

# Superpoint-guided Semi-supervised Semantic Segmentation of 3D Point Clouds

Shuang Deng, Qiulei Dong\*, Bo Liu, and Zhanyi Hu

**Abstract**—3D point cloud semantic segmentation is a challenging topic in the computer vision field. Most of the existing methods in literature require a large amount of fully labeled training data, but it is extremely time-consuming to obtain these training data by manually labeling massive point clouds. Addressing this problem, we propose a superpoint-guided semi-supervised segmentation network for 3D point clouds, which jointly utilizes a small portion of labeled scene point clouds and a large number of unlabeled point clouds for network training. The proposed network is iteratively updated with its predicted pseudo labels, where a superpoint generation module is introduced for extracting superpoints from 3D point clouds, and a pseudo-label optimization module is explored for automatically assigning pseudo labels to the unlabeled points under the constraint of the extracted superpoints. Additionally, there are some 3D points without pseudo-label supervision. We propose an edge prediction module to constrain features of edge points. A superpoint feature aggregation module and a superpoint feature consistency loss function are introduced to smooth superpoint features. Extensive experimental results on two 3D public datasets demonstrate that our method can achieve better performance than several state-of-the-art point cloud segmentation networks and several popular semi-supervised segmentation methods with few labeled scenes.

## I. INTRODUCTION

3D point cloud semantic segmentation draws increasing attention in the field of computer vision. In recent years, a large number of Deep Neural Networks (DNNs) [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13] for point cloud semantic segmentation have been proposed. Although these methods have a great ability to obtain the semantic features of point clouds, most of them require a large number of accurately labeled 3D scenes, and manually labeling point clouds is time and labor-intensive.

Recently, some weakly supervised segmentation methods [14], [15], [16], [17], [18], [19], [20], [21], [22], [23],

[24] for 3D point clouds have been proposed, which could be roughly divided into two groups according to two different kinds of training datasets: (1) weakly supervised methods whose training dataset contains a small portion of labeled points sampled from each 3D training scene; (2) weakly supervised methods (also called semi-supervised methods) whose training dataset contains a small portion of labeled 3D scenes. The former group of methods [15], [18], [21], [22], [23], [24] require point sampling for all 3D scenes, and the point clouds sampled from some dense 3D scenes will still be somewhat dense, and the labor costs of assigning point labels will not be reduced too much. Compared with the former group of methods, the semi-supervised methods [14], [16], [17], [19], [20] are able to significantly reduce labeling costs. Hence, we focus on the semi-supervised point cloud segmentation problem in this paper.

For solving the semi-supervised semantic segmentation problem for 3D point clouds, some methods [14], [16], [19], [20] introduce additional information of point clouds. Expert knowledge is utilized in [14], [19], [20] and Mei et al. [16] considers the consistency of scans stream. Besides, the point clouds used by the methods [14], [17] are CAD models, which are much simpler than 3D scenes. In addition, there are some methods [25], [26], [27], [28] to solve the semi-supervised segmentation problem for 2D images. However, since 3D point cloud is an unordered and irregular structure, these methods cannot be applied to 3D point clouds directly.

It is noted that some existing works [5], [8], [18], [19], [20], [21], [29] for 3D point cloud semantic segmentation utilize superpoints to improve their performances. In [5], [18], [21], [29], the point clouds are geometrically partitioned by minimizing a global energy function. Liu et al. [20] employs spectral clustering for superpoint generation. These methods do not consider the color information of 3D point clouds, where some classes of objects are only different in color from the surrounding objects (*i.e.* window and board). And minimizing the global energy function is time-consuming. Landrieu et al. [8] formulates superpoints generation as a deep metric learning problem. But this partition method requires semantic information of the 3D point clouds.

Addressing the aforementioned issues, we propose a superpoint-guided semi-supervised segmentation network for 3D point clouds. The labeled and unlabeled point clouds will be processed in different ways. We use the ground truth labels to supervise the labeled point clouds. And the pseudo labels predicted from unlabeled point clouds are used for self-training. Since the pseudo labels are not completely accurate, we utilize the superpoints to optimize pseudo labels.

This work was supported by the National Natural Science Foundation of China (Grant Nos. U1805264 and 61991423), the Strategic Priority Research Program of the Chinese Academy of Sciences (XDB32050100), Beijing Science and Technology Program (Grant No. Z211100011021004), and the Foundation of the Lab of Space Optoelectronic Measurement & Perception (Grant No. 502K0019118).

S. Deng, and Q. Dong are with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, and also with the School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing 100049, China, and also with the Center for Excellence in Brain Science and Intelligence Technology, Chinese Academy of Sciences, Beijing 100190, China (e-mail: [shuang.deng, qldong]@nlpr.ia.ac.cn).

B. Liu, and Z. Hu are with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, and also with the School of Future Technology, University of Chinese Academy of Sciences, Beijing 100049, China (e-mail: liubo2017@ia.ac.cn, huzy@nlpr.ia.ac.cn).

\*Corresponding author.

