# Interpretable Edge Enhancement and Suppression Learning for 3D Point Cloud Segmentation

Haoyi Xiu, Xin Liu, Weimin Wang, Kyoung-Sook Kim, Takayuki Shinohara, Qiong Chang, and Masashi Matsuoka

*Abstract*—**3D point clouds can flexibly represent continuous surfaces and can be used for various applications; however, the lack of structural information makes point cloud recognition challenging. Recent edge-aware methods mainly use edge information as an extra feature that describes local structures to facilitate learning. Although these methods show that incorporating edges into the network design is beneficial, they generally lack interpretability, making users wonder *how* exactly edges help. To shed light on this issue, in this study, we propose the Diffusion Unit (DU) that handles edges in an interpretable manner while providing decent improvement. Our method is interpretable in three ways. First, we theoretically show that DU learns to perform task-beneficial edge enhancement and suppression. Second, we experimentally observe and verify the edge enhancement and suppression behavior. Third, we empirically demonstrate that this behavior contributes to performance improvement. Extensive experiments performed on challenging benchmarks verify the superiority of DU in terms of both interpretability and performance gain. Specifically, our method achieves state-of-the-art performance in object part segmentation using ShapeNet part and scene segmentation using S3DIS. Our source code will be released at https://github.com/martianxiu/DiffusionUnit.**

*Index Terms*—**3D point clouds, Edge enhancement and suppression, Edge awareness, Diffusion.**

## I. INTRODUCTION

A 3D point cloud is a basic and flexible shape representation where the surface is represented as a set of discrete points in the 3D space. Powered by the recent advancement of cost-effective sensor technology, numerous large-scale datasets are publicly available to research communities, facilitating deep learning–based point cloud understanding. The application field of such research can be autonomous driving [1]–[3], remote sensing [4]–[7], and so on.

Among various types of deep neural networks (DNNs), Convolutional Neural Network (CNN) is the most effective algorithm for computer vision. However, the unstructured and unordered nature of point clouds prevents researchers from directly applying CNN to the data. Early works often project point clouds to regular grids before applying CNN [8]–[11] or utilizing special data structures [12]–[15]. However, most of these methods may suffer from information loss caused by projections. This issue is resolved by PointNet [16], which consists of permutation-invariant operations and can handle point clouds in a lossless manner. Further, PointNet++ extends PointNet to achieve hierarchical learning [17]. Based on the PointNet++ framework, varieties of convolution methods [18]–[23] and transformers [24]–[26] for point clouds are proposed.

On the other hand, point clouds naturally lack structures; hence, incorporating edge information into the model design can recover the structural information, leading to an improved performance [27]–[29]. Edges of a point are often constructed by explicitly connecting that point with its local neighbors. The resulting edges are used, for instance, as extra features describing the local structure [30], [31] or as similarity measures reflecting spatial consistency [32], [33]. Although these edge-aware methods show that the incorporation of edges improve performance, they are not interpretable, making users wonder *how* exactly edges contribute. This lack of interpretability can be problematic because it makes it difficult for users diagnose the problem of methods when they fail.

To resolve this issue, in this study, we propose the *Diffusion Unit (DU)* that handles edges in an interpretable manner while providing decent performance gain. Built on the diffusion theory [34], [35], DU is designed to learn how to enhance task-beneficial edges (e.g., part boundaries of objects) and suppress irrelevant discontinuities. The behavior of DU is interpretable in three ways: 1) we theoretically figure out which component of DU is responsible for edge enhancement and suppression and how it works; 2) we experimentally observe and verify the edge enhancement and suppression behavior; 3) we empirically demonstrate that this behavior contributes to performance improvement. The effect of DU is shown in Fig. 1.

Learning to enhance/suppress edges is useful in segmentation tasks, in which smooth intra-region predictions and accurate recognition of boundaries are crucial. To this end, we design *DU-Net*, a network tailored to point cloud segmentation. We perform extensive experiments to verify the effectiveness of DU-Net as well as the edge enhancement/suppression behavior of DU. In particular, we achieve state-of-the-art performance in object part segmentation using ShapeNet part [36] and scene segmentation using S3DIS [37].

Our main contributions are summarized as follows:

- We propose DU that performs automatic edge enhancement and suppression learning in an interpretable manner.
- We theoretically analyze and experimentally demonstrate the edge enhancement and suppression behavior of DU.
- We design DU-Net for point cloud segmentation and perform extensive experiments to verify its state-of-the-art performance and design choices.

The remainder of this paper is organized as follows. Section II revisits the related works concerning deep learning methods for 3D point clouds, especially, the edge-aware ones, and diffusion methods. Section III elaborates on the derivation and theoretical analysis of DU regarding edge enhancement and suppression learning, and the design of DU-Net. Section IV verifies the empirical interpretability of edge enhancement and suppression behavior of DU. Moreover,
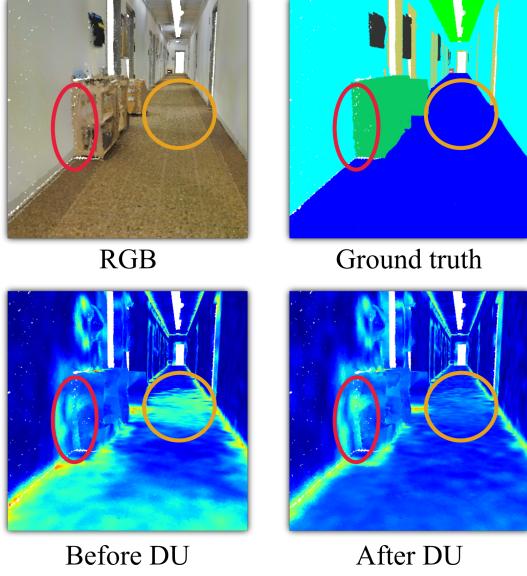
Fig. 1. The effect of DU. The heat map indicates the local feature smoothness. A low response indicates a smooth region while a high response indicates a non-smooth region. Red and orange circles indicate enhanced and suppressed features after applying DU. As can be seen, task-beneficial edges (in this case, object boundaries) are enhanced/preserved whereas intra-region discontinuities are suppressed.

experiment results on part and scene segmentation tasks are also reported. Then, extensive ablation studies are also made to validate the design of DU. We conclude the work in Section V.

## II. RELATED WORK

In this section, we first provide a brief review of deep learning methods for point cloud analysis. Then, we review approaches that facilitate point cloud understanding with the explicit awareness of edge information. Finally, we briefly revisit the diffusion methods on which we build our DU.

### A. Deep Learning for 3D Point Clouds

*a) Projection-Based Methods:* Owing to the unstructured and unordered nature of point clouds, early methods that deal with point clouds are projection-based methods. Those methods project point clouds to regular grids (e.g., 2D pixels [8], [9], [38] and 3D voxels [10]–[13], [39], [40]) to enable matured regular convolution to operate on point clouds. The problem of these methods is that, however, they lose fine-grained geometric details through projections.

*b) MLP-Based Methods:* MLP-based methods operate directly on the raw point clouds, thereby preventing information loss. The representative method is PointNet [16] whose main operations consist of shared MLPs and symmetric functions. Both operations are permutation-invariant, and thus the irregularity of point clouds is well-resolved. To obtain hierarchical feature representation, PointNet++ [17] applies PointNets to local subsets of points grouped according to different scales or resolutions. Afterwards, various advancement are made based on PointNet and PointNet++ [23], [41]–[43].

*c) Point Convolution–Based Methods:* The point convolution–based method defines a convolution operation on unstructured point clouds. For instance, some studies construct a pseudo grid on which input points are projected [20], [30], [44]. Since the constructed grid is structured, a standard convolution operation becomes applicable. On the other hand, the other methods predict convolution filters using sub-networks by dynamically transforming positional features [18], [19], [45]. Though effective, these methods do not explicitly consider the local structures characterized by edges which can describe the underlying surface more accurately.

