

# SegGroup: Seg-Level Supervision for 3D Instance and Semantic Segmentation

An Tao<sup>1</sup>Yueqi Duan<sup>2</sup>Yi Wei<sup>1</sup>Jiwen Lu<sup>1</sup>Jie Zhou<sup>1</sup><sup>1</sup>Tsinghua University<sup>2</sup>Stanford University

## Abstract

*Most existing point cloud instance and semantic segmentation methods rely heavily on strong supervision signals, which require point-level labels for every point in the scene. However, such strong supervision suffers from large annotation costs, arousing the need to study efficient annotating. In this paper, we discover that the locations of instances matter for 3D scene segmentation. By fully taking the advantages of locations, we design a weakly supervised point cloud segmentation algorithm that only requires clicking on one point per instance to indicate its location for annotation. With over-segmentation for pre-processing, we extend these location annotations into segments as seg-level labels. We further design a segment grouping network (SegGroup) to generate pseudo point-level labels under seg-level labels by hierarchically grouping the unlabeled segments into the relevant nearby labeled segments, so that existing point-level supervised segmentation models can directly consume these pseudo labels for training. Experimental results show that our seg-level supervised method (SegGroup) achieves comparable results with the fully annotated point-level supervised methods. Moreover, it also outperforms the recent weakly supervised methods given a fixed annotation budget.*

## 1. Introduction

Recent years have witnessed significant progress on analyzing different 3D geometric data structures, including point cloud [43, 45], mesh [16, 47], voxel grid [44, 58], multi-view [9, 50] and implicit function [40, 38, 11]. Due to the popularity of varying scanning devices, 3D point cloud data is easy to obtain and thus arouses more and more attention. Recently, many deep learning methods have been proposed to directly operate on point clouds and have achieved encouraging performance [43, 45, 57, 54, 51].

Point cloud instance and semantic segmentation are two fundamental but challenging tasks in 3D understanding. Given a point cloud, instance segmentation aims to find all existing objects and mark each object with a unique instance label and a semantic class, while semantic seg-

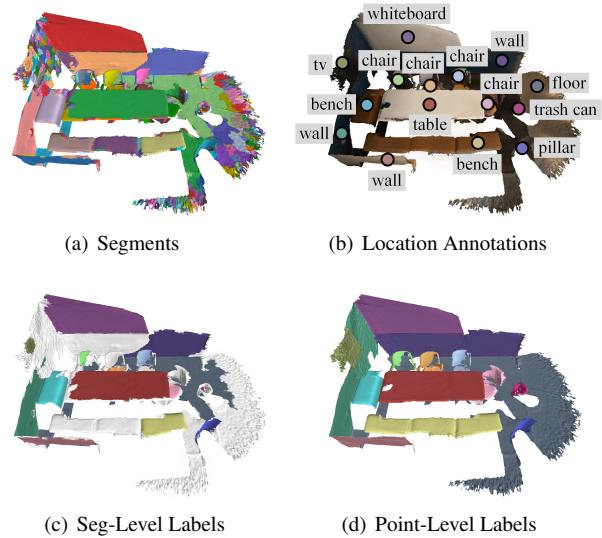


Figure 1. Illustration of annotation process to give instance locations. (a) Given a scene, we first perform over-segmentation to get segments as the data pre-processing step. (b) We click one point on the most representative segment of each instance to give the location and annotate it with semantic classes and instance IDs. (c) Then, we extend location annotations into segments to obtain seg-level labels. (d) Compared with point-level labels which require 22.3 minutes per scene to annotate all points, seg-level labels only require 1.93 minutes per scene to give location information of each instance by clicking on 0.03% points.

mentation only predicts semantic classes. Over the past few years, strong point-level supervisions that annotate every point in the scene have derived rapid performance improvement on point cloud instance and semantic segmentation tasks [53, 63, 17, 61, 26, 7]. However, as each scene may contain a large number of points, it is highly time-consuming to annotate all the points. For example, ScanNet [8] is a widely used large-scale real-world indoor dataset. It contains 1,613 scenes and each scene has 150,000 points on average. Even though ScanNet adopts over-segmentation to reduce the annotation workload, it still needs about 22.3 minutes to annotate all the segments in one scene. With a growing number of unlabeled 3D point cloud

data in real-world applications, a natural question is raised: is it necessary to label all points in a point cloud scene?

In this paper, we discover that *the locations of instances matter for 3D scene segmentation*. Compared with the 2D images which lack depth dimension, locations of instances can be more precisely in 3D scenes. To fully take the advantages of locations, we only need to click on one point per instance to indicate its location for segmentation, rather than labeling every point in the scene. With these locations, we design a weakly supervised point cloud segmentation task to explore the importance of locations.

Figure 1 illustrates the annotation process to give instance locations. Following the labeling strategy of point-level labels in ScanNet [8], we first perform over-segmentation to obtain segments as pre-processing. Unlike 2D images which usually face occlusions and lightening variances, 3D data structures are suitable for over-segmentation for their apparent boundaries between different simple geometry parts. Then, we click one point on the most representative<sup>1</sup> segment of each instance to give location annotations and extend them into corresponding segments as seg-level labels. Finally, the seg-level labels are adopted as supervision signals for point cloud segmentation. According to our manual annotations, seg-level labels only require 1.93 minutes per scene to click 0.03% points compared with the strong point-level labels. However, thank to the over-segmentation, these annotated segments contain 29.42% points of the whole scene.

To learn point cloud instance and semantic segmentation under seg-level supervision, we design a two-stage approach. We first propose a segment grouping network (SegGroup) to generate pseudo point-level labels from seg-level labels for the remaining unannotated points on the training set. Then, we adopt an existing point-level supervised point cloud segmentation model for standard training. The two stages are trained separately, and the evaluation of the segmentation performance is conducted on the existing point-level supervised model. In the SegGroup network, we propagate label information by grouping unlabeled segments into the relevant nearby labeled segments and conduct the grouping operation hierarchically. After all grouping operations, all points in the point cloud scene are assigned with labels that are considered as pseudo point-level labels.

Experimental results show that our seg-level supervised method (SegGroup) achieves comparable results with the fully annotated point-level supervised methods. It also outperforms the recent weakly supervised methods [18, 55] given a fixed annotation budget. These results validate that annotating location information is a low-cost but high-yield labeling manner for 3D scene segmentation.

<sup>1</sup>In this work, we consider the largest segment of the instance as the most representative segment to indicate the instance location.

## 2. Related Work

**Point Cloud Segmentation.** Approaches on point cloud semantic segmentation can be mainly classified into two categories: voxel-based [14, 7, 6] and point-based [43, 45, 54, 51, 57]. Voxel-based approaches voxelize point clouds into 3D grids in order to apply powerful 3D CNNs, while point-based approaches directly design models on point clouds to learn per-point local features. Recent methods on point-based scheme include neighbouring feature pooling [23, 66, 65, 20], graph message passing [48, 52, 34], attention-based aggregation [59, 62], and kernel-based convolution [51, 49, 57]. For point cloud instance segmentation, there are two common strategies to find instances in 3D scenes: detection-based [61, 63, 17] and segmentation-based [53, 41, 33, 26, 12, 15]. Detection-based approaches first extract 3D bounding boxes using object detection techniques, and then inside each box find the object mask. By contrast, segmentation-based approaches first predict semantic labels for each point with a semantic segmentation framework, and then group points into different objects.

While most existing methods heavily rely on strong point-level labels, only a few works have studied weakly supervised point cloud semantic segmentation. Wei *et al.* [55] proposed scene-level labels to indicate the semantic classes appearing in a point cloud scene, and subcloud-level labels to indicate the semantic classes appearing in a spherical subcloud sampled from a given point cloud. Hou *et al.* [18] and Xu *et al.* [60] studied point cloud segmentation under a fraction of labeled points. Although these weak label forms reduce annotation effort significantly, they cannot provide any instance-specific location information, which is a very important supervision signal in the 3D scene. In contrast, our seg-level labels are annotated per instance and easy to be adopted by annotators when they face the need for labeling.

