. In downstream task experiments, we decay the learning rate by 0.5, and other settings follow the original implementation. The full detail settings are provided in the Appendix.

4.2 Fully unsupervised semantic segmentation

We evaluate our pre-training model on fully unsupervised semantic segmentation tasks using the method in Section 3.5 to validate the effectiveness of object representation. As shown in Table 2, our method surpasses previous unsupervised methods by a huge margin and is closer to the weakly-supervised method, despite Point-GCC being not specifically designed for unsupervised downstream tasks. With the same backbone PointNet++, Point-GCC surpasses previous work SL3D [9] by +9.8% mIoU, and +7.8% mIoU compared with more powerful Point Transformer on ScanNetV2 dataset. The result proves that Point-GCC has learned rich object representation in unsupervised pre-training.

Fine-tuning semantic segmentation. Additionally, we fine-tune the pre-training model for semantic segmentation to verify the consistent improvement of our method. With supervised fine-tuning, the model gains significant improvements by +5.4% mIoU on ScanNetV2 dataset, which proves that our method has learned intrinsic representations of the point cloud.

4.3 Transfer learning on downstream tasks

3D Object detection. For 3D object detection task, we pre-train the PointNet++ [30] backbone for VoteNet [29], VoteNet+FF [37] and GroupFree-3D [25] and the MinkResNet [13] backbone for TR3D [37], TR3D+FF [37] respectively. Table 1 shows the results on ScanNetV2, SUN-RGBD, and S3DIS datasets. Our method gains stable and significant improvements for various settings. Compared with previous 3D self-supervised methods with the common baseline model VoteNet, our method achieves higher AP₅₀ than the previous highest model Ponder [21] by +3.1% on ScanNetV2

Method	ScanNetV2			S3DIS			
	AP	AP ₅₀	AP ₂₅	AP	AP ₅₀	Prec ₅₀	Rec ₅₀
Supervised Only							
PointGroup [22]	34.8	56.7	71.3	-	57.8	61.9	62.1
HAISt [11]	43.5	64.4	75.6	-	-	71.1	65.0
SoftGroup [†] [39]	45.8	67.6	78.9	51.6	66.1	73.6	66.6
Self-supervised Pre-training							
TD3D [23]	46.2	71.1	81.3	48.6	65.1	74.4	64.8
+ Point-GCC	47.3 (+1.1)	71.3 (+0.2)	81.6 (+0.3)	<u>50.5</u> (+1.9)	<u>65.4</u> (+0.3)	<u>75.5</u> (+1.1)	<u>65.9</u> (+1.1)
TD3D [†] [23]	-	-	-	52.1	67.2	75.2	68.7
+ Point-GCC [†]	-	-	-	53.6 (+1.5)	68.4 (+1.2)	76.6 (+1.4)	69.5 (+0.8)

Table 3: 3D instance segmentation results on ScanNetV2 and S3DIS dataset. The overall best results are **bold**, and the best results with the same baseline model are underlined. + means fine-tuning with pre-training on the corresponding dataset. [†] means with extra training dataset ScanNetV2.

Point-level		Object-level		Unsupervised Segmentation	Object Detection	
Contra.	Recon.	Cluster.	Contra.	mIoU	AP ₂₅	AP ₅₀
✓	✓	✓	✓	18.27	65.3	44.1
✓	✓	✓	✗	16.07	65.0	43.6
✓	✓	✗	✗	-	64.8	43.0
✓	✗	✗	✗	-	64.4	42.8
✗	✓	✗	✗	-	63.3	42.7
✗	✗	✗	✗	-	62.3	40.8

Table 4: Ablation study of the hierarchical supervision. - means the model can't perform the unsupervised segmentation task due to the lack of the pseudo-label.

and +1.1% on SUN RGB-D. For more recent models, our model also significantly boosts VoteNet+FF, GroupFree-3D, TR3D, TR3D+FF on multiple datasets and achieves new state-of-the-art results with 69.7% AP₂₅, 54.0% AP₅₀ on SUN RGB-D and 75.1% AP₂₅, 56.7% AP₅₀ on S3DIS datasets.