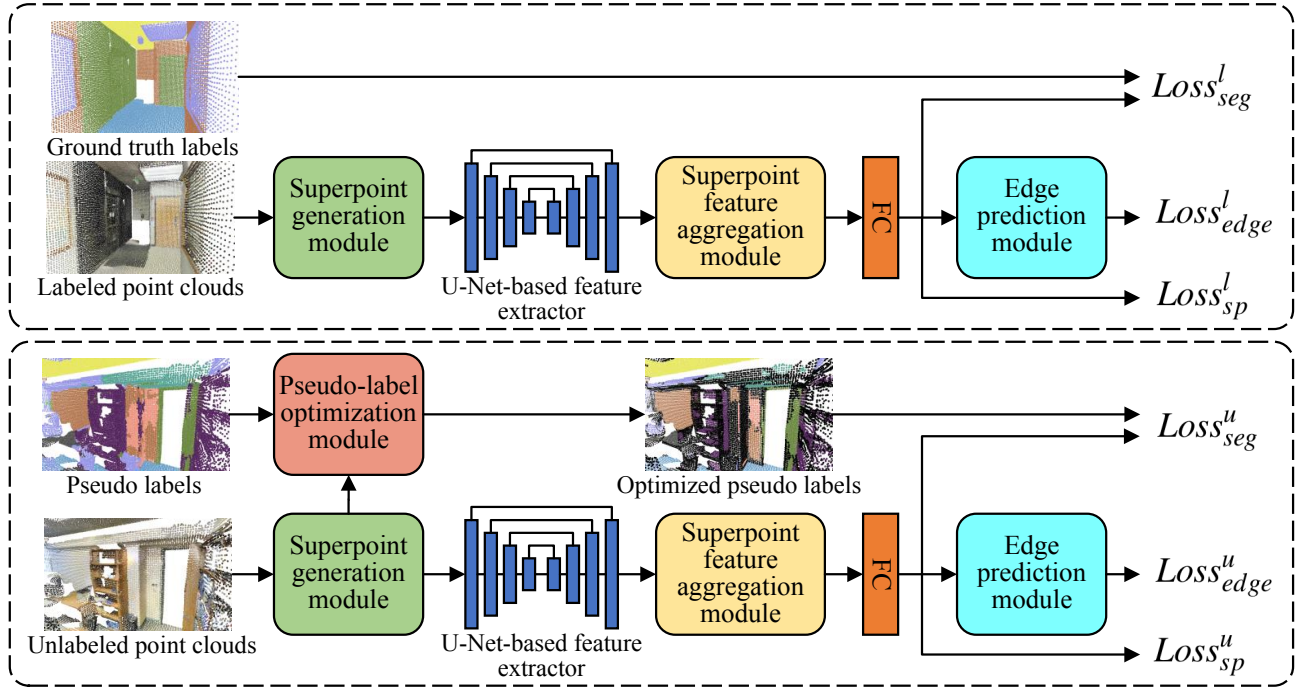


Fig. 1. Architecture of the proposed network. The structures in two dashed boxes are for labeled and unlabeled point clouds respectively. FC represents the fully-connected layer.

Specifically, we propose a superpoint generation module, named as SPG module, to combine the superpoints produced by geometry-based and color-based Region Growing algorithms [30], and a pseudo-label optimization module, named as PLO module, to modify and delete pseudo labels with low confidence in each superpoint. There are some 3D points without pseudo-label supervision. We propose an edge prediction module, named as EP module, to constrain the features from edge points of geometry and color. A superpoint feature aggregation module, named as SPFA module, and a superpoint feature consistency loss function are introduced to smooth the point features in each superpoint.

In sum, the main contributions of this paper include:

- For solving the semi-supervised semantic segmentation problem of 3D point clouds effectively and efficiently, we utilize the superpoints generated by combining geometry-based and color-based Region Growing algorithms to optimize pseudo labels predicted from unlabeled point clouds.
- We propose an edge prediction module, a superpoint feature aggregation module and a superpoint feature consistency loss function for constraining point features without pseudo labels.
- We propose the superpoint-guided semi-supervised segmentation network for 3D point clouds. The experimental results on two 3D public datasets show that the proposed method outperforms several state-of-the-art segmentation networks and several popular semi-supervised methods with few labeled scenes.

## II. SUPERPOINT-GUIDED SEMI-SUPERVISED SEGMENTATION NETWORK

In this section, we propose the superpoint-guided semi-supervised segmentation network for 3D point clouds. Firstly, we introduce the architecture of the proposed network. Secondly, we describe the details of the superpoint generation module (SPG module), the pseudo-label optimization module (PLO module), the edge prediction module (EP module), the superpoint feature aggregation module (SPFA module) and the superpoint feature consistency loss function respectively. Lastly, we end up with the final training loss of the network.

### A. Architecture

As shown in Fig. 1, our end-to-end superpoint-guided semi-supervised segmentation network consists of two branches. The inputs of one branch are labeled point clouds and their ground truth labels, and the other branch are unlabeled point clouds and their pseudo labels. The pseudo labels are predicted by our network from unlabeled point clouds. Both branches consist of a superpoint generation module (SPG module), a feature extractor based on U-Net [31], a superpoint feature aggregation module (SPFA module), a fully connected layer (FC), and an edge prediction module (EP module). And their parameters are shared. For the branch of unlabeled point clouds, there is a pseudo-label optimization module (PLO module) to optimize the pseudo labels. The U-Net-based feature extractor consists of four encoder layers and four decoder layers. The encoder layers are Local Feature Aggregation layers in RandLA-Net [9], and the decoder layers are MLPs.

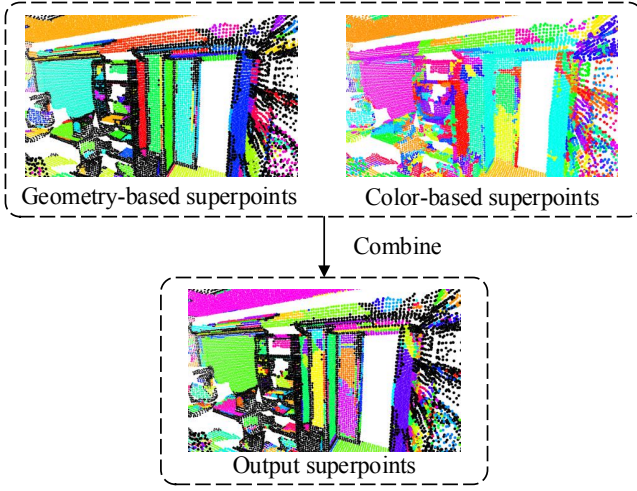


Fig. 2. The process of combining superpoints produced by geometry-based and color-based Region Growing algorithms. The black points are not clustered as superpoints due to the curvature threshold in the geometry-based Growing Region algorithm.