**Weakly Supervised Image Segmentation.** Many works have been proposed for weakly supervised image segmentation, where image-level supervision [2, 24, 42, 39, 67, 1, 68, 3] and bounding box supervision [10, 32, 28, 19, 28] are two major lines in both instance and semantic segmentation. Image-level labels indicate the semantic classes appearing in an image. A common strategy for image-level supervised segmentation is to train a classification model to recover class activation maps [67, 30, 24, 56] and then use them as pseudo labels to train a segmentation model. In contrast, bounding boxes frame every object with semantic labels, alleviating the need to estimate class activation maps. Segmentation masks can further be refined by heuristic cues [28, 46] or mean-field inference [31, 39]. There are also some other weak supervision forms. For example, Bearman *et al.* [4] proposed point labels to annotate each object. Lin *et al.* [37] proposed scribble labels by dragging the cursor in the center of the objects.

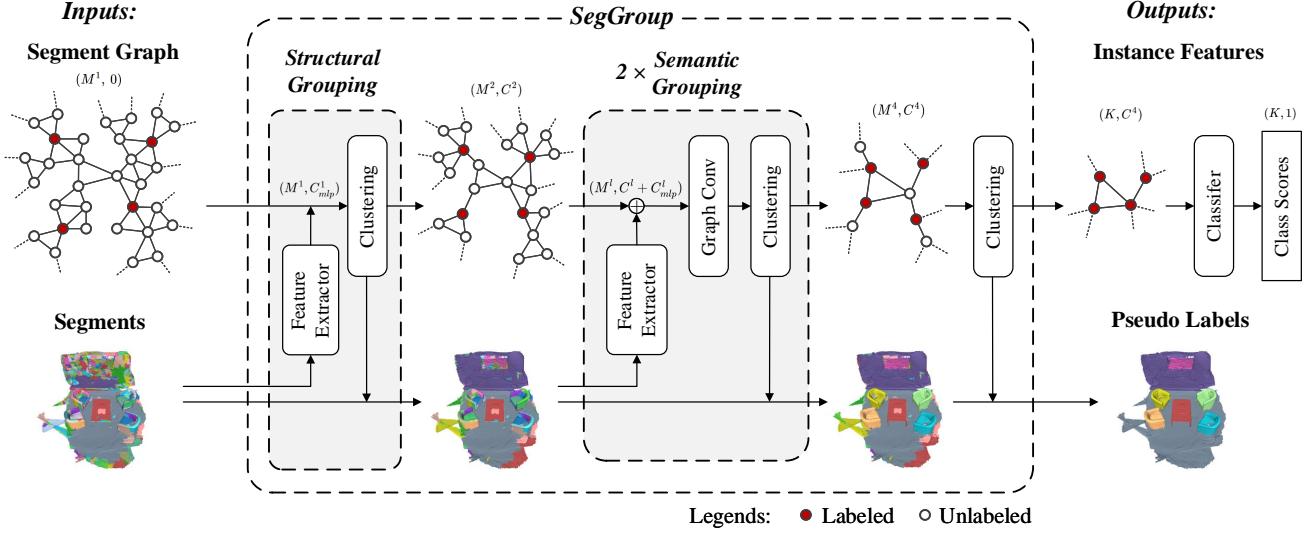


Figure 2. The Structure of our SegGroup network. The inputs of the SegGroup network consist of a segment graph and the corresponding point cloud segments. Each node in the graph is attached with a feature vector, which describes the corresponding point cloud segment. Label information is extended from segments into nodes where the red color indicates labeled nodes. Edges in the graph denote the connected segments that are adjacent in the scene and the lengths of edges represent the similarity of segments. We design three grouping layers followed by a clustering step to group segments into instances by clustering nodes in the graph hierarchically. The classifier in the last part of the framework is used for network training.

### 3. Approach

In this section, we first describe the seg-level annotations for point cloud segmentation. Then, we detail our SegGroup network to generate pseudo point-level labels from the annotated seg-level labels on the training set. Finally, we introduce the implementation details.

#### 3.1. Seg-Level Annotation

Following the labeling strategy of point-level labels in ScanNet [8], we first employ a normal-based graph cut method [13, 27] to perform over-segmentation on the mesh to obtain segments as pre-processing. Unlike 2D images that may suffer from occlusions and lightning variances, 3D data structures usually have clear boundaries between different simple geometry parts thereby easier to perform over-segmentation. The results of over-segmentation remain unchanged in all the subsequent operations, including the manual annotation process and SegGroup network learning. Although some adjacent similar parts that belong to different instances may be mistakenly segmented into one by the over-segmentation step, during the annotation process the annotators found that these effects are small and ignorable.

To annotate seg-level labels, we design a WebGL annotation tool in the browser. Different from point-level labels that annotate every point in the scene, we click on one point per instance to indicate its location and give its se-

mantic class. In order to annotate more precise locations, we ask the annotators to click the point on the most representative segment of the instance.<sup>2</sup> With the help of over-segmentation, these location annotations are automatically extended into segments as seg-level labels. In this paper, we consider the largest segment of an instance as the most representative segment. If it is hard to distinguish the largest segment, e.g. the sizes of some large segments for an instance are almost equal, we allow the annotators to click the point on the most central segment. See supplementary for more annotation details.

#### 3.2. SegGroup

To learn point cloud segmentation models with seg-level labels on the training set, we propose a two-stage approach. We first design a segment grouping network (SegGroup) to generate pseudo labels for the remaining unlabeled points on the training set. Then, we adopt an existing point-level supervised point cloud segmentation model (such as Point-Group [26]) to consume the generated pseudo point-level labels for standard training and evaluate its performance on the testing set. The two stages are trained separately.

The structure of our SegGroup network is shown in Figure 2. To fully utilize the locations of instances in 3D scene to generate pseudo labels, we assume that all segments of a single instance are interconnected so that segments can

<sup>2</sup>If an object instance is divided into several unconnected large parts, we click one point on each of them.

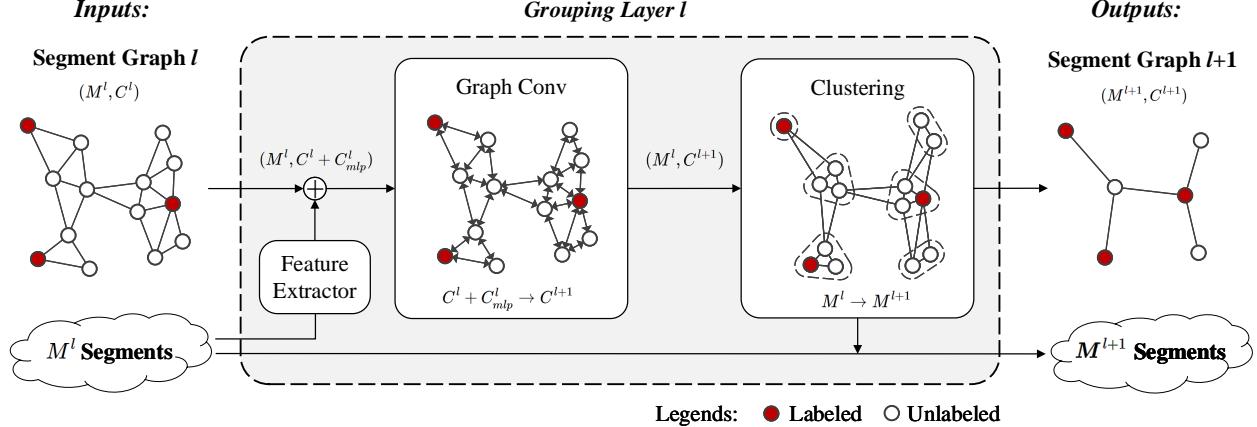


Figure 3. The Structure of the semantic grouping layer. The inputs consist of a segment graph in layer  $l$  as well as the corresponding point cloud segments. Same as the settings in Figure 2, each node in the graph represents a unique point cloud segment. The node features are first concatenated with segment features from point cloud segments by a feature extractor, and then updated by a graph convolution network. Finally, a node clustering algorithm groups similar nodes into new nodes. With the clustering result, similar segments are merged into large segments.

be gathered into instances according to their neighbor relationship under the guidance of instance locations. Therefore, given an over-segmented point cloud scene, we build a segment graph where each node denotes a unique segment and the edge shows the adjacency between two neighboring segments in the scene. We extend the label information of each segment to the corresponding node in the graph, where the red color indicates labeled nodes. As the labeled nodes show the locations of instances, unlabeled nodes of the same instance are gradually grouped into labeled nodes through edges in the graph with our SegGroup network. The final outputs of the SegGroup are one node for each instance and all nodes are labeled with semantic classes and different instance IDs. When the nodes are merged, the segments are also merged into larger segments until all unlabeled segments are merged into nearby labeled segments. All points in the 3D scene are therefore labeled with semantic classes and instance IDs, and we take these labels as pseudo point-level labels. The additional classifier in the last part of the framework is used for network training.

