

Window Normalization: Enhancing Point Cloud Understanding by Unifying Inconsistent Point Densities

Qi Wang^{1,2}, Sheng Shi^{1*}, Jiahui Li¹, Wuming Jiang³, Xiangde Zhang¹

¹Northeastern University

²Key Laboratory of Data Analytics and Optimization for Smart Industry (Northeastern University)

³Beijing Eyecool Technology Co., Ltd.

Abstract

Downsampling and feature extraction are essential procedures for 3D point cloud understanding. Existing methods are limited by the inconsistent point densities of different parts in the point cloud. In this work, we analyze the limitation of the downsampling stage and propose the pre-abstraction group-wise window-normalization module. In particular, the window-normalization method is leveraged to unify the point densities in different parts. Furthermore, the group-wise strategy is proposed to obtain multi-type features, including texture and spatial information. We also propose the pre-abstraction module to balance local and global features. Extensive experiments show that our module performs better on several tasks. In segmentation tasks on S3DIS (Area 5), the proposed module performs better on small object recognition, and the results have more precise boundaries than others. The recognition of the sofa and the column is improved from 69.2% to 84.4% and from 42.7% to 48.7%, respectively. The benchmarks are improved from 71.7%/77.6%/91.9% (mIoU/mAcc/OA) to 72.2%/78.2%/91.4%. The accuracies of 6-fold cross-validation on S3DIS are 77.6%/85.8%/91.7%. It outperforms the best model PointNet++-XL (74.9%/83.0%/90.3%) by 2.7% on mIoU and achieves state-of-the-art performance. The code and models are available at <https://github.com/DBDXSS/Window-Normalization.git>.

1. Introduction

3D point clouds have been widely researched. It has rich applications, such as autonomous driving, defect detection, and robotics. A point cloud is a set of coordinates of a real object in the 3D Cartesian coordinate system. The obtained representations include color, reflection intensity, and some others. Different from 2D images, point clouds are sparse and disordered. There are huge differences between differ-

ent observations of the same object. These characteristics differ from 2D images, so point clouds and images are processed differently.

In recent years, many networks have been proposed for processing point clouds. It can be divided into projection-based, voxel-based, and point-based methods. Projection-based methods [5, 18, 21, 23, 35, 37] benefit from image-based methods. They transfer point clouds to 2D images via a projection-based model. Voxel-based methods [6, 13, 28, 33, 34] transform point clouds into regular dense structures and use image methods to further process the dense structures. However, projection-based and voxel-based methods inevitably suffer from the loss of spatial information. Point-based methods [7, 12, 15, 16, 22, 25, 30, 31, 36, 43–45, 47, 48] directly process the raw point cloud. Point-based methods have developed rapidly since PointNet++ [31] proposed the farthest point sampling (FPS) method and hierarchical processing network.

There are still some problems with point-based methods. As shown in Figure 1, the different parts of the point cloud have different point densities. The local neighborhoods obtained by existing methods have different volumes. They share the same network weights, which is adverse to the convergence of the network. Moreover, only using the max-pooling function (MaxPool) can not balance local and global features well. We propose a pre-abstraction group-wise window-normalization (PAGWN) module for this situation to improve the performance of point-based semantic segmentation networks. A better balance between local and global features can be obtained. With the proposed method, it can improve sampling accuracy and reduce the information loss caused by the sampling stage, enabling point-based methods to perform better.

Extensive experiments show that the PAGWN module performs better on several tasks. The proposed module performs better on small object recognition, and the results have more precise boundaries than others. In segmentation tasks on S3DIS (Area 5), the recognition of the sofa and the column is improved from 69.2% to 84.4%



Figure 1. Different parts of the point cloud have different point densities. (a) Whole scene. (b) High-density part. (c) Low-density part.

and from 42.7% to 48.7%, respectively. The benchmarks are improved from 71.7%/77.6%/91.9% (mIoU/mAcc/OA) to 72.2%/78.2%/91.4% and achieve state-of-the-art. The accuracies of 6-fold cross-validation on S3DIS are 77.6%/85.8%/91.7% (mIoU/mAcc/OA). It outperforms the best model PointNeXt-XL [32] (74.9%/83.0%/90.3% (mIoU/mAcc/OA)) by 2.7% on mIoU and achieves state-of-the-art performance. Overall, our contributions are summarized below:

- We study the downsampling stage in a novel way. A general feature extractor PAGWN is proposed for point cloud understanding. It can better balance local and global information and reduce information loss.
- The proposed module performs better on small object recognition, and the segmentation results have more precise boundaries than others. Many existing networks can achieve better performances by using the pre-abstraction group-wise window-normalization module without making any other changes.
- The benchmarks of segmentation tasks on S3DIS (Area 5) are improved from 71.7%/77.6%/91.9% (mIoU/mAcc/OA) to 72.2%/78.2%/91.4%. The accuracies of 6-fold cross-validation on S3DIS are 77.6%/85.8%/91.7%. It outperforms the best model PointNeXt-XL [32] (74.9%/83.0%/90.3%) by 2.7% on mIoU and achieves state-of-the-art performance.

2. Related Work

The methods of processing point clouds can be divided into the following types: Convert the point clouds into 2D images by multi-view projection; Perform regular grid division on the point clouds and transform them to voxels or lattices; Process the point clouds directly based on points. The point-based approaches are based on the advantage of not doing too much processing on the point clouds. They try to keep the original features of the point clouds and avoid information losses caused by data processing. The point-based networks can be divided into point-based multi-layer

perceptron (MLP) networks, point-based convolution networks, RNN-based networks, and graph-based networks. Among the point-based MLP networks, the methods used in local feature aggregation can be classified into pooling-based, attention-based, and transformer-based methods.

Pooling-based Methods The first point-based method is proposed by PointNet [30]. It directly processes point clouds with symmetric functions such as MLP networks and pooling methods. PointNet++ [31] proposes farthest point sampling. Moreover, a hierarchical extraction network is constructed based on PointNet. It gradually expands the receptive field of the network and improves the ability to aggregate local features. The proposal of FPS and the construction of hierarchical networks have brought new inspiration to point-based methods. Point-based methods have developed rapidly since PointNet++. PointSIFT [17] encodes information in eight spatial directions through a three-stage ordered convolutional model. The multi-scale features are spliced to adapt to different scales. PointWeb [47] builds a fully connected graph by exploring the pairing relationships of points within a local region. PointMLP [26] normalizes different samples, and local features are obtained by the concatenation method.

Self-attention-based Methods Attention mechanism [41] is introduced to enhance the local feature extraction ability. Yang et al. [45] propose a grouped random attention model to construct relationships between points. Gumbel Subset Sampling is proposed to replace FPS. This sampling method is insensitive to contours and can select more representative points. Local Spatial Aware layer [4] and Attention-based Score Refinement module [46] are proposed to learn spatial weights of local structure in point clouds. RandLA-Net [15] proposes random sampling for fast sampling and improves sampling efficiency. It uses Cartesian coordinates and points' features for stitching to learn spatial weights, and pooling is used for local feature aggregation.