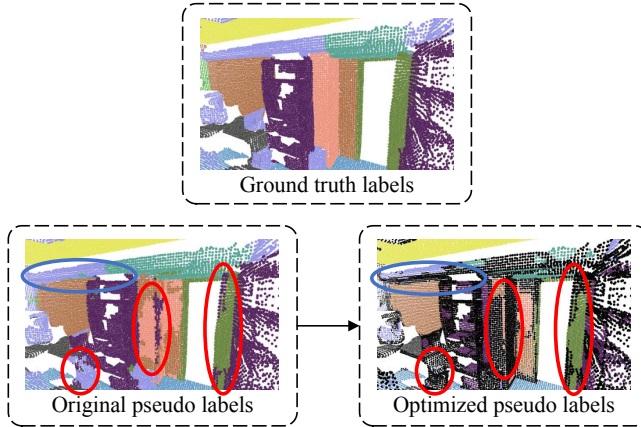


Fig. 3. The process of optimizing pseudo labels. Some pseudo labels inside the red and blue circles are incorrect. The black points have no pseudo labels after optimizing.

When a labeled 3D point cloud  $\mathbf{P}^l = \{p_1^l, p_2^l, \dots, p_{N^l}^l\} \in \mathbb{R}^{N^l \times 6}$  and its one-hot ground truth labels  $\mathbf{Y}^l = \{y_1^l, y_2^l, \dots, y_{N^l}^l\} \in \mathbb{R}^{N^l \times C}$ , and an unlabeled point cloud  $\mathbf{P}^u = \{p_1^u, p_2^u, \dots, p_{N^u}^u\} \in \mathbb{R}^{N^u \times 6}$  and its one-hot pseudo labels  $\mathbf{Y}^u = \{y_1^u, y_2^u, \dots, y_{N^u}^u\} \in \mathbb{R}^{N^u \times C}$  are given, where  $N^l$  and  $N^u$  are the number of points, 6 denotes the XYZ dimensions and RGB dimensions, and  $C$  is the number of semantic classes, we firstly send  $\mathbf{P}^l$  and  $\mathbf{P}^u$  to the SPG module to generate their superpoint collections  $\mathbf{S}^l = \{\mathbf{S}_1^l, \mathbf{S}_2^l, \dots, \mathbf{S}_{M^l}^l\}$  and  $\mathbf{S}^u = \{\mathbf{S}_1^u, \mathbf{S}_2^u, \dots, \mathbf{S}_{M^u}^u\}$ , where  $M^l$  and  $M^u$  are the number of superpoints. For the  $i^{th}$  superpoint in  $\mathbf{S}^l$ ,  $\mathbf{S}_i^l = \{p_{i_1}^l, p_{i_2}^l, \dots, p_{i_n}^l\} \in \mathbb{R}^{n \times 6}$ , where  $n$  is the number of points in this superpoint, similar in  $\mathbf{S}^u$ . Secondly, we send  $\mathbf{P}^l$  and  $\mathbf{P}^u$  to the U-Net-based feature extractor to construct their high-level representations  $\mathbf{F}^l = \{f_1^l, f_2^l, \dots, f_{N^l}^l\} \in \mathbb{R}^{N^l \times C_h}$  and  $\mathbf{F}^u = \{f_1^u, f_2^u, \dots, f_{N^u}^u\} \in \mathbb{R}^{N^u \times C_h}$ , where  $C_h$  is the dimension of high-level features. Then we send  $\mathbf{F}^l$  and  $\mathbf{F}^u$  to the SPFA

module to get feature maps  $\mathbf{G}^l = \{g_1^l, g_2^l, \dots, g_{N^l}^l\} \in \mathbb{R}^{N^l \times C_h}$  and  $\mathbf{G}^u = \{g_1^u, g_2^u, \dots, g_{N^u}^u\} \in \mathbb{R}^{N^u \times C_h}$  for feature smoothing in superpoints. After a FC layer, we obtain the final feature maps  $\mathbf{X}^l = \{x_1^l, x_2^l, \dots, x_{N^l}^l\} \in \mathbb{R}^{N^l \times C}$  and  $\mathbf{X}^u = \{x_1^u, x_2^u, \dots, x_{N^u}^u\} \in \mathbb{R}^{N^u \times C}$ .

### B. Superpoint Generation Module

We propose a novel superpoint generation module, named as SPG, to produce superpoints effectively and efficiently. The geometric and color characteristics of classes of objects in scenes are different. Some classes of objects are different in geometry and color from the surrounding objects (*i.e.* chair and table). But there are also some classes of objects are only different in geometry from the surrounding objects (*i.e.* beam and column), or only different in color from the surrounding objects (*i.e.* window and board). Some existing superpoint generation methods [5], [18], [20], [21], [29] can only geometrically partition the 3D point clouds, which ignore the color information. The proposed SPG module combines geometry-based and color-based superpoints produced by the Region Growing algorithm [30], which has low computational complexity.