**Grouping Layer.** We design three grouping layers followed by a final clustering process to group segments into instances by clustering nodes in the graph hierarchically. The first layer is a structural grouping layer, whose objective is to group similarly structured segments into one segment to reduce the computational costs of the subsequent grouping layers. The second and the third layer are semantic grouping layers. Because the semantic and structural grouping layers are almost the same except that the semantic grouping layer has an additional graph convolution network, we only introduce the semantic grouping layer in the following content.

The structure of the semantic grouping layer is shown in Figure 3. The inputs of the  $l$ -th layer consist of a segment graph with  $M^l$  nodes as well as its corresponding  $M^l$  point cloud segments. Each node is represented by a  $C^l$  dimensional vector which is considered as the node feature. The output of the grouping layer is a new graph with  $M^{l+1}$  nodes as well as its corresponding  $M^{l+1}$  point cloud segments, where  $M^{l+1} \leq M^l$ . The node features of the output new graph are in  $C^{l+1}$  dimensional.

**Node Feature Acquisition.** We first adopt a shared EdgeConv [54] network to obtain  $C_{mlp}^l$  dimensional segment features for each input point cloud segment of the grouping layer individually, which serves as a local feature learning module to extract semantic information. The EdgeConv network is depicted as “Feature Extracor” in Figure 2 and 3. During the forward-propagation process of the SegGroup network, small segments are gradually merged into large segments, so that the receptive field of the feature extractor in the deeper grouping layer also becomes larger to extract more macroscopic information. The input  $C^l$  dimensional node features are concatenated with their corresponding newly extracted  $C_{mlp}^l$  dimensional segment features to form new node features in  $C^l + C_{mlp}^l$  dimensional.

Then, we update the node features with a graph convolution network (GCN) [29], which is depicted as “Graph Conv” in Figure 2 and 3. Through the neighbor relationship of nodes in the graph, the goal of GCN is to semantically narrow down the difference between nodes belonging to the same instance and expand the difference between nodes belonging to different instances. Given a node feature  $\vec{h}_i^l$  and its adjacent neighbor node feature  $\vec{h}_j^l$  in the  $l$ -th grouping layer, the similarity coefficient  $e_{ij}^l$  is computed according

to the distance between  $\vec{h}_i^l$  and  $\vec{h}_j^l$ :

$$e_{ij}^l = \exp(-\lambda \|\vec{h}_i^l - \vec{h}_j^l\|_2) \quad (1)$$

where  $\lambda$  is a positive parameter to control the slope. A small distance indicates similar node features, which leads to a large coefficient. If  $i = j$ , the coefficient  $e_{ii}^l = 1$ . To make coefficients easily comparable across different nodes, we normalize them across all choices of  $j$ :

$$a_{ij}^l = \frac{e_{ij}^l}{e_{ii}^l + \sum_{k \in \mathcal{N}_i^l} e_{ik}^l} \quad (2)$$

where  $\mathcal{N}_i^l$  is the neighborhood of node  $i$  in the graph. The normalized similarity coefficients are used to compute a linear combination of adjacent segment features. After a shared linear transformation parametrized by a weight matrix  $\mathbf{W}^l$  and a nonlinearity  $\sigma$ , the output node feature  $\vec{h}_i^{l'}$  is finally computed as:

$$\vec{h}_i^{l'} = \sigma \left( a_{ii}^l \mathbf{W}^l \vec{h}_i^l + \sum_{k \in \mathcal{N}_i^l} a_{ik}^l \mathbf{W}^l \vec{h}_k^l \right) \quad (3)$$

In grouping layer  $l$ , the dimension of  $\vec{h}_i^l$  is  $C^l + C_{mlp}^l$ , and the dimension of  $\vec{h}_i^{l'}$  is  $C^{l+1}$ .

**Node Clustering.** In grouping layer  $l$ , we further design a straightforward but effective node clustering algorithm to group similar neighboring nodes into one node. For all neighboring pairwise nodes, if they do not belong to different instances and the distance between node features is below a given threshold  $e_\tau^l$ , we merge the two nodes into a new node. Finally, the algorithm produces  $M^{l+1}$  nodes from  $M^l$  nodes in grouping layer  $l$ . The node features of the new nodes are obtained by a max-pooling operation on their merged node features.

After three grouping layers, most of the nodes are merged. For the few remaining unlabeled nodes, we further perform another node clustering algorithm to group all unlabeled nodes into nearby relevant nodes to generate final instance proposals. Different from the node clustering algorithm in the three grouping layers, this node clustering algorithm traverses all unlabeled nodes. This algorithm adopts a greedy strategy that obtains a locally optimal solution by merging each unlabeled node into the most similar neighbor node in the graph by comparing distances between node features. Finally, all the nodes in the graph are labeled and each instance in the 3D scene corresponds to a unique node. The point cloud segments become our propagated instance proposals which are considered as point-level pseudo labels. Both the two clustering algorithms are depicted as "Clustering" in Figure 2 and 3. More details of the two algorithms are given in the supplementary.

**Network Training.** The SegGroup network outputs a graph containing  $K$  nodes, each of which is attached with a label and a node feature. Because these  $K$  nodes correspond to  $K$  different instances, we extend the label information into instances and consider the node features as instance features. To train the SegGroup network, we adopt a classifier network to obtain a score of labeled semantic class for each instance. A cross-entropy loss is further computed for backward-propagation. In essence, our network training scheme alternates between two steps: (i) forward-propagation: with the network parameters fixed, the SegGroup propagates labels to unlabeled segments as pseudo labels by clustering. The output instance features are further processed by the classifier to obtain semantic scores. (ii) backward-propagation: with the pseudo labels fixed, we use the cross-entropy loss computed by the semantic scores to optimize the parameters of the SegGroup and the classifier. As the network parameters are optimized, the instance features which are gathered by the node features in all grouping layers can capture more semantic information. Similar ideas of this EM-like algorithm are widely-used to effectively refine the generated pseudo labels in weakly- and semi-supervised learning, such as [10, 37, 39].

Although the clustering algorithm may make mistakes and spread to the ambiguous parts with the hierarchical grouping operations, the generated pseudo labels are relatively accurate around the locations of instances that are indicated by seg-level labels. During the training process, the feature extractor and GCN can gradually learn semantically related features for each instance according to the 3D structures around the instance locations. As the features of nodes belonging to the same instance contain more semantic information and become closer in feature space, the clustering algorithm can gather them more correctly with the help of the given instance location and the SegGroup can output better pseudo labels. After such a virtuous cycle in network training, we finally obtain the well-trained pseudo labels and they are used in the strong supervised instance and semantic point cloud segmentation.

### 3.3. Implementation Details

The feature extractor in the first semantic layers extracts 64-dim features with one shared multi-layer perceptron (MLP) layer (64), while the extractor in the second layer adopts two shared MLP layers (64, 64) to obtain 64-dim features. In all GCNs, we set the output feature dimension the same as the input dimension and  $\lambda = 1/8$ . The threshold  $e_\tau^l$  for clustering is set as 6, 2, 2 in each grouping layer. As in every layer the node features are concatenated with new segment features, the output feature dimension of the SegGroup network is  $128 + 64 + 64 = 256$ . Then, the classifier obtains class scores with two fully-connected layers. In the network training, we use SGD [5] to optimize

Table 1. The class-specific semantic IoUs (%) of the generated pseudo labels on ScanNet training set.