3D Instance segmentation. For 3D instance segmentation task, we pre-train the MinkResNet backbone for TD3D [23] on ScanNet and S3DIS datasets. Table 3 shows the results on ScanNetV2 and S3DIS validation sets. Downstream models gain remarkable performance by +1.1% AP on ScanNetV2, +1.9% on S3DIS and +1.5% on S3DIS with extra train data, demonstrating our method's general improvement across multiple settings.

Interestingly, the improvements for the PointNet++ backbone widely surpass the MinkResNet backbone. We guess that sparse convolution architecture implicitly aligns the color information from features and the geometry information from fine-grained sparse voxel operation. It may be a kind of explanation for why 3D sparse convolution has better performance and efficiency on various tasks.

4.4 Ablation study And Discussion

To analyze the effectiveness of our approach, we further explore additional experiments to measure the contribution of each component to the final representation quality. For efficiency, all ablation experiments are implemented with VoteNet setting on pre-training and object detection.

Hierarchical supervision. To further explore the improvement of our hierarchical supervision, we conduct ablation studies with different components. Table 4 shows the unsupervised semantic segmentation results with pre-training and object detection results with fine-tuning. The results show that both contrastive learning and reconstructive learning in point-level contribute to the final results. Even though just with point-level supervision, our method has achieved higher AP₂₅ and AP₅₀ than the previous best model Ponder by +1.2% and +2.0%. Furthermore, the swapped prediction and object-level contrastive learning also provide remarkable improvements for AP₅₀ and AP₂₅, especially AP₅₀. Intuitively, the improvement of AP₅₀ is more significant than AP₂₅ demonstrating that object-level supervision improves the model with a more precise view of objects.

Color		Geometry		Object Detection		Object Picking		Object Detection	
Branch	Branch	AP ₂₅	AP ₅₀	Threshold	AP ₂₅	AP ₅₀	Threshold	AP ₂₅	AP ₅₀
✓	✓	64.8	43.0	only max score	64.5	43.0			
✗	✓	62.4	40.9	2.0 / class num.	65.3	44.1			
✓	✗	62.5	39.4	1.8 / class num.	64.4	43.7			
✗	✗	62.3	40.8	1.5 / class num.	64.2	43.2			

Table 5: Ablation study of the Geometry-Color Contrast approach.

Table 6: Ablation study of the Object sampling strategy.



Figure 3: The visualization of reconstruction results from Point-GCC. Note that we decrease the point size in geometry reconstruction to avoid the block from noisy points.

Geometry-Color Contrast. To verify the importance of our Geometry-Color Contrast approach, we compare the results with a single reconstruction branch setting. Table 5 shows object detection results with different pre-training branches. The results show that the performance with a single branch of whether geometry or color reconstruction obviously declines, which proves our Geometry-Color Contrast plays an essential role in the significant performance.

Object sampling strategy. The result in table 4 shows that object-level supervision provides the most obvious boost for AP₅₀. We compare the results with different object sampling strategies to analyze the object samples used in object-level contrastive learning. The results in table 6 show that the more confident object samples are, the greater performance we achieve. However, only using the maximum score sample, the performance decays because of over-fitting.

4.5 Visualization

Figure 3 shows the visualization of geometry and color reconstruction results from our method. The results show that our method can generate high-quality complement from one type of information in the point cloud consistently. The method may contain potential applications such as depth estimation and texture generation.

5 Conclusions

In this paper, we propose a new universal self-supervised 3D scene pre-training framework via **Geometry-Color Contrast (Point-GCC)**, which utilizes an architecture-agnostic Siamese network with hierarchical supervision. Extensive experiments show that Point-GCC significantly improves performance on unsupervised tasks without fine-tuning and a wide range of downstream tasks, especially achieving new state-of-the-art results on multiple datasets.

To the best of our knowledge, Point-GCC is the first study to explore the self-supervised paradigm that can better utilize the relations of different point cloud information, hence we elaborately design our plug-and-play pre-training framework to help improve various existing downstream methods, instead of directly designing a new architecture. We hope our work could attract more attention about the discriminative information of point cloud, which may inspire future point cloud representation learning works.