Transformer-based Methods Following the success of the Transformer structure in the vision domain [2, 3, 8–10,

24, 27, 38–40, 42], many studies [11, 14, 20, 48] have begun to use it for processing point clouds. Engel et al. [11] utilize SortNet and Global Feature Generation to extract local and global features. The global and local features are integrated using local-global attention. Guo et al. [14] use Transformer for feature extraction based on PointNet++. Offset-Attention is proposed to form a residual structure to improve the accuracy further. Zhao et al. [48] propose using positional encoding twice to improve the Transformer effect. Lai et al. [20] propose a StratifiedFormer to obtain long-range contextual information and achieve better performance.

To reduce the information loss, we propose the PAGWN module to unify the point densities. With the PAGWN module, the original network can perform better and maintain its structure unchanged.

3. Methods

3.1. Overview

Downsampling and feature extraction are essential procedures for 3D point cloud understanding. Existing methods are limited by the inconsistent point densities of different parts in the point cloud. The downsampling stage aims to enlarge the receptive field by reducing the point cloud size without losing too much information. There are two steps in the downsampling stage. In the first step, select a certain number of points to represent the entire point cloud, which commonly uses random sampling, FPS, and other methods. The second step is performing local feature aggregation on the sampled point cloud. It aims to enhance the information of the point cloud. Generally, the k -nearest neighbor (KNN) method or ball query (BQ) [31] is used to select neighbor points. Then MLP and MaxPool are used for local feature aggregation. Figure 2 illustrates the different methods: KNN and BQ.

With the development of hierarchical networks, the downsampling stage has been unavoidable in many existing methods to enlarge the receptive field effectively. Suppose P is the given point cloud. $P = \{(c_i, x_i) | i = 1, \dots, N\} \in R^{N \times (3+n)}$. Here c_i is the 3D Cartesian coordinate of the i -th point. x_i is the n -dimensional point cloud feature of the i -th point. N is the number of points contained in the point cloud. In PointNet++ [31], FPS is used for point cloud downsampling. BQ is used for local point cloud extraction. MLP and MaxPool are used for local feature aggregation:

$$\begin{aligned} \hat{x}_i &= \text{MaxPool}\{\text{MLP}([x_{i,j}])\}; \\ \forall x_i \in \text{FPS}(P); x_{i,j} &\in \text{BQ}(x_i), j = 1, \dots, K, \end{aligned} \quad (1)$$

where MLP comprises a fully-connected layer, a batch normalization layer, and an activation function. There are K neighbor points in total. x_i is the i -th point in the sampled point cloud. $x_{i,j}$ is the j -th neighbor of the i -th point.

PointNet++ [31] proposes that the local region obtained by BQ with a fixed scale is beneficial to learn local patterns. However, as shown in Figure 3, BQ cannot solve the problem of different point densities in different parts of the point cloud. If there are too many points in the local region, the amount of computation will be enormous. Furthermore, if there are too few points, the estimate of the center point will be easily biased. So multi-scale grouping (MSG) and multi-resolution grouping (MSR) are further proposed to extract multi-scale features and improve accuracy.

With the increment in model complexity and data scale, MSG and MSR can no longer meet the requirements for computational efficiency. However, using single-scale grouping will reduce the sampling accuracy. So, Point Transformer [48] and StratifiedFormer [20] use KNN to replace BQ:

$$\begin{aligned} \hat{x}_i &= \text{MaxPool}\{\text{MLP}([x_{i,j}])\}; \\ \forall x_i \in \text{FPS}(P); x_{i,j} &\in \text{KNN}(x_i), j = 1, \dots, K. \end{aligned} \quad (2)$$

However, the local region sampled by KNN has no fixed scale. Different parts with different volumes share the same network weights, which is adverse to the convergence of the network. Therefore, it is necessary to propose a more efficient sampling method based on KNN, which motivates us to conduct further studies to improve the accuracy of the downsampling stage. We wish to develop a general downsampling scheme that works for segmentation networks.

3.2. Window-normalization

Let us consider the problem of different scales of local regions obtained by KNN. Existing normalization methods include Layernorm and Batchnorm. The Batchnorm balances the scale between different samples, and the Layer-norm normalizes the features of a single instance. They both can not solve the problem. So, we propose a new normalization method—window normalization with prior knowledge.

For a point cloud P with N points, a point cloud Q ($Q \subset P$) with M ($M < N$) points is firstly obtained by the sampling stage. Then, the K nearest points to $x \in Q$ are selected from the original point cloud P as the local region to improve the sampling accuracy. The feature aggregation is performed on the local region to obtain a new sampled point \hat{x} . As shown in Figure 4, the local region with K points is obtained by the KNN method. Due to the different densities of different parts in the point cloud, the volumes occupied by the K points are quite different. Most existing aggregation methods use the MLP network to learn local features directly and use the max-pooling function for feature aggregation. However, different local regions sharing the same weight will produce aggregation loss.

Neither Batchnorm nor Layernorm can solve this problem. Batchnorm unifies the scales of different samples, which cannot solve this problem. Layernorm normalizes

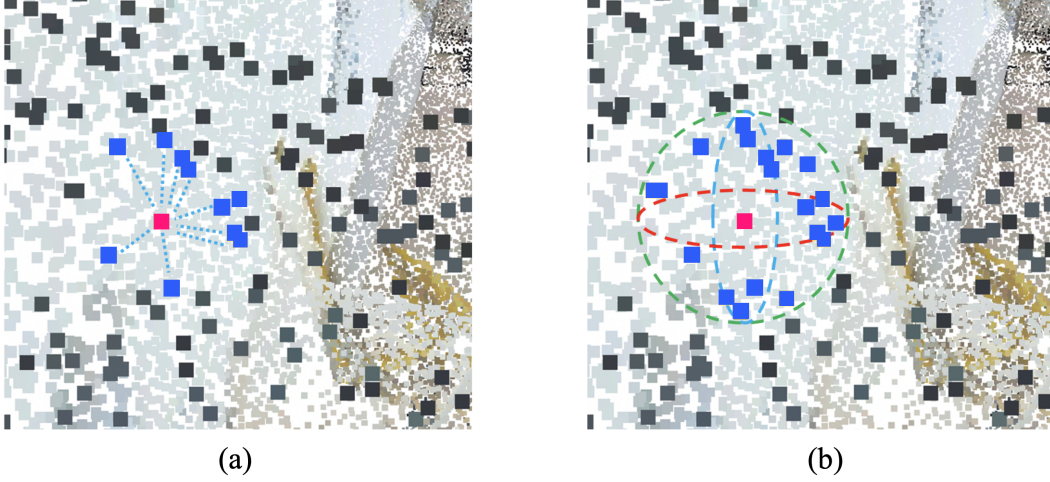


Figure 2. Different methods to obtain neighbor points. The neighborhood obtained by the KNN method has a fixed number of points but cannot limit the neighborhood volume. The neighborhood obtained by the BQ method has a fixed volume but no fixed number. (a) uses the k -nearest neighbor method to gain local region, while (b) uses ball query.