*d) Transformer-Based Methods:* Inspired by the success of Transformers [46] in NLP field, many studies have been tried to construct transformer-based methods for point cloud data [25], [47]. A detailed review of transformers for 3D point clouds is presented in [48]. Here, we focus on transformers that are designed for point cloud segmentation. Since self-attention is a set operator, it can be straightforwardly applied to point clouds. PAT [24] applies the original self-attention directly to point clouds. PointASNL [49], on the other hand, combines a point convolution with self-attention so that the network considers local and global characteristics simultaneously. Point Transformer [33] applies vector attention to the local subset of points to better model the local structure. Point-BERT [50] enhances the performance by masked point modeling. Stratified Transformer [26] achieves efficient modeling of long-range dependencies using the stratified strategy.

### B. Edge-Aware Methods

To better model the local structure, edge-aware methods incorporate edge information into network design in various ways. EdgeConv [28] is the pioneering method to achieve edge awareness. It performs convolution on the edge embedding to better model the local geometric structure. The idea of EdgeConv is adopted in numerous works (e.g., [29]–[31], [51]). Although performance is improved by incorporating edge information, it is difficult to analyze rigorously how it is achieved. Furthermore, local edge features are simply mixed with global ones by MLPs, which may be suboptimal since they are significantly distinct features. On the other hand, some methods regard edge information as a measure of semantic distance. For instance, [32] uses edges as a descriptor indicating the relationships between a center point and its neighbors. This idea is sometimes combined with the attention mechanism [33], [52], [53], which converts edges into spatial weights. However, such a forced conversion may undermine the rich structural information contained in the edges. On the other hand, edges information may be used as an extra source of supervision. [54] uses edge to maintain spatial consistency. CBL [55] enhances the prediction at boundaries by designing a dedicated loss function. In [56], a joint learning framework in which edge detection and other tasks are combined is proposed. However, excessive reliance on high-quality labels may affect their effectiveness when labels are noisy.

In contrast, DU offers significantly better interpretability both theoretically and empirically compared with existing

edge-aware methods. Moreover, DU automatically learn to enhance task-beneficial edges and suppress irrelevant ones without additional supervision.

### C. Diffusion Methods

In essence, diffusion methods, which are motivated by the diffusion equation, model the smoothing of data (e.g., images) as diffusion processes. The core of the diffusion methods is the diffusivity [34], [57], which is often defined as a function of edge. As a result, diffusion methods remove small edges while preserving significant edges, achieving edge-preserving smoothing [58]. Therefore, such techniques are extensively studied in image processing [34], [35], [59], [60] and by other communities (e.g., computer graphics [61]–[64]). Recently, although several works model the global information propagation [65]–[68], the probabilistic point cloud generation [69], [70] or completion [71] as diffusion processes, extending the diffusion equation for adaptive edge enhancement and suppression for 3D point cloud segmentation, which we exclusively cope with in this study, remains unexplored.

### III. METHOD

In this section, we first introduce the origin of DU by providing a brief background about the diffusion equation. Then, we give the continuous definition of DU and provide a theoretical analysis to reveal the mechanism of edge enhancement and suppression learning. We subsequently introduce the discretized DU that eases the implementation on modern machines. Lastly, we design DU-Net for point cloud segmentation.

### A. Preliminary

The diffusion equation describes the movement of diffusive substances from regions of higher concentration to lower concentration without creating or destroying mass [35]. For instance, when hot water is poured into cold water, heat diffuses until the temperature of the water becomes the same everywhere. Let $u(p, t)$ denote the concentration at the position $p$ and time $t$. The amount of substances that flow through per unit area per unit time (the flux) is described by Fick's law:

$$s = -g \cdot \nabla u \, , \tag{1}$$

where $\nabla$ denotes the gradient operator and $g$ denotes the diffusivity. The fact that diffusion processes do not create or destroy mass is expressed by the continuity equation:

$$\partial_t u = -\text{div}(s) \, . \tag{2}$$

The continuity equation indicates that the change of concentration is caused only by the flux, which is measured by the divergence operator (div). Finally, the diffusion process is described by combining the above two equations:

$$\partial_t u = \text{div}(g \cdot \nabla u) \quad t \geq 0 \, , \tag{3}$$

with the initial condition $u(p, 0) = u_0(p)$ and the boundary condition as appropriate.

### B. Diffusion Unit (DU)

Inspired by the diffusion equation, we propose the diffusion unit (DU) that facilitates edge enhancement and suppression learning for 3D point clouds. Suppose a continuous spatial-temporal multi-channel point cloud $\mathbf{u} = \mathbf{u}(\mathbf{p}, t) = (u_1(\mathbf{p}, t), u_2(\mathbf{p}, t), ..., u_d(\mathbf{p}, t))$, where $d$ is the number of channels, $t$ is time, $\mathbf{p}$ denotes the position vector, and the initial condition is $\mathbf{u}(\mathbf{p}, 0) = \mathbf{h}$. DU is defined as:

$$\partial_t \mathbf{u} = \text{div}(\phi(\nabla \mathbf{u})), \quad t \geq 0 \, , \tag{4}$$

where $\partial_t \mathbf{u} \in \mathbb{R}^d$ is the output of DU at time $t$, $\nabla \mathbf{u}$ encodes channel-wise spatial gradient. Note that the choice of diffusivity $g$ in Eq. (3) has a significant impact on performance. Finding the appropriate $g$ often requires domain knowledge and involves numerous trials and errors [57]. In our definition, we replace the handcrafted $g$ with a *trainable filter* $\phi : \mathbb{R}^d \to \mathbb{R}^d$ so that the diffusivity function can be learned w.r.t data and task.

In practice, we use `1D-Conv` to implement $\phi$ so that the diffusivity function can consider the information in all feature channels jointly. This is important, for instance, when the information of interest is encoded only in the subspace of the entire feature space.

### C. Edge Enhancement and Suppression Learning

In this part, we explain how DU performs edge enhancement and suppression learning. For simplicity, we consider a step edge (i.e., an abrupt change in feature values) convolved by a Gaussian. Ideally, such a structure can be detected by the spatial gradient $\nabla \mathbf{u}$. Without loss of generality, we assume that the edge is aligned with $x$ axis ($\mathbf{u}_y = \mathbf{u}_z = \mathbf{0}$). The profile of the step edge and its derivatives are illustrated in Fig. 2. Now we focus on a single output channel $i$ for brevity.

In this context, Eq. (4) can be simplified as

$$(\partial_t \mathbf{u})_i = \frac{\partial}{\partial x} (\phi_i(\mathbf{u}_x)) = \nabla \phi_i \cdot \mathbf{u}_{xx} \tag{5}$$

$$= \sum_{j=1}^{d} (\phi_i')_j \cdot (\mathbf{u}_{xx})_j \, , \quad i = 1, 2, ..., d \, . \tag{6}$$

We are interested in the evolution of the edge over time, i.e., how $\mathbf{u}_x$ changes during applications of DU. We can derive that:

$$(\partial_t \mathbf{u}_x)_i = \left( \frac{\partial}{\partial t} \frac{\partial \mathbf{u}}{\partial x} \right)_i = \frac{\partial}{\partial x} \left( \frac{\partial \mathbf{u}}{\partial t} \right)_i \tag{7}$$

$$= \frac{\partial}{\partial x} \left( \sum_{j=1}^{d} (\phi_i')_j \cdot (\mathbf{u}_{xx})_j \right) \tag{8}$$

$$= \sum_{j=1}^{d} (\phi_i'')_j \cdot (\mathbf{u}_{xx})_j^2 + (\phi_i')_j \cdot (\mathbf{u}_{xxx})_j \, . \tag{9}$$

As shown in Fig. 2, at the inflection point $\mathbf{u}_{xx} = 0$ and $\mathbf{u}_{xxx} < 0$. Therefore, the contribution of a particular input channel $j$ to the output channel $i$ (the sign of Eq. (9)) is determined by the sign of $(\phi_i')_j$. Specifically, channel $j$ has a
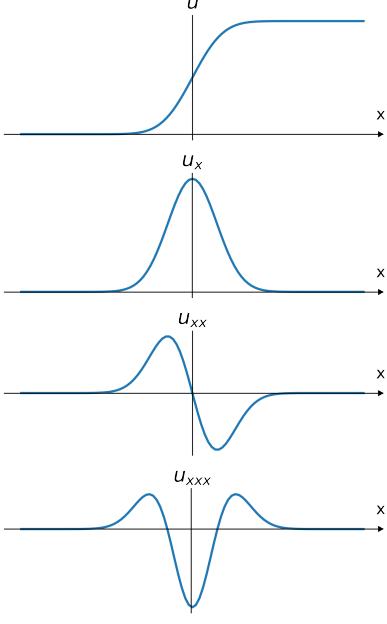
Fig. 2. Profiles of the smoothed step edge (top) and its derivatives. At the inflection point, $\mathbf{u}_{xx} = 0$ while $\mathbf{u}_{xxx} < 0$, which determines the behavior of DU (see Eq. (9)).

positive impact if $(\phi'_i)_j < 0$, whereas it has a negative impact if $(\phi'_i)_j > 0$.