References

- [1] Iro Armeni, Ozan Sener, Amir R. Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2016. 2, 7
- [2] Yuki Markus Asano, Christian Rupprecht, and Andrea Vedaldi. Self-labelling via simultaneous clustering and representation learning. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020. 3
- [3] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. *Advances in neural information processing systems*, 32, 2019. 3
- [4] Hangbo Bao, Li Dong, Songhao Piao, and Furu Wei. Beit: BERT pre-training of image transformers. In *Int. Conf. Learn. Represent. (ICLR)*. OpenReview.net, 2022. 3
- [5] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European conference on computer vision (ECCV)*, pages 132–149, 2018. 3
- [6] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In Vittorio Ferrari, Martial Hebert, Cristian Sminchisescu, and Yair Weiss, editors, *Computer Vision - ECCV 2018 - 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XIV*, volume 11218 of *Lecture Notes in Computer Science*, pages 139–156. Springer, 2018. 3
- [7] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020. 3
- [8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. 3, 5, 13
- [9] Fernando Julio Cendra, Lan Ma, Jiajun Shen, and Xiaojuan Qi. SL3D: self-supervised-self-labeled 3d recognition. *CoRR*, abs/2210.16810, 2022. 2, 7
- [10] Jianlong Chang, Lingfeng Wang, Gaofeng Meng, Shiming Xiang, and Chunhong Pan. Deep adaptive image clustering. In *Proceedings of the IEEE international conference on computer vision*, pages 5879–5887, 2017. 3
- [11] Shaoyu Chen, Jiemin Fang, Qian Zhang, Wenyu Liu, and Xinggang Wang. Hierarchical aggregation for 3d instance segmentation. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 15447–15456. IEEE, 2021. 2, 8
- [12] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. In *Proc. Int. Conf. Mach. Learn. (ICML)*, volume 119 of *Proceedings of Machine Learning Research*, pages 1597–1607. PMLR, 2020. 3
- [13] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. 2, 7, 13
- [14] MMDetection3D Contributors. MMDetection3D: OpenMMLab next-generation platform for general 3D object detection. <https://github.com/open-mmlab/mmdetection3d>, 2020. 7, 13
- [15] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. 5, 15
- [16] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. 2, 7

- [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. 3
- [18] Runpei Dong, Zekun Qi, Linfeng Zhang, Junbo Zhang, Jianjian Sun, Zheng Ge, Li Yi, and Kaisheng Ma. Autoencoders as cross-modal teachers: Can pretrained 2d image transformers help 3d representation learning? In *Int. Conf. Learn. Represent. (ICLR)*, 2023. 1, 3
- [19] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders are scalable vision learners. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022. 3
- [20] Ji Hou, Benjamin Graham, Matthias Nießner, and Saining Xie. Exploring data-efficient 3d scene understanding with contrastive scene contexts. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 15587–15597. Computer Vision Foundation / IEEE, 2021. 1, 2, 6
- [21] Di Huang, Sida Peng, Tong He, Xiaowei Zhou, and Wanli Ouyang. Ponder: Point cloud pre-training via neural rendering. *CoRR*, abs/2301.00157, 2023. 1, 6, 7, 14
- [22] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. Pointgroup: Dual-set point grouping for 3d instance segmentation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 4866–4875. Computer Vision Foundation / IEEE, 2020. 2, 8
- [23] Maksim Kolodiaznyi, Danila Rukhovich, Anna Vorontsova, and Anton Konushin. Top-down beats bottom-up in 3d instance segmentation. *CoRR*, abs/2302.02871, 2023. 2, 8, 13
- [24] Haotian Liu, Mu Cai, and Yong Jae Lee. Masked discrimination for self-supervised learning on point clouds. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 4
- [25] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 2929–2938. IEEE, 2021. 2, 6, 7, 13, 14, 16
- [26] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *International conference on machine learning*, pages 97–105. PMLR, 2015. 3
- [27] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *Int. Conf. Comput. Vis. (ICCV)*, pages 2886–2897. IEEE, 2021. 2
- [28] Yatian Pang, Wenxiao Wang, Francis E. H. Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked autoencoders for point cloud self-supervised learning. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 1, 2, 3
- [29] Charles R. Qi, Or Litany, Kaiming He, and Leonidas J. Guibas. Deep hough voting for 3d object detection in point clouds. In *Int. Conf. Comput. Vis. (ICCV)*, pages 9276–9285. IEEE, 2019. 2, 6, 7, 13, 14
- [30] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Adv. Neural Inform. Process. Syst. (NIPS)*, pages 5099–5108, 2017. 2, 7, 13
- [31] Zekun Qi, Runpei Dong, Guofan Fan, Zheng Ge, Xiangyu Zhang, Kaisheng Ma, and Li Yi. Contrast with reconstruct: Contrastive 3d representation learning guided by generative pretraining. In *International Conference on Machine Learning (ICML)*, 2023. 1, 3, 4
- [32] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Abed Al Kader Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. 2
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proc. Int. Conf. Mach. Learn. (ICML)*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR, 2021. 3

