# PointAtrousGraph: Deep Hierarchical Encoder-Decoder with Point Atrous Convolution for Unorganized 3D Points

Liang Pan[1], Chee-Meng Chew[2] and Gim Hee Lee[3]

*Abstract*— **Motivated by the success of encoding multi-scale contextual information for image analysis, we propose our PointAtrousGraph (PAG) - a deep permutation-invariant hierarchical encoder-decoder for efficiently exploiting multi-scale edge features in point clouds. Our PAG is constructed by several novel modules, such as Point Atrous Convolution (PAC), Edge-preserved Pooling (EP) and Edge-preserved Unpooling (EU). Similar with atrous convolution, our PAC can effectively enlarge receptive fields of filters and thus densely learn multi-scale point features. Following the idea of non-overlapping max-pooling operations, we propose our EP to preserve critical edge features during subsampling. Correspondingly, our EU modules gradually recover spatial information for edge features. In addition, we introduce chained skip subsampling/upsampling modules that directly propagate edge features to the final stage. Particularly, our proposed auxiliary loss functions can further improve our performance. Experimental results show that our PAG outperform previous state-of-the-art methods on various 3D semantic perception applications.**

## I. INTRODUCTION

Owing to the effectiveness in capturing spatially-local correlations of convolution operations, deep convolution neural networks (CNNs) have yielded impressive results for many image-based tasks. In order to encode multi-scale contextual information, CNNs probe incoming image features with filters or pooling operations at multiple rates and multiple effective fields-of-view [1], such as Atrous Spatial Pyramid Pooling (ASPP) [2] and Pyramid Pooling Module [3], [4]. Both strategies effectively capture the contextual information at multiple scales, hence significantly improving the capabilities of CNNs. Unorganized point cloud is a simple and straight-forward representation of 3D structures, which is frequently applied by modern intelligent robotics applications, such as autonomous driving and human-robot interactions. However, the unorderedness and irregularity of 3D point clouds make the conventional convolution operation inapplicable. Despite novel filtering kernels are recently proposed, limited studies have been carried out on designing deep hierarchical encoder-decoder architectures to learn multi-scale point features for 3D semantic perception.

PointNet [5] is a pioneering ***permutation-invariant*** network that directly processes unordered point clouds by using many symmetric functions. It considers each 3D



Fig. 1. Point Atrous Convolution (better view in color). The figure on the left denotes the conventional method in selecting neighboring points $\{q_1, q_2, q_3, q_4, q_5\}$, and the right figure denotes our constructed neighborhood graph $\{q_2, q_4, q_6, q_8, q_{10}\}$ with sampling rate equals to 2. By adding the sampling rate parameter, our PAC can perform the convolution operation over a larger field of view without increasing computation load.

point independently and meanwhile overlooks local geometric details. PointNet++ [6] adaptively combines multi-scale local features by its proposed abstraction layer. There are nearly the same encoded local feature vectors if two selected controid points share the same local regions. These permutation-invariant networks [7], [8] that capture fine geometric structures from local neighborhoods provide better point cloud inference results, which is evident that features of local neighbors can improve deep learning on 3D points. Nonetheless, receptive fields of their filters are limited to small local regions by constructing local neighborhood graphs. Most recently, PAN [9] introduces a novel Point Atrous Convolution (PAC) module, which can effectively enlarge receptive fields of filters by introducing a sampling rate to equivalently sparsely sample the neighboring point features. PAN can densely exploit multi-scale edge features by using PAC modules with different sampling rates. However, PAN does not have a hierarchical encoder-decoder architecture, which becomes inefficient when dealing with high-dimensional dense point features.

In this paper, we propose the PointAtrousGraph (PAG) - a deep permutation-invariant hierarchical encoder-decoder to exploit multi-scale local geometric details with novel PAC modules for point cloud analysis. To address the **overlapped neighborhood graph problem**, we apply our edge-preserved pooling (EP) operation to preserve critical edge features during subsampling. Therefore, our PAG can exploit and preserve multi-scale local geometrical details hierarchically. In a similar fashion, our edge-preserved unpooling (EU) operation is applied to recover the spatial information of sparse high-dimensional point features. Furthermore, we introduce tailored chained skip subsampling/upsampling modules to directly propagate point features from each hierarchy. Additionally, we propose novel auxiliary loss functions, maximum mean discrepancy (MMD) loss and deeply supervised loss, which further increases our inference accuracy. Our PAG also requires less training memory consumption and shorter training time than most existing networks that highly rely on

[1]Liang Pan is with Advanced Robotics Centre, National University of Singapore `pan.liang@u.nus.edu`

[2]Chee-Meng Chew is with Department of Mechanical Engineering, National University of Singapore `chewcm@nus.edu.sg`

[3]Gim Hee Lee is with Computer Vision and Robotic Perception (CVRP) Lab, Department of Computer Science, National University of Singapore `dcslgh@nus.edu.sg`

neighborhood graphs for 3D points. Experiments show that our PAG achieves better performance than previous state-of-the-art methods in various point cloud applications, including 3D object classification, object-part segmentation and 3D semantic segmentation.

## II. RELATED WORK

**Unorganized Point Cloud Analysis.** *Due to the un-orderedness of 3D points, a point cloud with N 3D points has a total of N! permutations in the data feeding order. Hence, it is important for a network to maintain invariance of all possible permutations* [5]. The pioneering work PointNet [5] achieves permutation-invariance by applying symmetric functions. Many following works [6], [10], [7], [8] propose more complicated symmetric operations to exploit local geometrical details in 3D points. Semantic labeling on point cloud is more challenging than classification and object-part segmentation. SPG [11] and SGPN [12] both construct super point graphs to refine their semantic labeling results. PAN [9] proposes a novel PAC module to effectively exploit multi-scale local edge features. However, unlike many networks for semantic labeling tasks on images, they [11], [12], [13], [9] do not have hierarchical encoder-decoder architectures, which limits their performance.

**Hierarchical Encoder-Decoder.** Deep hierarchical encoder-decoder architectures are widely and successfully used for many image-based tasks, such as human pose estimation [14], [15], semantic segmentation [16], [17], [18], [19], [20], [21], [22], [23], [24], [25], optical flow estimation [26], [27], and object detection [28], [29], [30]. The encoder-decoder architecture, stacked hourglass module, is based on the successive steps of pooling and upsampling, which produces impressive results on human pose estimation [14]. Lin et al. [28] introduced the feature pyramid network for object detection. As for semantic segmentation tasks, U-Net [16] and DeconvNet [25] follow the symmetric encoder-decoder architectures, and they refine the segmentation masks by utilizing features in low-level layers. DeepLabv3+ [1] takes advantage of both the encoder-decoder architecture and the atrous convolution modules to effectively change the fields-of-view of filters to capture multi-scale contextual information, which provides new state-of-the-art performance on many semantic segmentation benchmarks. *Typically, deep hierarchical encoder-decoder architectures contain: (1) **an encoder module** that progressively reduces the feature resolution, enlarges the receptive fields of filters and captures higher semantic information; (2) **a decoder module** that gradually recovers the spatial information* [1].