As a result, enhancing $((\partial_t \mathbf{u}_x)_i > 0)$ or suppressing $((\partial_t \mathbf{u}_x)_i < 0)$ the edge can be adaptively learned by the filter $\phi$, by collectively using fragmentary information in all channels.

### D. Discretization

To deal with discrete point clouds, discretization of Eq. (4) is necessary. Let $\mathbf{u}_s, \mathbf{u}_n \in \mathbb{R}^d$ denote the feature of the center point and its spatial neighbors, respectively. Let $n \in \mathcal{N}_s$ represents $n$-th neighbor of the center point. First, we adopt the explicit scheme for time discretization and have:

$$\partial_t \mathbf{u}_s \approx \mathbf{u}_s^{t+1} - \mathbf{u}_s^t . \tag{10}$$

Second, we discretize the space using the finite difference and have:

$$\text{div}(\phi(\nabla \mathbf{u})) \approx \frac{1}{|\mathcal{N}_s|} \sum_{n \in \mathcal{N}_s} \phi\left(\mathbf{u}_n^t - \mathbf{u}_s^t\right) , \tag{11}$$

where $|\mathcal{N}_s|$ represents the number of neighbors. Combining Eqs. (4), (10), (11), we have

$$\mathbf{u}_s^{t+1} = \mathbf{u}_s^t + \frac{1}{|\mathcal{N}_s|} \sum_{n \in \mathcal{N}_s} \phi\left(\mathbf{u}_n^t - \mathbf{u}_s^t\right) . \tag{12}$$

Moreover, we incorporate Batch Normalization [72] and ReLU [73] activation function $\varphi = \texttt{ReLU} \circ \texttt{BatchNorm} : \mathbb{R}^d \to \mathbb{R}^d$ into the second term on the right-hand side of Eq. (12) to facilitate training and encourage sparsity. Finally, the discretized DU is defined as:

$$\mathbf{u}_s^{t+1} = \mathbf{u}_s^t + \varphi\left(\frac{1}{|\mathcal{N}_s|} \sum_{n \in \mathcal{N}_s} \phi\left(\mathbf{u}_n^t - \mathbf{u}_s^t\right)\right) . \tag{13}$$

Note that the neural network functions $\phi$ and $\varphi$ work together and are responsible for the learning of enhancing and suppressing edges (we provide qualitative verification in Section IV-A). The above definition can be efficiently computed and easily parallelized, fitting modern GPU-empowered deep learning frameworks. The detailed computation flow of DU is shown in Fig. 3 (middle).

### E. Network Architecture

Further, we construct DU-Net to tackle point cloud segmentation. An overview of the network architecture is shown in Fig. 3 (top). The network architecture adopts U-Net [74]–like encoder-decoder style.

The encoder is responsible for hierarchical feature abstraction. In each encoder stage, input point features are transformed by certain functions and subsequently downsampled. The downsampling ensures an efficient encoding of multi-resolution characteristics, leading to a more discriminative feature representation. After downsampling, point features are passed to the next stage. We construct a CNN-like encoder based on KPConv [20], which is a standard convolution operator adopted in many previous works for its excellent performance and ease of implementation [26], [56]. Specifically, we select the depthwise [75] version to reduce the overall complexity. Following prior works [17], [33], we adopt the farthest point sampling as the downsampling scheme.

The decoder reverses the process of the encoder by performing upsampling and DU stage by stage. An upsampling layer receives the output of the previous layer and recovers the resolution of points in the adjacent upper stage of the encoder. We adopt U-Net–style skip connection to assist feature reconstruction. The upsampled features are subsequently processed by DU so that task-beneficial edges of each resolution are kept sharpened while irrelevant discontinuities are smoothed out. We apply DUs in the decoder because features of each resolution in the encoder are passed to the decoder through skip connections. Therefore, backpropagated signals at a resolution level can simultaneously reach both layers in the encoder and decoder, facilitating optimization. We refer to our decoder as DU-decoder for clarity and brevity.

## IV. EXPERIMENTS

In this section, we conduct experiments to answer the following questions:

1) Does DU really perform edge enhancement/suppression as described by the theoretical analysis in Section III-C?
2) Can DU-Net compete with recent cutting-edge networks?
3) Is the design choices of DU and DU-Net reasonable?

To answer these questions, we use standard benchmarks of object part segmentation and scene segmentation tasks.

For object part segmentation, we use ShapeNet part dataset [17], [36], which contains 16,881 3D models that are classified to 16 object categories. Each model is annotated with several parts (less than 6) from 50 object part classes. For instance, an airplane category usually consists of parts including the airplane body, wings, tails, and engine classes.
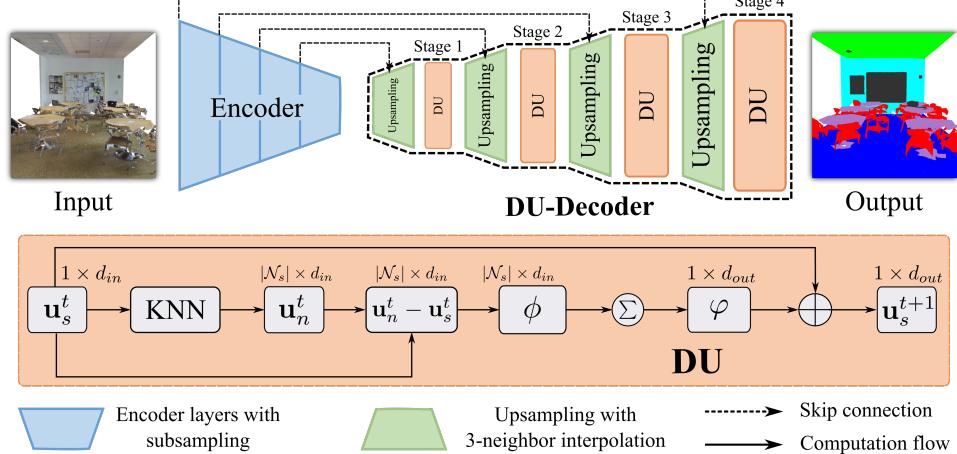
Fig. 3. The network architecture used in this study and the computation flow of the DU. Top: the network architecture. Middle: computation flow of the DU. $\sum$ and $\oplus$ denote neighborhood aggregation and element-wise sum, respectively.

For scene segmentation, we use Stanford large-scale 3D indoor spaces (S3DIS) [37]. The dataset contains 3D scans taken from six areas including 272 rooms. For instance, rooms such as lobby, office, and conference rooms are included. Each point is annotated with a class from 13 categories. A room, for example, in the dataset includes classes such as floor, wall, door, and table.

### A. Verifying the Behaviors of DU

We have theoretically analyzed that DU learns to enhance or suppress edges (see Section III-C). Here we experimentally verify the edge enhancement and suppression behaviors of DU. The features used for this experiment are taken from the last DU (stage 4 in Fig. 3) since it is the closest layer to the final classification layer. To verify the effect of DU, we examine *local feature smoothness* that summarizes how different the features are from their local neighbors. Formally, smoothness is defined as $||\sum_{n \in \mathcal{N}_s} \mathbf{f}_n - \mathbf{f}_s||$, where $\mathbf{f}_s$ and $\mathbf{f}_n$ denote the features of a center point and its neighbor from some layer. It can be observed that the smoothness value becomes closer to zero if the feature of the center point is similar to those of its neighbors while the value increases when they are dissimilar. Therefore, we can qualitatively grasp the edge enhancement/suppression behavior of DU by comparing smoothness values before and after applying DU.