- [34] Zhongzheng Ren, Ishan Misra, Alexander G. Schwing, and Rohit Girdhar. 3d spatial recognition without spatially labeled 3d. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [2](#), [7](#)
- [35] Dávid Rozenberszki, Or Litany, and Angela Dai. Language-grounded indoor 3d semantic segmentation in the wild. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. [15](#)
- [36] Danila Rukhovich, Anna Vorontsova, and Anton Konushin. FCAF3D: fully convolutional anchor-free 3d object detection. In Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part X*, volume 13670 of *Lecture Notes in Computer Science*, pages 477–493. Springer, 2022. [2](#), [6](#)
- [37] Danila Rukhovich, Anna Vorontsova, and Anton Konushin. TR3D: towards real-time indoor 3d object detection. *CoRR*, abs/2302.02858, 2023. [2](#), [6](#), [7](#), [13](#), [16](#)
- [38] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015. [2](#), [7](#)
- [39] Thang Vu, Kookhoi Kim, Tung Minh Luu, Thanh Nguyen, and Chang D. Yoo. Softgroup for 3d instance segmentation on point clouds. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 2698–2707. IEEE, 2022. [2](#), [8](#)
- [40] Yikai Wang, Xinghao Chen, Lele Cao, Wenbing Huang, Fuchun Sun, and Yunhe Wang. Multimodal token fusion for vision transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 12176–12185. IEEE, 2022. [1](#), [14](#)
- [41] Jiacheng Wei, Guosheng Lin, Kim-Hui Yap, Tzu-Yi Hung, and Lihua Xie. Multi-path region mining for weakly supervised 3d semantic segmentation on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4384–4393, 2020. [2](#), [7](#)
- [42] Junyuan Xie, Ross Girshick, and Ali Farhadi. Unsupervised deep embedding for clustering analysis. In *International conference on machine learning*, pages 478–487. PMLR, 2016. [3](#)
- [43] Saining Xie, Jiatao Gu, Demi Guo, Charles R. Qi, Leonidas J. Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part III*, volume 12348 of *Lecture Notes in Computer Science*, pages 574–591. Springer, 2020. [1](#), [2](#), [3](#), [4](#), [6](#)
- [44] Zhenda Xie, Zigang Geng, Jingcheng Hu, Zheng Zhang, Han Hu, and Yue Cao. Revealing the dark secrets of masked image modeling. *CoRR*, abs/2205.13543, 2022. [3](#)
- [45] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. ULIP: learning unified representation of language, image and point cloud for 3d understanding. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [3](#)
- [46] Hao Yang, Chen Shi, Yihong Chen, and Liwei Wang. Boosting 3d object detection via object-focused image fusion. *CoRR*, abs/2207.10589, 2022. [1](#)
- [47] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022. [1](#)
- [48] Hanlei Zhang, Hua Xu, Ting-En Lin, and Rui Lyu. Discovering new intents with deep aligned clustering. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(16):14365–14373, May 2021. [3](#), [6](#)
- [49] Junbo Zhang, Guofan Fan, Guanghan Wang, Zhengyuan Su, Kaisheng Ma, and Li Yi. Language-assisted 3d feature learning for semantic scene understanding. In *AAAI*, 2023. [1](#)
- [50] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Point-M2AE: Multi-scale masked autoencoders for hierarchical point cloud pre-training. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. [1](#), [2](#), [3](#)
- [51] Zaiwei Zhang, Rohit Girdhar, Armand Joulin, and Ishan Misra. Self-supervised pretraining of 3d features on any point-cloud. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 10232–10243. IEEE, 2021. [1](#), [6](#)

A Implementation Details

In this section, we provide the details and hyperparameters for pre-training and transfer learning.