## III. METHODS

Our PointAtrousGraph (PAG) is focused on learning multi-scale edge features by applying a deep hierarchical encoder-decoder architecture. To maintain the permutation-invariant property, our PAG is made up of symmetric functions, such as shared mlp, max-pooling and feature concatenation. In particular, the PAC module is applied as a fundamental
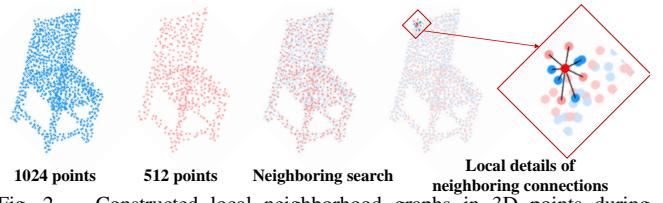


**Fig. 2.** Constructed local neighborhood graphs in 3D points during subsampling processes (better view in color).

building block, which can effectively learn multi-scale dense edge features in 3D points. Furthermore, we propose our edge-preserved pooling (EP) operation, which benefits in constructing deep hierarchical networks by preserving critical edge features during subsampling processes. Our EP operation also enlarges the receptive fields by decreasing 3D point feature density. In a similar manner, our edge-preserved unpooling (EU) operation gradually recovers the high-dimensional point feature density by considering 3D point spatial locations. We also directly propagate point features from different hierarchies to the final stage by tailored chained skip subsampling/upsampling modules. In addition, we propose novel auxiliary losses to further increase our inference accuracy.

### A. Point Atrous Convolution for Dense Point Feature

Two typical methods, ball query and k-nearest neighbors (kNN), are applied to exploit local geometric details in point clouds. However, the ball query algorithm applied by PointNet++ [6] always selects the first #K points in a specified search ball with a predefined radius, which cannot guarantee that closest points can be selected [31]. In order to exploit sufficient local contextual information in 3D points, conventional networks either construct large neighborhood graphs (large #K) [7], [8] or concatenate multi-scale local edge features (large #$C_f$) [32], [6]. Both above strategies, however, make their networks cumbersome and inefficient.

In contrast, Point Atrous Convolution (PAC) [9] can arbitrarily enlarge its receptive field in dense point features without increasing its computation volume (small #K and #$C_f$). Inspired by atrous convolution [33] for image analysis, PAC modules introduce an important sampling rate parameter $r$ to equivalently sparsely sample the neighboring point features. A PAC operation can be formulated as:

$$X'_p = g(H_\Theta(X_p, X_q^{1 \cdot r}), \ldots, H_\Theta(X_p, X_q^{k \cdot r})), \quad (1)$$

where $X_p$ is the feature of this centroid point $p$, $X_q^{k \cdot r}$ is the feature of the $(k \cdot r)^{th}$ nearest neighbor of point $p$, $r$ is the sampling rate, $k$ is the number of total searched neighboring points, $g(\cdot)$ denotes a max-pooling function, and $H_\Theta(\cdot)$ denotes the edge kernel $h_\theta(X_p \oplus (X_p - X_q^{k \cdot r}))$. $h_\theta$ is a shared mlp layer and $\oplus$ denotes feature concatenation operation. An example is also illustrated in Fig. 1.

### B. Edge-preserved Pooling

Pooling layers, especially ***non-overlapping max-pooling*** layers, are widely used in CNNs [34], [35], which summarizes the outputs of neighboring groups of neurons in the same kernel map in image domains [36]. Scherer et al. [37] report that the increment of step size of overlapping

Fig. 3. The overlapped neighborhood graph problem.

pooling windows deteriorates their recognition performance because *maxima in overlapping window regions are merely duplicated in the next layer and neighboring pixels are more correlated.* Cireşan et al. [38] replace subsampling layers with non-overlapping max-pooling layers in the CNNs [39], which achieves surprisingly rapid learning speed and better performance. Various subsampling methods on 3D point clouds [40], [6] have been proposed. Nonetheless, they either do not summarize local geometrical features, or ignore the problem caused by overlapped local neighborhood graphs.

An example for 3D points subsampling processes is illustrated in Fig. 2. The input point feature size is 1,024, and 512 point features are selected after a subsampling process. If only those features of "selected" points are propagated, local geometric details will thus be overlooked. Another strategy is to propagate edge features by considering the features of local neighboring points (also shown in Fig. 2). The searched neighboring point features often consist of both "selected" point features (shown in pink) and "discarded" point features (shown in blue) during subsampling. Two "selected" centroid points can often have overlapped neighborhood graphs, and even share the same local neighboring points. Hence, if we propagate the encoded neighboring point features, two neighboring "selected" points can have similar or even the same features. The respective information of each 3D point vanishes, especially when multiple subsampling operations are performed to construct a deep hierarchical network for 3D points. We entitle this as the **overlapped neighborhood graph problem** shown in Fig. 3.

In view of this, we propose our edge-preserved pooling (EP) module, which effectively captures local geometrical details while maintaining preserving respective features of each point. In line with the idea of non-overlapping max-pooling operations, we encode local edge features by considering neighboring point features in the original 1,024 points. Due to the absence of regular grids in 3D point clouds, we select neighboring point features by constructing neighborhood graphs in metric spaces. To preserve both distinctive individual point features and local geometrical edge features, our EP module is designed as:

$$X'_p = X_p \oplus g(X_q^1, ..., X_q^k), \tag{2}$$

where $X_p$ is the feature of selected centroid and $X_q^k$ is the point features of its $k$ nearest neighbors. Consequently, our EP operation explicitly propagates point features of each "selected" centroid and also summarizes its local point features, which is consistent with the idea behind non-overlapping max pooling for propagating image features.

## C. Edge-preserved Unpooling

To recover spatial information for image features, low-level image features from encoder are often applied to refine the high-level features in decoder, especially when

using a symmetric hierarchical encoder-decoder architecture [16], [25]. Our edge-preserved unpooling (EU) module also considers the point features of centroids and their local neighboring point features searched in metric spaces. Unlike PointNet++, our EU module does not need to consider the "d-dim coordinates" associated with each point:

$$X'_p = X_p^e \oplus w(X_q^1, ..., X_q^k), \tag{3}$$

where $X_p^e$ is the corresponding feature propagated from the encoder by a skip connection directly, $w(\cdot)$ is the inverse distance weighted average operation and $X_q^k$ is feature of the $k^{th}$ nearest neighbor of point $p$ in its previous hierarchy.

## D. Deep Hierarchical Encoder-Decoder

Based on our PAC, EP and EU modules, we construct the deep hierarchical encoder-decoder architecture PointAtrousGraph (PAG) to learn multi-scale features for 3D point classification and segmentation tasks (shown in Fig. 4). The same encoder architecture is applied by our classification and segmentation networks, which consists three hierarchies to gradually reduce the point feature density and meanwhile enlarge the receptive fields for learning higher semantic point features. Within each hierarchy, we successively lay out two PAC layers with increasing sampling rates to gradually exploit larger local geometric details. Our decoder architectures are designed differently with respect to different applications. In addition, we also propose different skip connection modules, chained skip subsampling and chained skip upsampling, for classification and segmentation, respectively.

**Classification Network.** Our classification network (enclosed by red dashed lines in Fig. 4) aims to encode a global point feature vector by exploiting multi-scale local geometrical details in a point cloud. The main stream is our network encoder, which unravels the multi-scale contextual information capturing problem by applying many PAC modules in a hierarchical fashion. In addition to the main stream which consecutively propagates features, we also propose the **chained skip subsampling** module (enclosed by blue dash-dotted lines in Fig. 4). The chained skip subsampling module progressively feeds forward features of each hierarchy to the final stage. In each hierarchy, we select the same set of centroid points with the corresponding EP module to construct local neighborhood graphs. However, we only propagate features of the neighbors in the chained skip subsampling operations, which is different from the EP module. After concatenating hierarchical point features from both streams, the global feature is obtained by applying a global max-pooling. Thereafter, two fully-connected (FC) layers are employed to yield the final classification results.