Fig. 4 shows some examples of the change in local feature smoothness after applying DU from the ShapeNet part dataset. In this case, DU-Net is trained to assign part categories to each point. As can be seen, DU-Net is able to perform edge enhancement and suppression simultaneously. For instance, the third row shows a laptop whose boundary between its screen and keyboard is enhanced while the feature of its keyboards is smoothed. Therefore, DU successfully enhances the task-beneficial edges and suppressed the unhelpful discontinuities. Furthermore, we observe that DU is capable of differentiating task-beneficial edges from other geometric edges (the top and last row). In the last row of Fig. 4, for example, the sharp edges on the boundary (but not on the part boundary) of the table are aggressively suppressed after applying DU. On the



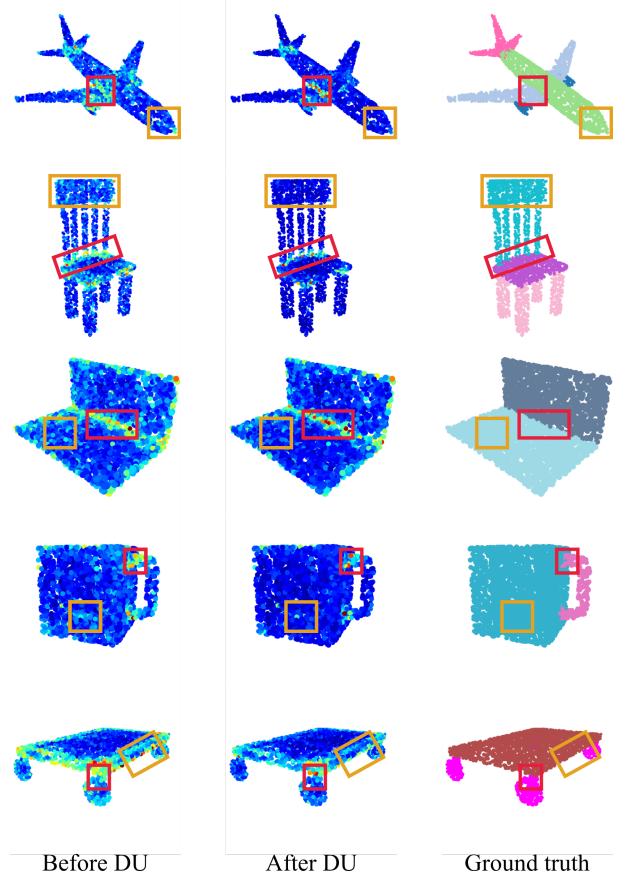Before DU     After DU     Ground truth

Fig. 4. Examples of local feature smoothness before and after applying DU using the ShapeNet part test set. A low response indicates a smooth region while a high response indicates a non-smooth region. Red and orange rectangles indicate enhanced and suppressed edges after DU. The task-beneficial edges (in this case, part boundaries) are enhanced (become brighter) while within-part features are suppressed (become darker).

other hand, the points near the part boundary are enhanced by DU. Thus, DU can understand the semantics of the objects and perform edge enhancement/suppression selectively.

Fig. 5 demonstrates the effect of DU in complex scenes. The

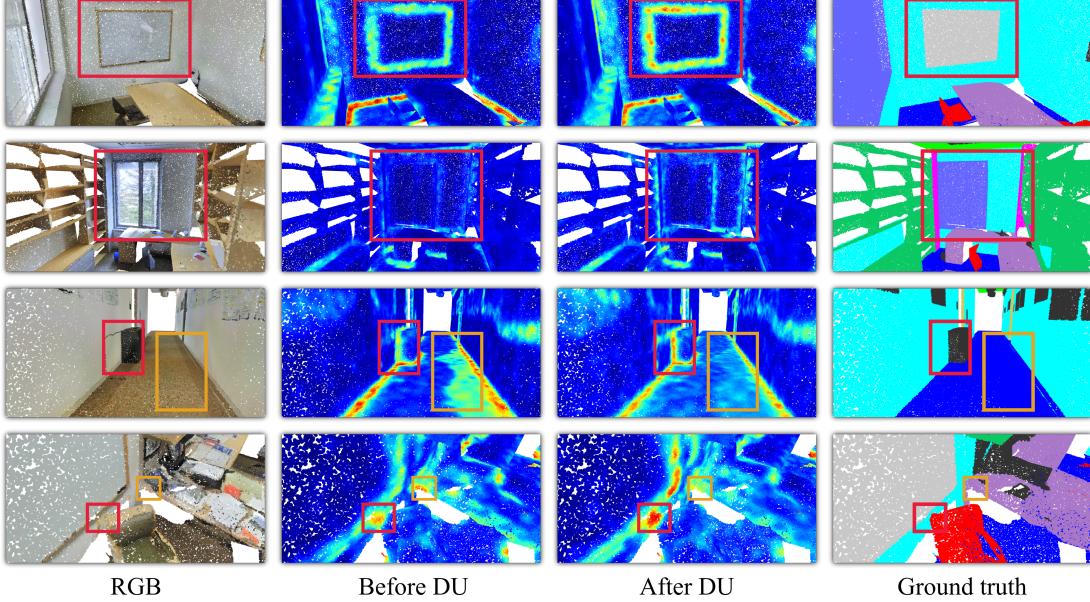|  |  |  |  |
|:---:|:---:|:---:|:---:|
| RGB | Before DU | After DU | Ground truth |

Fig. 5. Examples of local feature smoothness before and after applying DU using the S3DIS Area 5 test set. A low response indicates a smooth region while a high response indicates a non-smooth region. Red and orange rectangles indicate enhanced and suppressed features after DU. The task-beneficial edges (in this case, object boundaries) are enhanced (become brighter) while intra-region features are suppressed (become darker).

top row shows that DU successfully locates and enhances the boundary of a whiteboard. As a result, the boundary becomes more salient. In the second row, DU also manages to enhance the object boundary (e.g., the boundary between the window and wall) that is barely distinguishable before applying DU. Interestingly, only the boundaries between different categories are enhanced while other discontinuities such as the textures within the window and the window frame remain suppressed; therefore, the DU clearly understands and distinguishes between task-beneficial and other edges. The third row shows a case in which edge enhancement and suppression occur simultaneously. The discontinuous edges that appear on the floor are smoothed after the DU while the boundary between a box and wall is enhanced. In the last row, DU precisely enhances the point where the wall and the chair are attached, while nicely suppressing the edge of the table that is close to the chair but not attached to it.

Thus, the task-beneficial edge enhancement/suppression behavior of DU can be verified.

### B. Object Part Segmentation

We adopt the standard train-test split [17]. All available points are used for input. 3D coordinates along with surface normal information are used as the input features. We adopt the standard data augmentation strategy for point cloud analysis. Specifically, random anisotropic scaling and random translation are used for data augmentation. We train the model using one NVIDIA Tesla V100 GPU. The model is trained for 150 epochs. SGD is used for optimization with an initial learning rate of 0.1. The learning rate is decayed by 0.1 when the epoch reaches 90 and 120. Following the common practice [20], [29], [45], we use voting for post-processing. To assess the performance, we use widely adopted instance-wise mIoU (ImIoU) and category mIoU (CmIoU) defined in [16].

TABLE I
RESULT OF OBJECT PART SEGMENTATION (%). THE BOLD TEXT INDICATES THE BEST PERFORMANCE. * DENOTES METHODS THAT ADOPT VOTING FOR POST-PROCESSING.

| Method | ImIoU | CmIoU |
|:---|:---:|:---:|
| PointNet [16] | 83.7 | 80.4 |
| PointNet++ [17] | 85.1 | 81.9 |
| PointConv [19] | 85.7 | 82.8 |
| PointCNN [18] | 86.1 | 84.6 |
| RSCNN* [45] | 86.2 | 84.0 |
| KPConv* [20] | 86.4 | 85.1 |
| PAConv* [29] | 86.1 | 84.6 |
| Point Transformer [33] | 86.6 | 83.7 |
| Stratified Transformer* [26] | 86.6 | 85.1 |
| CurveNet* [31] | 86.8 | - |
| Ours (w/o DU) | 86.2 | 83.6 |
| Ours (w/ DU) | 86.7 | 84.5 |
| Ours (w/ DU)* | **87.1** | **85.2** |

The quantitative result is listed in Table I. DU-Net achieves state-of-the-art performance both in terms of ImIoU and CmIoU. This reveals that DU-Net not only performs well in the major classes but also achieves satisfactory accuracy for each class. Moreover, although the plain network (w/o DU in Table I) fails to compete with the cutting-edge networks, DU succeeds in improving its performance, achieving the best performance.