Sem IoU	Label	wall floor cab bed chair sofa table door wind bkshf pic cntr desk curt fridg showr toil sink bath ofurn	Mean
MPRM [55]	Scene	47.3 41.1 10.4 43.2 25.2 43.1 21.9 9.8 12.3 45.0 9.0 13.9 21.1 40.9 1.8 29.4 14.3 9.2 39.9 10.0	24.4
MPRM [55]	Subcloud	58.0 57.3 33.2 <b>71.8</b> 50.4 <b>69.8</b> 47.9 42.1 44.9 <b>73.8</b> 28.0 21.5 49.5 <b>72.0</b> 38.8 44.1 42.4 20.0 48.7 34.4	47.4
Layer 1	Seg	37.4 51.5 30.5 20.3 18.7 11.3 35.0 38.5 11.6 9.2 <b>47.6</b> 36.4 22.0 5.5 25.3 14.0 13.5 21.2 21.0 20.2	24.5
Layer 2	Seg	65.5 81.4 44.1 35.0 31.5 22.3 49.0 61.4 36.3 17.6 43.6 49.1 36.7 56.4 41.4 69.1 22.0 27.8 25.2 37.0	42.6
Layer 3	Seg	67.8 80.1 50.1 37.9 43.7 32.0 51.8 64.2 43.2 27.4 43.5 49.2 39.7 61.2 49.3 70.0 31.0 36.1 29.6 48.6	47.8
Layer 4	Seg	67.9 80.1 50.9 38.1 51.2 33.5 54.5 64.3 43.7 28.4 44.9 49.3 40.5 61.4 51.5 70.0 36.0 42.1 30.3 51.2	49.5
SegGroup	Seg	<b>71.0</b> <b>82.5</b> <b>63.0</b> 52.3 <b>72.7</b> 61.2 <b>65.1</b> <b>66.7</b> <b>55.9</b> 46.3 42.7 <b>50.9</b> <b>50.6</b> 67.9 <b>67.3</b> <b>70.3</b> <b>70.7</b> <b>53.1</b> <b>54.5</b> <b>63.7</b>	<b>61.4</b>

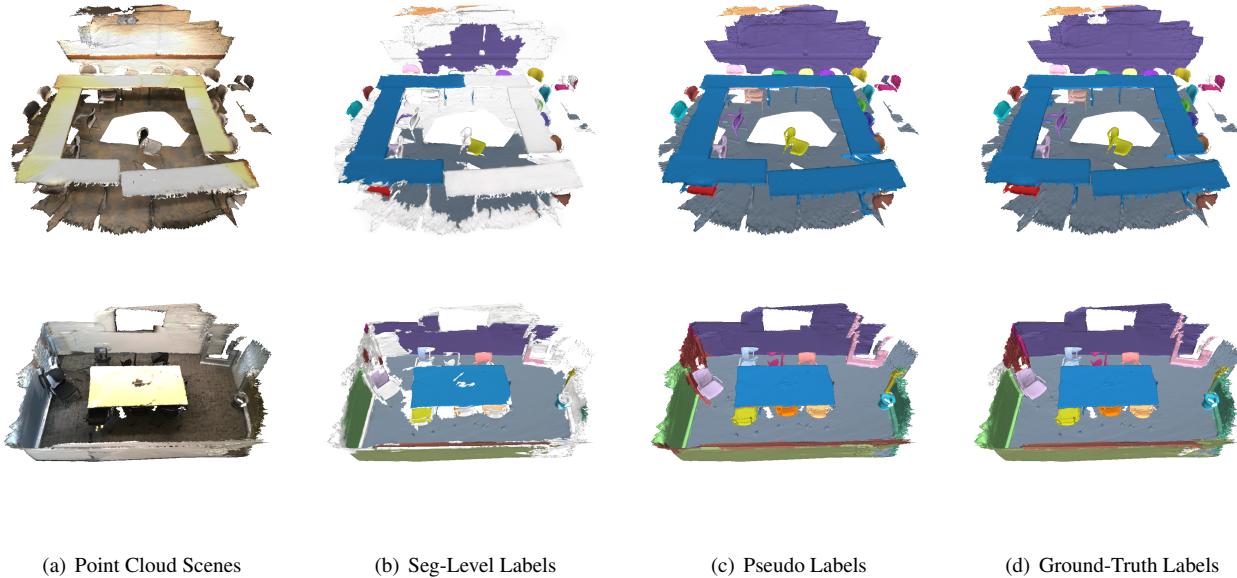


Figure 4. The qualitative visualization results of pseudo label generation. In (b), (c) and (d), various colors indicate different (pseudo) instance labels and points in white are unlabeled. With the benefit of seg-level labels, our SegGroup network can generate qualitative pseudo point-level labels from instance locations. (Best viewed in color.)

the SegGroup network with the learning rate as 0.1 and momentum as 0.9. See supplementary for more implementation details.

## 4. Experiments

In this section, we first evaluate the accuracy of the generated pseudo labels compared with ground truth labels. Then, we use the pseudo labels to train point cloud instance and semantic segmentation networks respectively and test their performance. Finally, we compare the annotation efficiency of different label forms with a fixed labeling time.

We adopted the ScanNet [8] dataset to conduct our experiment. ScanNet is a widely used large-scale real-world indoor 3D dataset, containing 1,201 training scenes, 312 validation scenes, and 100 hidden test scenes. The dataset has 40 object classes, and the evaluation of semantic segmentation is conducted on 20 classes. For instance segmen-

tation, the wall and floor classes are ignored and only 18 classes are used for evaluation. We annotated seg-level labels in all the 40 classes, where we trained the SegGroup network on these labeled classes. We report the results on 18 classes for instance segmentation tasks and 20 classes for semantic segmentation tasks.

### 4.1. Pseudo Label Evaluation

We first evaluated the generated pseudo labels by instance and semantic IoUs. Table 1 show the class-specific semantic IoU results of the pseudo labels on the ScanNet training set. We also tested the pseudo labels generated in every layer of the SegGroup network. As the label information is propagated from labeled segments into unlabeled segments with the number of grouping layers increases, more points in the point cloud scene are annotated with the pseudo labels until all points are annotated. Experiments show that the generated pseudo labels are very

close to the ground truth labels by annotating one point per instance. From the results on different classes, we observe that our method performs significantly better on the classes of *wall* and *floor* which have simple structures as well as *chair*, *sofa*, *table*, *shower*, *curtain* and *toilet* that are easy to identify and separated from a smooth surface. Table 1 also shows the evaluation results of the pseudo labels generated from scene-level labels and subcloud-level labels by MPRM [55], where our generated pseudo labels outperform them by a large gap. The mean IoU at Layer 1 of SegGroup is already higher than that of the scene-level labels by MPRM.

Besides quantitative results that compare our generated pseudo labels with the ground truth, we also show qualitative visualizations for a more intuitive illustration in Figure 4. In ScanNet, there are still a very small number of points remaining unlabeled for point-level annotation, so that we can also observe some white details in Figure 4(d). By only annotating one point per instance for its location, we obtain the labels for the corresponding segments with a number of points. In contrast to point-level labels, our labeling method is very cost-effective. With the benefit of seg-level labels, our SegGroup network can generate qualitative pseudo point-level labels from instance locations. From the visualization results, we observe that the pseudo labels match the ground-truth labels in almost all the areas.

## 4.2. Point Cloud Instance Segmentation

For the evaluation of point cloud instance segmentation task, we employed PointGroup [26] to train a point cloud instance segmentation model based on the generated pseudo point-level labels. Table 2 shows the recent point cloud instance segmentation results on the ScanNet testing set. AP averages the scores with IoU (Intersection of Union) threshold set from 50% to 95% with a step size of 5%, while AP<sub>25</sub> and AP<sub>50</sub> denote the AP scores with IoU threshold set as 25% and 50% respectively. We observe that our seg-level supervised method achieves competitive results with the recent strong supervised methods. However, the seg-level labels only require 1.93 minutes to annotate one scene in ScanNet on average, while the strong point-level labels need 22.3 minutes.