A.1 Experimental details

Pre-training Details Our pre-training settings are based on the default detection configs in MMDection3D [14]. For all settings, we use the AdamW optimizer with an initial learning rate of 0.001 and weight decay of 0.0001. More detailed training configurations are shown in Table 7.

Downstream Transferring Details For 3D object detection task, we fine-tune with the PointNet++ [30] backbone for VoteNet [29], VoteNet+FF [37] and GroupFree-3D [25] and the MinkResNet [13] backbone for TR3D [37], TR3D+FF [37] respectively. For 3D instance segmentation task, we fine-tune with the MinkResNet backbone for TD3D [23] on ScanNet and S3DIS datasets. For 3D semantic segmentation task, we fine-tune with PointNet++ backbone for PointNet++(SSG) [30]. We follow the default setting in the corresponding downstream model. More detailed training configurations are shown in Table 7.

Config	VoteNet	GroupFree-3D	PointNet++(SSG)	TR3D	TD3D
backbone	PointNet2SASSG	PointNet2SASSG	PointNet2SASSG	MinkResNet	MinkResNet
Pre-training Details					
optimizer	AdamW	AdamW	AdamW	AdamW	AdamW
learning rate	1e-3	1e-3	1e-3	1e-3	1e-3
weight decay	1e-4	1e-4	1e-4	1e-4	1e-4
training epochs	400	400	400	200	200
lr scheduler	cosine	cosine	cosine	cosine	cosine
batch size	4x8	4x8	4x8	4x4	4x4
augmentation	default	default	default	default	default
Downstream Transferring Details					
optimizer	AdamW	AdamW	Adam	AdamW	AdamW
learning rate	4e-3	3e-3	1e-3	1e-3	1e-3
weight decay	1e-2	5e-4	1e-2	1e-4	1e-4
training epochs	72	160	200	24	66
lr scheduler	step	step	cosine	step	step
batch size	8x8	8x4	16x2	16x1	6x1
augmentation	default	default	default	default	default
GPU device	RTX 2080Ti	RTX 2080Ti	RTX 2080Ti	RTX 3090	RTX 3090

Table 7: Training recipes for pretraining and downstream fine-tuning.

A.2 Implementation details of Deep Clustering

We provide a pseudo-code for Deep Clustering via Swapped Prediction training loop in Pytorch style as follows. All of the hyperparameters are the same as the previous works [8] in 2D. The temperature is set to 0.1 and the Sinkhorn regularization parameter is set to 0.05 for all runs.

```
# C: prototypes (DxK) i.e., linear layer
# feat: feature from backbone + projection head
# temp: temperature

def swapped_prediction(feats_geo, feats_color):

    # normalize prototypes
    with torch.no_grad():
        C = normalize(C, dim=0, p=2)

    # normalize features
    norm_geo = normalize(feats_geo, dim=-1, p=2)
```

```

norm_color = normalize(feats_color, dim=-1, p=2)

# compute scores
scores_geo = mm(norm_geo, C)
scores_color = mm(norm_color, C)

# compute assignments
with torch.no_grad():
    q_geo = sinkhorn(scores_geo)
    q_color = sinkhorn(scores_color)

# convert scores to probabilities
p_geo = Softmax(scores_geo / temp)
p_color = Softmax(scores_color / temp)

# swap prediction problem
loss = - 0.5 * mean(q_geo * log(p_color) + q_color * log(p_geo))

# Sinkhorn-Knopp
def sinkhorn(scores, eps=0.05, niters=3):
    Q = exp(scores / eps).T
    Q /= sum(Q)
    K, B = Q.shape
    u, r, c = zeros(K), ones(K) / K, ones(B) / B
    for _ in range(niters):
        u = sum(Q, dim=1)
        Q *= (r / u).unsqueeze(1)
        Q *= (c / sum(Q, dim=0)).unsqueeze(0)
    return (Q / sum(Q, dim=0, keepdim=True)).T

```

B Additional Experiments

B.1 Additional comparison

Some baseline models [25, 29] do not use color information. Table 8 shows the additional baseline model with color information for a fair comparison from our reproduction and other works [21, 40], which get slight improvement and even decrease. The results prove what we mentioned in the main paper: directly concatenating all information can not adapt the model to discriminately learn different aspects of point clouds, demonstrating our work’s necessity.