Moreover, we compare the predictions of DU-Net and the plain network to provide some insights into how DU improves performance. The results are shown in Fig. 6. First, we find that DU-Net produces smoother predictions compared with the plain counterpart. For instance, the example of an airplane in the top row shows that the plain network has difficulty in distinguishing between the tail and body, resulting in ragged predictions. In contrast, DU-Net is able to produce smooth
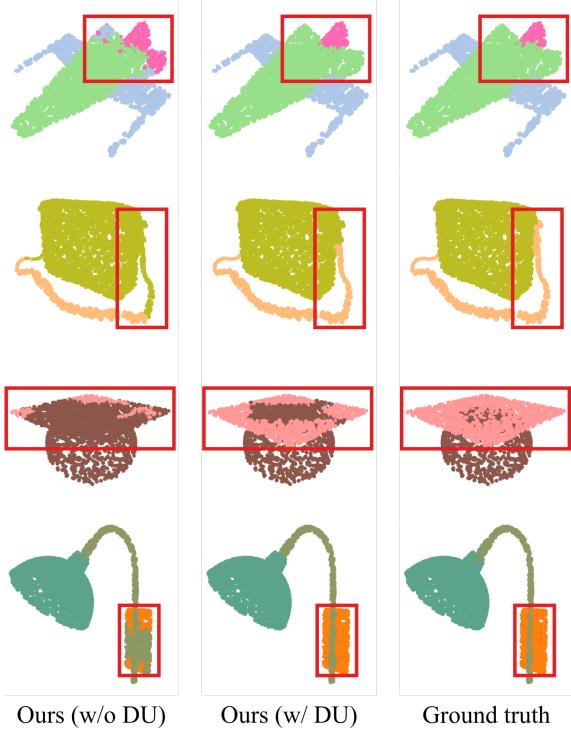
Fig. 6. Qualitative results of object part segmentation. The red rectangles show the improvements brought by DU. As can be seen, DU encourages the model to produce smoother and more boundary-aware predictions.

predictions even though the shape of the tail is similar to that of the body. We believe that the adaptability of DU in handling edges successfully suppresses the unhelpful edges, which in turn produces smoother predictions. Second, it is found that DU-Net tends to be more boundary-aware. In the second row, although both networks misclassify some of the points on the bag handle as bag body points, the result of DU-Net is much more accurate than the one of the plain counterpart, revealing that DU-Net has a better localization of the boundary between the bag body and handle. We conjecture that the ability to enhance task-beneficial edges makes features near boundaries more discriminative. As a result, DU-Net can detect part boundaries more precisely, thereby facilitating smooth predictions within boundaries and mitigating cross-boundary misclassifications (the third and last row).

### C. Scene Segmentation

Similar to [76], we advocate using Area five for testing and others for training. Following the common practice [20], [33], [42], we first grid-subsample each room with a grid size of 4cm for training. For testing, all points are evaluated. We form an input point cloud by taking at most 24,000 points from a room. We use 3D coordinates and color information as input features. Random vertical rotation, random anisotropic scaling, random jitter, random color drop, and random color auto-contrasting are used for data augmentation, following [26], [33]. Following the previous work [33], we train the model for 76,500 iterations. For this dataset, four NVIDIA Tesla V100 GPUs are used for training. For optimization, AdamW [77]

algorithm is used with an initial learning rate of 0.01. The learning rate is decayed based on the cosine annealing schedule. We use mIoU to quantitatively assess the performance.

The results are reported in Table II. DU-Net achieves state-of-the-art performance in terms of mIoU. We find that DU-Net performs especially well in classifying classes such as door and bookcase (book. in Table II). Since these two classes on often appear on the wall, one possible reason for the high performance is that the edge enhancement capability of DU makes the boundary feature more discriminative, easing the detection of subtle but crucial boundaries. As can be seen, the plain network (w/o DU) fails to compete with the recent strong transformer-based methods, while incorporating DUs significantly improves the performance, reaching the state of the art in terms of mIoU. As for the class-wise performance, we can verify that DUs successfully improve the performance by around 0.1–8.3, demonstrating the usefulness of edge enhancement and suppression learning in scene segmentation.

To intuitively understand how DU improves performance, we qualitatively analyze the effect of DU. The results are shown in Fig. 7. We observe that DU-Net (w/ DU in Fig. 7) generally better discriminates between similar classes compared with the plain counterpart (w/o DU in Fig. 7). For instance, the example in the top row shows that the plain network is greatly confused with the window, wall, and column classes. DU-Net, on the other hand, correctly identifies those classes and produces much more accurate and smoother predictions. We believe that the adaptability of DU leads to an enhanced sensitivity to object boundaries. Consequently, DU-Net can detect boundaries more accurately, which in turn facilitates spatial consistency within the detected boundaries. Further, the example in the second row shows that DU makes the model more shape-aware. The plain network shows significant confusion in differentiating between pipes and the ceiling as the pipes are attached to the ceiling. In contrast, DU-Net shows much better predictions. We conjecture that explicit consideration of edges enables DU-Net to obtain more shape-aware representations, resulting in better discrimination. Moreover, we observe that DU-Net can make effective use of semantic information when geometric information is inadequate. In the third row, the plain net completely fails to identify the door likely because the door is embedded in the wall, making it geometrically difficult to differentiate the door from the wall. On the contrary, DU-Net succeeds in producing far better predictions. Similarly, in the last row, although the plain net manages to detect the centers of drawings, it fails to identify precise boundaries. In contrast, DU-Net produces fine predictions with accurate boundaries, demonstrating its advanced capability to utilize semantic information.

### D. Design Analysis

We further study the design of DU and DU-Net. All experiments are conducted on the S3DIS dataset and the experiment settings are the same as described in Section IV-C.

*a) Ablation Study on DU:* The results are shown in Table III. First, the effect of functions $\phi$ and $\varphi$ are investigated. The performance drops significantly by 2.8 when $\phi$ is

TABLE II
RESULT OF SCENE SEGMENTATION (%). THE BOLD TEXT INDICATES THE BEST PERFORMANCE.

| Method | mIoU | ceil. | floor | wall | beam | column | window | door | chair | table | book. | sofa | board | clutter |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PointNet [16] | 41.1 | 88.8 | 97.3 | 69.8 | 0.5 | 3.92 | 46.3 | 10.8 | 52.6 | 58.9 | 40.3 | 5.9 | 26.4 | 33.2 |
| SegCloud [76] | 48.9 | 90.1 | 96.1 | 69.9 | 0.0 | 18.4 | 38.4 | 23.1 | 75.9 | 70.4 | 58.4 | 40.9 | 13.0 | 41.6 |
| PointCNN [18] | 57.3 | 92.3 | 98.2 | 79.4 | 0.0 | 17.6 | 22.8 | 62.1 | 80.6 | 74.4 | 66.7 | 31.7 | 62.1 | 56.7 |
| PointWeb [52] | 60.3 | 92.0 | 98.5 | 79.4 | 0.0 | 21.1 | 59.7 | 34.8 | 88.3 | 76.3 | 69.3 | 46.9 | 64.9 | 52.5 |
| KPConv [20] | 67.1 | 92.8 | 97.3 | 82.4 | 0.0 | 23.9 | 58.0 | 69.0 | 91.0 | 81.5 | 75.3 | 75.4 | 66.7 | 58.9 |
| JSENet [56] | 67.7 | 93.8 | 97.0 | 83.0 | 0.0 | 23.2 | 61.3 | 71.6 | 89.9 | 79.8 | 75.6 | 72.3 | 72.7 | 60.4 |
| CBL [55] | 69.4 | 93.9 | 98.4 | 84.2 | 0.0 | 37.0 | 57.7 | 71.9 | 91.7 | 81.8 | 77.8 | 75.6 | 69.1 | 62.9 |
| Point Transformer [33] | 70.4 | 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 |
| Stratified Transformer [26] | 72.0 | 96.2 | 98.7 | 85.6 | 0.0 | 46.1 | 60.0 | 76.8 | 84.5 | 92.6 | 75.2 | 77.8 | 78.1 | 64.0 |
| Ours (w/o DU) | 69.7 | 94.8 | 98.5 | 85.0 | 0.0 | 31.7 | 58.6 | 75.7 | 80.8 | 91.8 | 76.1 | 75.3 | 79.5 | 58.8 |
| Ours (w/ DU) | **72.2** | 95.6 | 98.6 | 85.2 | 0.0 | 40.0 | 60.7 | 82.7 | 83.1 | 90.8 | 83.5 | 78.5 | 75.9 | 64.1 |



