

# Multi Point-Voxel Convolution (MPVConv) for Deep Learning on Point Clouds

Wei Zhou<sup>1</sup>, Xin Cao<sup>1</sup>, Xiaodan Zhang<sup>1</sup>, Xingxing Hao<sup>1</sup>, Dekui Wang<sup>1</sup>, and Ying He<sup>2</sup>

<sup>1</sup>School of Information Science and Technology, Northwest University, Xi'an, China

<sup>2</sup>School of Computer Science and Engineering, Nanyang Technological University, Singapore

The existing 3D deep learning methods adopt either individual point-based features or local-neighboring voxel-based features, and demonstrate great potential for processing 3D data. However, the point based models are inefficient due to the unordered nature of point clouds and the voxel-based models suffer from large information loss. Motivated by the success of recent point-voxel representation, such as PVCNN, we propose a new convolutional neural network, called Multi Point-Voxel Convolution (MPVConv), for deep learning on point clouds. Integrating both the advantages of voxel and point-based methods, MPVConv can effectively increase the neighboring collection between point-based features and also promote independence among voxel-based features. Moreover, most of the existing approaches aim at solving one specific task, and only a few of them can handle a variety of tasks. Simply replacing the corresponding convolution module with MPVConv, we show that MPVConv can fit in different backbones to solve a wide range of 3D tasks. Extensive experiments on benchmark datasets such as ShapeNet Part, S3DIS and KITTI for various tasks show that MPVConv improves the accuracy of the backbone (PointNet) by up to 36%, and achieves higher accuracy than the voxel-based model with up to 34× speedups. In addition, MPVConv outperforms the state-of-the-art point-based models with up to 8× speedups. Notably, our MPVConv achieves better accuracy than the newest point-voxel-based model PVCNN (a model more efficient than PointNet) with lower latency.

*Index Terms*—MPVConv, 3D Deep Learning, Point Clouds, Point-Voxel, Semantic Segmentation.

## I. INTRODUCTION

3D deep learning for point clouds has received much attention in both industry and academia thanks to its potential for a wide range of applications, such as autonomous driving and robots. The main technical challenges are due to the sparse and irregular nature of point clouds.

The existing 3D deep learning methods can be roughly divided into voxel- and point-based methods according to the representations of point clouds. The voxel-based methods convert the irregular and sparse point clouds into regular 3D grids so that the widely studied convolutional neural networks (CNN) can be applied directly [1], [2], [3]. Since their performance heavily depends on the voxelization resolution, the voxel-based methods often suffer from large information loss when the resolution is low, as multiple adjacent points are quantized into the same grid, which are indistinguishable. Conversely, a high resolution volume would preserve the fine-detailed information, but requires significant amount of GPU memory and computation time due to the cubic complexity of volumes. In contrast, the point-based methods can handle high-resolution models, since they process the raw points in a local and separate manner [4], [5], [6], [7], [8]. Taking advantage of the sparse representation of point clouds, the point-based methods consume much less GPU memory than the voxel-based methods. However, due to lack of regularity, they suffer from expensive random memory access and dynamic kernel computation during the point and its nearest neighbor searching [9].

Motivated by the merits and limitations of each type of methods, several researchers proposed mixed representations to overcome the challenges of high accuracy demand and limited computational resources available on GPUs recently [9], [10], [11], [12], [13], [14]. However, most of these point-voxel combined methods are only for solving a specific task with a point-voxel based framework. For example, PV-RCNN [14] and PV-RCNN++ [11] focus on 3D object detection, and Pvdeconv [10] targets 3D auto-encoding CAD construction, FusionNet [12] is for semantic segmentation. Tang *et al.* proposed a sparse point-voxel convolution for efficient 3D architecture searching [13]. To our knowledge, there is no general point-voxel based method for solving different kind of tasks.

To design effective deep neural networks for 3D analysis, one must take into consideration the performance, efficiency, as well as generality and flexibility for various tasks. In this paper, we propose Multi Point-Voxel Convolution (MPVConv) which takes the advantages of both the voxel- and point-based methods, and can work with different backbones for a wide range of 3D tasks. In the previous 3D CNN models, the point-based features are individual and the voxel-based features are based on local neighborhood. In contrast, our MPVConv conducts 3D CNN and MLP on both points and voxels by multi Point-Voxel modules. As a result, it increases not only the neighboring collection for point-based features, but also the independence among voxel-based features. Extensive experiments show that MPVCNN outperforms the state-of-the-art 3D CNN models in terms of accuracy and efficiency.

## II. RELATED WORK

**Voxel-based 3D learning.** Inspired by the success of CNN on 2D images [15], [16], [17], researchers transfer point cloud

Corresponding author: Wei Zhou (email: mczhouwei12@gmail.com), Xin Cao (email: xin\_cao@163.com), Xiaodan Zhang (email: xiaodanzhang@nwu.edu.cn).