Besides the experiments above, we also compared the efficiency of different labeling strategies given a fixed annotation time budget on the training set. Table 3 shows point cloud instance segmentation results under different supervisions on the ScanNet validation set with different numbers of training scenes. Given the labeling time for the whole training set (1201 scenes) with seg-level labels, only 104 scenes can be annotated with point-level labels. Compared with the Contrastive Scene Contexts (CSC) [18] method which proposes an active labeling strategy (denoted as init+act.) to annotate points, our annotations are per in-

Table 2. Point cloud instance segmentation results (%) on the ScanNet testing set compared with the point-level supervised methods.

Method	Conference	AP	AP <sub>50</sub>	AP <sub>25</sub>
Point-Level Supervision:				
OccuSeg [15]	CVPR'20	48.6	67.2	78.8
PointGroup [26]	CVPR'20	40.7	63.6	77.8
3D-MPA [12]	CVPR'20	35.5	61.1	73.7
MTML [33]	ICCV'19	28.2	54.9	73.1
3D-BoNet [61]	NeurIPS'19	25.3	48.8	68.7
3D-SIS [17]	CVPR'19	16.1	38.2	55.8
GSPN [63]	CVPR'19	15.8	30.6	54.4
SGPN [53]	CVPR'18	4.9	14.3	39.0
Seg-Level Supervision:				
SegGroup (PointGroup)	-	24.6	44.5	63.7

Table 3. Point cloud instance segmentation results (%) on the ScanNet validation set with a fixed labeling time budget on the training set.

Method	Scenes	AP	AP <sub>50</sub>	AP <sub>25</sub>
Point-Level Supervision:				
PointGroup [26]	104	9.4	19.9	36.4
Init+Act. Point Supervision:				
CSC-20 (PointGroup) [18]	1201	-	27.2	-
CSC-50 (PointGroup) [18]	1201	-	35.7	-
Seg-Level Supervision:				
SegGroup (PointGroup)	1201	<b>23.4</b>	<b>43.4</b>	<b>62.9</b>

stance and indicate the locations. In seg-level labels, we annotate 41.2 points on average for each 3D scene.<sup>3</sup> The results show that our seg-level labels obtain much better performance than both point-level labels and init+act. point annotations given the same annotation budget for the instance segmentation task.

## 4.3. Point Cloud Semantic Segmentation

For the point cloud semantic segmentation task, we trained KPConv [51] with the generated pseudo labels for evaluation. Table 4 shows the results of both the recent point-level and subcloud-level methods on the ScanNet testing set. We use mIoU (mean Intersection of Union) as the evaluation metric. The results show that we achieve comparable results with the strong point-level supervised methods. Considering the huge time cost of the point-level labels, our seg-level supervised method is very competitive. Moreover, we outperform the subcloud-level supervised method with the same KPConv network for semantic segmentation, which shows the effectiveness of the seg-level labels.

We also compared the efficiency of different labeling strategies on point cloud semantic segmentation following

<sup>3</sup>Because the CSC method only provides its performance under 20 and 50 point annotations, we list both of them in Table 3 and 5 for reference.

Table 4. Point cloud semantic segmentation results (%) on the ScanNet testing set compared with the point-level and subcloud-level supervised methods.

Method	Conference	mIoU
Point-Level Supervision:		
OccuSeg [15]	CVPR'20	76.4
MinkowskiNet [7]	CVPR'19	73.4
JSENet [21]	ECCV'20	69.9
FusionNet [64]	ECCV'20	68.8
KPConv [51]	ICCV'19	68.4
DCM-Net [47]	CVPR'20	65.8
HPEIN [25]	ICCV'19	61.8
SegGCN [35]	CVPR'20	58.9
TextureNet [22]	CVPR'19	56.6
3DMV [9]	ECCV'18	48.8
PointCNN [36]	NeurIPS'18	45.8
PointNet++ [45]	NeurIPS'17	33.9
ScanNet [8]	CVPR'17	30.6
Subcloud-Level Supervision:		
MPRM (KPConv) [55]	CVPR'20	41.1
Seg-Level Supervision:		
SegGroup (KPConv)	-	61.1

the settings in Sec. 4.2. Table 5 shows the point cloud semantic segmentation results under different supervisions on the ScanNet validation set with a fixed annotation budget on the training set. Given the labeling time for the whole training set with seg-level labels, only 773 scenes and 104 scenes can be annotated with subcloud-level labels and point-level labels, respectively. Compared with the recent method CSC [18], we achieve very competitive performance by using less point annotations (41.2 vs. 50) with a more efficient backbone (KPConv [51] vs. MinkowskiNet [7]). Nevertheless, our SegGroup still outperforms all other methods under the same annotation cost. Combining the results in both Table 3 and 5, our seg-level labels provide an optimal tradeoff between annotation time and segmentation accuracy, which shows that precise instance locations are crucial for 3D scene understanding. Moreover, annotating instance location as seg-level labels can be a low-cost but high-yield labeling manner for 3D instance and semantic segmentation.

#### 4.4. Ablation Study

We also compare the annotation quality of locations between manual labeling by annotators and mechanical labeling from the ground-truth point-level labels. Following the settings in Sec. 4.2, we employed PointGroup [26] to train a point cloud instance segmentation model based on the generated pseudo point-level labels. Table 6 shows the point cloud instance segmentation results under various annotation manners on the ScanNet validation set. For mechanical seg-level labels generated from Top-N segments, we choose to annotate one point randomly in the range of its Top-N

Table 5. Point cloud semantic segmentation results (%) on the ScanNet validation set with a fixed labeling time budget on the training set.

Method	Scenes	mIoU
Point-Level Supervision: KPConv [51]	104	53.2
Subcloud-Level Supervision: MPRM (KPConv) [55]	773	41.5
Init+Act. Point Supervision: CSC-20 (MinkowskiNet) [18]	1201	53.8
CSC-50 (MinkowskiNet) [18]	1201	<b>62.9</b>
Seg-Level Supervision: SegGroup (KPConv)	1201	62.4

Table 6. Point cloud instance segmentation results (%) on the ScanNet validation set supervised by seg-level labels obtained with different annotation manners.

Manner	AP	AP <sub>50</sub>	AP <sub>25</sub>
Manual: Top-1 Segment	23.4	43.4	62.9
Mechanical: Top-1 Segment	24.8	46.2	64.6
Top-2 Segment	25.0	43.9	62.6
Top-3 Segment	24.2	43.0	62.4

largest segments. As human annotators may not accurately choose the largest segment for every instance, the results are similar with the Top-3 segments selected by machines. Moreover, for the labels generated mechanically, the performance is similar for Top-1 and Top-3 labeling strategies. Therefore, seg-level labels are practical for annotators to click on one point as the location on the most representative segment of an instance in 3D scene.

## 5. Conclusion

In this paper, we have explored the importance of instance locations for 3D scene segmentation and design a weakly supervised point cloud segmentation task for full exploitation. More specifically, we click on one point per instance to indicate its location and extend these location annotations into segments as seg-level labels by over-segmentation. We further design a segment grouping network (SegGroup) to generate pseudo point-level labels by grouping seg-level annotated segments into instances hierarchically, so that existing strong-supervised methods can directly consume the pseudo labels for training. Experimental results on both instance and semantic segmentation show that SegGroup effectively generates high-quality pseudo point-level labels from the locations of instances given the seg-level labels, which well balances the annotation cost and segmentation accuracy.