RGB      Ours (w/o DU)      Ours (w/ DU)      Ground truth

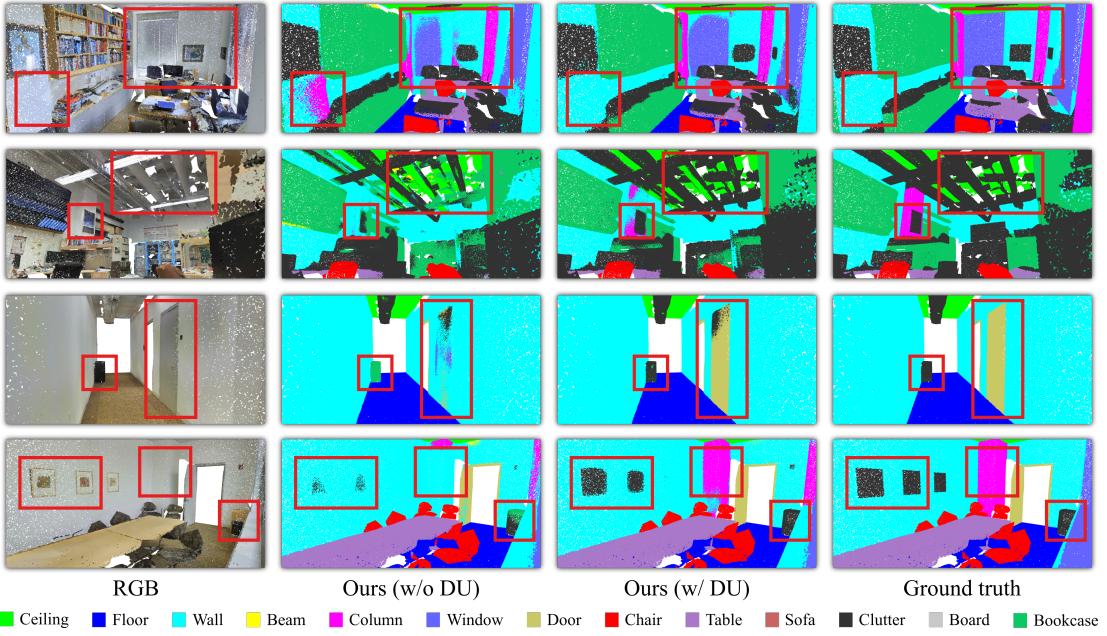■ Ceiling ■ Floor ■ Wall ■ Beam ■ Column ■ Window ■ Door ■ Chair ■ Table ■ Sofa ■ Clutter ■ Board ■ Bookcase

Fig. 7. Qualitative results of scene segmentation using S3DIS. The red rectangles show the improvements brought by DU. In general, DU-Net (w/ DU) tends to better discriminate between similar classes (the first row) and produce smoother and boundary-aware predictions (the second row). In some cases, DU-Net correctly identifies objects that the plain counterpart (w/o DU) almost completely failed to recognize (the third and fourth rows).

TABLE III
ABLATION STUDY ON DU. MODEL 1 INDICATES THE DEFAULT DESIGN OF DU. NUMBERS IN PARENTHESES DENOTE THE PERFORMANCE REDUCTION COMPARED WITH THE DEFAULT DU. ALL EXPERIMENTS ARE PERFORMED ON THE S3DIS DATASET.

| Model | #neighbors ($|\mathcal{N}_s|$) | $\phi$ | $\varphi$ | mIoU |
|---|---|---|---|---|
| 1 (default) | 16 | ✓ | ✓ | 72.2 |
| 2 | 16 | - | ✓ | 69.4 $(-2.8)$ |
| 3 | 16 | ✓ | - | 70.8 $(-1.4)$ |
| 4 | 16 | - | - | 69.8 $(-2.4)$ |
| 5 | 8 | ✓ | ✓ | 70.9 $(-1.3)$ |
| 6 | 24 | ✓ | ✓ | 71.2 $(-1.0)$ |
| 7 | 28 | ✓ | ✓ | 70.6 $(-1.6)$ |

removed. This is reasonable since the adaptability of DU is mostly provided by $\phi$. Therefore, the importance of $\phi$ in the design of DU is verified. Next, we observe that removing $\varphi$ from DU results in a reduction of 1.4. Recall that $\varphi = \texttt{ReLU} \circ \texttt{BatchNorm}$ is mainly applied to ease optimization. Thus, we

believe that removing $\varphi$ increases the optimization difficulty, which adversely affects the performance. On the other hand, the reduction of the performance is less severe compared with removing $\phi$; hence, expressiveness is relatively more important than optimization in the design of DU. As expected, removing both functions results in a significant reduction in performance. Interestingly, the reduction of performance caused by removing both functions is less severe than the one that removes $\phi$. We conjecture that the direct application of $\varphi$ without any transformation may lead to information loss since $\varphi$ contains $\texttt{ReLU}$. Then, we vary $|\mathcal{N}_s|$ from 8 to 28 to analyze the effect of neighborhood size in the design of DU. As can be observed, the lower scores are obtained when $|\mathcal{N}_s|$ takes too small (8) or too large values (28). On the other hand, better performance is achieved when $|\mathcal{N}_s|$ takes intermediate values. Therefore, we set $|\mathcal{N}_s|$ to 16 throughout this study.

*b) Influence of DUs at Different Stages:* Since DUs are applied to various stages of the decoder, we investigate the individual influence of DU at each stage in terms of

TABLE IV
INFLUENCE OF DUS AT DIFFERENT STAGES. STAGES INDICATE DIFFERENT
RESOLUTIONS (SEE FIG. 3). THE NUMBERS OF PARAMETERS (M) AND
FLOPS (G) ARE ALSO REPORTED. THE NUMBERS IN PARENTHESES
DENOTE THE PERFORMANCE REDUCTION COMPARED WITH THE DEFAULT
DU-NET (DU AT ALL STAGES).

| Model | #param. (M) | FLOPs (G) | mIoU |
|---|---|---|---|
| DU at all stages (Ours) | 8.07 | 11.94 | 72.2 |
| DU at stage 1 | 7.80 | 10.05 | 70.3 (−1.9) |
| DU at stage 2 | 7.18 | 10.06 | 70.6 (−1.6) |
| DU at stage 3 | 7.03 | 10.07 | 71.0 (−1.2) |
| DU at stage 4 | 6.99 | 10.07 | 70.3 (−1.9) |
| w/o DU | 6.98 | 9.43 | 69.7 (−2.5) |

TABLE V
COMPARISON OF DIFFERENT DECODERS. A DECODER IS CONSTRUCTED
BY STACKING THE DECODER LAYERS IN A WAY SIMILAR TO THE ONE
DESCRIBED IN SECTION III-E. ALL DECODERS ARE PAIRED WITH THE
ENCODER OF DU-NET. TYPE DENOTES THE MAIN OPERATION IN EACH
DECODER LAYER. THE VALUES IN THE PARENTHESES INDICATE THE
PERFORMANCE REDUCTION COMPARED WITH DU-DECODER.