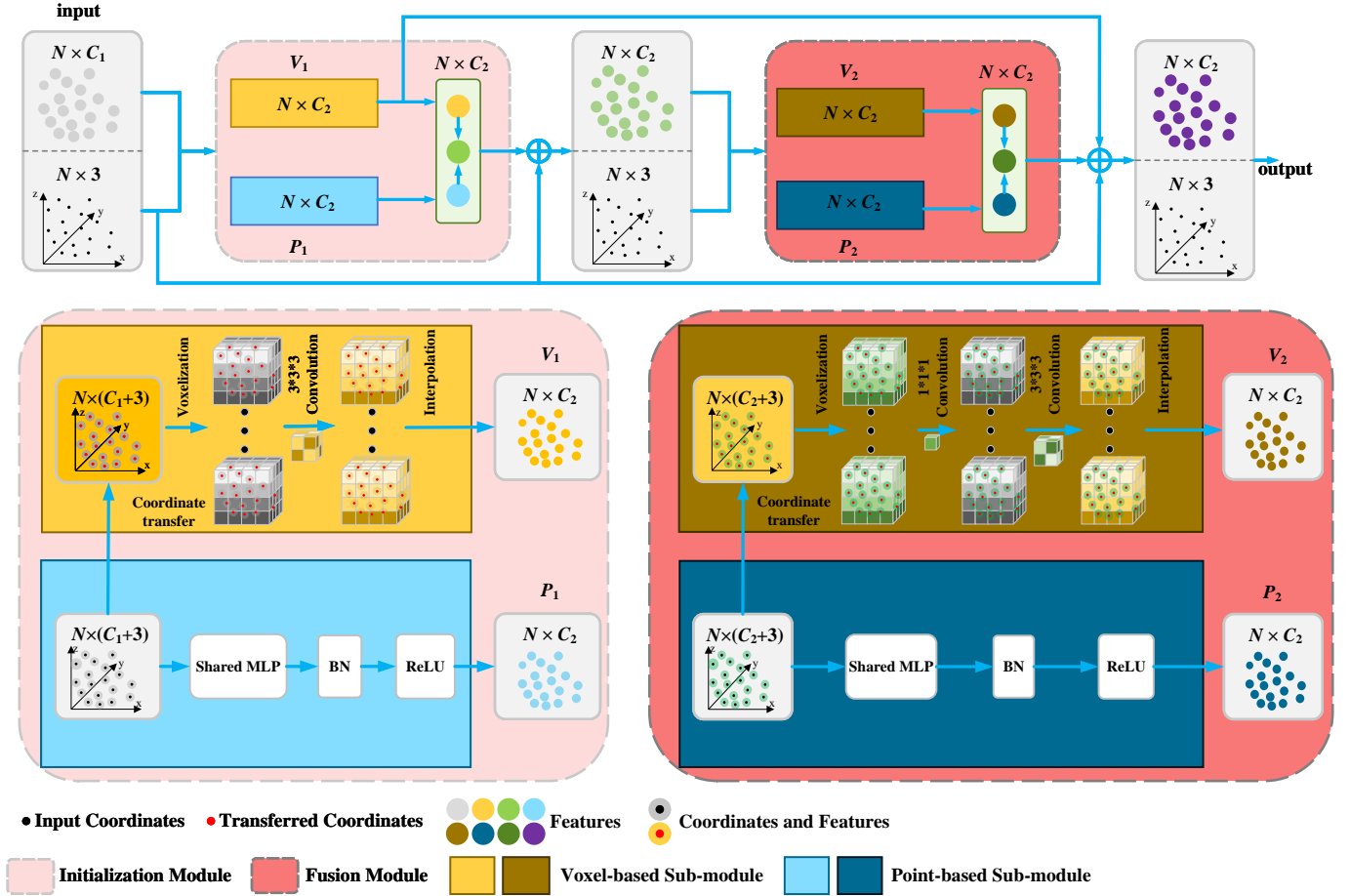


Fig. 1. MPVConv increases neighborhood collection in point-based features and independence in voxel-based features by adopting 3D CNNs and MLPs on both types of features.

representation to volumetric representation and attempt to adopt convolution over it [18], [1], [19], [20], [21]. As is known to all, the larger the voxel resolution is, the more detailed feature information is contained, and the calculation of 3D volumetric convolution increases exponentially. To slow down this problem, researchers adopt octree to construct efficient convolutional architectures to increase the efficiency of computation with high voxel resolutions [22], [2], [22]. State-of-the-art researches have demonstrated that volumetric representation can also be applied in 3D shape classification [23], [24], [20], 3D point cloud segmentation [25], [26], [27] and 3D object detection [3]. Although the voxel-based methods have a great advantage in data structuring for 3D CNN, its computation efficiency is still greatly limited by the sizes of voxel resolution.

**Point-based 3D learning.** PointNet [6], the first deep neural network that directly consumes 3D points, adopts simple symmetric function (maxpooling) to achieve permutation invariance. Since PointNet does not take local geometry into consideration, one typical way to improve it is to adopt hierarchical aggregation for extracting local features [7], [4]. PointCNN [5], SpiderCNN [28] and Geo-CNN [29] dynamically generate local geometric structures to capture points' neighboring features. RSNet [30] adopts a lightweight local dependency module to efficiently model local structures of point clouds. 3D-GCN [31],

DGCNN [8], Grid-GCN [32], SpecConv [33], SPGraph [34], GAC [35] adopt graph convolutional networks to conduct 3D point cloud learning, while RS-CNN [36], PCNN [37], SCN [38] and KCNet [39] make use of geometric relations for point cloud analysis. Furthermore, SPLATNet [40] uses sparse bilateral convolutional layers to build the network, and SO-Net [41] proposes permutation invariant architectures for learning with unordered point clouds. SSNet [42] combines Morton-order curve and point-wise to conduct self-supervised learning. SPNet [43] uses a self-prediction for 3D instance and semantic segmentation of point clouds. RandLA-Net [44] introduces a local feature aggregation module to preserve geometric details. As pointed out in [9], random memory access and dynamic kernel computation are the performance bottleneck of point-based methods.

**Point-voxel-based 3D learning.** Point-voxel-based methods combine the voxel- and point-based approaches. Aiming at 3D object detection, PV-RCNN [14] attracts the multi-scaled 3D voxel CNN features by PointNet++ [7], then the voxel- and point-based features are aggregated to a small set of representative points. Based on PV-RCNN, PV-RCNN++ [11] has improved the computation efficiency of point-based parts. Like PV-RCNN and PV-RCNN++ are specific to 3D object detection, other point-voxel-based 3D learning methods [10],

[12], [13] are specially designated only for 3D construction, 3D semantic segmentation or 3D structuring. PVCNN [9] adopts interpolation to obtain the voxel-based CNN features, and then aggregates the point-based features and the interpolated voxel-based features as output. Comparing with the existing point-voxel-based methods [9], [10], [11], [12], [13], [14], our method can deal with different kind of 3D tasks by applying MPVConv in various backbones with replacing the corresponding convolution module of backbones with MPVConv. Despite the resemblance of both combining voxel-based and point-based features, our work differs from [9], which integrates both the advantages of the voxel- and point-based methods both into the voxel and point-based features, as we can effectively increase the neighboring collection between point-based features and also promote the independence among voxel-based features. Moreover, we use shared MLPs to process the voxel-based features and adopt 3D CNN on the point-based features, thus improving the performances of various 3D deep learning tasks.

### III. MULTI POINT-VOXEL CONVOLUTION

State-of-the-art 3D deep learning methods are based on either voxel-based CNN methods or point-based network. Generally speaking, 3D voxel-based CNN method owns good data locality and regularity for low-resolution point neighbors (voxel grids), while the point-based approaches are independent for each point, so it could capture high-resolution information.

Inspired by point-voxel-based methods and convolution operation, we propose a new method, called Multi Point-Voxel Convolution (MPVConv), to conduct convolution on point clouds. Our method has two features that distinguishes itself from the existing point-voxel methods. First, the existing methods such as PV-RCNN [14], PV-RCNN++ [11], Pvdeconv [10] and FusionNet [12] are designed for a specific task, while our method can be applied to different tasks by replacing the corresponding convolution in backbones with MPVConv. Second, the previous methods (e.g. [14], [11], [9]) apply 3D volumetric convolution to voxels and MLP operation to points separately, and then fuse neighboring voxel features and independent point features together as output. In contrast, MPVConv applies 3D CNN and MLP to both points and voxels by multi point-voxel modules, yielding *independent* voxel features and local-neighboring point features.

As illustrated in Figure 1, MPVConv consists of an initialization module (left) and a fusion module (right). Initialization module generates initial individual point-based features and initial local-neighboring voxel-based features. Fusion module attaches neighboring and independent attributes to the output point- and voxel-based features of initialization module respectively. In each module, there are point-based and voxel-based sub-modules. The point-based sub-modules use shared MLPs to extract features for each point in an independent manner, hence they only consume a small amount of memory even for high-resolution point clouds. The voxel-based sub-modules aggregate features using volumetric CNNs and tri-linear interpolation to map voxel features back to points. Due to low voxel resolution adopted in MPVConv, the voxel sub-modules also have small memory consumption.

#### A. Initialization Module

The initialization module generates initial individual point-based features and initial local-neighboring voxel-based features from the input 3D data  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\} = \{(\mathbf{p}_1, \mathbf{f}_1), \dots, (\mathbf{p}_n, \mathbf{f}_n)\} \subseteq \mathbb{R}^{3+C_1}$ , where  $\mathbf{p}_i$  is the 3D point coordinates  $\mathbf{p}_i = (x_i, y_i, z_i)$ , it may also contain further information, e.g. RGB color and normal;  $\mathbf{f}_i$  is the output feature of previous layer which includes  $C_1$  channels.

##### 1) Voxel-based Sub-module

3D CNN on voxel grid is a popular selection for state-of-the-art 3D deep learning researches. Due to its high regularity and efficient structuring with 3D CNN, we adopt the voxel-based sub-module to capture the initial neighboring information for initialization module.