**Segmentation Network.** A segmentation task can be regarded as a per-point classification. Hence, our segmentation network (enclosed by golden dashed lines in Fig. 4) densely learns multi-scale edge features for each input 3D point. Accordingly, we propose a hierarchical decoder to progressively recover the high-dimensional point feature density. Our decoder for segmentation has a similar architecture as our encoder, which also has three hierarchies. Likewise, we lay out two PAC layers with decreasing sampling rates to
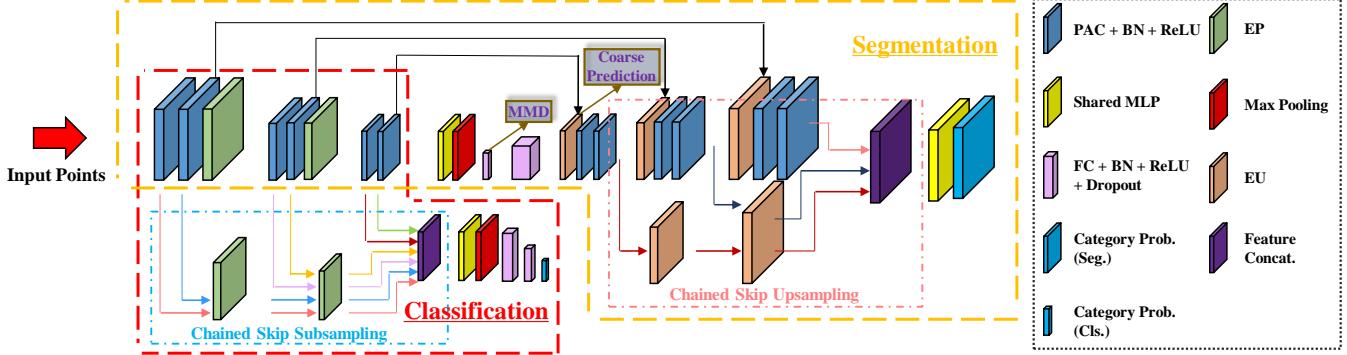
Fig. 4. Our PointAtrousGraph (PAG) architecture (better view in color). Our classification and segmentation networks have the same designed network encoder architecture. Our classification network is enclosed by red dashed lines, and our segmentation network is enclosed by golden dashed lines. The chained skip subsampling module is applied in our classification network, which is enclosed by blue dash-dotted lines. Likewise, our chained skip upsampling module is enclosed by pink dash-dotted lines in our segmentation network.

gradually aggregate features for each point. Particularly, we apply the regressed global features, which largely increases our performance. We also employ the **chained skip up-sampling** module to consecutively upsample point features of each hierarchy (enclosed by pink dash-dotted lines in Fig. 4). At each stage, we apply the same upsampling strategy with the EU module with the same constructed local neighborhood graphs. Unlike EU module, we do not concatenate the centroid features. Finally, we concatenate all the hierarchical features for each 3D point to perform pre-point predictions.

### E. Auxiliary Loss Functions

We also introduce auxiliary losses, maximum mean discrepancy (MMD) and deeply supervised losses, which are mainly applied for segmentation tasks.

MMD criterion [41] that is commonly used in variational auto-encoder architectures is performed over our embedded global point features. The MMD loss quantifies the similarity between two distributions by comparing all their moments. By applying the kernel trick, the MMD loss is defined as:

$$L_{mmd}(q\|p) = E_{q(z),q(z\prime)}[k(z,z\prime)] + E_{p(z),p(z\prime)}[k(z,z\prime)]$$
$$-2E_{q(z),p(z\prime)}[k(z,z\prime)], \qquad (4)$$

where $L_{mmd} \geq 0$, $q(z)$ denotes our embeded feature distribution and $p(z)$ denotes a prior Gaussian distribution (we use $N(\mu = 0, \sigma = 1.0)$). $L_{mmd} = 0$ if and only if $q = p$.

Similar with image-based segmentation networks [3], [42], we add a deeply supervised (cross-entropy) loss $L_{ds}$ at the first stage of our segmentation network decoder. Along with the master branch (cross-entropy) loss $L_{master}$, auxiliary loss functions also pass through all previous layers. Consequently, we train our segmentation network by minimizing the following joint loss function:

$$L_{all} = L_{master} + w_{mmd}L_{mmd} + w_{ds}L_{ds}, \qquad (5)$$

where $w_{mmd}$ and $w_{ds}$ are designed weights to balance corresponding auxiliary losses.

### F. Discussion

Many networks [5], [8], [40] that respect the permutation-invariant property are focused on designing convolution kernels for unordered 3D points. Previous studies [6], [7], [8], [43], [10] reveal effectiveness of considering local geometric details. However, limited attentions have been received on designing deep hierarchical encoder-decoder architectures for capturing multi-scale local contexts. In contrast, our PAG - a deep hierarchical encoder-decoder, adaptively learns multi-scale edge features by varying the field of view for filters in each layer. Motivated by the success of exploiting multi-scale contextual information by operations for image feature learning [17], [33], [44], [25], [3], such as atrous convolution, non-overlapping max-pooling and unpooling, we apply novel modules, including PAC, EP and EU, for our multi-scale edge feature learning in unorganized 3D points. Our PAC module effectively enlarges its receptive field in dense point features without increasing training parameters. With our EP modules, our PAG is capable of propagating high-dimensional features while preserving local edge features and respective information of each point. To recover the density for point features, we apply interpolation operations by considering features of those spatial neighboring points. The proposed auxiliary loss functions further increase our network performance. Moreover, our hierarchical architecture consists of multiple subsampling operations, which significantly decreases its computation amount and thus largely increases the efficiency.

## IV. EXPERIMENTS

Our PAG is evaluated on three point cloud analysis tasks, including shape classification, object-part segmentation and semantic segmentation. Without additional processes, our PAG outperform previous state-of-the-art methods.

### A. Implementation Details

Our networks are implemented with TensorFlow [45] on an NVIDIA GTX1080Ti GPU. We report our results with respect to different input and training strategies to achieve fair comparisons for shape classification. For segmentation tasks, we follow the same training and evaluation setting in [5]. To improve efficiency, we fix #K = 10 for all our neighborhood graphs. We apply the farthest point sampling (FPS) algorithm [6] for subsampling points by considering their metrics (static) or features (dynamic). Our code and models are publicly available on the project website[1].

---

[1]https://github.com/paul007pl/PointAtrousGraph

TABLE I

SHAPE CLASSIFICATION RESULTS ON MODELNET40 [46].

| | Method | Input | OA | Mem. | Time |
|---|---|---|---|---|---|
| With Up-axis Rotation | PointNet [5] | 1,024 pts | 89.2 | 2.4GB | 3-6h |
| | PointNet++ [6] | 1,024 pts | 90.7 | 11.1GB | ≃20h |
| | Wang(41-spec-cp) [47] | 1,024 pts | 91.5 | - | ≃12h |
| | MRTNet(kd-tree) [48] | 4,000 pts | 91.7 | | |
| | DGCNN [7] | 1,024 pts | 92.2 | 8.9GB | 11.4h |
| | PAN [9] | 1,024 pts | 92.2 | 9.2GB | 8.7h |
| | PAG (ours) | 1,024 pts | **92.2** | 2.4GB | 4.4h |
| | PointNet++ [6] | 5,000 pts + n | 91.9 | | |
| | SpiderCNN [8] | 1,024 pts + n | 92.4 | 11.3GB | |
| | Wang(41-spec-cp) [47] | 2,048 pts + n | 92.1 | - | ≃20h |
| | PAN [9] | 1,024 pts + n | 92.6 | 9.2GB | 8.8h |
| | PAG (ours) | 1,024 pts + n | **92.7** | 2.6GB | 4.5h |
| W/o Up-axis Rotation | KCNet [10] | 1,024 pts | 91.0 | | |
| | Atzmon et al. [49] | 1,024 pts | 92.3 | - | ≃25h |
| | PointCNN• [50] | 1,024 pts | 92.5 | | |
| | SO-Net [43] | 2,048 pts | 90.9 | - | ≃3h |
| | A-CNN [31] | 1,024 pts | 92.6 | | |
| | PAN [9] | 1,024 pts | 93.1 | 9.2GB | 8.7h |
| | PAG (ours) | 1,024 pts | **93.1** | 2.4GB | 4.4h |
| | SO-Net [43] | 5,000 pts + n | 93.4 | | |
| | PAN [9] | 5,000 pts + n | 93.4 | | |
| | PAG (ours) | 5,000 pts + n | **93.8** | | |