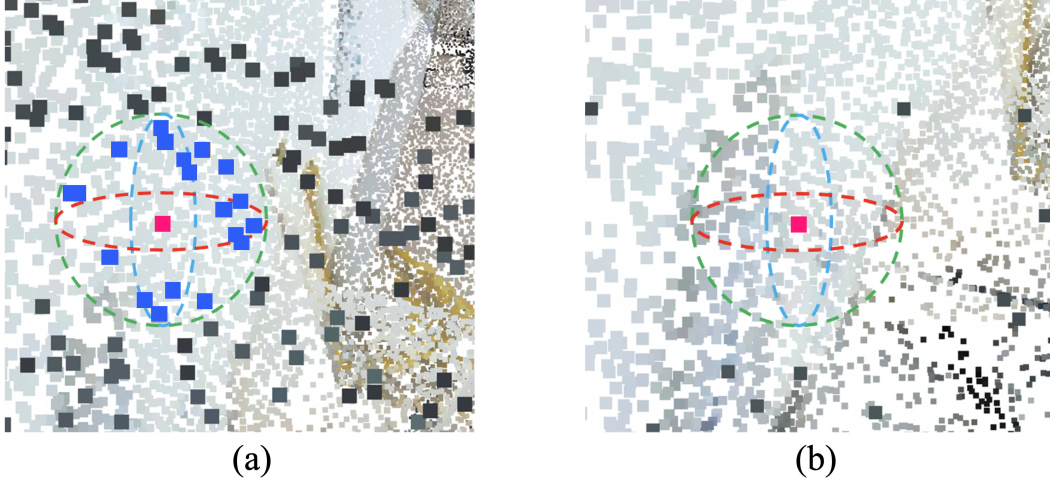


Figure 3. Different parts with different point densities have a different number of neighbor points with BQ on a fixed scale. (a) The BQ method will produce more neighbor points in high-density parts. (b) If the density is too low, there will be fewer, even no, points in the local region.

individual samples. Local regions need to be normalized as a single sample to solve regional scale inconsistency.

In addition to the errors caused by different neighborhood volumes, the normalization method also needs further improvement. Affected by KNN, the center of the sampled region is probably not the sampling point. Expectations will be biased. As shown in Figure 5, the K neighbor points sampled from the point cloud are all on one side of the center point. They are not evenly distributed around the center point. The features obtained by center normalization will deviate from the center point.

Inspired by Layernorm and Batchnorm, a window normalization method that uses prior knowledge is proposed.

It is used to normalize the features of the local region. It can accelerate the learning process and improve the network’s accuracy. We normalize the neighborhoods obtained by different sampling regions and achieve the unification of neighborhood features. Based on prior knowledge, the center point is taken as the mean point instead of using the mean value of point features in the neighborhood. It can maintain global features and reduce the offset problem caused by sampling. As shown in Figure 6, x_i is the center point. x'_i is the mean of the sampling points. $x_{i,k}$ is the k -th neighbor point of the i -th center point. Window-normalization (windownorm) follows the normalization method, where the features are subtracted from the expectation and divided by the

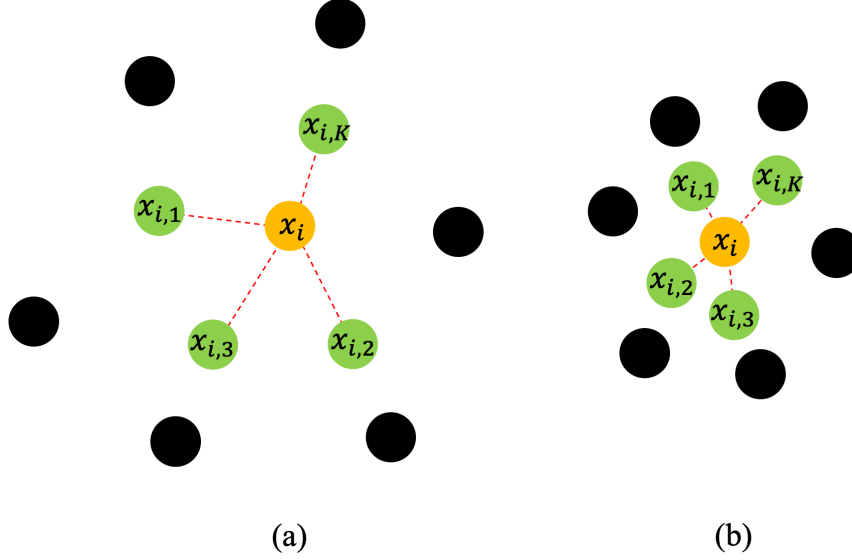


Figure 4. Apply KNN on point clouds with different point densities. Neighbors obtained in low-density regions have a larger volume but the volume is smaller in high-density regions. (a) Low-density region. (b) High-density region.

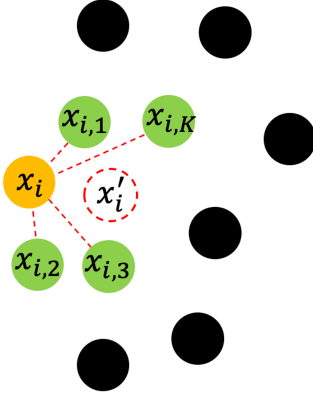


Figure 5. The expectation of the local region by KNN is different from the center point.

standard deviation. The expectation within the local region should be the feature of the sampling point. The variance used by the normalization is the standard deviation of all features in the local region. The normalization formula can be expressed as:

$$\hat{x}_{i,j} = \frac{x_{i,j} - x_i}{\sigma + \varepsilon}, \quad \sigma = \sqrt{\frac{1}{K \times n - 1} \sum_{j=1}^K (x_{i,j} - x_i)^2}; \quad (3)$$

$$\forall x_i \in FPS(P); x_{i,j} \in KNN(x_i), j = 1, \dots, K,$$

where σ is the standard deviation. $\varepsilon = 1e - 5$ is a small constant used to guarantee the stability of the calculation.

n is the feature dimension. x_i is the center point. $x_{i,j}$ is the j -th neighbor of x_i . Window normalization solves the problem of different window scales without any learnable parameters. The window itself determines the different window expectations and variances. The calibration formulae can be described as follows:

$$\begin{aligned} x_{i,j}^* &= \hat{x}_{i,j} + x_i = \frac{x_{i,j} - x_i}{\sigma + \varepsilon} + x_i \\ &= \frac{x_{i,j}}{\sigma + \varepsilon} + \left(1 - \frac{1}{\sigma + \varepsilon}\right) \cdot x_i \\ &= \lambda \cdot x_{i,j} + (1 - \lambda) \cdot x_i, \quad \lambda = \frac{1}{\sigma + \varepsilon}; \end{aligned} \quad (4)$$

$$\begin{aligned} E(x_{i,j}^*) &= \frac{E(x_{i,j}) - x_i}{\sigma + \varepsilon} + x_i \\ &= \lambda \cdot E(x_{i,j}) + (1 - \lambda) \cdot x_i; \end{aligned} \quad (5)$$

$$\begin{aligned} Var(x_{i,j}^*) &= \left(\frac{1}{\sigma + \varepsilon}\right)^2 \cdot Var(x_{i,j}) \\ &= \lambda^2 \cdot Var(x_{i,j}). \end{aligned} \quad (6)$$

where $x_{i,j}^*$ is the rectified point of $x_{i,j}$.