**Point Transformation.** Before we implement the voxel-based sub-module, we conduct point transformation to eliminate the influences from the translation and the scale variations of 3D point coordinates. Firstly, we calculate the mean point  $\bar{\mathbf{p}}$  of the input 3D data, and translate each point with  $\bar{\mathbf{p}}$  by  $\mathbf{p}_i = \mathbf{p}_i - \bar{\mathbf{p}}$ . Then we search the farthest point  $\|\mathbf{p}\|_{max}$  as the radius, and transform the whole points into the unit sphere:

$$\hat{\mathbf{p}}_i = \frac{\mathbf{p}_i}{2 * \|\mathbf{p}\|_{max}} + 0.5, \quad (1)$$

where  $\hat{\mathbf{p}}_i$  is the transformed point coordinates. During the experiments, we find that feature transformation reduces the accuracy of our model. Therefore, we transform only point coordinates, and mark the transformed points as  $\{\hat{\mathbf{p}}_i, \hat{\mathbf{f}}_i\}$ .

**3D Voxel CNN.** After the point transformation, the coordinates range from 0 to 1. To match the voxel resolution, we scale the coordinates by a factor  $\{r - 1\}$ , and mark the scaled coordinates as  $\hat{\mathbf{p}}_{ri}(\hat{x}_{ri}, \hat{y}_{ri}, \hat{z}_{ri})$ . Then the points  $\hat{\mathbf{p}}_{ri}$  and its corresponding features  $\hat{\mathbf{f}}_i$  are voxelized into the voxels with low spatial resolution of  $r \times r \times r$ , and the features of voxels are the mean features of all inside points:

$$\mathbf{F}(u, v, w) = \frac{\sum_{i=1}^n [(\hat{x}_{ri}, \hat{y}_{ri}, \hat{z}_{ri}) \cdot \hat{\mathbf{f}}_i]}{n}, \quad (2)$$

where  $\mathbf{F}(u, v, w)$  is the feature of voxel  $(u, v, w)$ ,  $n$  is the number of interior points inside the voxel, and  $[\cdot]$  is the floor function. We discuss the choice of voxelization resolution  $r$  in Section IV-E. Next, we adopt a series of  $3 \times 3 \times 3$  3D CNNs to aggregate the neighboring feature information of voxels. To capture the neighborhood information more accurately, we increase the feature channels to  $C_2$  in 3D CNN implementation. **Voxel Feature Interpolation.** Then we interpolate the voxel features into the common domain of point cloud, as we need to aggregate the information of voxel-based sub-module and point-based sub-module. The common operation for interpolation is the tri-linear interpolation and nearest-neighbor interpolation. Similar to the conduction in PVCNN, we adopt tri-linear interpolation to transform the voxel features into point features  $\mathbf{V}_1 = \{V_{11}, \dots, V_{1n}\} \subseteq \mathbb{R}^{C_2}$ .

##### 2) Point-based Sub-module

Although the voxel-based sub-module aggregates the neighboring feature information for the input 3D data, its extracted information is in a coarse manner with low voxel resolution. In order to make up for this defect, we adopt point-based

shared MLPs (1D convolution with kernel of 1) to make efficient learning on each point for the input 3D data. For the sake of convenient aggregation of voxel-based sub-module and point-based sub-module, the output feature channels of point-based sub-module are consistent with the output channels of voxel-based sub-module. The output features are implemented with batch normalization [45] and ReLU activation function [46], and these point-based features are marked as  $\mathbf{P}_1 = \{P_{11}, \dots, P_{1n}\} \subseteq \mathbb{R}^{C_2}$ .

When both the feature information of voxel-based and point-based sub-module are obtained, the initialization module will output  $\mathbf{V}_1$  and  $\mathbf{P}_1$  to the fusion module to strengthen the point-based and voxel-based features.

### B. Fusion Module

The fusion module is used to strengthen the point- and voxel-based features from the initialization module. The enhancement of feature information is reflected in two aspects: (1) attaching the neighboring collection for the individual point-based features  $\mathbf{P}_1$ ; (2) increasing the independence of the neighboring voxel-based features  $\mathbf{V}_1$ .

#### Neighboring collection for individual point-based feature.

Fusion module obtains the input features from the output aggregation feature information of initialization module. In terms of network structures, the point-based sub-module in fusion module is the same as the one in initialization module, the voxel-based sub-module got a little differences. The voxel-based and point-based information ( $\mathbf{P}_1, \mathbf{V}_1$ ) from the initialization module have aggregated together before inputting to the fusion module, but as we know, the voxel-based features  $\mathbf{V}_1$  have carried out information exchanging between different channels by 3D CNN, and yet the point-based features  $\mathbf{P}_1$  haven't exchanged. To enhance the neighboring collection for the point-based features  $\mathbf{P}_1$  from the initialization module with 3D CNN in the voxel-based sub-module of the current fusion module, we firstly fuse  $\mathbf{P}_1$  and  $\mathbf{V}_1$  to obtain the fused features  $\{\mathbf{V}_1 + \mathbf{P}_1\} \subseteq \mathbb{R}^{C_2}$ , then voxelize  $\{\mathbf{V}_1 + \mathbf{P}_1\}$  with low spatial resolution  $r \times r \times r$ . In order to further strengthen the neighboring collection for  $\mathbf{P}_1$  and exchange the collection between different channels of  $\mathbf{V}_1$ , we set a 3D CNN with kernel of  $1 \times 1 \times 1$  in the voxel-based sub-module of the current fusion module before the regular  $3 \times 3 \times 3$  3D CNN operations, thus we obtain the enhanced features  $\mathbf{V}_2 = \{V_{21}, \dots, V_{2n}\} \subseteq \mathbb{R}^{C_2}$ .

#### Independence for the neighboring voxel-based features.

To attach the independent attribute for voxel-based features  $\mathbf{V}_1$ , and increase the fine granularity of point-based features  $\mathbf{P}_1$  by the way, we adopt shared MLP to the fused features  $\{\mathbf{V}_1 + \mathbf{P}_1\} \subseteq \mathbb{R}^{C_2}$ . The parameters and network structure are consistent with the point-based sub-module of initialization module, and here we obtain the enhanced individual features  $\mathbf{P}_2 = \{P_{21}, \dots, P_{2n}\} \subseteq \mathbb{R}^{C_2}$ . After the completion of point-based sub-module and voxel-based sub-module in the fusion module, next is to output the aggregating information. To output the final features, in addition to the fusion module's aggregation information, the initialization module's voxel-based features are also aggregated to the output, and thus we output  $\{\mathbf{V}_1 + \mathbf{V}_2 + \mathbf{P}_2\} \subseteq \mathbb{R}^{C_2}$  as the final features. See Section IV-E for the discussion on output aggregating features.

## IV. EXPERIMENTAL RESULTS

This section presents the implementation details and compares our method with the state-of-the-art frameworks on benchmark datasets for various tasks, including ShapeNet Parts (object part segmentation) [47], S3DIS (indoor scene segmentation) [48], [49] and KITTI (3D object detection) [50]. We also conduct ablation study to study the performance of MPVConv.

### A. Implementation Details

**Initialization module.** We implemented MPVConv in Py-Torch. In the initialization module, the voxel-based sub-module contains two  $3 \times 3 \times 3$  3D CNNs with stride 1 and padding 1. Each 3D CNN is followed by 3D batch normalization [45] and the leaky ReLU activation function [51]. The point-based sub-module contains shared MLPs which is a 1D CNN with kernel of 1 which converts the feature channels to make it consistent with the output feature channels of voxel-based sub-module, next to the 1D CNN is 1D batch normalization [45] and ReLU activation function [46].

**Fusion module.** The fusion module follows the initialization module and takes its aggregation-fused feature as input. Its voxel-based sub-modules contain one  $1 \times 1 \times 1$  3D CNN with stride 1 and padding 0, and the rest part is the same as the one in the voxel-based sub-modules of the initialization module. The point-based sub-module also keeps the output feature channels consistent with its voxel-based sub-module.