## References

- [1] Jiwoon Ahn, Sunghyun Cho, and Suha Kwak. Weakly supervised learning of instance segmentation with inter-pixel relations. In *CVPR*, pages 2209–2218, 2019. 2
- [2] Jiwoon Ahn and Suha Kwak. Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In *CVPR*, pages 4981–4990, 2018. 2
- [3] Aditya Arun, CV Jawahar, and M Pawan Kumar. Weakly supervised instance segmentation by learning annotation consistent instances. In *ECCV*, pages 254–270, 2020. 2
- [4] Amy Bearman, Olga Russakovsky, Vittorio Ferrari, and Li Fei-Fei. What’s the point: Semantic segmentation with point supervision. In *ECCV*, pages 549–565, 2016. 2
- [5] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *COMPSTAT*, pages 177–186, 2010. 5
- [6] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Fast point R-CNN. In *ICCV*, pages 9775–9784, 2019. 2
- [7] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4D spatio-temporal convnets: Minkowski convolutional neural networks. In *CVPR*, pages 3075–3084, 2019. 1, 2, 8
- [8] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3D reconstructions of indoor scenes. In *CVPR*, pages 5828–5839, 2017. 1, 2, 3, 6, 8, 12
- [9] Angela Dai and Matthias Nießner. 3DMV: Joint 3D-multi-view prediction for 3D semantic scene segmentation. In *ECCV*, pages 452–468, 2018. 1, 8
- [10] Jifeng Dai, Kaiming He, and Jian Sun. BoxSup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In *ICCV*, pages 1635–1643, 2015. 2, 5
- [11] Yueqi Duan, Haidong Zhu, He Wang, Li Yi, Ram Nevatia, and Leonidas J. Guibas. Curriculum DeepSDF. In *ECCV*, pages 51–67, 2020. 1
- [12] Francis Engelmann, Martin Bokeloh, Alireza Fathi, Bastian Leibe, and Matthias Nießner. 3D-MPA: Multi-proposal aggregation for 3D semantic instance segmentation. In *CVPR*, pages 9031–9040, 2020. 2, 7
- [13] Pedro F. Felzenszwalb and Daniel P. Huttenlocher. Efficient graph-based image segmentation. *IJCV*, 59(2):167–181, 2004. 3
- [14] Benjamin Graham, Martin Engelcke, and Laurens Van Der Maaten. 3D semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, pages 9224–9232, 2018. 2
- [15] Lei Han, Tian Zheng, Lan Xu, and Lu Fang. OccuSeg: Occupancy-aware 3D instance segmentation. In *CVPR*, pages 2940–2949, 2020. 2, 7, 8
- [16] Rana Hanocka, Amir Hertz, Noa Fish, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. MeshCNN: a network with an edge. *TOG*, 38(4):1–12, 2019. 1
- [17] Ji Hou, Angela Dai, and Matthias Nießner. 3D-SIS: 3D semantic instance segmentation of RGB-D scans. In *CVPR*, pages 4421–4430, 2019. 1, 2, 7
- [18] Ji Hou, Benjamin Graham, Matthias Nießner, and Saining Xie. Exploring data-efficient 3D scene understanding with contrastive scene contexts. In *CVPR*, 2021. 2, 7, 8
- [19] Cheng-Chun Hsu, Kuang-Jui Hsu, Chung-Chi Tsai, Yen-Yu Lin, and Yung-Yu Chuang. Weakly supervised instance segmentation using the bounding box tightness prior. In *NeurIPS*, volume 32, pages 6586–6597, 2019. 2
- [20] Qingsong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. RandLA-Net: Efficient semantic segmentation of large-scale point clouds. In *CVPR*, pages 11108–11117, 2020. 2
- [21] Zeyu Hu, Mingmin Zhen, Xuyang Bai, Hongbo Fu, and Chiew-lan Tai. JSENet: Joint semantic segmentation and edge detection network for 3D point clouds. In *ECCV*, pages 222–239, 2020. 8
- [22] Jingwei Huang, Haotian Zhang, Li Yi, Thomas Funkhouser, Matthias Nießner, and Leonidas J. Guibas. Texturenet: Consistent local parametrizations for learning from high-resolution signals on meshes. In *CVPR*, pages 4440–4449, 2019. 8
- [23] Qiangui Huang, Weiyue Wang, and Ulrich Neumann. Recurrent slice networks for 3D segmentation of point clouds. In *CVPR*, pages 2626–2635, 2018. 2
- [24] Zilong Huang, Xinggang Wang, Jiasi Wang, Wenyu Liu, and Jingdong Wang. Weakly-supervised semantic segmentation network with deep seeded region growing. In *CVPR*, pages 7014–7023, 2018. 2
- [25] Li Jiang, Hengshuang Zhao, Shu Liu, Xiaoyong Shen, Chi-Wing Fu, and Jiaya Jia. Hierarchical point-edge interaction network for point cloud semantic segmentation. In *ICCV*, pages 10433–10441, 2019. 8
- [26] Li Jiang, Hengshuang Zhao, Shaoshuai Shi, Shu Liu, Chi-Wing Fu, and Jiaya Jia. PointGroup: Dual-set point grouping for 3D instance segmentation. In *CVPR*, pages 4867–4876, 2020. 1, 2, 3, 7, 8
- [27] Andrej Karpathy, Stephen Miller, and Li Fei-Fei. Object discovery in 3D scenes via shape analysis. In *ICRA*, pages 2088–2095, 2013. 3
- [28] Anna Khoreva, Rodrigo Benenson, Jan Hosang, Matthias Hein, and Bernt Schiele. Simple does it: Weakly supervised instance and semantic segmentation. In *CVPR*, pages 876–885, 2017. 2
- [29] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, pages 1–10, 2017. 4
- [30] Alexander Kolesnikov and Christoph H Lampert. Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In *ECCV*, pages 695–711, 2016. 2
- [31] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *NeurIPS*, pages 109–117, 2011. 2
- [32] Viveka Kulharia, Siddhartha Chandra, Amit Agrawal, Philip Torr, and Ambrish Tyagi. Box2Seg: Attention weighted loss and discriminative feature learning for weakly supervised segmentation. In *ECCV*, pages 290–308, 2020. 2
- [33] Jean Lahoud, Bernard Ghanem, Marc Pollefeys, and Martin R. Oswald. 3D instance segmentation via multi-task metric learning. In *ICCV*, pages 9256–9266, 2019. 2, 7