It can be seen that if the original feature is unbiased, the corrected feature will still be unbiased. If the standard deviation is larger, the local point cloud would be more likely to be sparse. Moreover, the rectified neighbor point features will be more unbiased and practical. Furthermore, if there are few points around the center point, less distinguishable information can be provided in the neighborhood of the center point, leading to difficulties in classification. The proposed method can improve the unbiasedness and

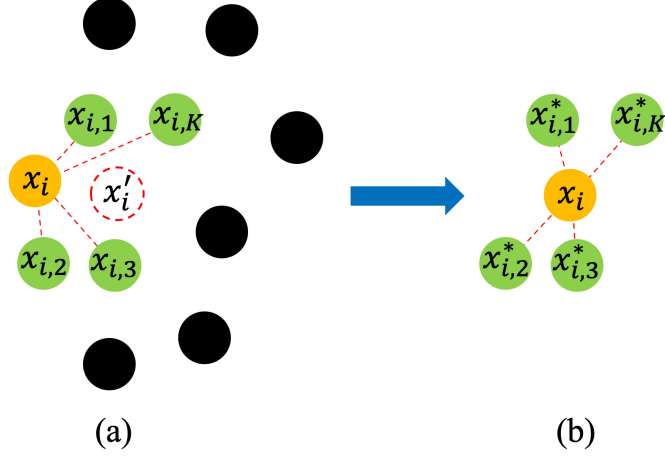


Figure 6. We use the window-normalization method with prior knowledge to calibrate local features. (a) Original expectation by KNN. (b) Calibrated expectation by proposed method

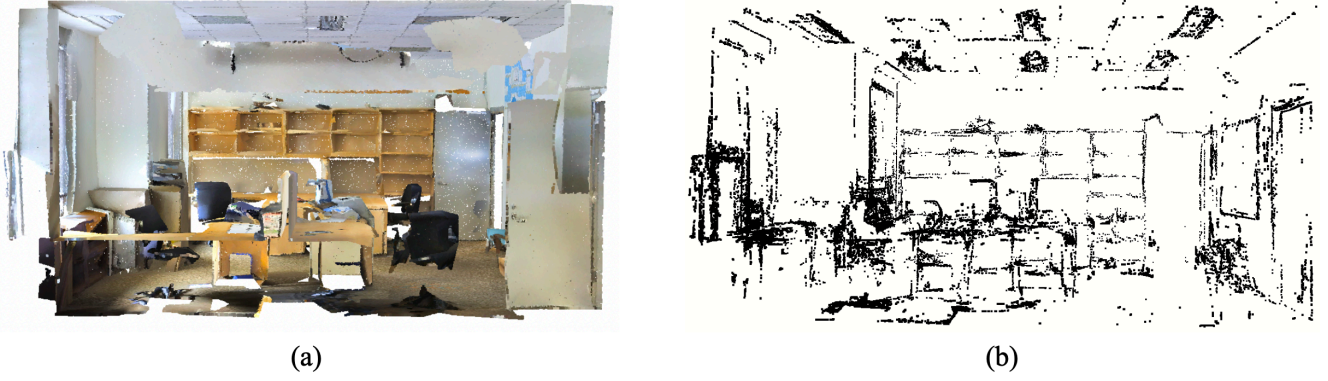


Figure 7. The standard deviation of different regions in the point cloud. The center points with the standard deviation of neighbors larger than 1 are displayed. (a) A point cloud with coordinates and RGB features of a scene in S3DIS (Area 5). (b) Points with the standard deviation of neighbors larger than 1.

effectiveness of such regions, which is beneficial for segmenting complex objects.

To show the sparsity in natural scenes, we calculate the standard deviation of different regions in the point cloud. The center points with the standard deviation larger than 1 are displayed in Figure 7. As can be seen, the extracted points are almost boundaries, which is difficult for semantic segmentation.

3.3. Group-wise Window-normalization

As shown in Figure 8, we argue that points in the neighborhood express two different kinds of feature information: texture information that enhances the local features of the center point and spatial information that increases the receptive field. Therefore, we divide the neighbor points into two groups for windownorm. The m neighbors closest to the center point are used as the first set of points. Besides, the remaining neighbor points are used as the second set

of points. Windownorm is used for the two sets of points respectively.

With the group-wise Window-normalization module, the scales of two sets of points are unified. The first set of points enhances the texture information. And the second set provides a larger receptive field and richer spatial information. The formula is:

$$\begin{aligned} \hat{x}_{i,j} &= \frac{x_{i,j} - x_i}{\sigma + \varepsilon}; \forall x_i \in FPS(P); \\ x_{i,j} &\in KNN(x_i), j = 1, \dots, K; \\ \sigma &= \begin{cases} \sqrt{\frac{1}{m \times n - 1} \sum_{j=1}^m (x_{i,j} - x_i)^2}; \\ \sqrt{\frac{1}{(K-m) \times n - 1} \sum_{j=m+1}^K (x_{i,j} - x_i)^2}. \end{cases} \end{aligned} \quad (7)$$

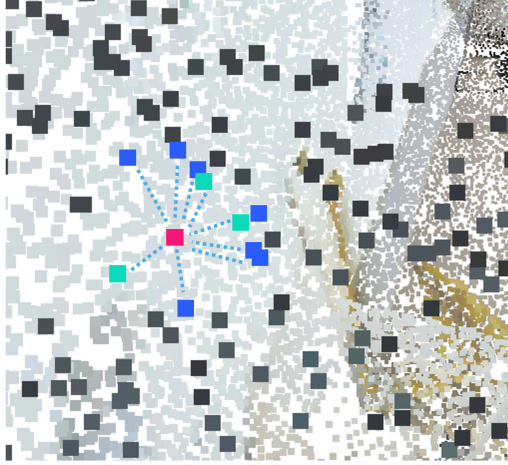


Figure 8. In the local region, green points are closer to the center point, providing texture information that enhances the local features. Blue points are far from the center point, providing spatial information with a bigger receptive field.

3.4. Pre-Abstraction Group-wise Window-normalization Module

When solving semantic segmentation problems, we usually supplement the spatial information of points with position encoding. However, position encoding and sampling points are both used as spatial information to enhance the features of the center point. The position encoding makes the feature dimension of neighbors larger than that of the center point. It leads the model learning to be more biasedness toward position information. We hope the sampled features include local and global information in the sampling stage.

To balance local and global features, we propose a pre-processing module. It can pre-aggregate the features of location encoding and neighborhood points. The model can be expressed as:

$$\hat{x}_{i,j} = LB(GWN([c_{i,j}, x_{i,j}])); \forall x_i \in FPS(P); \quad (8)$$

$$(c_{i,j}, x_{i,j}) \in KNN(x_i), j = 1, \dots, K,$$

where GWN is the group-wise window-normalization module, $[\cdot]$ is the concatenation operation. LB consists of a linear layer and a Batchnorm layer. The input and output feature dimension of the linear layer is $n+3$ and n , respectively. The pre-processing module efficiently reduces the spatial feature dimension to n . So, the neighborhood feature dimension is the same as the center point.

