

SPNet: Multi-Shell Kernel Convolution for Point Cloud Semantic Segmentation

Yuyan Li*, Chuanmao Fan*, Xu Wang*, and Ye Duan

University of Missouri, USA
{y1235, cf7b6, xwf32, duanye}@umsystem.edu

Abstract. Feature encoding is essential for point cloud analysis. In this paper, we propose a novel point convolution operator named Shell Point Convolution (SPConv) for shape encoding and local context learning. Specifically, SPConv splits 3D neighborhood space into shells, aggregates local features on manually designed kernel points, and performs convolution on the shells. Moreover, SPConv incorporates a simple yet effective attention module that enhances local feature aggregation. Based upon SPConv, a deep neural network named SPNet is constructed to process large-scale point clouds. Poisson disk sampling and feature propagation are incorporated in SPNet for better efficiency and accuracy. We provided details of the shell design and conducted extensive experiments on challenging large-scale point cloud datasets. Experimental results show that SPConv is effective in local shape encoding, and our SPNet is able to achieve top-ranking performances in semantic segmentation tasks.

Keywords: Point Cloud · Semantic Segmentation · Attention · Deep Learning

1 Introduction

Deep learning has achieved great success in image classification [9,10,12], semantic segmentation [20,35] and object detection [5,4,8]. However, deep learning based point cloud analysis is still a challenging topic. One major reason is that point clouds are non-uniformly sampled from a large, continuous 3D space, making them lack of regular grid structure. To tackle this, one straightforward approach is to voxelize point cloud into 3D regular grids and utilize standard 3D Convolutions [28,16]. But the voxelization approach has a major limitation, the discretization step inevitably loses geometric information. To address this problem, many researchers have proposed approaches to directly process point clouds. One of the seminal works is PointNet proposed by Qi et al. [23]. It uses Multi Layer Perceptron (MLP) and global pooling to preserve permutation invariance and gather a combination of local and global feature presentation. This work is further improved in their follow-up work PointNet++ [24] which adds the local geometry pooling/sampling over local neighborhood. Later works seek

* Equal contribution

other ways to enhance local feature aggregation. For example, Pointweb [37] constructs a dense fully-linked web, ShellNet [36] conducts convolution based on statistics from concentric spherical shells. Besides point-based methods, several graph-based methods are proposed to capture 3D shape and structures of point clouds. Methods such as [15,14] treat point clouds as nodes in a graph whose edges carry learnable affinity/similarity between adjacent points.

Recently, there has been another thread of research [34,31,32,29,17] that proposes learnable kernel functions which define convolutional kernels on a continuous space. One of these approaches is KPConv [29] introduced by Thomas et al. KPConv [29] proposes kernel point operator that consists of a set of local point filters which simulate 2D image convolution processes. Features from unordered point clouds are aggregated on either rigid or deformable kernel points. Using structural kernel points makes convolution feasible on a continuous space.

Following this line of work, we propose a novel multi-shell kernel point convolution named SPConv. Our SPConv operator partitions local 3D space into shells, each shell contains a set of rigid kernel points which aggregate local supporting point features. We perform kernel point convolution on each shell individually, then integrate the output features by an additional 1D convolution operation. The last convolution learns the contributions from shells and enhances shell correlation. An illustration of our SPConv is shown in Figure 1. Comparing to deformable KPConv [29], our method has an enhanced structure learning module, and does not require additional regularization during training. Furthermore, we find that incorporating low-level features such as color, normal, etc. in all layers for local feature re-weighting can be very effective for improving network performance. We propose two different approaches to accomplish the task, (1)Gaussian function based and (2)learning based. The first approach is hand-crafted and does not necessarily add GPU computation. The second approach has learnable weights and shows more robustness.