- [34] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *CVPR*, pages 4558–4567, 2018. 2
- [35] Huan Lei, Naveed Akhtar, and Ajmal Mian. SegGCN: Efficient 3D point cloud segmentation with fuzzy spherical kernel. In *CVPR*, pages 11611–11620, 2020. 8
- [36] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhua Di, and Baoquan Chen. PointCNN: Convolution on X-transformed points. In *NeurIPS*, pages 828–838, 2018. 8
- [37] Di Lin, Jifeng Dai, Jiaya Jia, Kaiming He, and Jian Sun. Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In *CVPR*, pages 3159–3167, 2016. 2, 5
- [38] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *CVPR*, pages 4460–4470, 2019. 1
- [39] George Papandreou, Liang-Chieh Chen, Kevin P. Murphy, and Alan L. Yuille. Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In *ICCV*, pages 1742–1750, 2015. 2, 5
- [40] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *CVPR*, pages 165–174, 2019. 1
- [41] Quang-Hieu Pham, Thanh Nguyen, Binh-Son Hua, Gemma Roig, and Sai-Kit Yeung. JSIS3D: joint semantic-instance segmentation of 3D point clouds with multi-task pointwise networks and multi-value conditional random fields. In *CVPR*, pages 8827–8836, 2019. 2
- [42] Pedro O. Pinheiro and Ronan Collobert. From image-level to pixel-level labeling with convolutional networks. In *CVPR*, pages 1713–1721, 2015. 2
- [43] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *CVPR*, pages 652–660, 2017. 1, 2
- [44] Charles R. Qi, Hao Su, Matthias Nießner, Angela Dai, Mengyuan Yan, and Leonidas J. Guibas. Volumetric and multi-view CNNs for object classification on 3D data. In *CVPR*, pages 5648–5656, 2016. 1
- [45] Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, pages 5105–5114, 2017. 1, 2, 8
- [46] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. “GrabCut” interactive foreground extraction using iterated graph cuts. *TOG*, 23(3):309–314, 2004. 2
- [47] Jonas Schult, Francis Engelmann, Theodora Kontogianni, and Bastian Leibe. DualConvMesh-Net: Joint geodesic and euclidean convolutions on 3D meshes. In *CVPR*, pages 8612–8622, 2020. 1, 8
- [48] Yiru Shen, Chen Feng, Yaoqing Yang, and Dong Tian. Mining point cloud local structures by kernel correlation and graph pooling. In *CVPR*, pages 4548–4557, 2018. 2
- [49] Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. SPLATNet: Sparse lattice networks for point cloud processing. In *CVPR*, pages 2530–2539, 2018. 2
- [50] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3D shape recognition. In *ICCV*, pages 945–953, 2015. 1
- [51] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Fleuret, and Leonidas J. Guibas. KPConv: Flexible and deformable convolution for point clouds. In *ICCV*, pages 6411–6420, 2019. 1, 2, 7, 8
- [52] Lei Wang, Yuchun Huang, Yaolin Hou, Shenman Zhang, and Jie Shan. Graph attention convolution for point cloud semantic segmentation. In *CVPR*, pages 10296–10305, 2019. 2
- [53] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. SGPN: Similarity group proposal network for 3D point cloud instance segmentation. In *CVPR*, pages 2569–2578, 2018. 1, 2, 7
- [54] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph CNN for learning on point clouds. *TOG*, 38(5):1–12, 2019. 1, 2, 4
- [55] 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 *CVPR*, pages 4384–4393, 2020. 2, 6, 7, 8
- [56] Yunchao Wei, Jiashi Feng, Xiaodan Liang, Ming-Ming Cheng, Yao Zhao, and Shuicheng Yan. Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In *CVPR*, pages 1568–1576, 2017. 2
- [57] Wenxuan Wu, Zhongang Qi, and Li Fuxin. PointConv: Deep convolutional networks on 3D point clouds. In *CVPR*, pages 9621–9630, 2019. 1, 2
- [58] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Ligang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D shapenets: A deep representation for volumetric shapes. In *CVPR*, pages 1912–1920, 2015. 1
- [59] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional ShapeContextNet for point cloud recognition. In *CVPR*, pages 4606–4615, 2018. 2
- [60] Xun Xu and Gim Hee Lee. Weakly supervised semantic point cloud segmentation: Towards 10x fewer labels. In *CVPR*, pages 13706–13715, 2020. 2
- [61] Bo Yang, Jianan Wang, Ronald Clark, Qingyong Hu, Sen Wang, Andrew Markham, and Niki Trigoni. Learning object bounding boxes for 3D instance segmentation on point clouds. In *NeurIPS*, pages 1–10, 2019. 1, 2, 7
- [62] Jiancheng Yang, Qiang Zhang, Bingbing Ni, Linguo Li, Jinxian Liu, Mengdie Zhou, and Qi Tian. Modeling point clouds with self-attention and gumbel subset sampling. In *CVPR*, pages 3323–3332, 2019. 2
- [63] Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas J. Guibas. GSPN: Generative shape proposal network for 3D instance segmentation in point cloud. In *CVPR*, pages 3947–3956, 2019. 1, 2, 7
- [64] Feihu Zhang, Jin Fang, Benjamin Wah, and Philip Torr. Deep fusionnet for point cloud semantic segmentation. In *ECCV*, pages 644–663, 2020. 8
- [65] Zhiyuan Zhang, Binh-Son Hua, and Sai-Kit Yeung. ShellNet: Efficient point cloud convolutional neural networks us-

- ing concentric shells statistics. In *ICCV*, pages 1607–1616, 2019. 2
- [66] Hengshuang Zhao, Li Jiang, Chi-Wing Fu, and Jiaya Jia. PointWeb: Enhancing local neighborhood features for point cloud processing. In *CVPR*, pages 5565–5573, 2019. 2
- [67] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *CVPR*, pages 2921–2929, 2016. 2
- [68] Yanzhao Zhou, Yi Zhu, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. Weakly supervised instance segmentation using class peak response. In *CVPR*, pages 2209–2218, 2018. 2

## Appendix

### A. Seg-Level Label Annotation

In this paper, we design a WebGL annotation tool in the browser to annotate seg-level labels. Figure 5 shows the interface of our annotation tool, which includes a scene display window on the left and a control panel on the right. The annotator can rotate and pan the scene to browse and annotate seg-level labels by mouse clicking.

To annotate seg-level labels, we design an annotation interface that requires the annotator to label both the semantic class and the instance ID of the location of an instance (denoted as “Annotating from Scratch” in Figure 5(a)). In this annotation interface, the scene is displayed with original scanned colors at the beginning. The annotator needs to choose a semantic class before annotating the location of each instance. If the next instance to annotate shares the same semantic class with the last annotated instance, the annotator only needs to click on “Add” in the control panel to use the last chosen semantic class. To show the over-segmentation results to facilitate annotation, the segment corresponding to the mouse cursor is displayed in red. When a instance location is annotated, the color of the segment that corresponds to the instance location changes from red to a new color to indicate the segment is annotated. Different colors of the annotated segments indicate they belong to different instances. The mouse cursor location (XYZ) on the surface of the scene mesh and the ID of the segment that corresponds to the cursor location are shown at the bottom of the interface. In Figure 5(a), the annotator is preparing to annotate the pillow.

Because scenes in ScanNet [8] dataset have ground-truth point-level labels, in this paper we choose to annotate our seg-level labels based on the ground-truth labels to reduce the annotation difficulty (denoted as “Annotating with GT Labels” in Figure 5(b)). Compared with annotating from scratch in Figure 5(a), in this annotation mode the annotator do not need to annotate the semantic class of each instance location. In Figure 5(b), different colors indicate different instances, and the white color indicates instances are labeled. The scene is displayed with non-white colors at the beginning of the annotation process. The annotator needs to annotate on every instance to make the scene become white in all areas. Same as the interface in Figure 5(a), in Figure 5(b) the segment corresponding to the mouse cursor is displayed in red, and the status of the mouse cursor is shown at the bottom of the interface. Because the ground-truth labels of the scene are given in this annotation mode, the instance information of the segment that corresponds to the mouse cursor location is also displayed at the bottom of the interface. When a instance location is annotated by mouse clicking, the color of the segment that corresponds

to the instance location changes from red to black. At the same time, the color of the instance last labeled turns white to indicate this instance is already labeled. In Figure 5(b), the annotator has just annotated the cabinet and is preparing to annotate the pillow. The annotation status of Figure 5(a) and 5(b) is the same. Compared with annotating from scratch, in this annotation mode only when the annotator annotates a new instance, the annotation of the last instance is completed and the instance number is added in the annotation progress. Therefore, the instance number in Figure 5(b) is one less than the instance number in Figure 5(a).

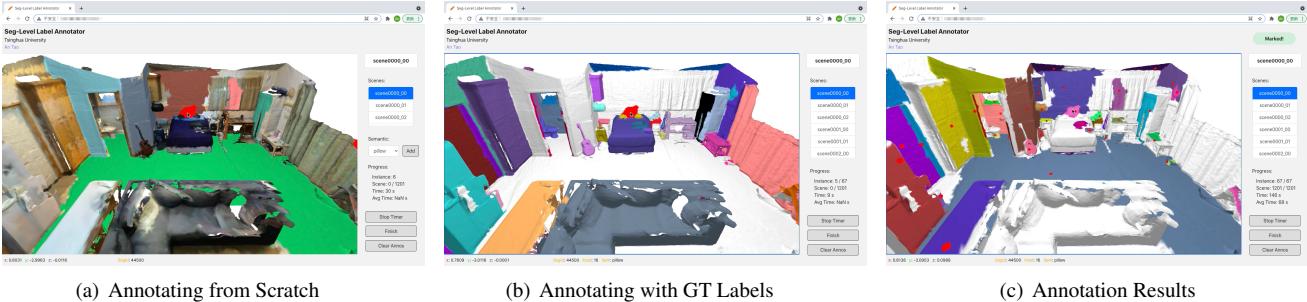
After annotation, the WebGL annotation tool provides an interface to display the annotation results (denoted as “Annotation Results” in Figure 5(c)). The annotation results include location annotations and seg-level labels. In the display window, the positions of red balls indicate the location annotations of instances. For seg-level labels, different colors indicate they belong to different instances. The white areas of the scene are unlabeled.

For annotating the seg-level labels with ground-truth labels, the average annotation time of a scene in ScanNet is 68s. Because we do not annotate semantic class for each instance location, we need to add the semantic annotation time for the standard annotation process from scratch in Figure 5(a). The average time to annotate a semantic class is 1.5s. Since there are 32 instances in average per scene, the total average annotation time is calculated as  $68 + 32 \times 1.5 = 116$ s (1.93 minutes).