The structure of the PAGWN module is shown in Figure 9. The following formula can express the overall structure

of the proposed sampling method:

$$x_{i,j}^* = LB_2([LB_1(GWN([c_{i,j}, x_{i,j}]))], x_i])$$

$$\hat{x}_i = ReLU\left(\text{Maxpool}\left\{x_{i,j}^*\right\}_{j=1,\dots,K}\right); \forall x_i \in FPS(P); \quad (9)$$

$$(c_{i,j}, x_{i,j}) \in KNN(x_i), j = 1, \dots, K,$$

where the structure of LB_1 and LB_2 is the same as that of LB . The input feature dimension of LB_1 is $n+3$, and the output feature dimension is n . LB_1 is the pre-processing of local spatial information. The input and output feature dimensions of LB_2 are $2n$. LB_2 is the aggregation of center and neighbor point information, which is used to aggregate local and global information.

4. Experiments

4.1. Semantic Segmentation on S3DIS (Area 5)

We evaluate the effectiveness of the proposed PAGWN module on several networks and tasks. The Point Transformer [48] and StratifiedFormer [20] are used as the backbone, respectively.

Datasets The S3DIS dataset contains 271 rooms from three different buildings in six regions with 13 categories [1]. We test our networks on Area 5 and train on other areas. One RTX 3080Ti GPU is used to train networks. Mean intersection over union (mIoU), mean of class-wise accuracy (mAcc), and overall pointwise accuracy (OA) are used as evaluation metrics.

Setup We use Point Transformer and StratifiedFormer as the backbone, respectively. All settings of the dataset are the same as the backbone. We replace its downsampling module with our PAGWN module. The rest parameters are set according to the repositories [19, 29].

Results We re-experiment Point Transformer and StratifiedFormer to ensure the results are compared under the same conditions. After that, we experiment with the Point Transformer first. The accuracies are improved from 70.6%/77.1%/90.5% (mIoU/mAcc/OA) to 71.4%/77.9%/91.1%. It surpasses PointNeXt [32] to reach second. Furthermore, based on the StratifiedFormer, the recognition of the sofa and the column are improved from 69.2% to 84.4% and from 42.7% to 48.7%, respectively. The accuracies on StratifiedFormer are improved from 71.7%/77.6%/91.9% (mIoU/mAcc/OA) to 72.2%/78.2%/91.4% and achieve state-of-the-art. Figure 10 compares the quality of the semantic segmentation results of Point Transformer, StratifiedFormer, and our model. The specific results are shown in Table 1. Extensive experiments show that networks can achieve better performance with the PAGWN module without modifying the other network structure, parameters, and training strategy. The proposed module performs better on small object recognition, and the results have more precise boundaries than others.

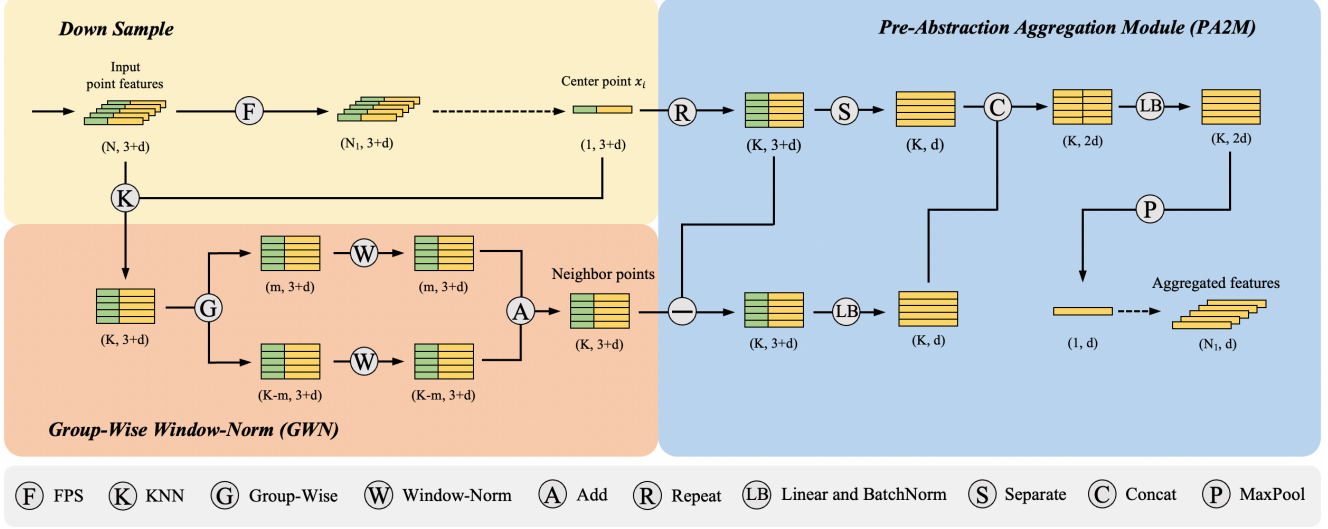


Figure 9. Illustration of our pre-abstraction group-wise window-normalization module. GWN module normalizes the volume of the local region and provides a variety of information. PAGWN module pre-aggregates the features of location encoding and neighborhood points to balance local and global features.

Table 1. Results of different models on S3DIS (Area 5, %).

Method	mIoU	mAcc	OA	ceiling	floor	wall	beam	column	window	door	table	chair	sofa	bookcase	board	clutter
PointNet [30]	41.1	49.0	-	88.8	97.3	69.8	0.1	3.9	46.3	10.8	59.0	52.6	5.9	40.3	26.4	33.2
PAT [45]	60.1	70.8	-	93.0	98.5	72.3	1.0	41.5	85.1	38.2	57.7	83.6	48.1	67.0	61.3	33.6
PointWeb [47]	60.3	66.6	87.0	92.0	98.5	79.4	0.0	21.1	59.7	34.8	76.3	88.3	46.9	69.3	64.9	52.5
KPConv [38]	67.1	72.8	-	92.8	97.3	82.4	0.0	23.9	58.0	69.0	81.5	91.0	75.4	75.3	66.7	58.9
Point Transformer [48]	70.4	76.5	90.8	94.0	98.5	86.3	0.0	38.0	63.4	74.3	89.1	82.4	74.3	80.2	76.0	59.3
PointNeXt [32]	70.5	76.8	90.6	94.2	98.5	84.4	0.0	37.7	59.3	74.0	83.1	91.6	77.4	77.2	78.8	60.6
StratifiedFormer [20]	72.0	78.1	91.5	-	-	-	-	-	-	-	-	-	-	-	-	-
Point Transformer (ours)	70.6	77.1	90.5	93.9	97.6	85.7	0.0	44.7	61.0	78.4	82.8	90.3	75.1	72.2	77.7	58.6
+ PAGWN	71.4	77.9	91.1	94.7	98.3	87.2	0.0	37.2	62.4	75.9	82.8	91.0	81.9	74.0	81.5	60.7
StratifiedFormer (ours)	71.7	77.6	91.9	95.5	98.7	87.3	0.0	42.7	63.6	74.4	85.7	91.0	69.2	79.0	80.3	65.0
+ PAGWN	72.2	78.2	91.4	95.0	98.1	85.9	0.0	48.7	62.0	70.0	82.7	92.1	84.4	78.3	76.8	64.0

It can be seen that the results of our improved module are closer to the ground truth. From all three scenarios, we can find that the PAGWN module performs better in recognizing large contiguous regions such as the ceiling, wall, and window. In addition, our model has apparent advantages in recognizing small objects. With the proposed PAGWN module, Point Transformer can recognize the smoke detector on the ceiling, which could not be recognized before. Furthermore, only our proposed module can locate the clock, cabinet, and bookcase in all tested networks. We can also find that our proposed module has improved the recognition of object boundaries. The recognized objects have more precise boundaries and are closer to the ground truth, which can be seen from the frame in the second scene and the pillow in the third scene. All this shows that the PAGWN module can better balance and utilize local and global information.