## B. Shape Classification

We evaluate the performance of our network for shape classification on the ModelNet40 dataset [46]. ModelNet40 benchmark contains 13,834 CAD models from 40 categories, and it is split into a training (9,843 models) and a test set (2,468 models). Most 3D models from ModelNet40 are pre-aligned to the common up-direction and horizontal-facing direction. To better approximate the scenarios in real world applications, the pre-aligned direction can be ignored by applying random up-axis rotations. For fair comparisons, we report our classification network performance in both training settings, provided Table I. The column "OA" provides the overall classification accuracy (percentage). The column "Mem." and "Time" denotes the required training memory consumption and time, respectively. According to Table I, our PAG yields the best classification accuracy for all the training settings. Furthermore, our PAG requires smaller memory footprint and shorter time for training than most existing networks.

**Ablation Study.** We report the ablation study results in Table II. Without applying our PAC modules, the classification accuracy dropped due to insufficient contextual information. Different from previous works [7], [8], setting the number of selected nearest neighbors as 20 does not increase our classification accuracy. Furthermore, setting a smaller number of selected nearest neighbors decreases the training memory requirement and also improves the training efficiency. Additionally, the proposed chained skip subsampling (CSS) module can further improve the classification performance. In our EP operations, if we only propagate "centroid" or "neighbor" features, our classification accuracy will drop significantly.

**Robustness to Sampling Density Variation.** In real applications, point clouds are often partially captured and thus become irregular and incomplete. Our PAG heavily relies on kNN, which cannot guarantee a fixed scale for exploited geometric details, especially for non-uniformly distributed
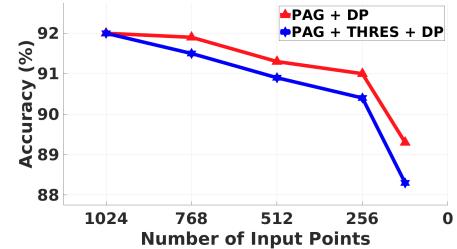
Fig. 5. Curve showing our advantage of dealing with non-uniform distributed point inputs. DP denotes randomly dropout input points, and THRES denotes additional radii thresholds.

TABLE II

ABLATION STUDY OF THE CLASSIFICATION NETWORK.

| K | PAC | Centroid/Neighbors | CSS | OA |
|---|---|---|---|---|
| 10 | × | both | ✓ | 92.0 |
| 20 | ✓ | both | ✓ | 91.9 |
| 10 | ✓ | both | × | 92.2 |
| 10 | ✓ | centroid | ✓ | 91.7 |
| 10 | ✓ | neighbors | ✓ | 92.0 |
| 10 | ✓ | both | ✓ | 92.7 |

3D points. To this end, we revise the PAC module by applying two searching radii to bound selected neighboring points in a specified searching field ($r_{min}$, $r_{max}$). Those searching radii denote the distance between a centroid point and its neighboring point. In line with [6], we randomly drop points (DP) to imitate non-uniform and sparse input points. In our experiments, we observe that applying additional radii negatively impact our performance (shown in Fig. 5). The reasons can be: 1) our PAG can exploit multi-scale point features to yield robust prediction; 2) it is difficult to manually set the predefined searching thresholds, which cannot be self-adaptive to irregularly incomplete 3D points.

## C. Object-part Segmentation

Object-part segmentation is demonstrated on the ShapeNet-part dataset [51], which contains 16,881 shapes represented as separate point clouds from 16 categories with per-point annotation. Following previous work [5], we split the dataset into a training (14,034 objects) and a test (2,847 objects) set. We apply the mean Intersection over Union (IoU) metric to evaluate our method on each point. For each instance, we obtain the mean IoU by averaging IoUs for all part types in the corresponding object category. The overall instance IoU ("pIoU") is computed by averaging IoUs over all the tested instances. Category-wise IoU is computed as the average of all the instances under the corresponding category. Mean category IoU ("mpIoU") is thus computed as the mean of all the category-wise IoUs. In Table III, our model presents the best overall mean IoU at **86.4%**. The qualitative segmentation results are shown in Fig. 6.

**Ablation Study.** Table IV provides the results for ablation study on the object-part segmentation task. The results

Fig. 6. Qualitative results of the part segmentation task.

TABLE III

SEGMENTATION COMPARISONS ON SHAPENET [51] AND S3DIS [52].

| Method | ShapeNet | | k-fold S3DIS | | Area-5 S3DIS | |
|---|---|---|---|---|---|---|
| | mpIoU | pIoU | OA | mIoU | OA | mIoU |
| SGPN● [12] | 82.8 | 85.8 | 80.8 | 50.4 | - | - |
| RSNet● [13] | 81.4 | 84.9 | - | 56.5 | - | - |
| SPG● [11] | - | - | 85.5 | 62.1 | 86.4 | 58.0 |
| SegCloud● [53] | - | - | - | - | - | 48.9 |
| PointCNN● [50] | **84.6** | 86.1 | 88.1 | 65.4 | 85.9 | 57.3 |
| PCCN [54] | - | - | - | - | - | 58.3 |
| PointNet [5] | 80.4 | 83.7 | 78.5 | 47.6 | - | 41.1 |
| PointNet++ [6] | 81.9 | 85.1 | - | - | - | - |
| SO-Net [43] | 81.0 | 84.9 | - | - | - | - |
| SPLATNET$_{3D}$ [55] | 82.0 | 84.6 | - | - | - | - |
| Atzmon et. al [49] | 81.8 | 85.1 | - | - | - | - |
| SpiderCNN [8] | 82.4 | 85.3 | - | - | - | - |
| DGCNN [7] | 82.3 | 85.2 | 84.1 | 56.1 | - | - |
| PAN [9] | 82.6 | 85.7 | 85.9 | 61.4 | - | - |
| A-CNN [31] | 84.0 | 86.1 | 87.3 | 62.9 | - | - |
| PAG (ours) | 84.0 | **86.4** | **88.1** | **65.9** | **86.8** | **59.3** |

illustrate that the proposed chained skip upsampling (CSU) module and the extracted global features are beneficial to increasing the segmentation accuracy. Our proposed PAC module further increases our segmentation accuracy. Setting the number of selected nearest neighbors as 20 does not increase our segmentation accuracy, but largely increases its training memory consumption. In our EP operations, ignoring either centroid point features or neighbor point features can lead to a drop of its accuracy. By adding global point features (GF) and auxiliary loss functions (Aux.L), our segmentation network can yield better results.

### D. Semantic Segmentation

We evaluate our model on large-scale 3D semantic segmentation (see Fig. 7) on the Stanford 3D Indoor Semantic Dataset (S3DIS) [52]. This dataset contains 3D RGB point clouds of 271 rooms from 3 different buildings split into 6 areas. To achieve fair comparison, we apply the same setting of the training strategy with PointNet [5]. We randomly sample 4,096 points in each block for training, and all the points are used for testing. The semantic segmentation results (k-fold and Area-5) are provided in Table III. Those methods [50], [12], [13], [11], [56] (denoted by ●) obtain unfair advantages by applying specific preprocessing and/or postprocessing procedures, such as super point graphs or recurrent neural networks. Moreover, they do not maintain the important permutation-invariant property, which makes their results unreliable and sensitive with respect to the points feeding order. Without additional processes, our permutation-invariant PAG provides the best results on the OA (**88.1%** and **86.8%**) and mIoU (**65.9%** and **59.3%**).