Using SPConv as our building block, we build a deep architecture SPNet. Similar to standard CNNs which utilize downsampling and upsampling strategy to reduce computation cost as well as to enlarge receptive field, we use Poisson disk sampling (PDS) for downsampling, and feature propagation (FP) for upsampling. In section 4, we evaluate the effectiveness of our methods on the most competitive indoor segmentation datasets. Notably, experimental results show that we achieve top-ranking performances. Our main contribution is summarized as follows:

- We propose a multi-shell kernel convolution operator that shows powerful local shape encoding ability.
- We introduce a simple yet effective attention mechanism for local neighbor feature re-weighting. This attention module improves performance and speeds up convergence.
- We present a comprehensive architecture design which outperforms stage-of-the-arts on challenging large-scale indoor datasets.

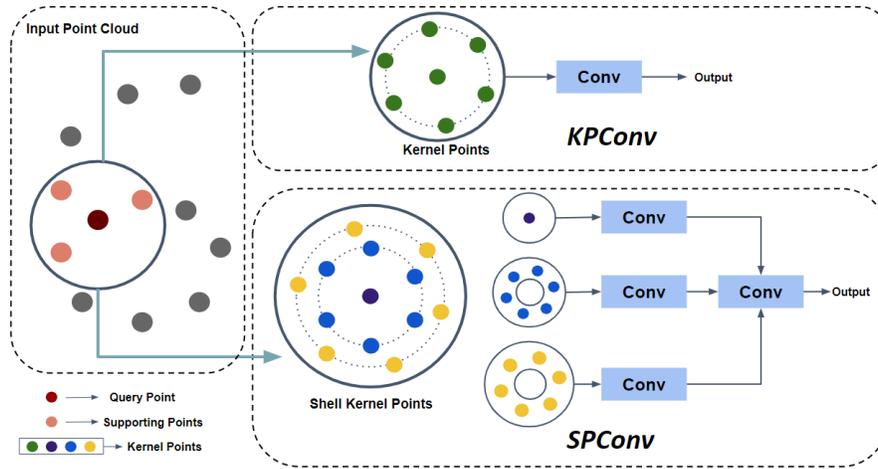


Fig. 1. Comparison between SPConv and KPCConv. For a query point, a range search is performed to locate supporting points. KPCConv defines a set of kernel points to aggregate local features and performs point convolution. Our SPConv has multiple shells, each shell contains one set of kernel points. Point convolutions are conducted for shells individually to encode distinctive geometric information. An additional convolution layer is used to fuse shell outputs together, as an enhancement of structure correlation.

2 Related Work

View-based and Voxel-based Methods. One classic category of point cloud representations is multi-view representation. MVCNN [27] renders 3D shape into images from various viewpoints and combines features from CNNs to predict point labels. However, these methods suffer from surface occlusion and density variation, which make it difficult to capture the internal structure of the shape. Another strategy is to convert point cloud into a 3D voxel structure which can be processed by standard 3D convolutions. VoxNet [21] and subsequent work [33,19] discretize point cloud into 3D volumetric grids. To improve efficiency on processing high resolution 3D voxels, recent researches [6,2] process volumetric data only on non-empty voxels.

Point-wise MLP Methods. Point-based methods receive great attention since PointNet [23] was proposed. In PointNet [23], points go through shared MLPs to obtain high dimensional features followed by a global max-pooling layer. In order to capture local neighborhood context, PointNet++ [24] is developed by hierarchically applying pointnet in local regions. There are extensive works based on PointNet++. For example, PointWeb [37] builds a dense fully connected web to explore local context, and uses an Adaptive Feature Adjustment module for feature refinement. ShellNet [36] proposes a ShellConv operator with concentric spherical shells to capture representative features.