The geometry-based Region Growing algorithm iteratively produces superpoints. In each iteration, there is one point with a minimum curvature value in the unsegmented points to be selected as a seed and added to a seeds set and a superpoint. Then, the following three steps are repeated until there are no point in the seeds set: (1) Finding the neighbouring points of seeds and testing their angles between their normals, these neighbouring points will be added to the current superpoint if the angles are less than the threshold value  $t_{ang}$ . (2) If the curvatures of these neighbouring points are less than the threshold value  $t_{cvr}$ , then these points are added to the seeds set. (3) Current seeds are removed from the seeds set. If there are no unsegmented points whose curvatures are smaller than the threshold value  $t_{cvr}$  in the scene, the process of iterations is terminated. Due to the curvature threshold  $t_{cvr}$ , some points will not be clustered to superpoints.

The color-based Region Growing algorithm is similar to the geometry-based ones. There are two main differences in the color-based algorithm. The first one is that it uses color instead of normals. The color threshold value is  $t_{clr}$ . The second one is that it uses the merging algorithm for segmentation control. Two neighbouring clusters with a small difference between average color are merged together. The color-based Region Growing algorithm has no curvature threshold, so every point can be clustered.

After obtaining the superpoints from the geometry-based and color-based Region Growing algorithms, we over-segment every geometry-based superpoint based on the color-based superpoints, which can be seen in Fig. 2. It is noted that the geometric edge points will not be clustered as superpoints due to the curvature threshold  $t_{cvr}$ . The final merged superpoints  $\mathbf{S}^l$  and  $\mathbf{S}^u$  could be used by the PLO module, the SPFA module, and the superpoint feature consistency loss function.

### C. Pseudo-label Optimization Module

Since the pseudo labels  $\mathbf{Y}^u$  predicted by the network are not completely accurate, and the points in same superpoint should have same semantic labels in most cases, we utilize the superpoints to optimize pseudo labels. We propose a novel pseudo-label optimization module, named as PLO module, to modify and delete pseudo labels with low confidence.

As shown from the red circle areas in the second row of Fig. 3, incorrect pseudo labels generally have no geometric and color rules. So we can constrain pseudo labels by the geometry and color-based superpoints. Specifically, for a superpoint  $\mathbf{S}_i^u (i = 1, 2, \dots, M^u)$  with  $n$  points, we first count the number of points contained in each semantic category  $n_j (\sum_{j=1}^C n_j = n)$ . Then we find the category  $c_i$  that contains the most points, which can be formulated as:

$$c_i = \arg \max_j (n_j). \quad (1)$$

If  $n_{c_i} > t_{plo} \times n$ , where  $t_{plo}$  is a ratio parameter, we modify all the pseudo labels in superpoint  $\mathbf{S}_i^u$  to  $c_i$ , otherwise all the pseudo labels in this superpoint will be deleted. We also delete the pseudo labels of points which are not clustered as superpoints in the geometry-based Growing Region algorithm [30]. After above operations being done on all superpoints in the unlabeled point clouds, the optimized pseudo labels  $\tilde{\mathbf{Y}}^u = \{\tilde{y}_1^u, \tilde{y}_2^u, \dots, \tilde{y}_{N^u}^u\} \in \mathbb{R}^{N^u \times C}$  are shown in the second row of Fig. 3.

### D. Edge Prediction Module

The geometry-based Region Growing algorithm [30] does not contain edge points due to the curvature threshold setting. And the predicted pseudo labels of edge points are usually unstable, which can be seen from the area inside the blue circle in the second row of Fig. 3. So the pseudo labels of many edge points are deleted after PLO module. We design an edge prediction module, named as EP module, to constrain the features of edge points in another way. We consider not only geometric edge points, but also color edge points. The geometric edge points are composed of points that are not clustered by the geometry-based region growing algorithm. The color edge points are those points whose neighboring points do not belong to the same color-based superpoint.

The EP module consists of two FC layers, which reduce the number of feature channels to two, to predict whether the point is a geometric or color edge point. The activation function of the first FC layer is Leaky ReLU (LReLU) [32]. The activation function of the second FC layer is Sigmoid. For the features of unlabeled point cloud  $\mathbf{X}^u$ , the outputs of the EP module are  $\mathbf{E}^u = \{e_1^u, e_2^u, \dots, e_{N^u}^u\} \in \mathbb{R}^{N^u \times 2}$ , which can be formulated as:

$$e_i^u = \text{Sigmoid}(\text{FC}(\text{LReLU}(\text{FC}(x_i^u)))) \quad (2)$$

where  $e_i^u$  is the  $i$ -th element of  $\mathbf{E}^u$ . The labels of EP module for the unlabeled point cloud  $\mathbf{P}^u$  are  $\hat{\mathbf{E}}^u = \{\hat{e}_1^u, \hat{e}_2^u, \dots, \hat{e}_{N^u}^u\} \in \mathbb{R}^{N^u \times 2}$ , where the values of edge points are 1, otherwise 0.

So the edge prediction loss function for the unlabeled point cloud  $\text{Loss}_{edge}^u$  is:

$$\text{Loss}_{edge}^u = \frac{1}{N^u} \sum_{i=1}^{N^u} \sum_{c=1}^2 -\hat{e}_{i,c}^u \log(e_{i,c}^u) - (1 - \hat{e}_{i,c}^u) \log(1 - e_{i,c}^u) \quad (3)$$

where  $e_{i,c}^u$  is the  $c$ -th channel of  $e_i^u$ . The edge prediction loss function for the labeled point cloud  $\text{Loss}_{edge}^l$  is obtained by the same way.

### E. Smoothing Superpoint Features

In the PLO module, the pseudo labels of some superpoints are deleted, the features in these superpoints are not constrained. Besides, the points within same superpoint should have similar semantic features in most cases. So we propose a superpoint feature aggregation module, named as SPFA module, and a superpoint feature consistency loss function to smooth superpoint features.