## V. ANALYSIS AND FUTURE DIRECTIONS

**Permutation-Invariance Property.** Our proposed modules are made up of symmetric functions, such as shared mlp, max-pooling and feature concatenation. Consequently, our PAG is a permutation-invariant network.

**Sampling Neighboring Points.** Ball query introduced by PointNet++ is another strategy to sample neighboring points. However, PointNet++ only select the first #K points found within the radius, which cannot guarantee that nearest points
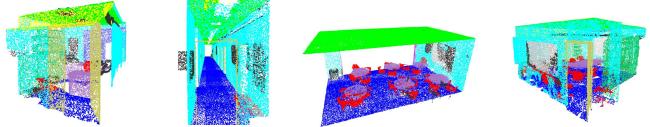


Fig. 7. Qualitative results on the semantic segmentation task.

TABLE IV

ABLATION STUDY OF THE SEGMENTATION NETWORK.

| K | PAC | Centroid/Neighbors | CSU | GF | Aux.L | pIoU |
|---|---|---|---|---|---|---|
| 10 | ✓ | both | ✓ | ✗ | ✗ | 85.3 |
| 10 | ✗ | both | ✓ | ✓ | ✗ | 85.7 |
| 20 | ✓ | both | ✓ | ✓ | ✗ | 85.7 |
| 10 | ✓ | both | ✗ | ✓ | ✗ | 85.6 |
| 10 | ✓ | centroid | ✓ | ✓ | ✗ | 85.8 |
| 10 | ✓ | neighbors | ✓ | ✓ | ✗ | 85.7 |
| 10 | ✓ | both | ✓ | ✓ | ✗ | 86.1 |
| 10 | ✓ | both | ✓ | ✓ | ✓ | 86.4 |

can be selected. Furthermore, PointNet++(MSG) search for many neighboring points for 3 times (e.g. #K={16,32,128}) and then concatenate all exploited features together (e.g. #C$_f$=64+128+128) in each hierarchy, which makes it computationally expensive [6]. In contrast, our PAC module equivalently sparsely sample the neighboring points in feature spaces with a sampling rate parameter "r". In our PAG, we consecutively lay out PAC modules in series to aggregate multi-scale edge features without increasing the selected neighboring points (#K) and the point feature size (#C$_f$).

**Hierarchical Encoder-Decoder Architecture.** Inspired by the success of encoding rich contextual image information at multiple scales for image-based applications, we focus on capturing multi-scale local geometrical details in 3D points by a hierarchical encoder-decoder architecture. Our PAG contains: (1) an encoder module that progressively decreases point feature density, enlarges the field of view of filters and learns higher semantic edge features; (2) a decoder module that gradually recovers the density for high-dimensional point features. Deep hierarchical CNNs commonly use the non-overlapping max-pooling operations, which decrease the image feature resolution while summarizing local-spatial contextual information. To resolve the **overlapped neighborhood graph problem**, we propose our EP operation, which concatenates both "centroid" and "neighbors" features for hierarchically propagating edge-aware features. In our ablation studies, we observe that the "neighbors" features benefit in generating global features for classification tasks. The "centroid" features are beneficial for propagating respective information of each individual point, hence leading to better segmentation results. Following spatial pyramid pooling operations [57], [58], we can also perform our EP operations at several scales for future studies.

## VI. CONCLUSION

We proposed the novel PointAtrousGraph (PAG) architecture to capture multi-scale local geometric details for hierarchical point features learning. Many novel modules are proposed and then evaluated in our experiments, which significantly increases our evaluation accuracy on classification and segmentation tasks. In particular, our network consumes much smaller training memory than previous state of the art models. Hence, our architectures can be conveniently applied for more complicated tasks of large-scale 3D scenes.

## Supplementary Material

### A. Overview

Our PointAtrousGraph (PAG) is focused on designing a deep permutation-invariant encoder-decoder architecture to learn hierarchical edge features for unorganized point cloud analysis. Consequently, we emphasize and elaborate those most related and important factors in our main paper: 1) **unorganized point cloud analysis**; 2) **permutation-invariance**; 3) **hierarchical encoder-decoder structure**. In this documents, more discussions, analysis and qualitative results are provided.

In Sec. *B*, we discuss and compare other methods and/or representations for organizing 3D scenes. In Sec. *C*, we elaborate the advantages and disadvantages for the two categories of deep learning methods for unorganized point cloud analysis: 1) employing symmetric functions; 2) learning canonical orderings. The details and ablation studies of our network architectures are provided in Sec. *D*, which also reveals the benefits and effectiveness of learning multi-scale local geometric details with both our PAC and edge-preserved pooling operations. The Sec. *E* provides more analysis on our point atrous convolution (PAC). Thereafter, we illustrate the intuition and effectiveness of applying k-nearest neighbors searches in feature spaces and metric spaces in Sec. *F*. In the end, more qualitative results are visualized in Sec. *G*.

### B. Organized 3D Structures vs. Unorganized Point Clouds

Besides the point cloud representation, 3D scenes can also be represented by volumetric grids and scalable indexing structures.

*1) Volumetric Grid:* Many works [59], [60], [46], [61], [62] use volumetric grids to regularize point clouds for deep learning. VoxNet [59] applies 3D CNN over a volumetric occupancy grid, but expensive computational and memory requirements limit its performance. VoxelNet [61] encodes a point cloud as a descriptive volumetric representation, and then it is connected with a region proposal network to generate detections. PointGrid [62] incorporates the 3D points within each volumetric grid cell. This allows the network to learn higher order local features of a 3D point cloud with a smaller memory footprint. SSCN [63] introduces submanifold sparse convolutions that maintain the sparsity of the voxelized input over many layers by fixing the location of active sites.

*2) Scalable Indexing Structure:* Scalable indexing structures, e.g. kd-tree and octree, are used in many existing works [64], [65], [66]. OctNet [64] hierarchically partitions the 3D space into a set of unbalanced octrees, which leads to a significant reduction in computational and memory resources. O-CNN [65] stores the octant information and CNN features into the graphics memory by designing an octree data structure. 3D CNN operations are only performed on those octants occupied by the 3D shape surfaces. Its computation time and memory requirement grow quadratically with the depth of the designed octree. Kd-Net [66] uses a kd-tree structure to construct the computational graph

with shared learnable parameters. Similar with CNNs, Kd-Net computes a sequence of hierarchical representations in a feed-forward bottom-up fashion.

*3) Comparison and Discussion:*

*a) Comparison with the Volumetric Grid Representation.:* The main advantages of using volumetric grids are as follows: 1) volumetric grid is the most straight-forward and understandable method to organize irregular 3D observations; 2) volumetric grids are regular and uniformly distributed; 3) a occupation grid or a distance field can be used as a mean of data representation with the help of a volumetric grid.

However, its disadvantages are also apparent: 1) volumetric grid is notorious for its memory and computationally inefficiency; 2) volumetric grid does not scale well for dense 3D data, which is incapable of exploiting the rich and detailed geometry of original 3D data; 3) the inference results can be sensitive to the size of the volumetric grid as well as the 3D grid resolution.