| Decoder layer | Type | mIoU |
|---|---|---|
| DU (Ours) | Edge-aware | 72.2 |
| Feature propagation [17] | MLP | 69.7 (−2.5) |
| PointDeconv [19] | Deconvolution | 70.0 (−2.2) |
| Vector attention [33] | Attention/Edge-aware | 70.3 (−1.9) |

performance and computational complexity. The results are listed in Table IV. From stages 1 to 4, the resolution increases because of upsamplings while the feature dimension decreases. First, it is found that a single DU can immediately improve performance. As can be seen, DU at each stage improves upon the plain network (w/o DU) by 0.6–1.3. Therefore, although the amount of improvement varies, each DU can make effective use of information at the corresponding stage. Second, we observe that the benefits brought by DU at each stage can be accumulated. Specifically, applying DUs at all stages achieves the best performance while the effect of a single DU is limited. We speculate that the adaptability of the DUs allows them to be tailored to each resolution, resulting in a cumulative improvement. Third, we find that because of the lightweight nature of the DUs, each DU produces improvement with only a slight increase in spatial and time complexity. For example, applying DU at stage 3 improves mIoU by 1.3, but only increases #param. and FLOPs by 0.7% and 6.7%, respectively. Thus, an acceptable cost-performance tradeoff can be achieved when computational resources are limited.

*c) Effectiveness of DU-decoder:* We compare DU with different decoder layers to validate its effectiveness. Several decoder layers of different types (second column in Table V) are selected for comparison. A decoder is constructed by stacking a decoder layer in a way similar to the one described in Section III-E and Fig. 3. The encoder of DU-Net is used as the default encoder for all decoders. The results are listed in Table V. The first baseline is the feature propagation layer, which is the core layer of PointNet++ decoder [17]. This decoder is a standard decoder that has been adopted by numerous previous works (e.g., [20], [26], [29], [45]). As

can be seen, DU-decoder outperforms this baseline by 2.5, demonstrating its effectiveness. The second baseline is the decoder based on PointDeconv, a deconvolution layer based on PointConv [19]. We can see that DU-decoder is superior to this one with a performance gain of 2.2. The third baseline is vector attention [33], which is the core component of the decoder of Point Transformer [33]. In addition, this baseline also uses edges as a similarity measure for generating attention weights. We can see that our decoder again outperforms this one by 1.9. Therefore, it is verified that DU-decoder is more effective than various types of decoders in segmentation.

## V. CONCLUSION

Edge information has been mainly used as an extra feature that enhances performance. While edges have proven useful, it is still unclear how exactly it improves performance. In this study, we propose DU that handles edges in an interpretable manner. The interpretability of DU is shown in three ways. First, we theoretically figure out that DU performs edge enhancement and suppression. Second, we verify the result of theoretical analysis by intuitive visualizations. Third, we validate that the incorporation of DUs into the network design improves performance. DU-Net is constructed to tackle point cloud segmentation. Specifically, DU-Net achieves the state-of-the-art performance in object part segmentation and scene segmentation. We believe that the interpretability of DU can provide users with a clear idea of how to use this method, what results to expect, and how to diagnose problems when it fails, which is valuable since most of the methods based on DNNs operate in a black-box manner.

In an application point of view, it will be interesting to apply DU to other domains where the modeling of edges is crucial. An example is damage detection using point clouds [78]. Damage is often characterized by non-smooth distributions; thus, it would be interesting to see how enhanced adaptability in handling edges influences the damage detection result.

## REFERENCES

[1] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, "Frustum pointnets for 3d object detection from rgb-d data," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 918–927.

[2] Y. Cui, R. Chen, W. Chu, L. Chen, D. Tian, Y. Li, and D. Cao, "Deep learning for image and point cloud fusion in autonomous driving: A review," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[3] Y.-L. Jin, Z.-Y. Ji, D. Zeng, and X.-P. Zhang, "Vwp: An efficient drl-based autonomous driving model," *IEEE Transactions on Multimedia*, 2022.

[4] Y. Chen, Y. Guo, Y. Wang, D. Wang, C. Peng, and G. He, "Denoising of hyperspectral images using nonconvex low rank matrix approximation," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 55, no. 9, pp. 5366–5380, 2017.

[5] X. Jiang, G. Li, X.-P. Zhang, and Y. He, "A semisupervised siamese network for efficient change detection in heterogeneous remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–18, 2021.

[6] N. Wambugu, Y. Chen, Z. Xiao, K. Tan, M. Wei, X. Liu, and J. Li, "Hyperspectral image classification on insufficient-sample and feature learning using deep neural networks: A review," *International Journal of Applied Earth Observation and Geoinformation*, vol. 105, p. 102603, 2021.

[7] X. Wang, J. Ma, J. Jiang, and X.-P. Zhang, "Dilated projection correction network based on autoencoder for hyperspectral image super-resolution," *Neural Networks*, vol. 146, pp. 107–119, 2022.

[8] A. Kanezaki, Y. Matsushita, and Y. 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*, 2018, pp. 5010–5019.

[9] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 945–953.

[10] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2015, pp. 922–928.

[11] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4490–4499.

[12] P.-S. Wang, Y. Liu, Y.-X. Guo, C.-Y. Sun, and X. Tong, "O-cnn: Octree-based convolutional neural networks for 3d shape analysis," *ACM Transactions On Graphics (TOG)*, vol. 36, no. 4, pp. 1–11, 2017.

[13] M. Tatarchenko, A. Dosovitskiy, and T. Brox, "Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2088–2096.

[14] R. Klokov and V. Lempitsky, "Escape from cells: Deep kd-networks for the recognition of 3d point cloud models," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 863–872.

[15] P.-S. Wang, C.-Y. Sun, Y. Liu, and X. Tong, "Adaptive o-cnn: A patch-based deep representation of 3d shapes," *ACM Transactions on Graphics (TOG)*, vol. 37, no. 6, pp. 1–11, 2018.

[16] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.

[17] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[18] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: Convolution on $\chi$-transformed points," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 828–838.

[19] W. Wu, Z. Qi, and L. Fuxin, "Pointconv: Deep convolutional networks on 3d point clouds," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9621–9630.

[20] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. J. Guibas, "Kpconv: Flexible and deformable convolution for point clouds," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 6411–6420.

[21] J. Zhang, C. Zhu, L. Zheng, and K. Xu, "Fusion-aware point convolution for online semantic 3d scene segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 4534–4543.

[22] H. Zhou, Y. Feng, M. Fang, M. Wei, J. Qin, and T. Lu, "Adaptive graph convolution for point cloud analysis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 4965–4974.

[23] G. Qian, Y. Li, H. Peng, J. Mai, H. A. A. K. Hammoud, M. Elhoseiny, and B. Ghanem, "Pointnext: Revisiting pointnet++ with improved training and scaling strategies," *arXiv preprint arXiv:2206.04670*, 2022.

[24] J. Yang, Q. Zhang, B. Ni, L. Li, J. Liu, M. Zhou, and Q. Tian, "Modeling point clouds with self-attention and gumbel subset sampling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3323–3332.

[25] Z. Qin, H. Yu, C. Wang, Y. Guo, Y. Peng, and K. Xu, "Geometric transformer for fast and robust point cloud registration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 143–11 152.

[26] X. Lai, J. Liu, L. Jiang, L. Wang, H. Zhao, S. Liu, X. Qi, and J. Jia, "Stratified transformer for 3d point cloud segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8500–8509.

[27] M. Simonovsky and N. Komodakis, "Dynamic edge-conditioned filters in convolutional neural networks on graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3693–3702.

[28] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *Acm Transactions On Graphics (tog)*, vol. 38, no. 5, pp. 1–12, 2019.

[29] M. Xu, R. Ding, H. Zhao, and X. Qi, "Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds," *arXiv preprint arXiv:2103.14635*, 2021.

[30] Z. Liu, H. Hu, Y. Cao, Z. Zhang, and X. Tong, "A closer look at local aggregation operators in point cloud analysis," in *European Conference on Computer Vision*. Springer, 2020, pp. 326–342.