It can be found from Table 1 that the mIoU and mAcc of StratifiedFormer have increased while the OA has decreased. From each category's results, our module has significantly improved the recognition of sofa from 69.2% to 84.4% and column from 42.7% to 48.7%. However, the recognition of the board and door is not suitable. Considering that sofa and column are more complex categories to be recognized, the PAGWN module is better for hard-to-classify items than for easy-to-classify items. So, it leads to the rise of mAcc and the decline of OA. As shown in Figure 6 and Figure 10, our method performs better in the boundary region where the σ is larger, which is essential for point cloud understanding.

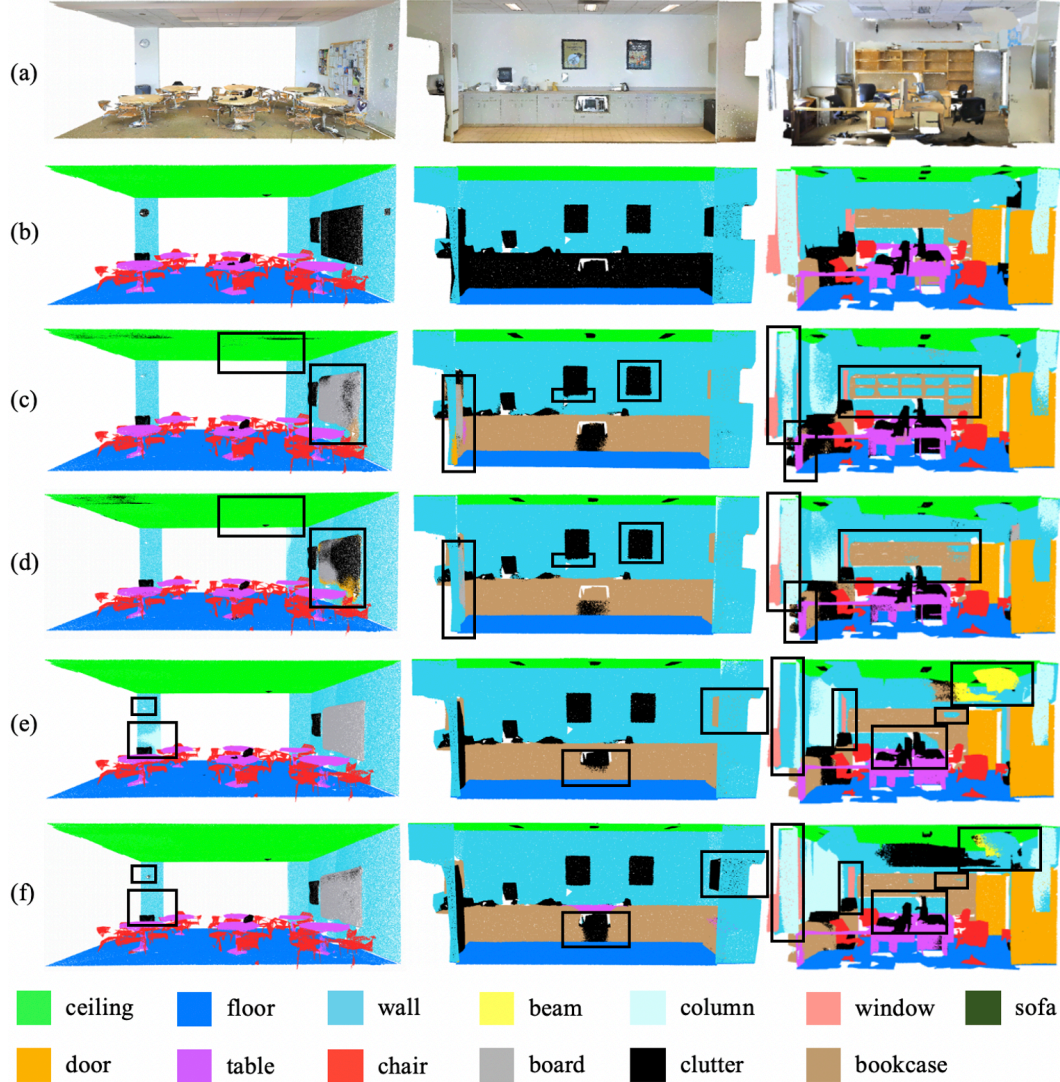


Figure 10. Quality comparison on three different scenes of S3DIS (Area 5). There are thirteen different kinds of categories colored by different colors. We use black boxes to highlight where our model is better than the original. The proposed module performs better on small object recognition, and the results have more precise boundaries than others. Only our proposed module can locate the clock, cabinet, and bookcase in all tested networks. (a) Point clouds with coordinates and RGB features. (b) Ground truth Semantic labels. (c) Point Transformer. (d) Point Transformer with PAGWN. (e) StratifiedFormer. (f) StratifiedFormer with PAGWN.

4.2. Semantic Segmentation on S3DIS (6-fold Cross-validation)

The 6-fold cross-validation experiments on S3DIS are conducted to prove our module’s effectiveness further. Point Transformer [48] and StratifiedFormer [20] are also kept as the backbone. We keep all settings and networks remained. Table 2 shows the details. As can be seen, with the proposed PAGWN module, the performances have been raised to 77.6%/85.8%/91.7% (mIoU/mAcc/OA). It outperforms the best model PointNeXt-XL [32] (74.9%/83.0%/90.3% (mIoU/mAcc/OA)) by 2.7% on mIoU

and achieves state-of-the-art performance.

Figure 11 illustrates the effects brought in different areas. Table 3 shows the exact results. It can be find that among all six areas, our method brings improvement in four of them. These four areas are the four with lower accuracy. As the difficulty of recognition increases, the improvement brought by our method also becomes larger. The mIoU of area 2 is originally 59.9%, which is the lowest among all six areas, and our method brings an obvious improvement of 10.8%. This shows that our method is more effective for complex sample recognition, which is consistent with our previous conclusions.

Table 2. Results of different models on S3DIS (6-fold cross-validation).