*b) Comparison with the Scalable Indexing Structure Representation.:* Owing to the inefficiency of organizing 3D points with volumetric grids, scalable indexing structures, such as kd-tree and octree, are utilized to organize 3D point clouds. It forms the computational graph to apply 3D convolutions level by level, where the convolution kernels often share learnable parameters. Consequently, the sequence of hierarchical representations can be computed in a feedforward bottom-up fashion [62]. The most important advantage on applying these structures is their capability on exploiting the sparsity of 3D data, and thus adaptively allocate computational and memory resources with respect to the data density. However, these data structures are more complicated, which makes it a very difficult task to implement the networks.

*c) Discussion.:* Point cloud is the most straight-forward representation and simple representation of 3D structures. Designing networks directly on 3D points can be much more flexible and convenient in comparison with applying the two organization methods above. Furthermore, observed 3D scenes in real applications, are usually irregularly and unpredictably distributed and occluded, which can lead to very unreliable understandings if they are organized by regular structures. The major difficult of applying the point cloud representation is the absence of 3D point ordering as well as the irregular distribution. To mitigate this limitation, two categories of networks are proposed. The methods belonging to the first category are following the pioneering network, PointNet [5], which applies symmetric functions to maintain the important permutation-invariant property. The other category is trying to learn a canonical ordering in 3D points, and the representative network is PointCNN [50]. We are going to discuss and compare the two categories of methods in detail in the next section.

### C. Learning Canonical Ordering vs. Utilizing Symmetric Functions

Due to the unorderedness of 3D points, a point cloud with N 3D points has a total of $N!$ permutations in the data

| Network Configuration | Hierarchy | K | Subsampling Rate | Acc. |
|---|---|---|---|---|
| ([32, 1],  [64, 2],  [128, 2],  [256, 4],  [512, 8]) | 1 | 10 | None | 92.4 |
| ([32, 1],  [64, 2],  [128, 4],  [256, 8]) | 1 | 10 | None | 92.3 |
| ([32, 2];  [64, 2];  [128, 2];  [256, 2];  [512, 2]) | 5 | 10 | 2 | 92.2 |
| ([32, 2];  [64, 2];  [128, 2];  [256, 2];  [512, 2]) | 5 | 5 | 2 | 92.1 |
| ([32, 1];  [64, 1];  [128, 2];  [256, 2];  [512, 2]) | 5 | 10 | 2 | 92.2 |
| ([32, 1], [32, 2];  [64, 1], [64, 2];  [128, 1], [128, 2];  [256, 1], [256, 2]) | 4 | 10 | 2 | 92.3 |
| ([32, 1], [64, 2];  [128, 1], [256, 2]) | 2 | 10 | 2 | 92.3 |
| ([64, 1], [64, 2];  [128, 1], [128, 2]) | 2 | 10 | 2 | 92.1 |
| ([64, 1], [64, 2];  [256, 1], [256, 2]) | 2 | 10 | 2 | 92.4 |
| ([64, 1], [64, 2];  [256, 1], [256, 2]) | 2 | 10 | 4 | 92.3 |
| ([32, 1], [32, 2];  [64, 1], [64, 2];  [128, 1], [128, 2]) | 3 | 10 | 2 | 92.1 |
| ([64, 2], [64, 4];  [128, 2], [128, 4];  [256, 2], [256, 4]) | 3 | 5 | 2 | 92.0 |
| ([64, 1], [64, 2];  [128, 1], [128, 2];  [256, 1], [256, 2]) | 3 | 10 | 2 | 92.3 |
| ([32, 1], [32, 2];  [64, 1], [64, 2];  [128, 1], [128, 2]) | 3 | 5 | 4 | 91.5 |
| ([64, 2], [64, 4];  [128, 2], [128, 4];  [256, 2], [256, 4]) | 3 | 10 | 4 | 91.8 |
| ([64, 1], [64, 2];  [256, 1], [256, 2];  [512, 1], [512, 2]) | 3 | 10 | 4 | 91.7 |
| ([64, 1], [64, 2];  [128, 1], [128, 2];  [256, 1], [256, 2]) | 3 | 10 | 4 | 92.7 |

feeding order. We hope that our inference results can be stable, reliable and also invariant with the order for input points. Two strategies are thus proposed: 1) preserving the permutation-invariance by only using symmetric functions; 2) learning a canonical ordering. PointNet [5] and its following works [6], [8], [7], [10], [43] maintain the permutation-invariant property by only applying symmetric functions, which is effectively applied and validated by many applications. Belonging to this category, all the operations in our PAG are symmetric, which makes our PAG a permutation-invariant architecture.

PointCNN [50] is recently released, which aims to learn an $\chi$-transformation from the input points to permutate themselves into canonical orders. Ideally, their inference results can be invariant to the input point ordering with the help of this $\chi$-transformation. ***However, their learned $\chi$-transformations are far from ideal, especially with the permutation equivalence aspect.*** In other words, they did not manage to resolve the canonical ordering learning problem with their $\chi$-transformations. Although they provided state-of-the-art results on various benchmarks, their results can be dependent and sensitive to the 3D points feeding orders. Furthermore, they do not provide the covariance analysis, which is important and necessary for analyzing the influence of the point feeding orderings on the inference accuracy. In addition, PointCNN prepare their own datasets for training and testing, and they also apply many specific strategies during training. All these operations provide additional and unfair advantages for PointCNN. Consequently, the results provided by PointCNN cannot be directly compared with those permutation-invariant networks that follow the idea behind PointNet. As an early-stage work, a rigorous understanding of PointCNN is still an open and unsolved problem.

### D. Network Architecture Details

Our classification and segmentation networks share the same encoder architecture. The encoder configuration that provides the best results is Encoder([64, 1], [64, 2]; [128, 1], [128, 2]; [256, 1], [256, 2]), where $[C, r]$ denotes the feature channel size $C$ and sampling rate $r$, respectively. The subsampling rate between each two hierarchies is set as 4.

For example, if the input has 1,024 points, then 256 points will be selected after an EP process.

After the encoder, we add a shared mlp layer to project each point feature to high-dimensional feature space. The fully-connected (FC) layers in our classification network have the following feature size: 512, 256 and the number of total classes in the evaluated dataset. As for our segmentation network, we set the FC layers with the following feature size: 512 and 1024. The decoder of our segmentation network has a similar architecture with our encoder, but all the configurations are set inversely: Decoder([256, 2], [256, 1]; [128, 2], [128, 1]; [64, 2], [64, 1]). Likewise, the upsampling rate is also set as 4. Each layer in the chained skip subsampling/upsampling layers has the same configuration with the corresponding layer in the main stream.

The configuration of our decoder for classification network is selected and fixed in advance, and our decoder for segmentation network follows the setting of our network encoder. Therefore, it is important to find the best configuration for our network encoder. We have evaluated multiple different encoder architectures on the classification task, as reported in Table V. In our experiments, we observe that the idea of learning edge features at different scales of receptive fields indeed improves our inference accuracy as well as convergence speed. Only apply the PAC modules (hierarchy equals to 1) can improve the inference accuracy, but the training can be inefficient, especially when dealing with dense high-dimensional point features. Applying hierarchical architectures without the PAC modules can lead to better efficiency but limited receptive fields in each network hierarchy. Therefore, our PAC and EP modules benefit each other for improving the learning capabilities on exploiting multi-scale local geometric details. **On the one hand, PAC can enlarge the field-of-view for filters without decreasing point feature density. On the other hand, EP increases the receptive fields as well as the network efficiency.** There can be some fluctuations in training performance from time to time. However, our selected configuration often provides relatively steady results, especially when applying the non-rotation training strategy.