We first introduce the SPFA module. For the  $i$ -th clustered point in the unlabeled point cloud  $p_i^u$ , we randomly sample  $K$  points  $p_{i_1}^u, p_{i_2}^u, \dots, p_{i_K}^u$  within the same superpoint as  $p_i^u$ , and thier high-level features  $f_{i_1}^u, f_{i_2}^u, \dots, f_{i_K}^u$ . The aggregated feature  $g_i^u$  for the point  $p_i^u$  is obtained by:

$$g_i^u = \frac{(f_i^u + \sum_{k=1}^K f_{i_k}^u)}{2}. \quad (4)$$

Obtaining  $g_i^l$  is in the same way.

Then we introduce the superpoint feature consistency loss functions  $\text{Loss}_{sp}^l$  and  $\text{Loss}_{sp}^u$ . We use the variance function as the metric criterion of smoothness. For the features of unlabeled point cloud  $\mathbf{X}^u$ , the loss function  $\text{Loss}_{sp}^u$  is formulated as:

$$\text{Loss}_{sp}^u = \frac{1}{N^u} \sum_{i=1}^{N^u} \sum_{c=1}^C w_i^u (x_{i,c}^u - \frac{\sum_{k=1}^K x_{i_k,c}^u}{K})^2 \quad (5)$$

where  $w_i^u$  is a boolean value whether  $p_i^u$  is within a superpoint.  $\text{Loss}_{sp}^l$  is obtained in the same way.

### F. Training Loss

We introduce the final training loss of the network. For the labeled point clouds, we calculate a multi-class cross-entropy loss  $\text{Loss}_{seg}^l$  between  $\mathbf{Y}^l$  and the Softmax of features  $\mathbf{X}^l$  as follows:

$$\text{Loss}_{seg}^l = -\frac{1}{N^l} \sum_{i=1}^{N^l} \sum_{c=1}^C y_{i,c}^l \log(\text{Softmax}(x_{i,c}^l)) \quad (6)$$

where  $y_{i,c}^l$  is the  $c$ -th channel of  $y_i^l$ . For the unlabeled point clouds, we calculate a weighted multi-class cross-entropy loss  $\text{Loss}_{seg}^u$  between  $\tilde{\mathbf{Y}}^u$  and features  $\mathbf{X}^u$  as follows:

$$\text{Loss}_{seg}^u = -\frac{1}{N^u} \sum_{i=1}^{N^u} \sum_{c=1}^C \tilde{w}_i^u \tilde{y}_{i,c}^u \log(\text{Softmax}(x_{i,c}^u)) \quad (7)$$

where  $\tilde{w}_i^u$  is a boolean value whether  $p_i^u$  has an optimized pseudo label after PLO module. The final loss function is formulated as:

$$\text{Loss} = \text{Loss}_{seg}^l + \text{Loss}_{seg}^u + \text{Loss}_{edge}^l + \text{Loss}_{edge}^u + \text{Loss}_{sp}^l + \text{Loss}_{sp}^u. \quad (8)$$

TABLE I  
SEMANTIC SEGMENTATION RESULTS (%) ON THE S3DIS DATASET  
(AREA-5).

|     | Methods            | mIoU         | mAcc         | OA           |
|-----|--------------------|--------------|--------------|--------------|
| 20% | RandLA-Net [9]     | 50.90        | 60.76        | 81.24        |
|     | GA-Net [13]        | 52.12        | 61.76        | 81.46        |
|     | SCF-Net [12]       | 51.78        | 61.19        | 81.61        |
|     | $\pi$ -Model [25]  | 51.58        | 59.46        | 82.09        |
|     | Mean Teacher [26]  | 51.44        | 62.27        | 81.70        |
|     | Pseudo-Labels [33] | 52.21        | 63.76        | 82.39        |
|     | Ours               | <b>55.49</b> | <b>65.45</b> | <b>83.55</b> |
| 10% | RandLA-Net [9]     | 45.64        | 58.58        | 79.08        |
|     | GA-Net [13]        | 43.85        | 52.60        | 78.41        |
|     | SCF-Net [12]       | 42.64        | 53.16        | 76.54        |
|     | $\pi$ -Model [25]  | 46.05        | 57.49        | 80.26        |
|     | Mean Teacher [26]  | 46.72        | 57.84        | 80.50        |
|     | Pseudo-Labels [33] | 47.78        | 61.40        | 81.13        |
|     | Ours               | <b>51.14</b> | <b>64.92</b> | <b>82.54</b> |

TABLE II  
SEMANTIC SEGMENTATION RESULTS (%) ON THE SCANNet DATASET.

|     | Methods            | mIoU         | mAcc         | OA           |
|-----|--------------------|--------------|--------------|--------------|
| 20% | RandLA-Net [9]     | 52.86        | 62.56        | 81.43        |
|     | GA-Net [13]        | 52.12        | 61.50        | 81.39        |
|     | SCF-Net [12]       | 52.05        | 61.32        | 81.31        |
|     | $\pi$ -Model [25]  | 53.07        | 62.78        | 81.52        |
|     | Mean Teacher [26]  | 52.98        | 62.65        | 81.48        |
|     | Pseudo-Labels [33] | 53.23        | 62.95        | 81.63        |
|     | Ours               | <b>55.12</b> | <b>63.61</b> | <b>82.43</b> |
| 10% | RandLA-Net [9]     | 49.34        | 58.20        | 79.66        |
|     | GA-Net [13]        | 49.03        | 58.05        | 79.29        |
|     | SCF-Net [12]       | 49.11        | 59.35        | 79.21        |
|     | $\pi$ -Model [25]  | 49.52        | 58.48        | 79.87        |
|     | Mean Teacher [26]  | 49.41        | 58.65        | 79.70        |
|     | Pseudo-Labels [33] | 50.25        | 59.37        | 79.92        |
|     | Ours               | <b>52.38</b> | <b>60.76</b> | <b>81.18</b> |

### III. EXPERIMENTS

In this section, we firstly introduce the details of experimental setup. Secondly, we evaluate the performances of proposed superpoint-guided semi-supervised segmentation network on two 3D public datasets with a few labeled 3D scenes. Thirdly, we explore the effect of  $t_{plo}$ . Finally, we end up with ablation analysis.