Method	Input	Object Detection	
		AP ₂₅	AP ₅₀
VoteNet [29]	xyz+height	58.6	33.5
VoteNet [21]	xyz+color+height	58.8	33.4
VoteNet*	xyz+height	62.3	40.8
VoteNet*	xyz+color	61.8	39.9
+ Point-GCC*	xyz+color	<u>65.3</u> (+3.0)	<u>44.1</u> (+3.3)
GroupFree-3D [25]	xyz	66.3	47.8
GroupFree-3D [40]	xyz+color	66.3	47.0
+ Point-GCC	xyz+color	<u>68.1</u> (+1.8)	<u>49.2</u> (+1.4)

Table 8: 3D Object detection results on ScanNetV2 validation set. * means the VoteNet with the stronger MMDetection3D implementation for a fair comparison.

B.2 Weakly positional embedding

To further explore the improvement of positional embedding, we conduct an additional ablation study with different settings. Table 9 shows the object detection results with different positional embedding. The results show that our model has stable effects under different position encoding conditions. This may be because we remove the positional embedding and adjust the input channel in the downstream fine-tuning stage so that the impact on the downstream task is reduced.

Positional Embedding	Object Detection	
	AP ₂₅	AP ₅₀
no pos	64.5	43.7
xy pos	64.7	43.8
norm pos	65.3	44.1

Table 9: Ablation study of the positional embedding. **no pos** means without positional embedding. **xy pos** means with the positional embedding of x and y, and the corresponding task aims only to reconstruct z, *i.e.*, height estimation task.

B.3 Pseudo-label Classes

We compare the results with different pseudo-label classes in deep clustering to analyze the effect of cluster distribution. The results in table 10 show that more pseudo-label classes degrade the performance of downstream tasks. We guess that finer-grained labels in the ScanNet dataset, such as matching to ScanNet200 [35], show more apparent characteristics of long-tailed data, contrary to our assumption of equal cluster distribution because we are agnostic to the ground truth labels in pre-training.

Pseudo-label Classes	Object Detection	
	AP ₂₅	AP ₅₀
20	65.3	44.1
40	64.3	43.7
200	63.9	42.6

Table 10: Ablation study of the pseudo-label classes in Deep Clustering.

B.4 Error elimination

To obtain statistically significant results, we train our models five times and evaluate each trained model five times independently on erratic downstream tasks. Table 11 shows the best and the average value (in brackets).

C Additional Visualization

C.1 Unsupervised semantic segmentation visualization

We provide additional visualization results of unsupervised semantic segmentation. Figure 4 shows that our method can clearly distinguish the main parts of different objects without supervision. However, for small or complex objects, the segment may be merged into others or ignored because of the equal partition cluster distribution from the Sinkhorn-Knopp algorithm [15]. The result demonstrates that our pre-training approach helps the model learn object representations to enhance performance on downstream tasks.

Method	ScanNetV2		SUN RGB-D		S3DIS	
	AP ₂₅	AP ₅₀	AP ₂₅	AP ₅₀	AP ₂₅	AP ₅₀
VoteNet+FF [37]	-	-	64.5 (63.7)	39.2 (38.1)	-	-
+ Point-GCC	-	-	<u>64.9</u> (64.2)	<u>41.3</u> (40.6)	-	-
GroupFree-3D [25]	66.3 (65.7)	47.8 (47.7)	-	-	-	-
+ Point-GCC	<u>68.1</u> (67.3)	<u>49.2</u> (48.7)	-	-	-	-
TR3D [37]	72.9 (72.0)	59.3 (57.4)	67.1 (66.3)	50.4 (49.6)	74.5 (72.1)	51.7 (47.6)
+ Point-GCC	<u>73.1</u> (72.2)	<u>59.6</u> (57.2)	<u>67.7</u> (66.1)	<u>51.0</u> (49.9)	74.9(72.6)	53.2(50.9)
+ Point-GCC [†]	-	-	-	-	<u>75.1</u> (73.5)	<u>56.7</u> (54.4)
TR3D+FF [37]	-	-	69.4 (68.7)	53.4 (52.4)	-	-
+ Point-GCC	-	-	<u>69.7</u> (69.0)	<u>54.0</u> (53.3)	-	-