Point Convolution Methods. Some recent works define explicit kernels for point convolution. KNet [25] develops a kernel correlation layer to compute affinities between each point’s K nearest neighbors and a predefined set of kernel points. Local features are acquired by graph pooling layers. SpiderCNN [34] designs a family of Taylor polynomial kernels to aggregate neighbor features. PointCNN [18] introduces X-transformation to exploit the canonical order of points. PCNN [31] builds a network using parametric continuous convolutional layers. SPH3D [17] uses spherical harmonic kernels during convolution on quantized space to identify distinctive geometric features. Our work is most related to KPConv [29], which defines rigid and deformable kernel points for feature aggregation. This convolution operator resolves point cloud ambiguity, alleviates varying density, and shows superior performances. Compared to KPConv [29], our SPConv enhances local structure correlation by incorporating shell-structured kernel points and learning on a larger neighborhood context.

3 Methodology

3.1 Review on Kernel Point Convolution

KPConv [29] effectively resolves the point cloud ambiguity by placing manually designed kernel points in a local neighborhood. This convolution simulates image-based convolutions. A typical image based 2D convolution with a $(2m + 1) \times (2m + 1)$ kernel at location $i, j \in \mathbb{Z}$ is defined as:

$$F * W = \sum_{x=-m}^m \sum_{y=-m}^m F(i-x, j-y)W(i, j) \quad (1)$$

where $x, y \in \{-m, \dots, m\}$, W is the learnable weight, $F(i, j)$ is the feature for pixel (i, j) . Image based convolution describes a one-to-one relationship between single kernel and single image pixel. Similarly, for a 3D point $p \in \mathbb{R}^3$ with a local neighborhood of radius R , point convolution can be defined as:

$$F * W = \sum_k^K F(p_k, p)W(k) \quad (2)$$

where $F(p_k, p)$ is the aggregated features on kernel point p_k . p_k carries learnable matrix $W_k \in \mathbb{R}^{C_{in} \times C_{out}}$. C_{in} and C_{out} are input and output feature channels respectively. With proper aggregation approach, the structure of supporting points can be well captured and learned by weight W . There are two key components of this convolution, placements of kernel points and aggregation function.

For a 3D point $x \in \mathbb{R}^3$ surrounded by neighboring points $x_j \in \mathbb{R}^3$ within a ball radius R , kernel points are distributed on the surface of a sphere with radius r , plus one point placed at center. Aggregated features $F(p_k, p)$ for kernel point p_k are computed as the sum of the features carried by neighboring points that fall into the influenced radius v . These neighboring features are weighted based

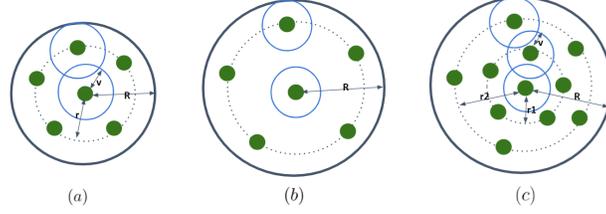


Fig. 2. A KPConv [29] operator is defined in (a) with kernel point radius r , neighborhood size R , kernel influence radius v . Kernel points have overlapping influence regions. When enlarging neighborhood size R to capture a larger context as shown in (b), kernel points become sparse and this may cause a loss of information for complex scenes with objects of different scales. Our SPConv (c) has multiple shells and keeps overlapping influence regions. r_1 and r_2 are kernel point radius for the second and third shell.

upon the Euclidean distance between p_k and p_j . An illustration of KPConv is shown in Figure 1.

$$F(p_k, p) = \sum_{p_j, \|p_j - p\| < R, \|p_k - p_j\| < v} F_{p_j} d(p_k, p_j) \quad (3)$$

where $d(p_k, p_j)$ denotes the correlation of kernel point p_k and a neighbor point p_j . This correlation can be calculated by a linear function:

$$d(p_k, p_j) = \max(0, 1 - \frac{\|p_k - p_j\|}{v}) \quad (4)$$

3.2 SPConv

We extend the work of KPConv and propose a new point convolution operator, SPConv. SPConv divides the local 3D space into a total of N shells. Each shell has one set of kernel points. Specifically, the innermost shell contains one central kernel point p_0 , for outer n^{th} ($n > 1$) shell, a set of kernel points $p_{1,m}$, $m \in M_n$ scatter on the surface of a sphere with radius r_n . Central kernel point impacts on a spherical region, and the n^{th} shell forms a ring-shaped influence area. As a result, all kernel points cover a spherical space of radius $(r_N + v)$. We perform kernel point convolution on N shells respectively, then stack shell features together along a new dimension. Finally, we use an additional convolution layer to further correlate shell structure. This convolution operator can be defined as follows:

$$(F(x) * W_1) * W_2 = \sigma(\sum_n^N \sigma(F(x) * W_1) W_{2,n}) \quad (5)$$

where σ refers to non-linear activation function. $W_1 \in \mathbb{R}^{K \times C_{in} \times C_{out}/2}$ is the learnable weight matrix for kernel point convolution, $W_2 \in \mathbb{R}^{N \times C_{out}/2 \times C_{out}}$ is

the weight matrix for shell correlation. To balance off efficiency and accuracy, we choose to use a total of 3 shells and 14 kernel points for the second and third shell respectively. An illustration of our SPConv is shown in Figure 1.

A detailed illustration of the influence area of SPConv kernel points is shown in Figure 2. The central kernel point encodes features from points that are spatially close to the query point. Kernel points located far from the center tend to encode more contextual information. Therefore, we learn the features by shells based on the distance from kernel point to center such that the encoded features can be representative for each shell. Furthermore, the 1D convolution layer applies weight matrix on the fused shell features, which enhances structure learning across shells. Our method aims to improve descriptive power of kernel points, so does KPConv [29] deformable version. Although deformable kernels provide more flexibility, regularization imposed on offsets is mandatory to account for mis-shifts. By contrast, our kernels are rigid and do not require regularization. In section 4.4, we compare our evaluation results with deformable KPConv [29]. Our method outperforms deformable KPConv with even less parameters.

3.3 Feature Attention

To further improve local feature encoding, we propose a feature attention module using low level features such as RGB or surface normal.

We propose two approaches for local feature attention. First approach is to apply a pre-defined Gaussian function:

$$\omega_k = \exp\left(-\frac{\|f(p) - f(s_k)\|}{2\sigma^2}\right) \quad (6)$$

where σ is a parameter that needs to be manually set. The second approach is to use sequential MLPs:

$$\omega_k = g(f(p) - f(s_k)), \quad (7)$$

where g is a sequence of MLPs activated by *ReLU*, except the last one which uses *sigmoid*. The final updated feature f' for point s_k can be calculated as follows with a residual connection:

$$f'(s_k) = \omega_k f(s_k) + f(s_k) \quad (8)$$

Both of the approaches improve performances and accelerate convergence speed. One issue for the first approach is that it is manually designed and not flexible. The second approach takes advantages of learnable weights and non-linear activations but adds more computation costs.

3.4 Network Architecture

We build a deep encoder-decoding network with point downsampling and upsampling to accomplish semantic segmentation task. A detailed SPNet architecture is shown in Figure 3(a).