#### A. Experimental Setup

The proposed superpoint-guided semi-supervised segmentation network is evaluated on two 3D public datasets, including S3DIS [34], and ScanNet [35]. In the geometry-based Region Growing algorithm [30], the curvature threshold  $t_{cvt}$  is 1, and the angle threshold  $t_{ang}$  is 3 degrees as defaulted in PCL Library [36]. In the color-based Region Growing

algorithm, the color threshold  $t_{clr}$  is 6 as defaulted in PCL Library [37]. In the PLO module, the ratio parameter  $t_{plo}$  is 0.8. The U-Net-based feature extractor parameters are consistent with the model before the FC layers in RandLA-Net [9], where  $C_h$  is 64. The output dimensionality of the first FC layer in EP module is 6. We train the network using the Adam optimizer with initial learning rate 0.01 and batchsize 6 for 100 epochs. In the first 50 epochs, we only optimize the network branch for labeled point clouds. And in the last 50 epochs, we train the whole network. The pseudo labels are updated after each epoch.

#### B. Evaluation on the S3DIS Dataset

The S3DIS dataset consists of 271 rooms in 6 different areas inside an office building. 13 semantic categories are assigned to each 3D point with XYZ coordinates and RGB features. Since the fifth area with 68 rooms does not overlap with other areas, experiments on Area-5 could better reflect the generalization ability of the framework. So we conducted our experiments on Area-5 validation. We randomly sample about 20% and 10% (40 and 20 rooms) of the 203 rooms respectively in the training set as labeled point clouds, and the remaining rooms in the training set are used as unlabeled point clouds. The evaluation metrics we use are mean class Intersection-over-Union (mIoU), mean class Accuracy (mAcc) and Overall Accuracy (OA).

We compare our superpoint-guided semi-supervised segmentation network to several state-of-the-art point cloud semantic segmentation methods including RandLA-Net [9], GA-Net [13], and SCF-Net [12], and several popular semi-supervised semantic segmentation methods based on RandLA-Net including  $\pi$ -Model [25], Mean Teacher [26], and Pseudo-Labels [33]. The RandLA-Net, GA-Net, and SCF-Net are only trained on the labeled data. In the  $\pi$ -Model and Mean Teacher, the dual inputs are the original point cloud and the point cloud after a random plane rotation and a random mirror transformation. In the Pseudo-Labels, the predicted labels are updated after each epoch. All the comparative methods utilize the same labeled and unlabeled splitting manner.

As seen from Table I,  $\pi$ -Model and Mean Teacher only improving mIoU by about 1% based on RandLA-Net indicates that the consistency between geometric transformed point clouds is not enough to constrain the unlabeled point cloud features. The results of Pseudo-Labels are worse than our method, indicating that there are some false-predicted pseudo labels which will affect the learning of network. Our method achieves best on all metrics due to its more effective use of unlabeled data. The results on 20% semi-supervised setting are better than on 10% semi-supervised setting, which may be attributed to more labeled point clouds.

#### C. Evaluation on the ScanNet Dataset

The ScanNet dataset contains 1,513 3D indoor scenes obtained by scanning and reconstruction, of which 1,201 are used for training and the remaining 312 are used for testing. 20 semantic categories are provided for evaluation.



TABLE III  
RESULTS OF DIFFERENT  $t_{plo}$  ON THE S3DIS DATASET (AREA-5).

|     | $t_{plo}$ values | mIoU         | mAcc         | OA           |
|-----|------------------|--------------|--------------|--------------|
| 20% | 0.70             | 53.93        | 64.80        | 82.92        |
|     | 0.75             | 54.48        | 65.20        | 83.23        |
|     | 0.80             | <b>55.49</b> | <b>65.45</b> | <b>83.55</b> |
|     | 0.85             | 54.13        | 64.79        | 82.81        |
|     | 0.90             | 53.38        | 64.02        | 82.53        |
| 10% | 0.70             | 50.59        | 64.56        | 81.93        |
|     | 0.75             | 50.97        | 64.78        | 82.21        |
|     | 0.80             | <b>51.14</b> | <b>64.92</b> | <b>82.54</b> |
|     | 0.85             | 50.67        | 63.56        | 82.19        |
|     | 0.90             | 49.57        | 60.07        | 82.15        |

TABLE IV  
ABLATION STUDY OF THE MODULES ON THE S3DIS DATASET (AREA-5).