### B. Part Segmentation

**Dataset.** We apply MPVConv to 3D part segmentation on ShapeNet Part [47], which consists of 16,881 3D shapes in 16 categories. We sample 2048 points from each shape as the training data. To make fair comparisons with others, we follow the same evaluation scheme as adopted in PVCNN [9], PointCNN [5], and SSCN [25] in our experiments.

**Architecture.** Using PointNet [6] as the backbone, we build MPVCNN by replacing its shared MLP layers with the proposed MPVConv layers.

**Training.** We use the cross-entropy loss function, set the batch size to 8, and adopt the ADAM optimizer [52] with learning rate of 0.001 for 200 epochs. We carried out the training process on a RTX 2080Ti GPU.

**Results.** We use the mean intersection-over-union (mIoU) as the evaluation metric. Same as the evaluation scheme in PointNet, we calculate the IoU of each shape by averaging the same shape parts' IoUs of the 2874 test models. We then compute the global mIoU by averaging the IoU of all shapes. We compare our model against the state-of-the-art point-based methods [4], [6], [41], [38], [40], [39], [30], [7], [8], [37], [28], [53], [31], [5], [54], [43], [55], [36], [56], voxel-based methods [1] and the recent point-voxel-based model [9]. To better balance time efficiency and accuracy, we also reduce the output feature channels to 50% and 25%, and call the downsized networks  $\text{MPVCNN}_{(0.5 \times Ch)}$  and  $\text{MPVCNN}_{(0.25 \times Ch)}$  respectively. Table I and Figure 2 show the results of MPVCNN on the ShapeNet part dataset. See also

the supplementary material for more visual results. MPVCNN outperforms the PointNet backbone with 3.3% increase of mIoU, even if we reduce the latency by 8.8%, we can still obtain 2.2% mIoU ahead of PointNet. Comparing with the voxel-based 3D-UNet and point-based SpiderCNN, our quarter version is  $34\times$  and  $8\times$  faster respectively, while achieving the better accuracy. In addition, MPVConv also outperforms other SOTA point-based methods (3D-GCN) with better accuracy by a large margin of 1.4%. Notably, even compared with the voxel-point-based PVCNN, we still gain higher mIoU by 0.3% with a small amount of time efficiency lost, and PVCNN is a state-of-the-art method aimed at speed improvement.

TABLE I  
EVALUATION RESULTS OF PART SEGMENTATION ON SHAPENET PART DATASET. P, V AND PV STAND FOR POINT-, VOXEL- AND POINT-VOXEL-BASED.

Method	Type	Batch Size	Input Size	mIoU (%)	GPU (GB)	Latency (ms)
Kd-Net [4]	P	8	4K	82.3	-	-
PointNet [6]	P	8	2K	83.7	<b>1.5</b>	21.7
3D-UNet [1]	V	8	96 <sup>3</sup>	84.6	8.8	682.1
SO-Net [41]	P	8	1K	84.6	-	-
SCN [38]	P	8	1K	84.6	-	-
SPLATNet [40]	P	-	-	84.6	-	-
KCNet [39]	P	8	2K	84.7	-	-
RSNet [30]	P	8	2K	84.9	0.8	74.6
PointNet++ [7]	P	8	2K	85.1	2.0	77.9
DGCNN [8]	P	8	2K	85.1	2.4	87.8
PCNN [37]	P	8	2K	85.1	-	-
SpiderCNN [28]	P	8	2K	85.3	6.5	170.7
GSNet [53]	P	-	-	85.3	-	-
3D-GCN [31]	P	8	2K	85.3	-	-
<b>MPVCNN</b> <sub>(0.25×Ch)</sub>	PV	8	2K	85.5	1.7	<b>19.8</b>
<b>MPVCNN</b> <sub>(0.5×Ch)</sub>	PV	8	2K	85.7	2.1	31.0
PointCNN [5]	P	8	2K	86.1	2.5	135.8
SPNet [43]	P	8	2K	86.2	-	-
CF-SIS [55]	P	8	2K	86.2	-	-
PVCNN [9]	PV	8	2K	86.2	1.6	50.7
RS-CNN [36]	P	8	2K	-	-	-
KPConv [54]	P	8	2K	86.4	-	-
<b>MPVCNN</b> <sub>(1×Ch)</sub>	PV	8	2K	<b>86.5</b>	2.8	81.9

### C. Indoor Scene Segmentation

**Dataset.** Stanford large-scale 3D Indoor Spaces (S3DIS) Dataset [48], [49] is a benchmark used for semantic indoor scene segmentation. S3DIS consists of 6 3D scanning indoor areas which totally includes 272 rooms. Similar to [9], [5], [57], we take Area-1, 2, 3, 4, 6 as the training set, and Area-5 (which is the only one that is not overlapping with others) as the testing set. We follow the data preprocessing and the evaluation criteria as in [5] before training and sample each block with 4096 points for training.

**Architecture.** In this task, we also use PointNet as the backbone for MPVCNN, and PointNet++ for MPVCNN++. Similar to the part segmentation experiments in Section IV-B, we also design a compact version by reducing the output feature channels to 12.5%, 25% and 50% in MPVCNN, and 50% in MPVCNN++.

**Training.** The batch size, optimizer, learning rate and criterion are set the same as what we have done in ShapeNet part experiments. The number of epoch is set to 50. It is also implemented on a single RTX 2080Ti GPU.

**Results.** Apart from mIoU, mean accuracy (mAcc) is also used to evaluate the performances of our proposed model. In addition to comparing with the state-of-the-art point-based [6], [8], [31], [30], [7], [58], [59], [42], [5], [32], [43], [55] and voxel-based methods [1], we also compare with the newest point-voxel-based model [9]. Table II shows the results of all methods on S3DIS dataset, and Figure 3 presents the visualization results, more visualization results of S3DIS are presented in the supplementary material. Specifically, the batch size of KPConv [54] is set to 1 since KPConv takes up too much memory to run. Our MPVCNN improves the mIoU of the backbone (PointNet) by more than **36%**, and MPVCNN++ outperforms its backbone (PointNet++) by a large margin in mIoU of more than **17%**. Notably, The compact MPVCNN of 25% feature channel outperforms the voxel-based 3D-UNet in accuracy with more than **11×** lower latency, and the point-based DGCNN by more than **15%** in mIoU with **3×** lower latency. Moreover, the full model of MPVCNN outperforms the state-of-the-art point-based model (Grid-GCN and PointCNN) with 2.4% increasing of the mIoU, and we increase the speed by nearly **4×** compared with PointCNN. At the same time, MPVCNN also outperforms the newest voxel-point-based method (PVCNN) both in mIoU by 4.5% and mAcc by 1.3%. Remarkably, the compact version of MPVCNN++ is faster than the extremely efficient voxel-point-based PVCNN++, and also outperforms it in accuracy by a large margin of 2%. The full MPVCNN++ can improve the performance more than **4%** in mIoU compared with PVCNN++. Overall, our MPVConv owns efficiency, small GPU memory requirements and good extensibility without sacrificing the performance of accuracy.

TABLE II  
EVALUATION RESULTS OF INDOOR SCENE SEGMENTATION ON S3DIS DATASET.