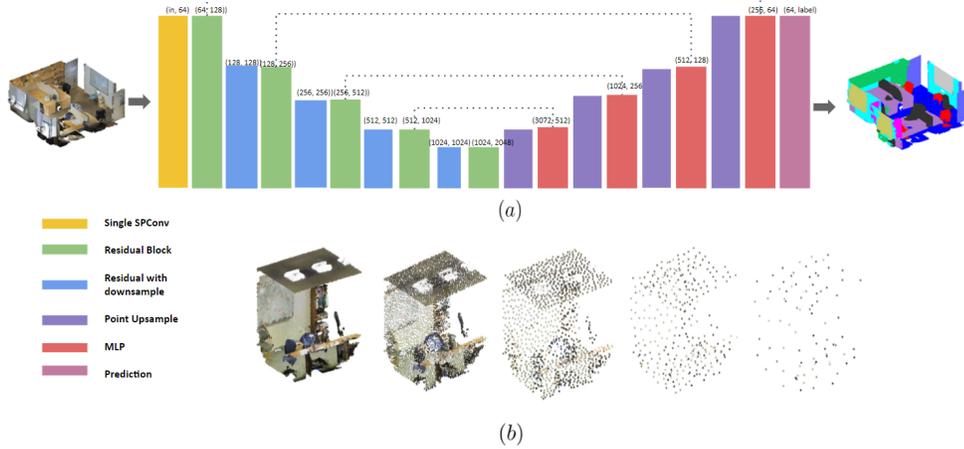


Fig. 3. (a) Illustration of the network architecture. (b) Downsampling process by PDS at each level.

Downsampling Strategy Similar to works [29,24], we adopt downsampling to reduce computation load as well as to increase receptive field. In our work, we favors Poisson disk sampling (PDS) strategy to deal with the varying density. PDS controls spatial uniformity by Poisson disk radius r_p , thus maintaining a minimal distance between points. Unlike grid sampling as used in [29] in which downsampling location are interpolated as the barycenter of a cell, PDS keeps the original locations of sub-sets and preserves shape patterns. Comparing to farthest point sampling (FPS) [24], PDS is faster when sampling large-scale points. A downsampling process by PDS is shown in Figure 3(b).

Upsampling Strategy With PDS, sampled points at each level are always a sub-set from input point sets. Therefore, we can accurately recover the point sampling patterns and gradually propagate features in decoder. We adopt feature propagation module proposed in [24]. For a point p_j at level j , its propagated features f are calculated as:

$$f = \sum_k^K w_k * f_k, w_k = \frac{d_k^2}{\sum_k^K d_k^2} \quad (9)$$

where d_k is the inverse Euclidean Distance between p_j and its k_{th} nearest neighbor at level $j - 1$.

4 Experiments

In this section, we evaluate the performance of our network on large-scale semantic segmentation datasets. We provide extensive ablation studies to justify

the effectiveness of our proposed methods. A comparison of scene segmentation results between existing methods and our is shown in Table 1.

4.1 Datasets

Stanford Large-Scale 3D Indoor Spaces (S3DIS) The S3DIS dataset [1] is a benchmark for large-scale indoor scene semantic segmentation. It consists of point clouds of six floors from three different buildings. Following the convention [23,18], We perform experiments on both 6-fold and Area 5 to evaluate our framework. For evaluation metrics, we use Overall point-wise accuracy (OA), and mean intersection over union (mIoU). The detailed results for individual class are listed in Table 2 and Table 3. We can see that our method has the highest scores for several challenging classes, such as door, wall and board.

Scannet The Scannet [3] dataset contains more than 1500 scanned scenes annotated with 20 valid semantic classes. It provides a 1,201/312 data split for training and testing. The Scannet dataset is reconstructed from RGB-D scanner. We report the per-voxel accuracy (OA) as evaluation metrics. As shown in Table 1, our framework achieves state-of-the-art performance.

4.2 Implementation Details

Parameter Setting SPNet uses residual block similar to [10]. Each block consists of one MLP for feature dimension reduction, one SPConv, and another MLP to increase feature dimension. SPNet consists of 5 encoding levels and 4 decoding levels, as shown in Figure 3. The kernel influence v_0 for the first encoding level is set to $0.04m$ for both S3DIS [1] and ScanNet [3]. For subsequent level l , kernel influence is increased to $v_l = 2^l v_0$. The rest of the parameters are adjusted according to v_l . For SPConv operator, we use a total of $K = 3$ shells. Kernel points of the second and the third shell are initialized on the surfaces of spheres with radius $r_2 = 1.5v_l, r_3 = 3v_l$ respectively. Query neighborhood radius R_l is set to $4v_l$, PDS radius is set to $0.75v$. For attention module, both color and normal are used to compute the attentional scores. For the Gaussian function, we set $\sigma = v_l$ at each level.