|     | Methods              | mIoU         | mAcc         | OA           |
|-----|----------------------|--------------|--------------|--------------|
| 20% | Baseline             | 50.90        | 60.76        | 81.24        |
|     | Baseline+SPFA        | 51.32        | 61.21        | 82.22        |
|     | Baseline+SPFA+PL     | 52.75        | 63.86        | 82.78        |
|     | Baseline+SPFA+PLO    | 53.95        | 64.57        | 83.04        |
|     | Baseline+SPFA+PLO+EP | 54.77        | 64.98        | 83.30        |
|     | Ours                 | <b>55.49</b> | <b>65.45</b> | <b>83.55</b> |
| 10% | Baseline             | 45.64        | 58.58        | 79.08        |
|     | Baseline+SPFA        | 46.05        | 59.62        | 80.38        |
|     | Baseline+SPFA+PL     | 47.86        | 61.59        | 81.38        |
|     | Baseline+SPFA+PLO    | 49.78        | 62.63        | 81.86        |
|     | Baseline+SPFA+PLO+EP | 50.45        | 63.25        | 82.17        |
|     | Ours                 | <b>51.14</b> | <b>64.92</b> | <b>82.54</b> |

We randomly sample about 20% and 10% (240 and 120 rooms) of the 1201 rooms in the training set as labeled scenes, and the remaining rooms in the training set are used as unlabeled scenes. The mIoU, mAcc, and OA are used as evaluation metrics.

The competitive methods we use following experiments on the S3DIS dataset. Table II shows the comparison results. As seen from Table II, the results of  $\pi$ -Model [25], Mean Teacher [26], and Pseudo-Labels [33] have a small improvement on the basis of RandLA-Net, which may be attributed to the fact that there are more semantic categories in ScanNet than S3DIS, which results in a small number of labeled points of some categories. It is not easy to learn the features of these categories by the DNNs. Our method achieves the state-of-the-art performance, probably due to the great pseudo-label filtering and feature constraining abilities.

#### D. Effect of $t_{plo}$

The ratio parameter  $t_{plo}$  in the PLO module affects the quality of the optimized pseudo labels, and results in affecting the final segmentation performances. Too small value of  $t_{plo}$  will result in pseudo labels with lower-confidence being assigned to superpoints, and too large value of  $t_{plo}$  will result in many correct pseudo labels being deleted.

Here we conduct experiments to analyze the effect of  $t_{plo}$  by setting different values  $\{0.7, 0.75, 0.8, 0.85, 0.9\}$ . We conduct experiments on Area-5 of the S3DIS dataset with the evaluation metrics mIoU, mAcc and OA. The results are listed in Table III. As seen from Table III, the results by setting  $t_{plo}$  to 0.8 achieve the best performance, so we use this value as  $t_{plo}$  in the PLO module.

#### E. Ablation Study

For ablation study, we stack the proposed sub-modules on the baseline step-to-step to prove the effectiveness of our method. Our baseline method employs a U-Net-based feature extractor from RandLA-Net [9], and is only trained on the labeled point clouds for 100 epochs. The comparing experiments are (1) baseline method, denoted as “Baseline”; (2) adding the SPG and SPFA modules on baseline and being trained on the labeled point clouds, denoted as “Baseline+SPFA”; (3) adding pseudo labels to unlabeled point clouds for supervision in the last 50 epochs based on (2), denoted as “Baseline+SPFA+PL”; (4) adding the PLO module on (3) for unlabeled point clouds, denoted as “Baseline+SPFA+PLO”; (5) adding the EP module on (4) for all point clouds, denoted as “Baseline+SPFA+PLO+EP”; and (6) adding the superpoint feature consistency loss functions  $Loss_{sp}^l$  and  $Loss_{sp}^u$  on (5), denoted as “Ours”. We conduct ablation study on Area-5 of the S3DIS dataset with the evaluation metrics mIoU, mAcc and OA. And 20% and 10% of the rooms in the training set are used for labeled point clouds.

As shown in Table IV, the performances on “Baseline+SPFA” being better than “Baseline” demonstrate the importance of smoothing the features in superpoints. “Baseline+SPFA+PL” achieves better than “Baseline+SPFA”, which may be attributed to the supervision of unlabeled point clouds. “Baseline+SPFA+PLO” performing better than “Baseline+SPFA+PL” indicates that the superpoints produced by combining geometry-based and color-based Region Growing algorithms [30] can help optimize pseudo labels effectively. The result of “Baseline+SPFA+PLO+EP” achieves better than “Baseline+SPFA+PLO”, which may be attributed to edge-point feature learning. “Ours” achieves best, which demonstrates that combining all these modules can reach the best results.

## IV. CONCLUSIONS

For using the large amount of unlabeled point clouds, we propose a superpoint-guided semi-supervised segmentation network for 3D point clouds. Specifically, we combine the superpoints produced by geometry-based and color-based Region Growing algorithms [30] to optimize the pseudo labels predicted by unlabeled point clouds. The features of points without pseudo labels are constrained by the superpoint feature aggregation module, the edge prediction module, and the superpoint feature consistency loss function. Our method can learn the discriminative features of unlabeled point clouds and achieve best performance on two 3D public datasets with a few number of labeled scenes in most cases.