Method	Type	Batch Size	Input Size	mIoU (%)	mAcc (%)	GPU (GB)	Latency (ms)
PointNet [6]	P	8	4K	42.97	82.54	<b>0.6</b>	<b>20.9</b>
DGCNN [8]	P	8	4K	47.94	83.64	2.4	178.1
KPConv [54]	P	1	4K	48.58	-	10.7	-
3D-GCN [31]	P	8	4K	51.90	84.60	-	-
RSNet [30]	P	8	4K	51.93	-	1.1	111.5
PointNet++ [7]	P	8	4K	52.28	-	-	-
TanConv [58]	P	8	4K	52.8	85.5	-	-
3D-UNet [1]	V	8	96 <sup>3</sup>	54.93	86.12	6.8	574.7
JSNet [59]	P	8	4K	54.5	87.7	-	-
SSNet [42]	P	-	-	55.00	61.20	-	-
<b>MPVCNN</b> <sub>(0.25×Ch)</sub>	PV	8	4K	55.30	86.91	2.3	51.9
PVCNN [9]	PV	8	4K	56.12	86.66	1.3	47.3
PointCNN [5]	P	16	2K	57.26	85.91	4.6	282.3
Grid-GCN [32]	P	8	4K	57.75	86.94	-	25.9
<b>MPVCNN</b> <sub>(1×Ch)</sub>	PV	8	4K	58.63	87.75	4.3	72.7
SPNet [43]	P	8	4K	58.80	65.90	-	-
CF-SIS [55]	P	8	4K	58.90	67.30	-	-
PVCNN++ [9]	PV	4	8K	58.98	87.12	0.8	69.5
<b>MPVCNN++</b> <sub>(0.5×Ch)</sub>	PV	4	8K	60.17	88.76	2.7	52.6
<b>MPVCNN++</b> <sub>(1×Ch)</sub>	PV	4	8K	<b>61.51</b>	<b>89.31</b>	4.3	72.7

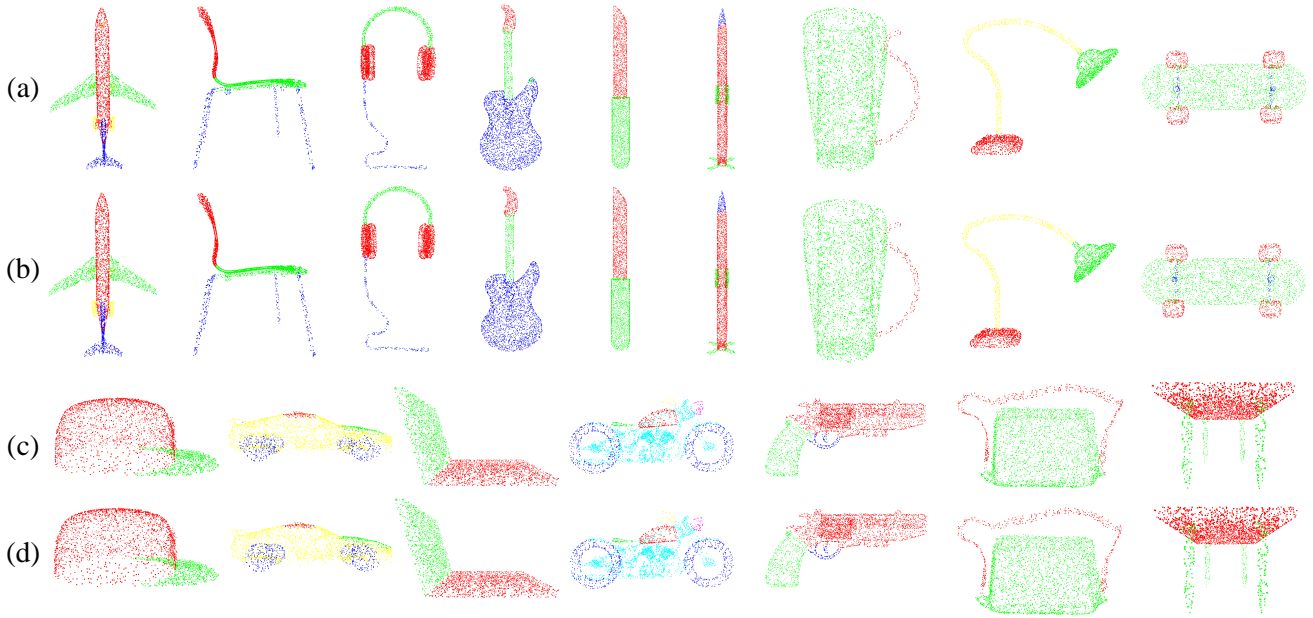


Fig. 2. Segmentation results on ShapeNet Part. (a)(c) Ground truth. (b)(d) Our results.

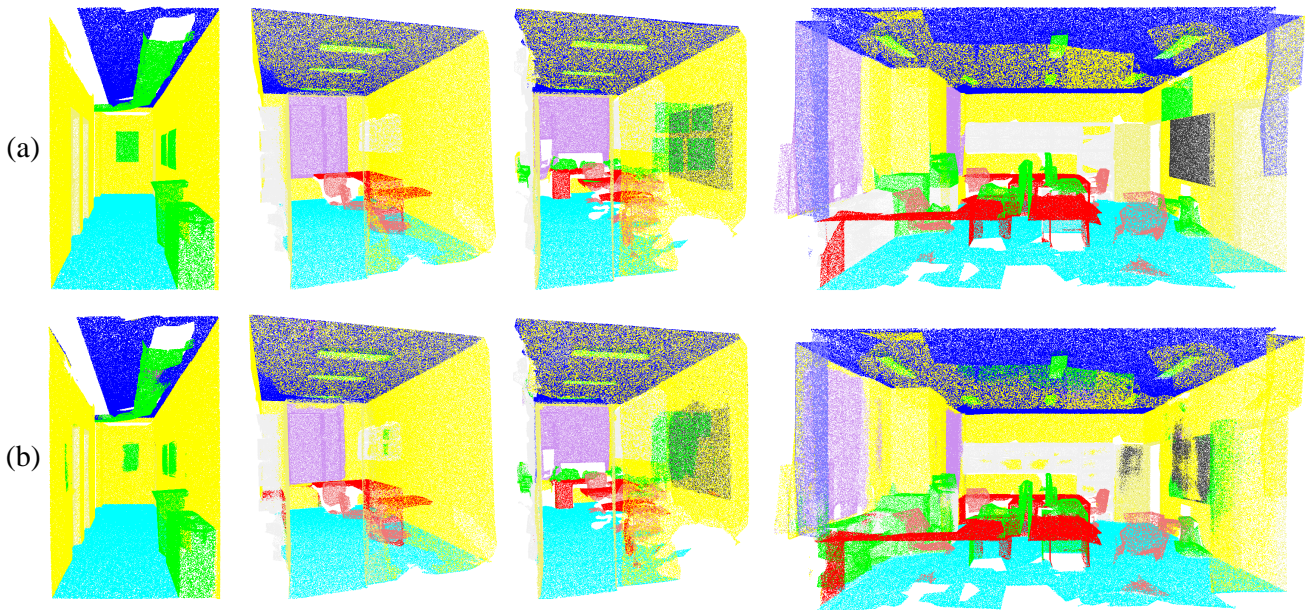


Fig. 3. Segmentation results on S3DIS Area-5. (a) Ground truth. (b) Our results.

#### D. 3D Detection

**Dataset.** We evaluate the performance of our model for 3D detection on KITTI [50], a benchmark for autonomous driving which contains 7,481 training and 7,518 test samples. To make fair comparison, we follow the experimental setup as Qi *et al.* [60]. We split the training samples to the *train* set with 3,712 samples and the *val* set with 3,769 samples.

**Architecture.** Without changing the whole network, we set F-PointNet [60] as the backbone by replacing the shared MLP layers of the instance segmentation network with our MPVConv to generate F-MPVCNN.

**Training.** The same setting for training as in part segmentation and indoor scene segmentation, except that epoch is set to 209 and batch size is set to 32.

**Results.** Mean average precision (mAP) is adopted to evaluate our model. We compare our model against F-PointNet, F-PointNet++ and F-PVCNN (whose backbone is also F-PointNet). Table III shows the experimental results on KITTI. Our F-MPVCNN outperforms all methods in all classes, and ours improves the mAP of F-PointNet (backbone) by up to **13.2%**. Compared with F-PointNet++, we improve the accuracy of it by a large margin up to 8.6% in the hard pedestrian class, and we are faster than it.