**Time and Memory Complexity.** Table VI reports the total

TABLE VI

MODEL SIZE AND INFERENCE TIME.

| Method | #Params(M) | Fwd.(ms) |
|---|---|---|
| PointNet(vanilla) [5] | 0.8 | 11.5 |
| PointNet [5] | 3.5 | 29.3 |
| PointNet++(msg) [6] | 1.7 | 98.1 |
| DGCNN [7] | 1.8 | 173.7 |
| SpiderCNN(3-layers) [8] | 3.2 | 275.4 |
| PAG (Ours) | 1.8 | 109.8 |

number of parameters that need to be trained, and the average inference time of a few representative networks. Compared with the other networks [7], [8] that highly rely on nearest neighborhood graphs, our model takes shorter time for a single forward process.

*E. Point Atrous Convolution Analysis*

Our point atrous convolution is inspired by the atrous convolution [33] in image domain. Conventional convolution is restricted in local regions due to applying the $3 \times 3$ convolution kernels. To enlarge the receptive fields of convolution kernels without increasing the overall computation amount, a sampling rate $r$ is introduced in conventional convolution kernel for selecting and sampling the neighboring image features. The most important advantage of atrous convolution is that it effectively enlarges the field of view for each operation kernel without sacrificing the feature density, which is evaluated by many image-based applications, especially the semantic segmentation tasks [33], [2], [1], [67].

In view of the success of atrous convolution in image domains, we extend the idea of atrous convolution for point cloud analysis. In a similar spirit, we apply the sampling rate $r$ to equivalently sparsely sample and select neighboring points to exploit and aggregate contextual information from a larger local neighborhood region. Hence, we introduce our point atrous convolution (PAC). Our PAC enables our network to exploit local contextual information in 3D points with flexible scales. However, it can be inefficient when we propagate point features to the high-dimensional feature space. In this work, we also propose our edge-preserved pooling (EP) operation, which is also able to enlarge the receptive fields for our convolution kernels. Nonetheless, the point feature density is decreased after an EP operation. In our experiments, **we observe that both operations benefit each other mutually.** More results with different network architectures are further elaborated in Sec. VI-D.

*F. Mixture of K-Nearest Neighbors Search in Feature Spaces and Metric Spaces*

In our PAG architecture, we exploit neighboring point features in feature spaces and metric spaces. Our PAC module is applied to learn edge features by considering local neighboring point features that are selected in feature spaces. In our experiments, we observe that constructing neighborhood graphs in feature spaces for edge feature extracting can lead to better inference results as well as faster convergence speed during training than constructing neighborhood graphs in metric spaces, as shown in Fig. 8.

However, as for constructing deep hierarchical networks, the proposed edge-preserved pooling (EP) and edge-
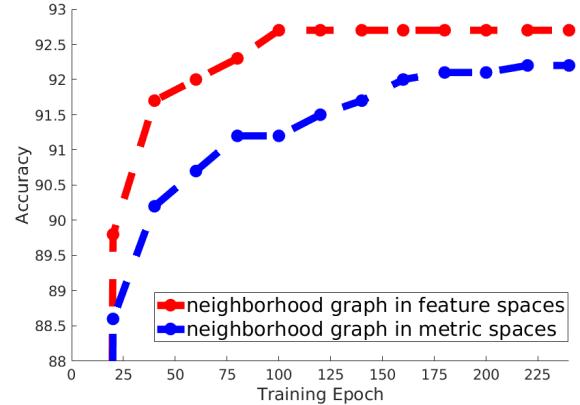


Fig. 8. Test accuracy (show the best) for PAC selecting neighboring points in feature spaces and metric spaces during training.

preserved unpooling (EU) both select neighboring points in the metric spaces. The intuitions of designing EP and EU modules are from the pooling and unpooling operations in image domains. Those operations for image analysis summarize/recover local contextual information with the help of regular image patterns. Due to the absence of regular grids in 3D points, the local geometric details can be summarized/recovered by constructing local neighborhood graphs in metric spaces instead. Moreover, the constructed neighborhood graphs in metric spaces are static, whereas the neighborhood graphs in feature spaces are dynamic. Accordingly, maintaining a stable hierarchical structure can benefit the convergence of the overall network structure, especially the EU operation, which is the second reason. Last but not the least, we discuss and analyze an important problem - the overlapped neighborhood graph problem, in the main paper. If we select neighboring point features in feature spaces, the overlapped neighborhood graph problem can be frequently encountered, which also deviates from our original intention - to capture characteristic local geometric details.

*G. More Qualitative Results*

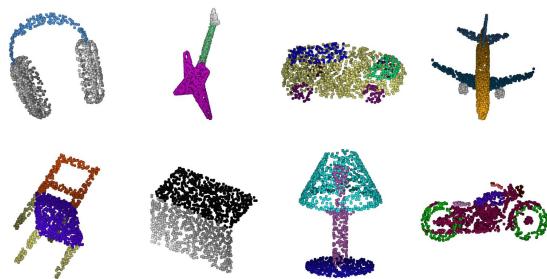More qualitative results are visualized in this section.



Fig. 9. More qualitative results on ShapeNet [51]

REFERENCES

[1] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 801–818.
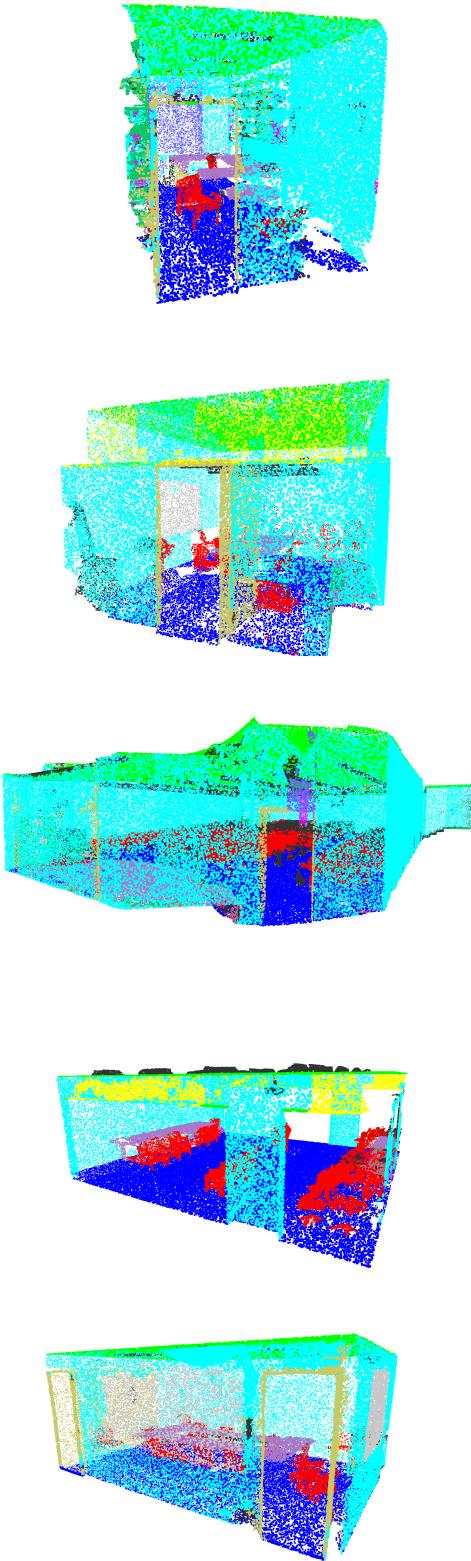
Fig. 10. More qualitative results for 3D semantic segmentation on S3DIS dataset [52].

[2] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017.

[3] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2881–2890.

[4] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.

[5] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, vol. 1, no. 2, p. 4, 2017.

[6] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *Advances in Neural Information Processing Systems*, 2017, pp. 5099–5108.

[7] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *arXiv preprint arXiv:1801.07829*, 2018.