### B. Node Clustering Algorithms

Algorithm 1 shows the detailed procedure to produce  $M^{l+1}$  nodes from  $M^l$  nodes in grouping layer  $l$ . Given a graph that has  $M^l$  nodes, we first initialize each node with a separate cluster and also assign the label information (labeled or unlabeled) to the clusters. Then, we gradually merge clusters according to conditions. More specifically, for each edge in the graph, we first check whether the connected two nodes are in the same cluster or their belonged clusters are both labeled. If neither of the two conditions is true, the two clusters do not belong to different instances and then we check whether the distance between the two node features is smaller than a threshold  $e_\tau^l$ . If the distance condition is met, we merge the two clusters and update label information. After finishing the edge traversal, we obtain  $M^{l+1}$  clusters and  $M^{l+1}$  corresponding labels. Then, we convert these clusters into new nodes and extend the label information to them. The node feature of each new node is obtained by a max-pooling operation on old node features of all nodes in each cluster. Before the max-pooling operation, the node features and node-wise distances are unchanged during the clustering process. Finally, we remove the edges inside each cluster to obtain a new graph.

Algorithm 2 shows the detailed procedure to produce  $K$



(a) Annotating from Scratch

(b) Annotating with GT Labels

(c) Annotation Results

Figure 5. The interfaces of our WebGL annotation tool. (a) We design a annotation interface that requires the annotator to label both the semantic class and the instance ID of each chosen segment. (b) Because scenes in ScanNet have ground-truth point-level labels, in this paper we choose to annotate our seg-level labels based on the ground-truth labels to reduce the annotation difficulty. (c) After annotation, the interface displays the annotation results including location annotations and seg-level labels.

nodes from  $M^l$  nodes after all grouping layers, where  $K$  is the number of instances in the 3D scene. For each unlabeled node, this algorithm merges it into the most similar neighbor node in the graph by comparing distances between node features. At each step when two nodes are merged, the label information is updated and the new node feature is computed by a max-pooling operation on the two old node features. The corresponding edge is also removed to form a new graph. In the point cloud scene, the two corresponding segments are also merged into one. The clustering process continues until no unlabeled nodes exist.

## C. Implementation Details

There are slight differences among the three feature extractors in the SegGroup network. For the extractor in the structural grouping layer, we first sample 64 points in each segment with farthest point sampling (FPS) and normalize them into a unit sphere. Therefore, the input 6-dim point cloud vectors (XYZ and RGB) only retain geometry and color information which is used to group similarly structured segments into one segment. After extracting local features for each point in EdgeConv, we conduct both a max-pooling and average-pooling on point clouds to obtain 128-dim features. The feature extractors in semantic grouping layers are standard EdgeConv. Besides XYZ and RGB, we also use an additional 3-dim centered coordinates based on its belonged segment to represent each point.

---

**Algorithm 1** Node Clustering in Layer  $l$ 


---

**Input:** Nodes  $\{n_i^l\}_{i=1}^{M^l}$ , node labels  $\{y_i^l\}_{i=1}^{M^l}$ ,  
 node features  $\{\vec{h}_i^l\}_{i=1}^{M^l}$ , edges  $\{o_s^l\}_{s=1}^{P^l}$ ,  
 distance threshold  $e_\tau^l$ .

**Output:** Nodes  $\{n_i^{l+1}\}_{i=1}^{M^{l+1}}$ , node labels  $\{y_i^{l+1}\}_{i=1}^{M^{l+1}}$ ,  
 node features  $\{\vec{h}_i^{l+1}\}_{i=1}^{M^{l+1}}$ , edges  $\{o_s^{l+1}\}_{s=1}^{P^{l+1}}$ .

```

1: for every node  $n_i^l$  do
2:    $C_i^l \leftarrow \{n_i^l\}$            // initialize clusters
3:    $z_i^l \leftarrow y_i^l$            // initialize cluster labels
4: for every edge  $o_s^l = (n_a^l, n_b^l)$  do
5:   Find  $n_a^l \in C_p^l, n_b^l \in C_q^l$ 
6:   if  $C_p^l == C_q^l$  then
7:     Continue           // belong to same cluster
8:   if  $z_p^l \neq \text{None}$  and  $z_q^l \neq \text{None}$  then
9:     Continue           // both clusters are labeled
10:    if  $\text{dist}(\vec{h}_a^l, \vec{h}_b^l) < e_\tau^l$  then
11:       $C_p^l \leftarrow C_p^l \cup C_q^l$            // merge clusters
12:      Delete  $C_q^l$ 
13:      if  $z_q^l \neq \text{None}$  then
14:         $z_p^l \leftarrow z_q^l$            // update cluster label
15:      Delete  $z_q^l$ 
16: Obtain  $M^{l+1}$  clusters and  $M^{l+1}$  labels
17: Sort clusters and labels into  $\{C_i^l\}_{i=1}^{M^{l+1}}$  and  $\{z_i^l\}_{i=1}^{M^{l+1}}$ 
18: for every cluster  $C_i^l$  do
19:    $n_i^{l+1} \leftarrow C_j^l$            // get node
20:    $y_i^{l+1} \leftarrow z_i^l$            // get node label
21:    $\vec{h}_i^{l+1} \leftarrow \text{maxpool}(\{\vec{h}_k^l\}_{n_k^l \in C_i^l})$  // get node feature
22: for every edge  $o_s^l = (n_a^l, n_b^l)$  do
23:   Find  $n_a^l \in C_p^l, n_b^l \in C_q^l$ 
24:   if  $C_p^l \neq C_q^l$  then
25:      $o_s^{l+1} \leftarrow (n_p^{l+1}, n_q^{l+1})$            // get edge
26: Obtain  $P^{l+1}$  edges
27: Sort edges into  $\{o_s^{l+1}\}_{s=1}^{P^{l+1}}$ 
28: return  $\{n_i^{l+1}\}_{i=1}^{M^{l+1}}, \{y_i^{l+1}\}_{i=1}^{M^{l+1}}, \{\vec{h}_i^{l+1}\}_{i=1}^{M^{l+1}},$   

 $\{o_s^{l+1}\}_{s=1}^{P^{l+1}}$            // return results for next layer

```

---



---

**Algorithm 2** Node Clustering in Final

---

**Input:** Nodes  $\{n_i^l\}_{i=1}^{M^l}$ , node labels  $\{y_i^l\}_{i=1}^{M^l}$ ,  
 node features  $\{\vec{h}_i^l\}_{i=1}^{M^l}$ , edges  $\{o_s^l\}_{s=1}^{P^l}$ .

**Output:** Nodes  $\{n_i^l\}_{i=1}^K$ , node labels  $\{y_i^l\}_{i=1}^K$ ,  
 node features  $\{\vec{h}_i^l\}_{i=1}^K$ , edges  $\{o_s^l\}_{s=1}^Q$ .

```

1: while exist  $y_i^l == \text{None}$  do
2:   for every node  $n_i^l$  that  $y_i^l == \text{None}$  do
3:     Find the neighborhood  $\mathcal{N}_i^l$  of node  $n_i^l$  by edges
4:     Find node  $n_k^l \in \mathcal{N}_i^l$  that minimize  $\text{dist}(\vec{h}_i^l, \vec{h}_k^l)$ 
5:      $n_k^l \leftarrow \{n_i^l, n_k^l\}$ 
6:     Delete node  $n_i^l$ 
7:     Delete node label  $y_i^l$ 
8:      $\vec{h}_k^l \leftarrow \text{maxpool}(\{\vec{h}_i^l, \vec{h}_k^l\})$ 
9:     Delete node feature  $\vec{h}_i^l$ 
10:    Delete edge  $o_s^l = (i, k)$ 
11: Obtain  $K$  nodes, node labels, and node features
12: Sort into  $\{n_i^l\}_{i=1}^K, \{y_i^l\}_{i=1}^K, \{\vec{h}_i^l\}_{i=1}^K$ 
13: Obtain  $Q$  edges
14: Sort into  $\{o_s^l\}_{s=1}^Q$ 
15: return  $\{n_i^l\}_{i=1}^K, \{y_i^l\}_{i=1}^K, \{\vec{h}_i^l\}_{i=1}^K, \{o_s^l\}_{s=1}^Q$ 

```

---