TABLE III  
KITTI

	Car			Pedestrian			Cyclist			Efficiency	
	Easy	Mod.	Hard	Easy	Mod.	Hard	Easy	Mod.	Hard	GPU Mem.	Latency
F-PointNet [60]	85.24	71.63	63.79	66.44	56.90	50.43	77.14	56.46	52.79	1.3GB	29.1ms
F-PointNet++ [60]	84.72	71.99	64.20	68.40	60.03	52.61	75.56	56.74	53.33	2.0GB	105.2ms
F-PVCNN [9]	85.25	72.12	64.24	70.60	61.24	56.25	78.10	57.45	53.65	1.4GB	58.9ms
F-MPVCNN	<b>85.66</b>	<b>72.63</b>	<b>64.62</b>	<b>71.12</b>	<b>62.34</b>	<b>57.13</b>	<b>79.85</b>	<b>58.28</b>	<b>54.62</b>	2.2GB	100.0ms

### E. Ablation Study

We conducted ablation experiments to understand the contributions of components to MPVConv. We trained the network and tested it on the ShapeNet part datasets [47], and considered half of the feature channels to output as it is a trade-off between accuracy and latency.

**Fusion module.** We tried to enlarge the voxel resolution and directly deepen 3D CNN network of voxel-based sub-module in initialization module before we think about constructing the fusion module for the whole MPVConv. As presented in Table IV (row 2), increasing the voxel resolution can slightly improve the performance of network in segmentation task, but at a high cost of GPU memory consumption and latency. Conversely, Table IV (row 3) shows that the accuracy of deepening network structure by 3D convolution is reduced on the contrary, and the GPU memory and latency increased also, but the extra computational cost can be neglected. So we tried another way to deepen the network, that is the fusion module. As shown in Table IV (row 4), the fusion module improves mIoU by 0.3% for MPVConv without consuming much more GPU memory and latency.

TABLE IV  
EFFECTS OF THE FUSION MODULE.

Model	mIoU	GPU Mem.	Latency
Init. module	85.50	1.97GB	22.3ms
Init. module (1.5×R)	85.55	2.35GB	37.1ms
Init. module (3×Conv3D)	85.33	2.08GB	25.1ms
Init. + Fusion modules	<b>85.76</b>	2.11GB	31.0ms

**Features.** Table V reveals the effects of various combinations of the four feature components of MPVConv. Model A, with only feature  $V_2$ , has performance downgrade, as the voxel-based feature  $V_2$  is not enough for 3D learning. The feature aggregations of D, F, G and H improve the performances of MPVConv, and the combination of  $V_1 + V_2 + P_2$  achieves the best performance.

**1×1×1 3D CNN.** As Table VI shows, 1×1×1 3D CNN can strengthen the information association and nonlinearity between different feature channels of point-based and voxel-based features, thus can improve the 3D learning ability of the network. With almost no increasing of latency and GPU memory, 1×1×1 3D CNN helps MPVConv to achieve better accuracy.

TABLE V  
EFFECTS OF DIFFERENT FEATURES FOR MPVCONV.

Model	$V_1$	$P_1$	$V_2$	$P_2$	mIoU
A			✓		85.39
B	✓	✓			85.50
C		✓	✓		85.43
D			✓	✓	85.58
E	✓	✓	✓		85.46
F		✓	✓	✓	85.54
G	✓		✓	✓	<b>85.76</b>
H	✓	✓	✓	✓	85.57

TABLE VI  
EFFECTS OF 1×1×1 3D CNN FOR MPVCONV.

Model	mIoU	GPU Mem.	Latency
MPVConv	85.72	2.024GB	31.2ms
MPVConv (1×1×1 CNN)	<b>85.76</b>	2.033GB	31.6ms

## V. CONCLUSION

We presented MPVConv, a 3D convolution neural network, for deep learning on point clouds. Combining both points and voxels, our method increases the neighboring collection between point-based features and promote the independence among voxel-based features with efficient convolutions. Experimental results on several benchmark datasets show that our method can improve the performance of common 3D tasks, such as segmentation and detection.

## REFERENCES

- [1] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3d u-net: learning dense volumetric segmentation from sparse annotation," in *MICCAI*, 2016.
- [2] G. Riegler, A. Osman Ulusoy, and A. Geiger, "Octnet: Learning deep 3d representations at high resolutions," in *CVPR*, 2017.
- [3] Y. Zhou and O. Tuzel, "Voxelnet: End-to-end learning for point cloud based 3d object detection," in *CVPR*, 2018.
- [4] R. Klokov and V. Lempitsky, "Escape from cells: Deep kd-networks for the recognition of 3d point cloud models," in *ICCV*, 2017.
- [5] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: Convolution on x-transformed points," in *NeurIPS*, 2018.
- [6] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *CVPR*, 2017.
- [7] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *NeurIPS*, 2017.
- [8] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," in *SIGGRAPH*, 2019.