Network Training Our network is implemented using PyTorch [22] on a single Nvidia Titan RTX for all experiments. We use a batch size of 8, initial learning rate of 0.001. Optimization is done with Adam optimizer ($\beta_1 = 0.9, \beta_2 = 0.999$) [13]. Learning rate decays by a factor of 0.3 every 50 epoch for S3DIS [1], and every 30 epoch for ScanNet [3].

4.3 Ablation Studies

To prove the effectiveness of our proposed method, we conduct a series of experiments on S3DIS [1], evaluate on Area 5. Our baseline employs kernel point convolution with 15 kernel points. As shown in Table 4, each time we add or replace a module while keeping the rest unchanged. First, combining feature propagation(FP) with PDS produces a +2.3% boost. An explanation is that PDS

Table 1. Comparative 3D scene segmentation scores on S3DIS [1], ScanNet [3] datasets. S3DIS [1] scores are reported in metric of mean Intersection over Union(mIoU) including Area5 and 6-fold cross validation. ScanNet [3] scores are reported as Overall Accuracy(OA) and mIoU. The symbol ‘-’ means the results are not available.

Methods	S3DIS(mIoU)		ScanNet
	Area5	6-fold	(OA)
PointNet [23]	41.1	47.6	-
PointNet++ [24]	-	54.5	84.5
DGCNN [26]	-	56.1	-
SPGraph [15]	58.0	62.1	-
ShellNet [36]	-	66.8	85.2
PointWeb [37]	60.3	66.7	85.9
GACNet [30]	62.9	-	-
RandLA-Net [11]	-	68.5	-
SPH3D-GCN [17]	59.5	68.9	-
Point2Node [7]	63.0	70.0	86.3
KPConv(R) [29]	65.4	69.6	-
KPConv(D) [29]	67.1	70.6	-
Minkowski [2]	65.4	-	-
Ours	69.9	73.7	89.5

Table 2. Semantic segmentation mIoU and OA scores on S3DIS [1] Area 5.

Method	mIoU	OA	ceil.	floor	wall	beam	col.	wind.	door	chair	table	book.	sofa	board	clut.
PointNet [23]	41.1	49.0	88.8	97.3	69.8	0.1	3.9	46.3	10.8	52.6	58.9	40.3	5.9	26.4	33.2
PointWeb [37]	60.3	87.0	91.9	98.5	79.4	0.0	21.1	59.7	34.8	76.3	88.3	46.9	69.3	64.9	52.5
Point2Node [7]	62.9	88.8	93.8	98.3	83.3	0.0	35.6	55.3	58.8	79.5	84.7	44.1	71.1	58.7	55.2
KPConv(R) [29]	65.4	-	92.6	97.3	81.4	0.0	16.5	54.5	69.5	90.1	80.2	74.6	66.4	63.7	58.1
KPConv(D) [29]	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
Ours	69.9	90.3	94.5	98.3	84.0	0.0	24.0	59.7	79.8	89.6	81.0	75.2	82.4	80.4	60.4

Table 3. Semantic segmentation mIoU and OA scores on S3DIS [1] 6-fold.