Method	mIoU (%)	mAcc (%)	OA (%)
PointNet [30]	47.6	66.2	78.5
PointWeb [47]	66.7	76.2	87.3
KPConv [38]	70.6	79.1	-
Point Transformer [48]	73.5	81.9	90.2
PointNeXt-XL [32]	74.9	83.0	90.3
Point Transformer + PAGWN	74.1	82.5	90.2
StratifiedFormer + PAGWN	77.6 (+2.7)	85.8 (+2.8)	91.7 (+1.4)

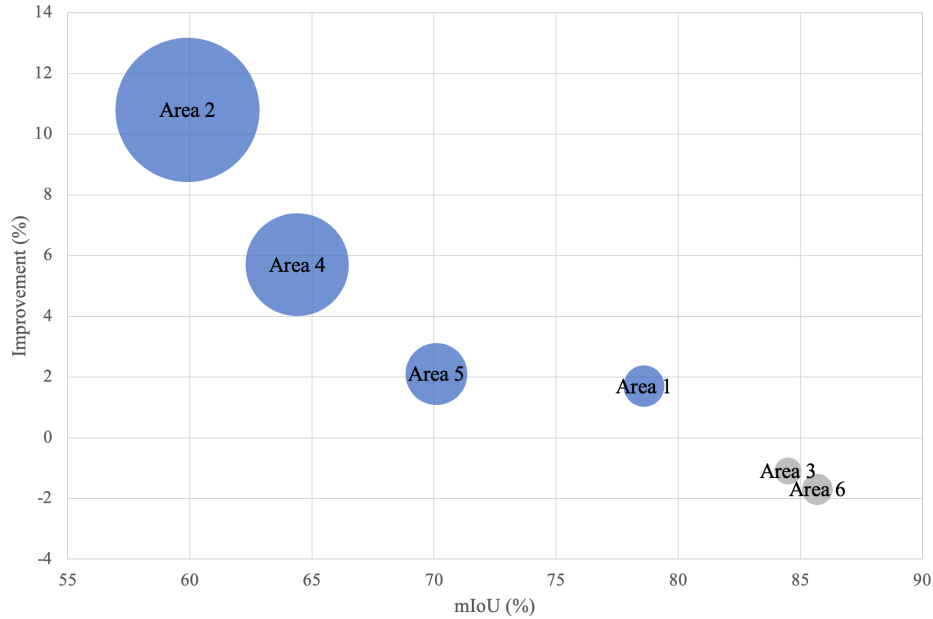


Figure 11. The effects we brought in different areas. There are six areas in the S3DIS dataset. Our method brings improvement to four of them whose accuracy is lower. As the difficulty of recognition increases, the improvement brought by our method also becomes larger. The mIoU of area 2 is originally 59.9%, which is the lowest among all six areas, and our method brings an obvious improvement of 10.8%. The blue ball means that we improve the accuracies and the gray one means the opposite. The bigger the ball is, the larger the effectiveness we bring.

Table 3. 6-fold cross-validation results (mIoU, %) on different areas in S3DIS.

Method	Area 1	Area 2	Area 3	Area 4	Area 5	Area 6
PointNeXt-XL [32]	78.6	59.9	84.5	64.4	70.1	85.7
StratifiedFormer + PAGWN	80.3 (+1.7)	70.7 (+10.8)	83.4 (-1.1)	70.1 (+5.7)	72.2 (+2.1)	84.2 (-1.5)

4.3. Ablation Study

We use Point Transformer as the backbone to conduct ablation studies for each module. The necessity of each module is analyzed. Table 4 shows the results of different components.

Window-normalization and Pre-Abstraction The role of windownorm is to map the neighborhood into the spherical unit space. The neighborhood volume obtained by its K

nearest points differs for different center points. Therefore, the scale and density of the point cloud greatly influence the neighborhood. If the point cloud features are directly learned without normalization, different neighborhoods will probably have different volumes. Since neighborhoods with different sizes share the same network weights, the model sacrifices the accuracy in each region to keep the efficiency of the whole point cloud. However, windownorm can ef-

Table 4. Ablation studies: Different components.

Component	mIoU (%)	mAcc (%)	OA (%)
Point Transformer	70.4	76.5	90.8
+ Window-Norm and Pre-Abstraction	70.9 (+0.5)	77.1 (+0.6)	91.1 (+0.3)
+ Group-wise	71.4 (+1.0)	77.9 (+1.4)	91.1 (+0.3)

fectively solve this problem. Besides, it can also solve the problem of sampling offset in the sampling stage. The sampling region can be effectively constrained around the center point.

Furthermore, it can be observed that the neighbor and center features have the same dimension. The Cartesian coordinates are used as a supplement to the spatial information of neighbor points, which makes the neighbor feature dimension larger than that of the center point. Neighbor information will have a greater weight in the network. In the downsampling stage, the network learns local features and maintains global features. The pre-abstraction module is used to pre-abstract the local feature, ensuring neighbor points have the same status as the center point. In this way, local and global features can achieve a better balance. With the constraints of the pre-abstraction module, windownorm can work better. The two modules cooperate to make the model perform better. As can be seen from Table 4, the model accuracies are improved from 70.4%/76.5%/90.8% (mIoU/mAcc/OA) to 70.9%/77.1%/91.1% with windownorm and pre-abstraction modules. The improvement of model accuracy shows that both windownorm and pre-abstraction modules are effective.

Group-wise The points closer to the center point are considered to provide more local texture information. Meanwhile, points farther away from the center provide spatial information with a larger receptive field. We take the m nearest points as the first set that provides local information. The remaining neighbor points serve as a second set that provides spatial information. Use group-wise windownorm for each of the two sets of points. A model that does not use a grouping strategy can be expressed as:

$$\begin{aligned}
 x_{i,j}^* &= LB_2([LB_1(WN([c_{i,j}, x_{i,j}]))], x_i) \\
 \hat{x}_i &= ReLU\left(\text{Maxpool}_{j=1,\dots,K}\{x_{i,j}^*\}\right); \forall x_i \in FPS(P); \quad (10) \\
 (c_{i,j}, x_{i,j}) &\in KNN(x_i), j = 1, \dots, K,
 \end{aligned}$$

where WN means windownorm. As seen from Table 4, the model performance has been further improved with the grouping strategy. The model accuracies are further improved from 70.9%/77.1%/91.1% (mIoU/mAcc/OA) to 71.4%/77.9%/91.1%.

In addition, we also conduct experiments and analyses on the number of grouping points m . m is used to determine the number of points involved in each type of information in the neighborhood. It should be noted that the number of grouping points m may have different effects for different models and tasks. We do experiments to find optimal m for our module. Table 5 shows the results of different m . It can be seen from Table 5 that the model accuracy is optimal with $m=3$.

Table 5. Ablation studies: The number of grouping points m .

m	mIoU (%)	mAcc (%)	OA (%)
1	70.1	76.5	90.8
2	69.8	76.3	90.8
3	71.4	77.9	91.1
4	70.2	77.0	90.7

4.4. Analysis of Sampling Complexity

The PAGWN module has one Linear layer to aggregate location information and neighbor features. Benefiting from parallel operations, the normalization process does not affect computational efficiency. In our proposed method, only a few learnable parameters are added. Compared to the original method, the PAGWN module has a small increment in memory usage and inference time. Table 6 compares complexity between different models on S3DIS.