- [9] Z. Liu, H. Tang, Y. Lin, and S. Han, "Point-voxel cnn for efficient 3d deep learning," in *NeurIPS*, 2019.
- [10] K. Cherenkova, D. Aouada, and G. Gusev, "Pvdeconv: Point-voxel deconvolution for autoencoding cad construction in 3d," in *ICIP*, 2020.
- [11] S. Shi, L. Jiang, J. Deng, Z. Wang, C. Guo, J. Shi, X. Wang, and H. Li, "Pv-rcnn+: Point-voxel feature set abstraction with local vector representation for 3d object detection," *arXiv preprint arXiv:2102.00463*, 2021.
- [12] F. Zhang, J. Fang, B. Wah, and P. Torr, "Deep fusionnet for point cloud semantic segmentation," in *ECCV*, 2020.
- [13] H. Tang, Z. Liu, S. Zhao, Y. Lin, J. Lin, H. Wang, and S. Han, "Searching efficient 3d architectures with sparse point-voxel convolution," in *ECCV*, 2020.
- [14] S. Shi, C. Guo, L. Jiang, Z. Wang, J. Shi, X. Wang, and H. Li, "Pv-rcnn: Point-voxel feature set abstraction for 3d object detection," in *CVPR*, 2020.
- [15] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *ICCV*, 2015.
- [16] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: towards real-time object detection with region proposal networks," *TPAMI*, 2016.
- [17] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *CVPR*, 2016.
- [18] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *IROS*, 2015.
- [19] 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 *CVPR*, 2016.
- [20] T. Le and Y. Duan, "Pointgrid: A deep network for 3d shape understanding," in *CVPR*, 2018.
- [21] D. Du, Z. Zhang, X. Han, S. Cui, and L. Liu, "Vipnet: A fast and accurate single-view volumetric reconstruction by learning sparse implicit point guidance," in *3DV*, 2020.
- [22] M. Tatarchenko, A. Dosovitskiy, and T. Brox, "Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs," in *ICCV*, 2017.
- [23] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *CVPR*, 2015.
- [24] 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," in *SIGGRAPH*, 2017.
- [25] B. Graham, M. Engelcke, and L. Van Der Maaten, "3d semantic segmentation with submanifold sparse convolutional networks," in *CVPR*, 2018.
- [26] H.-Y. Meng, L. Gao, Y.-K. Lai, and D. Manocha, "Vv-net: Voxel vae net with group convolutions for point cloud segmentation," in *ICCV*, 2019.
- [27] Z. Wang and F. Lu, "Voxsegnet: Volumetric cnns for semantic part segmentation of 3d shapes," *TVCG*, 2019.
- [28] Y. Xu, T. Fan, M. Xu, L. Zeng, and Y. Qiao, "Spiderncnn: Deep learning on point sets with parameterized convolutional filters," in *ECCV*, 2018.
- [29] S. Lan, R. Yu, G. Yu, and L. S. Davis, "Modeling local geometric structure of 3d point clouds using geo-cnn," in *CVPR*, 2019.
- [30] Q. Huang, W. Wang, and U. Neumann, "Recurrent slice networks for 3d segmentation of point clouds," in *CVPR*, 2018.
- [31] Z.-H. Lin, S. Y. Huang, and Y.-C. F. Wang, "Learning of 3d graph convolution networks for point cloud analysis," *TPAMI*, 2021.
- [32] Q. Xu, X. Sun, C.-Y. Wu, P. Wang, and U. Neumann, "Grid-gcn for fast and scalable point cloud learning," in *CVPR*, 2020.
- [33] C. Wang, B. Samari, and K. Siddiqi, "Local spectral graph convolution for point set feature learning," in *ECCV*, 2018.
- [34] L. Landrieu and M. Simonovsky, "Large-scale point cloud semantic segmentation with superpoint graphs," in *CVPR*, 2018.
- [35] L. Wang, Y. Huang, Y. Hou, S. Zhang, and J. Shan, "Graph attention convolution for point cloud semantic segmentation," in *CVPR*, 2019.
- [36] Y. Liu, B. Fan, S. Xiang, and C. Pan, "Relation-shape convolutional neural network for point cloud analysis," in *CVPR*, 2019.
- [37] M. Atzmon, H. Maron, and Y. Lipman, "Point convolutional neural networks by extension operators," in *SIGGRAPH*, 2018.
- [38] S. Xie, S. Liu, Z. Chen, and Z. Tu, "Attentional shapecontextnet for point cloud recognition," in *CVPR*, 2018.
- [39] Y. Shen, C. Feng, Y. Yang, and D. Tian, "Mining point cloud local structures by kernel correlation and graph pooling," in *CVPR*, 2018.
- [40] 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 *CVPR*, 2018.
- [41] J. Li, B. M. Chen, and G. H. Lee, "So-net: Self-organizing network for point cloud analysis," in *CVPR*, 2018, pp. 9397–9406.
- [42] A. Thabet, H. Alwassel, and B. Ghanem, "Self-supervised learning of local features in 3d point clouds," in *CVPR*, 2020.
- [43] J. Liu, M. Yu, B. Ni, and Y. Chen, "Self-prediction for joint instance and semantic segmentation of point clouds," in *ECCV*, 2020.
- [44] Q. Hu, B. Yang, L. Xie, S. Rosa, Y. Guo, Z. Wang, N. Trigoni, and A. Markham, "Randla-net: Efficient semantic segmentation of large-scale point clouds," in *CVPR*, 2020.
- [45] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015.
- [46] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics. JMLR Workshop and Conference Proceedings*, 2011, pp. 315–323.
- [47] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su *et al.*, "Shapenet: An information-rich 3d model repository," *arXiv preprint arXiv:1512.03012*, 2015.
- [48] 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 *CVPR*, 2016.
- [49] I. Armeni, S. Sax, A. R. Zamir, and S. Savarese, "Joint 2d-3d-semantic data for indoor scene understanding," *arXiv preprint arXiv:1702.01105*, 2017.
- [50] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [51] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *ICML*, 2013.
- [52] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [53] M. Xu, Z. Zhou, and Y. Qiao, "Geometry sharing network for 3d point cloud classification and segmentation," in *AAAI*, 2020.
- [54] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, and L. J. Guibas, "Kpconv: Flexible and deformable convolution for point clouds," in *ICCV*, 2019.
- [55] X. Wen, Z. Han, G. Youk, and Y.-S. Liu, "Cf-sis: Semantic-instance segmentation of 3d point clouds by context fusion with self-attention," in *ACM MM*, 2020.
- [56] K. Liu, Z. Gao, F. Lin, and B. M. Chen, "Fg-net: Fast large-scale lidar point clouds understanding network leveraging correlated feature mining and geometric-aware modelling," *arXiv preprint arXiv:2012.09439*, 2020.
- [57] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, and S. Savarese, "Segcloud: Semantic segmentation of 3d point clouds," in *3DV*, 2017.
- [58] M. Tatarchenko, J. Park, V. Koltun, and Q.-Y. Zhou, "Tangent convolutions for dense prediction in 3d," in *CVPR*, 2018.
- [59] L. Zhao and W. Tao, "Jsnet: Joint instance and semantic segmentation of 3d point clouds," in *AAAI*, 2020.
- [60] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas, "Frustum pointnets for 3d object detection from rgb-d data," in *CVPR*, 2018.



**Wei Zhou** received his B.E. degrees in Electronic and Information Engineering from Xidian University in 2013. From 2013 to 2018 he is doing 5-year-system Ph.D. candidate in Xi'an Institute of Optics and Precision Mechanics of CAS & University of Chinese Academy of Sciences. From 2016 to 2017, he was visiting Fraunhofer IGD & TU Darmstadt. Since 2019, he has been a lecturer with School of Information Science and Technology, Northwest University, Xi'an, China. His main research interests include image restoration, 3D semantic segmentation,

3D recognition, 3D feature matching and other point-cloud-based and image-based researches.





**Xin Cao** received the B.S. degrees in electronic engineering from Xidian University, Xi'an, China, in 2011 and the Ph.D. degree in pattern recognition and intelligent system from Xidian University, Xi'an, China, in 2016. Since 2016, he has been a lecturer with School of Information Science and Technology, Northwest University, Xi'an, China. His research interests include medical image analysis and optical molecular imaging.



**Xiaodan Zhang** received the B.Eng. degree in electronic information engineering and Ph.D. degree in intelligence information processing from Xidian University, Xi'an, China, in 2014 and 2019. She is currently a faculty member with Northwest University, Xi'an, China. Her current research interests include image aesthetic quality assessment, visual attention, and computationally modeling of human visual system.



**Xingxing Hao** received the B.S. degree in intelligent science and technology from Xidian University, Xi'an, China in 2014. Now, he is pursuing the Ph.D. degree in circuits and systems from the School of Artificial Intelligence, Xidian University. His research interests include image processing, combinatorial optimization and evolutionary algorithms.



**Dekui Wang** received the Ph.D. and B.S. degree in computer science and technology from Xidian University, Xi'an, China. He is working at School of Information Science and Technology, Northwest University, Xi'an, China. His current research interests include very large-scale integrated physical design automation with emphasis on synthesis, placement, routing, timing analysis and image processing.



**Ying He** received the BS and MS degrees in electrical engineering from Tsinghua University, China, and the PhD degree in computer science from Stony Brook University. He is currently an Associate Professor with School of Computer Science and Engineering, Nanyang Technological University, Singapore. His research interests fall into the general areas of visual computing and he is particularly interested in the problems which require geometric analysis and computation.

## APPENDIX A

In the supplementary materials, we show more comparable visualization results in part segmentation (Figure 4, Figure 5, Figure 6 and Figure 7) and indoor scene segmentation (Figure 8, Figure 9, Figure 10, Figure 11 and Figure 12) tasks with PointNet [6], PointNet++ [7], PVCNN [9] and our MPVCNN. Besides, we will make our code and models publicly available.