## REFERENCES

- [1] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 652–660.
- [2] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5099–5108.
- [3] B. Graham, M. Engelcke, and L. van der Maaten, "3d semantic segmentation with submanifold sparse convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 9224–9232.
- [4] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: Convolution on x-transformed points," in *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2018, pp. 820–830.
- [5] L. Landrieu and M. Simonovsky, "Large-scale point cloud semantic segmentation with superpoint graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4558–4567.
- [6] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *ACM Transaction on Graphics (TOG)*, vol. 38, pp. 1–12, 2019.
- [7] H. Zhao, L. Jiang, C.-W. Fu, and J. Jia, "Pointweb: Enhancing local neighborhood features for point cloud processing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5565–5573.
- [8] L. Landrieu and M. Bousaha, "Point cloud oversegmentation with graph-structured deep metric learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 7432–7441.
- [9] Q. Hu, B. Yang, L. Xie, S. Rosa, Y. Guo, Z. Wang, N. Trigoni, and A. Markham, "Randla-net: Efficient semantic segmentation of large-scale point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 11 105–11 114.
- [10] Q. Xu, X. Sun, C.-Y. Wu, P. Wang, and U. Neumann, "Grid-gcn for fast and scalable point cloud learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 5661–5670.
- [11] S. Deng, B. Liu, Q. Dong, and Z. Hu, "Rotation transformation network: Learning view-invariant point cloud for classification and segmentation," in *Proceedings of the IEEE International Conference on Multimedia and Expo (ICME)*, 2021, pp. 1–6.
- [12] S. Fan, Q. Dong, F. Zhu, Y. Lv, P. Ye, and F.-Y. Wang, "Scf-net: Learning spatial contextual features for large-scale point cloud segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 14 504–14 513.
- [13] S. Deng and Q. Dong, "Ga-net: Global attention network for point cloud semantic segmentation," *IEEE Signal Processing Letters (SPL)*, vol. 28, pp. 1300–1304, 2021.
- [14] Y. Wang, S. Asafi, O. van Kaick, H. Zhang, D. Cohen-Or, and B. Chen, "Active co-analysis of a set of shapes," *ACM Transactions on Graphics (TOG)*, vol. 31, pp. 1–10, 2012.
- [15] X. Xu and G. H. Lee, "Weakly supervised semantic point cloud segmentation: Towards 10x fewer labels," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 706–13 715.
- [16] J. Mei, B. Gao, D. Xu, W. Yao, X. Zhao, and H. Zhao, "Semantic segmentation of 3d lidar data in dynamic scene using semi-supervised learning," *IEEE Transactions on Intelligent Transportation Systems (TITS)*, vol. 21, pp. 2496–2509, 2020.
- [17] H. Li, Z. Sun, Y. Wu, and Y. Song, "Semi-supervised point cloud segmentation using self-training with label confidence prediction," *Neurocomputing*, vol. 437, pp. 227–237, 2021.
- [18] M. Cheng, L. Hui, J. Xie, and J. Yang, "Sspc-net: Semi-supervised semantic 3d point cloud segmentation network," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2021, pp. 1140–1147.
- [19] T.-H. Wu, Y.-C. Liu, Y.-K. Huang, H.-Y. Lee, H.-T. Su, P.-C. Huang, and W. H. Hsu, "Redal: Region-based and diversity-aware active learning for point cloud semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 15 510–15 519.
- [20] X. Shi, X. Xu, K. Chen, L. Cai, C. S. Foo, and K. Jia, "Label-efficient point cloud semantic segmentation: An active learning approach," *arXiv preprint: 2101.06931*, 2021.
- [21] Z. Liu, X. Qi, and C.-W. Fu, "One thing one click: A self-training approach for weakly supervised 3d semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 1726–1736.
- [22] J. Hou, B. Graham, M. Nießner, and S. Xie, "Exploring data-efficient 3d scene understanding with contrastive scene contexts," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 15 587–15 597.
- [23] Q. Hu, B. Yang, G. Fang, Y. Guo, A. Leonardis, N. Trigoni, and A. Markham, "Sqcn: Weakly-supervised semantic segmentation of large-scale 3d point clouds with 1000x fewer labels," *arXiv preprint: 2104.04891*, 2021.
- [24] Y. Zhang, Y. Qu, Y. Xie, Z. Li, S. Zheng, and C. Li, "Perturbed self-distillation: Weakly supervised large-scale point cloud semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 15 520–15 528.
- [25] S. Laine and T. Aila, "Temporal ensembling for semi-supervised learning," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017, pp. 1–13.
- [26] A. Tarvainen and H. Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," in *Proceedings of the Conference on Neural Information Processing Systems (NeurIPS)*, 2017, p. 1195–1204.
- [27] N. Souly, C. Spampinato, and M. Shah, "Semi supervised semantic segmentation using generative adversarial network," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 5689–5697.
- [28] X. Luo, J. Chen, T. Song, and G. Wang, "Semi-supervised medical image segmentation through dual-task consistency," in *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2021, pp. 8801–8809.
- [29] M. Cheng, L. Hui, J. Xie, J. Yang, and H. Kong, "Cascaded non-local neural network for point cloud semantic segmentation," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 8447–8452.
- [30] R. Adams and L. Bischof, "Seeded region growing," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 16, pp. 641–647, 1994.
- [31] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015, pp. 234–241.
- [32] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *Proceedings of the International Conference on Machine Learning (ICML)*, 2013, pp. 1–6.
- [33] D.-H. Lee, "Pseudo-label : The simple and efficient semi-supervised learning method for deep neural networks," in *Proceedings of the International Conference on Machine Learning Workshop (ICMLW)*, 2013, pp. 896–901.
- [34] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, "3d semantic parsing of large-scale indoor spaces," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 1534–1543.
- [35] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "ScanNet: Richly-annotated 3d reconstructions of indoor scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2432–2443.
- [36] S. Ushakov, "Region growing segmentation," [https://pcl.readthedocs.io/projects/tutorials/en/latest/region\\_growing\\_segmentation.html?highlight=region%20growing](https://pcl.readthedocs.io/projects/tutorials/en/latest/region_growing_segmentation.html?highlight=region%20growing).
- [37] S. Ushakov, "Color-based region growing segmentation," [https://pcl.readthedocs.io/projects/tutorials/en/latest/region\\_growing\\_rgb\\_segmentation.html?highlight=region%20growing](https://pcl.readthedocs.io/projects/tutorials/en/latest/region_growing_rgb_segmentation.html?highlight=region%20growing).