[8] Y. Xu, T. Fan, M. Xu, L. Zeng, and Y. Qiao, "Spidercnn: Deep learning on point sets with parameterized convolutional filters," *arXiv preprint arXiv:1803.11527*, 2018.

[9] L. Pan, P. Wang, and C. M. Chew, "Pointatrousnet: Point atrous convolution for point cloud analysis," *IEEE Robotics and Automation Letters*, 2019.

[10] Y. Shen, C. Feng, Y. Yang, and D. Tian, "Mining point cloud local structures by kernel correlation and graph pooling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, vol. 4, 2018.

[11] L. Landrieu and M. Simonovsky, "Large-scale point cloud semantic segmentation with superpoint graphs," *arXiv preprint arXiv:1711.09869*, 2017.

[12] W. Wang, R. Yu, Q. Huang, and U. Neumann, "Sgpn: Similarity group proposal network for 3d point cloud instance segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2569–2578.

[13] Q. Huang, W. Wang, and U. Neumann, "Recurrent slice networks for 3d segmentation of point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2626–2635.

[14] A. Newell, K. Yang, and J. Deng, "Stacked hourglass networks for human pose estimation," in *European Conference on Computer Vision*. Springer, 2016, pp. 483–499.

[15] G. Pavlakos, X. Zhou, K. G. Derpanis, and K. Daniilidis, "Coarse-to-fine volumetric prediction for single-image 3d human pose," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7025–7034.

[16] 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.

[17] V. Badrinarayanan, A. Kendall, and R. Cipolla, "Segnet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2481–2495, 2017.

[18] T. Pohlen, A. Hermans, M. Mathias, and B. Leibe, "Full-resolution residual networks for semantic segmentation in street scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4151–4160.

[19] G. Lin, A. Milan, C. Shen, and I. Reid, "Refinenet: Multi-path refinement networks for high-resolution semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1925–1934.

[20] C. Peng, X. Zhang, G. Yu, G. Luo, and J. Sun, "Large kernel matters– improve semantic segmentation by global convolutional network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4353–4361.

[21] M. Amirul Islam, M. Rochan, N. D. Bruce, and Y. Wang, "Gated feedback refinement network for dense image labeling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3751–3759.

[22] Z. Wojna, V. Ferrari, S. Guadarrama, N. Silberman, L.-C. Chen, A. Fathi, and J. Uijlings, "The devil is in the decoder," *arXiv preprint arXiv:1707.05847*, 2017.

[23] J. Fu, J. Liu, Y. Wang, J. Zhou, C. Wang, and H. Lu, "Stacked deconvolutional network for semantic segmentation," *IEEE Transactions on Image Processing*, 2019.

[24] Z. Zhang, X. Zhang, C. Peng, X. Xue, and J. Sun, "Exfuse: Enhancing feature fusion for semantic segmentation," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 269–284.

[25] H. Noh, S. Hong, and B. Han, "Learning deconvolution network for semantic segmentation," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1520–1528.

[26] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2758–2766.

[27] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "Flownet 2.0: Evolution of optical flow estimation with deep networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2462–2470.

[28] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2117–2125.

[29] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta, "Beyond skip connections: Top-down modulation for object detection," *arXiv preprint arXiv:1612.06851*, 2016.

[30] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg, "Dssd: Deconvolutional single shot detector," *arXiv preprint arXiv:1701.06659*, 2017.

[31] A. Komarichev, Z. Zhong, and J. Hua, "A-cnn: Annularly convolutional neural networks on point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7421–7430.

[32] M. Jiang, Y. Wu, and C. Lu, "Pointsift: A sift-like network module for 3d point cloud semantic segmentation," *arXiv preprint arXiv:1807.00652*, 2018.

[33] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2018.

[34] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[35] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[36] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[37] D. Scherer, A. Müller, and S. Behnke, "Evaluation of pooling operations in convolutional architectures for object recognition," in *International conference on artificial neural networks*. Springer, 2010, pp. 92–101.

[38] D. C. Ciresan, U. Meier, J. Masci, L. M. Gambardella, and J. Schmidhuber, "Flexible, high performance convolutional neural networks for image classification," in *Twenty-Second International Joint Conference on Artificial Intelligence*, 2011.

[39] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[40] F. Groh, P. Wieschollek, and H. Lensch, "Flex-convolution (deep learning beyond grid-worlds)," *arXiv preprint arXiv:1803.07289*, 2018.

[41] S. Zhao, J. Song, and S. Ermon, "Infovae: Information maximizing variational autoencoders," *arXiv preprint arXiv:1706.02262*, 2017.

[42] H. Zhao, Y. Zhang, S. Liu, J. Shi, C. Change Loy, D. Lin, and J. Jia, "Psanet: Point-wise spatial attention network for scene parsing," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 267–283.

[43] J. Li, B. M. Chen, and G. H. Lee, "So-net: Self-organizing network for point cloud analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9397–9406.

[44] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3431–3440.

[45] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "Tensorflow: a system for large-scale machine learning." in *OSDI*, vol. 16, 2016, pp. 265–283.

[46] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1912–1920.

[47] C. Wang, B. Samari, and K. Siddiqi, "Local spectral graph convolution for point set feature learning," *arXiv preprint arXiv:1803.05827*, 2018.

[48] M. Gadelha, R. Wang, and S. Maji, "Multiresolution tree networks for 3d point cloud processing," *arXiv preprint arXiv:1807.03520*, 2018.

[49] M. Atzmon, H. Maron, and Y. Lipman, "Point convolutional neural networks by extension operators," *arXiv preprint arXiv:1803.10091*, 2018.

[50] Y. Li, R. Bu, M. Sun, and B. Chen, "Pointcnn," *arXiv preprint arXiv:1801.07791*, 2018.

[51] L. Yi, V. G. Kim, D. Ceylan, I. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, L. Guibas *et al.*, "A scalable active framework for region annotation in 3d shape collections," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 6, p. 210, 2016.

[52] 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.

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

[54] S. Wang, S. Suo, W.-C. Ma, A. Pokrovsky, and R. Urtasun, "Deep parametric continuous convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2589–2597.

[55] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz, "Splatnet: Sparse lattice networks for point cloud processing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2530–2539.

[56] F. Engelmann, T. Kontogianni, A. Hermans, and B. Leibe, "Exploring spatial context for 3d semantic segmentation of point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 716–724.

[57] K. Grauman and T. Darrell, "Pyramid match kernels: Discriminative classification with sets of image features (version 2)," 2006.

[58] S. Lazebnik, C. Schmid, and J. Ponce, "Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2. IEEE, 2006, pp. 2169–2178.

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

[60] C. R. Qi, H. Su, M. Nießner, A. Dai, M. Yan, and L. J. Guibas, "Volumetric and multi-view cnns for object classification on 3d data," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5648–5656.

[61] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," *arXiv preprint arXiv:1711.06396*, 2017.

[62] T. Le and Y. Duan, "Pointgrid: A deep network for 3d shape understanding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9204–9214.

[63] B. Graham, M. Engelcke, and L. van der Maaten, "3d semantic segmentation with submanifold sparse convolutional networks," *arXiv preprint arXiv:1711.10275*, 2017.

[64] G. Riegler, A. O. Ulusoy, and A. Geiger, "Octnet: Learning deep 3d representations at high resolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, vol. 3, 2017.

[65] 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, p. 72, 2017.

[66] R. Klokov and V. Lempitsky, "Escape from cells: Deep kd-networks for the recognition of 3d point cloud models," in *Computer Vision (ICCV), 2017 IEEE International Conference on*. IEEE, 2017, pp. 863–872.

[67] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.