Method	mIoU	OA	ceil.	floor	wall	beam	col.	wind.	door	chair	table	book.	sofa	board	clut.
PointNet [23]	47.8	78.5	88.0	88.7	69.3	42.4	23.1	47.5	51.6	54.1	42.0	9.6	38.2	29.4	35.2
SPGraph [15]	62.1	85.5	89.9	95.1	76.4	62.8	47.1	55.3	68.4	69.2	73.5	45.9	63.2	8.7	52.9
PointCNN [18]	65.4	88.1	94.8	97.3	75.8	63.3	51.7	58.4	57.2	69.1	71.6	61.2	39.1	52.2	58.6
PointWeb [37]	66.7	87.3	93.5	94.2	80.8	52.4	41.3	64.9	68.1	71.4	67.1	50.3	62.7	62.2	58.5
KPConv(R) [29]	69.6	-	93.7	92.0	82.5	62.5	49.5	65.7	77.3	57.8	64.0	68.8	71.7	60.1	59.6
Point2Node [7]	70.0	89.0	94.1	97.3	83.4	62.7	52.3	72.3	64.3	75.8	70.8	65.7	49.8	60.3	60.9
KPConv(D) [29]	70.6	-	93.6	92.4	83.1	63.9	54.3	66.1	76.6	57.8	64.0	69.3	74.9	61.3	60.3
Ours	73.7	90.9	94.6	97.3	85.0	45.2	56.9	82.1	63.4	73.1	83.4	71.5	68.8	68.6	67.8

preserves shape patterns in every downsampling level, FP correctly retrieves the shape patterns by interpolating features from neighborhood at upsampling level. Grid sampling loses geometric information, and this incorrectness accumulates through multiple downsampling layers. Moreover, adding local feature attention module produces +2.9% gain. Finally, SPConv improves the performance by +3.4%, showing its great capability of local feature encoding. Our full pipeline exceeds baseline with grid sample by +4.5%, which has the state-of-the-art performance on S3DIS dataset [1].

Table 4. Ablation studies evaluated on Area 5 of S3DIS [1].

	mIoU	Gain Δ
Baseline + grid sampling	65.4	-
Baseline + grid sampling + FP	65.4	-
Baseline + PDS	66.0	+0.6
Baseline + PDS + FP	67.7	+2.3
SPConv + PDS + FP	68.8	+3.4
Baseline + Attention + PDS + FP	68.3	+2.9
SPConv + Attention + PDS + FP	69.9	+4.5

To illustrate the effectiveness of our proposed attention module, see Table 5. Learnable MLP-based method achieves better performance however it burdens the computation.

Table 5. Ablation studies with the proposed feature attention.

Approach	Input features	mIoU	Inference speed (iter/s)
Gaussian function	color	68.0	4.3
Gaussian function	normal	67.1	4.3
Gaussian function	color+normal	68.3	4.3
2layer MLP	color	68.7	4.1
2layer MLP	normal	68.2	4.1
2layer MLP	color+normal	68.7	3.0
3layer MLP	color+normal	69.9	3.0

5 Conclusions

In this work, we propose an architecture named SPNet for 3D point cloud semantic segmentation. We introduce a SPConv operator to effectively learn point cloud geometry. We demonstrate that with Poisson disk sampling as well as feature propagation, our network can go deep without losing much inherent shape patterns. Our framework outperforms many competing approaches proved by experimental results on public large-scale datasets. We will experiment our method on outdoor Lidar datasets and investigate more effective attention methods.

References

1. Armeni, I., Sax, S., Zamir, A.R., Savarese, S.: Joint 2d-3d-semantic data for indoor scene understanding. arXiv preprint arXiv:1702.01105 (2017)
2. Choy, C., Gwak, J., Savarese, S.: 4d spatio-temporal convnets: Minkowski convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3075–3084 (2019)