[31] T. Xiang, C. Zhang, Y. Song, J. Yu, and W. Cai, "Walk in the cloud: Learning curves for point clouds shape analysis," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021.

[32] L. Wang, Y. Huang, Y. Hou, S. Zhang, and J. Shan, "Graph attention convolution for point cloud semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 10 296–10 305.

[33] H. Zhao, L. Jiang, J. Jia, P. H. Torr, and V. Koltun, "Point transformer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 16 259–16 268.

[34] P. Perona and J. Malik, "Scale-space and edge detection using anisotropic diffusion," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 12, no. 7, pp. 629–639, 1990.

[35] J. Weickert, *Anisotropic diffusion in image processing*. Teubner Stuttgart, 1998, vol. 1.

[36] L. Yi, V. G. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas, "A scalable active framework for region annotation in 3d shape collections," *ACM Transactions on Graphics (ToG)*, vol. 35, no. 6, pp. 1–12, 2016.

[37] I. Armeni, O. Sener, A. R. Zamir, H. Jiang, I. Brilakis, M. Fischer, and S. Savarese, "3d semantic parsing of large-scale indoor spaces," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1534–1543.

[38] Y. Feng, Z. Zhang, X. Zhao, R. Ji, and Y. Gao, "Gvcnn: Group-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 264–272.

[39] B. Graham, M. Engelcke, and L. Van Der Maaten, "3d semantic segmentation with submanifold sparse convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 9224–9232.

[40] C. Choy, J. Gwak, and S. Savarese, "4d spatio-temporal convnets: Minkowski convolutional neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3075–3084.

[41] Z. Zhang, B.-S. Hua, and S.-K. Yeung, "Shellnet: Efficient point cloud convolutional neural networks using concentric shells statistics," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1607–1616.

[42] S. Lan, R. Yu, G. Yu, and L. S. Davis, "Modeling local geometric structure of 3d point clouds using geo-cnn," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 998–1008.

[43] S. Qiu, S. Anwar, and N. Barnes, "Geometric back-projection network for point cloud classification," *IEEE Transactions on Multimedia*, 2021.

[44] J. Mao, X. Wang, and H. Li, "Interpolated convolutional networks for 3d point cloud understanding," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1578–1587.

[45] Y. Liu, B. Fan, S. Xiang, and C. Pan, "Relation-shape convolutional neural network for point cloud analysis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8895–8904.

[46] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *arXiv preprint arXiv:1706.03762*, 2017.

[47] Z. Liu, Z. Zhang, Y. Cao, H. Hu, and X. Tong, "Group-free 3d object detection via transformers," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 2949–2958.

[48] D. Lu, Q. Xie, M. Wei, L. Xu, and J. Li, "Transformers in 3d point clouds: A survey," *arXiv preprint arXiv:2205.07417*, 2022.

[49] X. Yan, C. Zheng, Z. Li, S. Wang, and S. 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*, 2020, pp. 5589–5598.

[50] X. Yu, L. Tang, Y. Rao, T. Huang, J. Zhou, and J. Lu, "Point-bert: Pre-training 3d point cloud transformers with masked point modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 19 313–19 322.

[51] G. Li, M. Muller, A. Thabet, and B. Ghanem, "Deepgcns: Can gcns go as deep as cnns?" in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9267–9276.

[52] H. Zhao, L. Jiang, C.-W. Fu, and J. Jia, "Pointweb: Enhancing local neighborhood features for point cloud processing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5565–5573.

[53] H. Xiu, X. Liu, W. Wang, K.-S. Kim, T. Shinohara, Q. Chang, and M. Matsuoka, "Enhancing local feature learning for 3d point cloud processing using unary-pairwise attention," *arXiv preprint arXiv:2203.00172*, 2022.

[54] L. Jiang, H. Zhao, S. Liu, X. Shen, C.-W. Fu, and J. Jia, "Hierarchical point-edge interaction network for point cloud semantic segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 10 433–10 441.

[55] L. Tang, Y. Zhan, Z. Chen, B. Yu, and D. Tao, "Contrastive boundary learning for point cloud segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8489–8499.

[56] Z. Hu, M. Zhen, X. Bai, H. Fu, and C.-l. Tai, "Jsenet: Joint semantic segmentation and edge detection network for 3d point clouds," in *European Conference on Computer Vision*. Springer, 2020, pp. 222–239.

[57] M. J. Black, G. Sapiro, D. H. Marimont, and D. Heeger, "Robust anisotropic diffusion," *IEEE Transactions on image processing*, vol. 7, no. 3, pp. 421–432, 1998.

[58] H. Chen, J. Huang, O. Remil, H. Xie, J. Qin, Y. Guo, M. Wei, and J. Wang, "Structure-guided shape-preserving mesh texture smoothing via joint low-rank matrix recovery," *Computer-Aided Design*, vol. 115, pp. 122–134, 2019.

[59] J. Weickert, "Coherence-enhancing diffusion filtering," *International journal of computer vision*, vol. 31, no. 2, pp. 111–127, 1999.

[60] T. Brox, J. Weickert, B. Burgeth, and P. Mrázek, "Nonlinear structure tensors," *Image and Vision Computing*, vol. 24, no. 1, pp. 41–55, 2006.

[61] M. Desbrun, M. Meyer, P. Schröder, and A. H. Barr, "Implicit fairing of irregular meshes using diffusion and curvature flow," in *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*, 1999, pp. 317–324.

[62] U. Clarenz, U. Diewald, and M. Rumpf, *Anisotropic geometric diffusion in surface processing*. IEEE, 2000.

[63] C. L. Bajaj and G. Xu, "Anisotropic diffusion of surfaces and functions on surfaces," *ACM Transactions on Graphics (TOG)*, vol. 22, no. 1, pp. 4–32, 2003.

[64] Y. Li, C. Wang, J. Hong, J. Zhu, J. Guo, J. Wang, Y. Guo, and W. Wang, "Video vectorization via bipartite diffusion curves propagation and optimization," *IEEE Transactions on Visualization and Computer Graphics*, 2021.

[65] J. Atwood and D. Towsley, "Diffusion-convolutional neural networks," *Advances in neural information processing systems*, vol. 29, 2016.

[66] J. Klicpera, S. Weißenberger, and S. Günnemann, "Diffusion improves graph learning," *arXiv preprint arXiv:1911.05485*, 2019.

[67] J. Zhao, Y. Dong, M. Ding, E. Kharlamov, and J. Tang, "Adaptive diffusion in graph neural networks," *Advances in Neural Information Processing Systems*, vol. 34, 2021.

[68] B. Chamberlain, J. Rowbottom, M. I. Gorinova, M. Bronstein, S. Webb, and E. Rossi, "Grand: Graph neural diffusion," in *International Conference on Machine Learning*. PMLR, 2021, pp. 1407–1418.

[69] S. Luo and W. Hu, "Diffusion probabilistic models for 3d point cloud generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2837–2845.

[70] L. Zhou, Y. Du, and J. Wu, "3d shape generation and completion through point-voxel diffusion," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5826–5835.

[71] Z. Lyu, Z. Kong, X. XU, L. Pan, and D. Lin, "A conditional point diffusion-refinement paradigm for 3d point cloud completion," in *International Conference on Learning Representations*, 2022. [Online]. Available: https://openreview.net/forum?id=wqD6TfbYkrn

[72] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.

[73] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 315–323.

[74] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[75] L. Sifre and S. Mallat, "Rigid-motion scattering for texture classification," *arXiv preprint arXiv:1403.1687*, 2014.

[76] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, and S. Savarese, "Segcloud: Semantic segmentation of 3d point clouds," in *2017 international conference on 3D vision (3DV)*. IEEE, 2017, pp. 537–547.

[77] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," in *International Conference on Learning Representations*, 2019. [Online]. Available: https://openreview.net/forum?id=Bkg6RiCqY7

[78] H. Xiu, T. Shinohara, M. Matsuoka, M. Inoguchi, K. Kawabe, and K. Horie, "Collapsed building detection using 3d point clouds and deep learning," *Remote Sensing*, vol. 12, no. 24, p. 4057, 2020.