Table 11: 3D Object detection results on ScanNetV2 validation set. The main result is the best value, and the number within the bracket is the average value across 25 trials following previous works [25, 37]. [†] means with extra training dataset ScanNetV2.

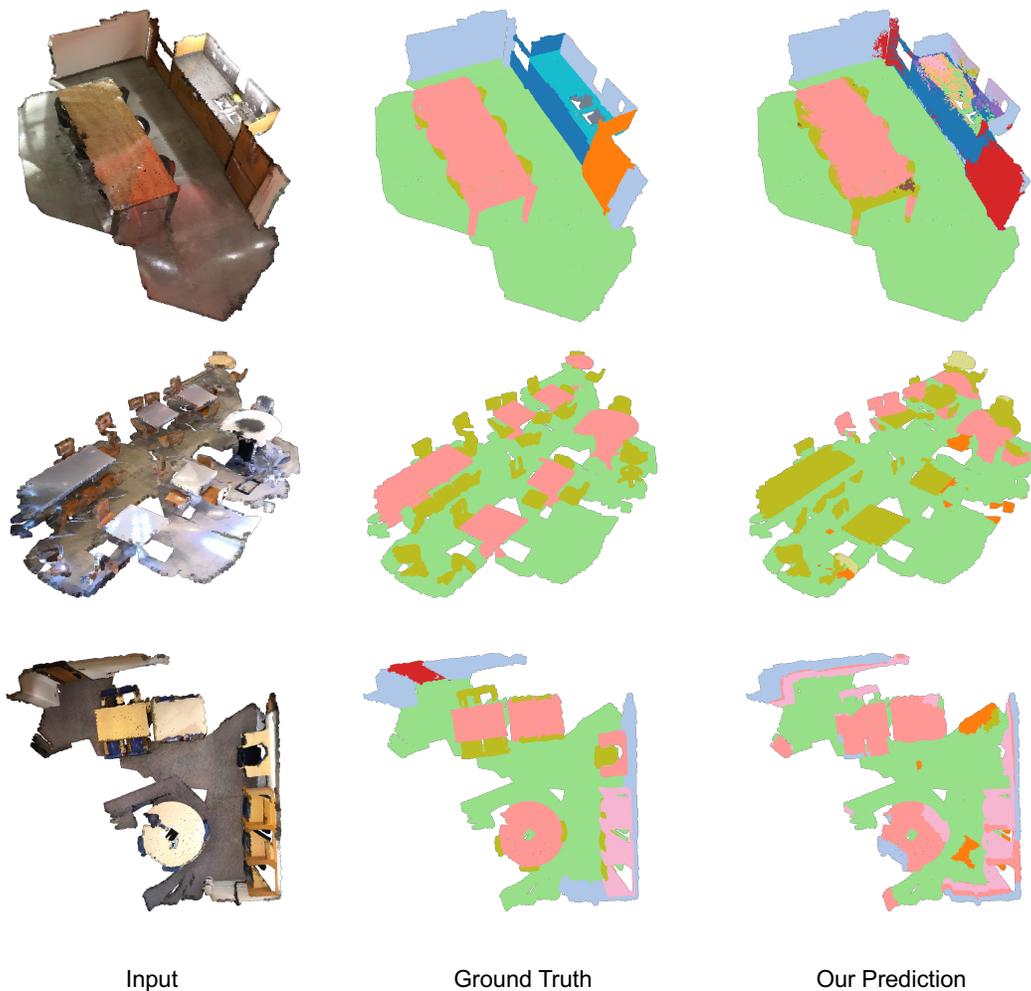


Figure 4: The visualization of unsupervised semantic segmentation results. For better visualization, we use the Hungarian matching alignment to project the pseudo-labels to ground-truth labels.