Table 6. Complexity comparison on S3DIS (Area 5).

Method	Point Transformer	PAGWN
Parameters	7,767,729	8,031,729
Video memory	10175m	11897m
Training time	34min10s	35min47s
Inference time	30min24s	31min45s

5. Conclusions

In this work, we try to improve the information acquisition rate in the downsampling stage. In particular, we propose the group-wise window-normalization module to

unify inconsistent point densities and obtain multi-type information, including texture and spatial information. Furthermore, the pre-abstraction strategy is leveraged to balance local and global features. With the proposed PAGWN module, the recognition of the sofa and column is improved from 69.2% to 84.4% and from 42.7% to 48.7%, respectively. The benchmarks of segmentation on S3DIS (Area 5) are improved from 71.7%/77.6%/91.9% (mIoU/mAcc/OA) to 72.2%/78.2%/91.4%. The accuracies of 6-fold cross-validation on S3DIS are 77.6%/85.8%/91.7%. It outperforms the best model PointNeXt-XL [32] (74.9%/83.0%/90.3%) by 2.7% on mIoU and achieves state-of-the-art performance. Extensive experiments show that the PAGWN module has the following advantages: 1) It is computationally efficient with a small number of learnable parameters. 2) The normalization of inconsistent point densities is essential to point cloud understanding. 3) Texture and spatial information are extracted to enhance the network performance. 4) Local and global information can be well balanced. 5) The proposed module performs better on small object recognition, and the results have more precise boundaries than others.

References

- [1] Iro Armeni, Ozan Sener, Amir R Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1534–1543, 2016. 7
- [2] Song Bai, Philip Torr, et al. Visual parser: Representing part-whole hierarchies with transformers. *arXiv preprint arXiv:2107.05790*, 2021. 2
- [3] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European conference on computer vision*, pages 213–229. Springer, 2020. 2
- [4] Lin-Zhuo Chen, Xuan-Yi Li, Deng-Ping Fan, Kai Wang, Shao-Ping Lu, and Ming-Ming Cheng. Lsanet: Feature learning on point sets by local spatial aware layer. *arXiv preprint arXiv:1905.05442*, 2019. 2
- [5] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017. 1
- [6] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. 1
- [7] Ruihang Chu, Yukang Chen, Tao Kong, Lu Qi, and Lei Li. Icm-3d: Instantiated category modeling for 3d instance segmentation. *IEEE Robotics and Automation Letters*, 7(1):57–64, 2021. 1
- [8] Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. *Advances in Neural Information Processing Systems*, 34:9355–9366, 2021. 2
- [9] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12124–12134, 2022. 2
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- [11] Nico Engel, Vasileios Belagiannis, and Klaus Dietmayer. Point transformer. *IEEE Access*, 9:134826–134840, 2021. 3
- [12] Qian Gao and Xukun Shen. Thickseg: Efficient semantic segmentation of large-scale 3d point clouds using multi-layer projection. *Image and Vision Computing*, 108:104161, 2021. 1
- [13] Benjamin Graham, Martin Engelcke, and Laurens Van Der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 9224–9232, 2018. 1
- [14] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021. 3
- [15] Qingyong 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11108–11117, 2020. 1, 2
- [16] 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 *European Conference on Computer Vision*, pages 222–239. Springer, 2020. 1
- [17] Mingyang Jiang, Yiran Wu, Tianqi Zhao, Zelin Zhao, and Cewu Lu. Pointsift: A sift-like network module for 3d point cloud semantic segmentation. *arXiv preprint arXiv:1807.00652*, 2018. 2
- [18] Asako Kanezaki, Yasuyuki Matsushita, and Yoshifumi Nishida. Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5010–5019, 2018. 1
- [19] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. <https://github.com/dvlab-research/Stratified-Transformer>, 2022. Accessed 24 May 2022. 7

- [20] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8500–8509, 2022. [3](#), [7](#), [8](#), [9](#)
- [21] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12697–12705, 2019. [1](#)
- [22] Huan Lei, Naveed Akhtar, and Ajmal Mian. Seggen: Efficient 3d point cloud segmentation with fuzzy spherical kernel. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11611–11620, 2020. [1](#)
- [23] Bo Li, Tianlei Zhang, and Tian Xia. Vehicle detection from 3d lidar using fully convolutional network. *arXiv preprint arXiv:1608.07916*, 2016. [1](#)
- [24] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10012–10022, 2021. [2](#)
- [25] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2949–2958, 2021. [1](#)
- [26] Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local geometry in point cloud: A simple residual mlp framework. *arXiv preprint arXiv:2202.07123*, 2022. [2](#)
- [27] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3164–3173, 2021. [2](#)
- [28] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 922–928. IEEE, 2015. [1](#)
- [29] Chunghyun Park, Junha Lee, Yang Heemin, and Kwonyoung Ryu. <https://github.com/POSTECH-CVLab/point-transformer>, 2021. Accessed 6 April 2022. [7](#)
- [30] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas 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*, pages 652–660, 2017. [1](#), [2](#), [8](#), [10](#)
- [31] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. [1](#), [2](#), [3](#)
- [32] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Abed Al Kader Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies. *arXiv preprint arXiv:2206.04670*, 2022. [2](#), [7](#), [8](#), [9](#), [10](#), [12](#)
- [33] G Riegler, AO Ulusoy, and GA OctNet. Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, page 3577. [1](#)
- [34] Shuran Song, Fisher Yu, Andy Zeng, Angel X Chang, Manolis Savva, and Thomas Funkhouser. Semantic scene completion from a single depth image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1746–1754, 2017. [1](#)
- [35] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015. [1](#)
- [36] Jingang Tan, Kangru Wang, Lili Chen, Guanghui Zhang, Jia-mao Li, and Xiaolin Zhang. Hcfs3d: Hierarchical coupled feature selection network for 3d semantic and instance segmentation. *Image and Vision Computing*, 109:104129, 2021. [1](#)
- [37] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3887–3896, 2018. [1](#)
- [38] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6411–6420, 2019. [2](#), [8](#), [10](#)
- [39] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. [2](#)
- [40] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 32–42, 2021. [2](#)
- [41] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [2](#)
- [42] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pvt v2: Improved baselines with pyramid vision transformer. *Computational Visual Media*, 8(3):415–424, 2022. [2](#)
- [43] Mutian Xu, Runyu Ding, Hengshuang Zhao, and Xiaojuan Qi. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3173–3182, 2021. [1](#)
- [44] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5589–5598, 2020. [1](#)

- [45] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3323–3332, 2019. 1, 2, 8
- [46] Chenxi Zhao, Weihao Zhou, Li Lu, and Qijun Zhao. Pooling scores of neighboring points for improved 3d point cloud segmentation. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 1475–1479. IEEE, 2019. 2
- [47] Hengshuang Zhao, Li Jiang, Chi-Wing Fu, and Jiaya Jia. Pointweb: Enhancing local neighborhood features for point cloud processing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5565–5573, 2019. 1, 2, 8, 10
- [48] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021. 1, 3, 7, 8, 9, 10