3. Dai, A., Chang, A.X., Savva, M., Halber, M., Funkhouser, T., Nießner, M.: Scannet: Richly-annotated 3d reconstructions of indoor scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5828–5839 (2017)
4. Girshick, R.: Fast r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 1440–1448 (2015)
5. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 580–587 (2014)
6. Graham, B., Engelcke, M., van der Maaten, L.: 3d semantic segmentation with submanifold sparse convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 9224–9232 (2018)
7. Han, W., Wen, C., Wang, C., Li, X., Li, Q.: Point2node: Correlation learning of dynamic-node for point cloud feature modeling. arXiv preprint arXiv:1912.10775 (2019)
8. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
9. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE international conference on computer vision. pp. 1026–1034 (2015)
10. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)
11. Hu, Q., Yang, B., Xie, L., Rosa, S., Guo, Y., Wang, Z., Trigoni, N., Markham, A.: Randla-net: Efficient semantic segmentation of large-scale point clouds. arXiv preprint arXiv:1911.11236 (2019)
12. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4700–4708 (2017)
13. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
14. Landrieu, L., Boussaha, M.: Point cloud oversegmentation with graph-structured deep metric learning. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 7440–7449 (2019)
15. Landrieu, L., Simonovsky, M.: Large-scale point cloud semantic segmentation with superpoint graphs. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4558–4567 (2018)
16. Le, T., Duan, Y.: Pointgrid: A deep network for 3d shape understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 9204–9214 (2018)
17. Lei, H., Akhtar, N., Mian, A.: Spherical kernel for efficient graph convolution on 3d point clouds. arXiv preprint arXiv:1909.09287 (2019)
18. Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B.: Pointcnn: Convolution on x-transformed points. In: Advances in neural information processing systems. pp. 820–830 (2018)
19. Li, Y., Pirk, S., Su, H., Qi, C.R., Guibas, L.J.: Fpnn: Field probing neural networks for 3d data. In: Advances in Neural Information Processing Systems. pp. 307–315 (2016)
20. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)

21. Maturana, D., Scherer, S.: Voxnet: A 3d convolutional neural network for real-time object recognition. In: 2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 922–928. IEEE (2015)
22. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch (2017)
23. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)
24. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in neural information processing systems. pp. 5099–5108 (2017)
25. Shen, Y., Feng, C., Yang, Y., Tian, D.: Mining point cloud local structures by kernel correlation and graph pooling. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4548–4557 (2018)
26. Simonovsky, M., Komodakis, N.: Dynamic edge-conditioned filters in convolutional neural networks on graphs. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3693–3702 (2017)
27. Su, H., Maji, S., Kalogerakis, E., Learned-Miller, E.: Multi-view convolutional neural networks for 3d shape recognition. In: Proceedings of the IEEE international conference on computer vision. pp. 945–953 (2015)
28. Tchapmi, L., Choy, C., Armeni, I., Gwak, J., Savarese, S.: Segcloud: Semantic segmentation of 3d point clouds. In: 2017 international conference on 3D vision (3DV). pp. 537–547. IEEE (2017)
29. Thomas, H., Qi, C.R., Deschaud, J.E., Marcotegui, B., Goulette, F., Guibas, L.J.: Kpconv: Flexible and deformable convolution for point clouds. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 6411–6420 (2019)
30. Wang, L., Huang, Y., Hou, Y., Zhang, S., Shan, J.: Graph attention convolution for point cloud semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 10296–10305 (2019)
31. Wang, S., Suo, S., Ma, W.C., Pokrovsky, A., Urtasun, R.: Deep parametric continuous convolutional neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2589–2597 (2018)
32. Wu, W., Qi, Z., Fuxin, L.: Pointconv: Deep convolutional networks on 3d point clouds. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9621–9630 (2019)
33. Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1912–1920 (2015)
34. Xu, Y., Fan, T., Xu, M., Zeng, L., Qiao, Y.: Spidercnn: Deep learning on point sets with parameterized convolutional filters. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 87–102 (2018)
35. Yang, M., Yu, K., Zhang, C., Li, Z., Yang, K.: Densespp for semantic segmentation in street scenes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3684–3692 (2018)
36. Zhang, Z., Hua, B.S., Yeung, S.K.: Shellnet: Efficient point cloud convolutional neural networks using concentric shells statistics. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1607–1616 (2019)
37. Zhao, H., Jiang, L., Fu, C.W., Jia, J.: Pointweb: Enhancing local neighborhood features for point cloud processing. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5565–5573 (2019)