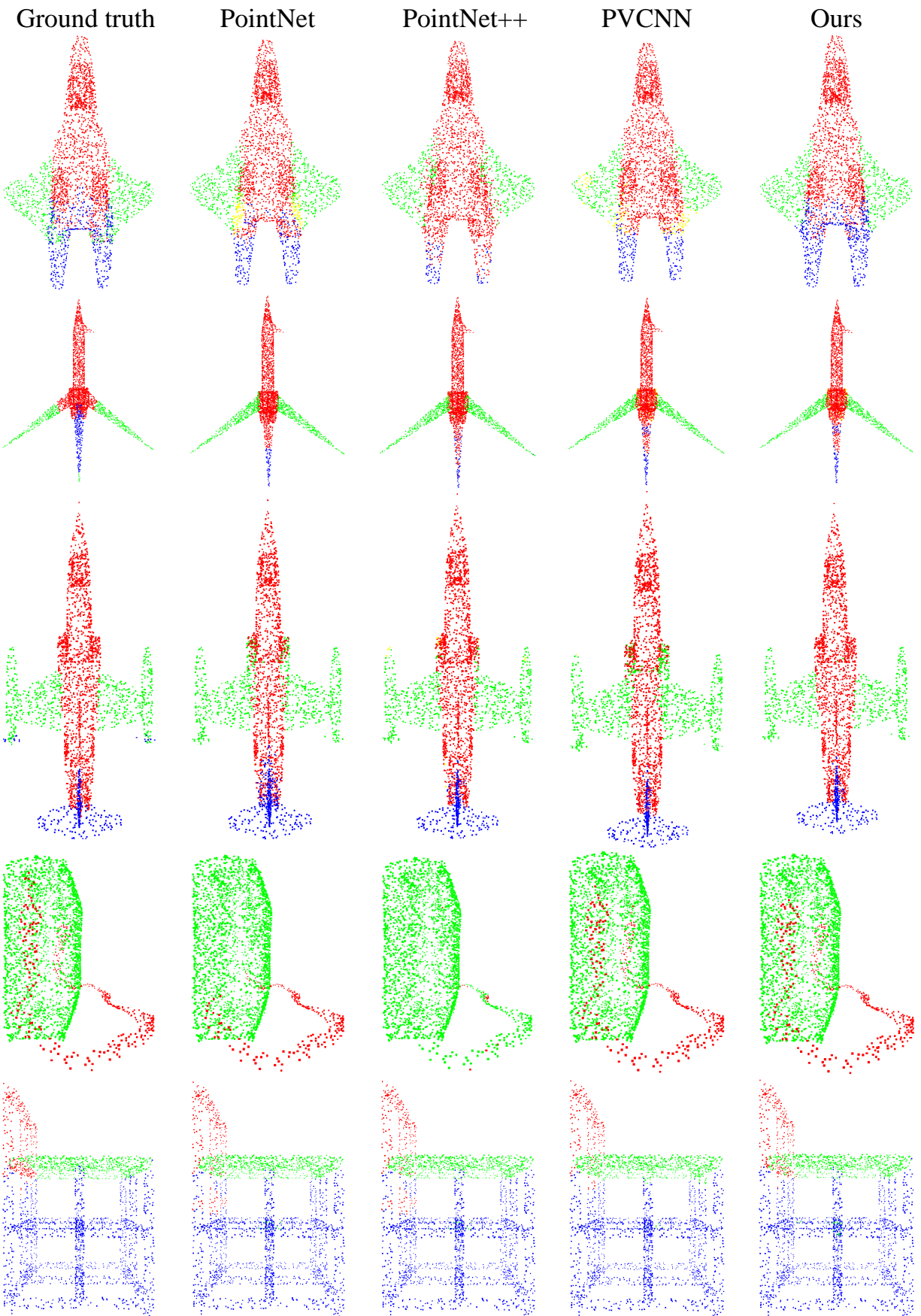


Fig. 4. Part segmentation results of PointNet [6], PointNet [7] PVCNN [9] and our MPVCNN on Shapenet Part.

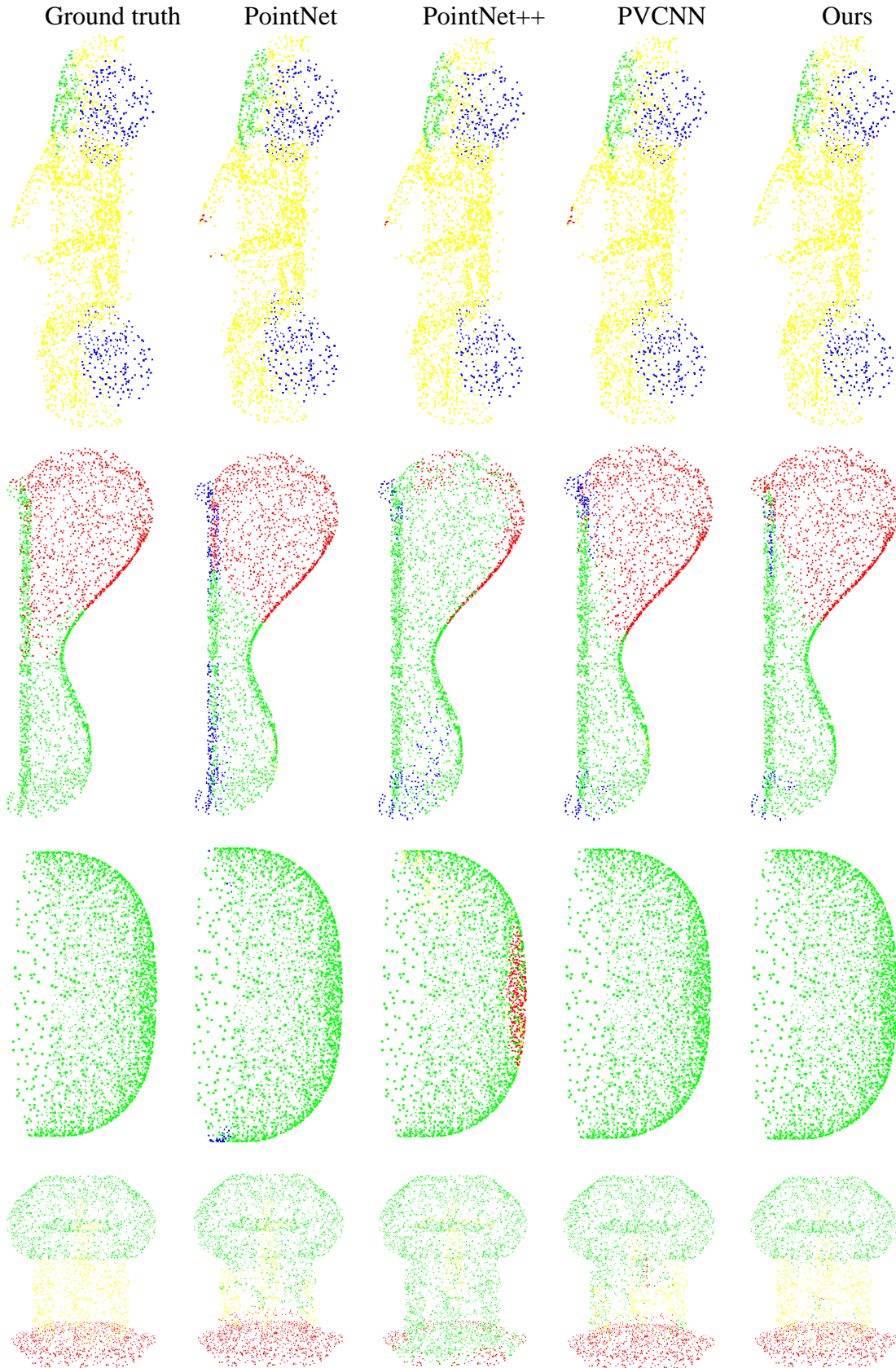


Fig. 5. Part segmentation results of PointNet [6], PointNet [7] PVCNN [9] and our MPVCNN on Shapenet Part.

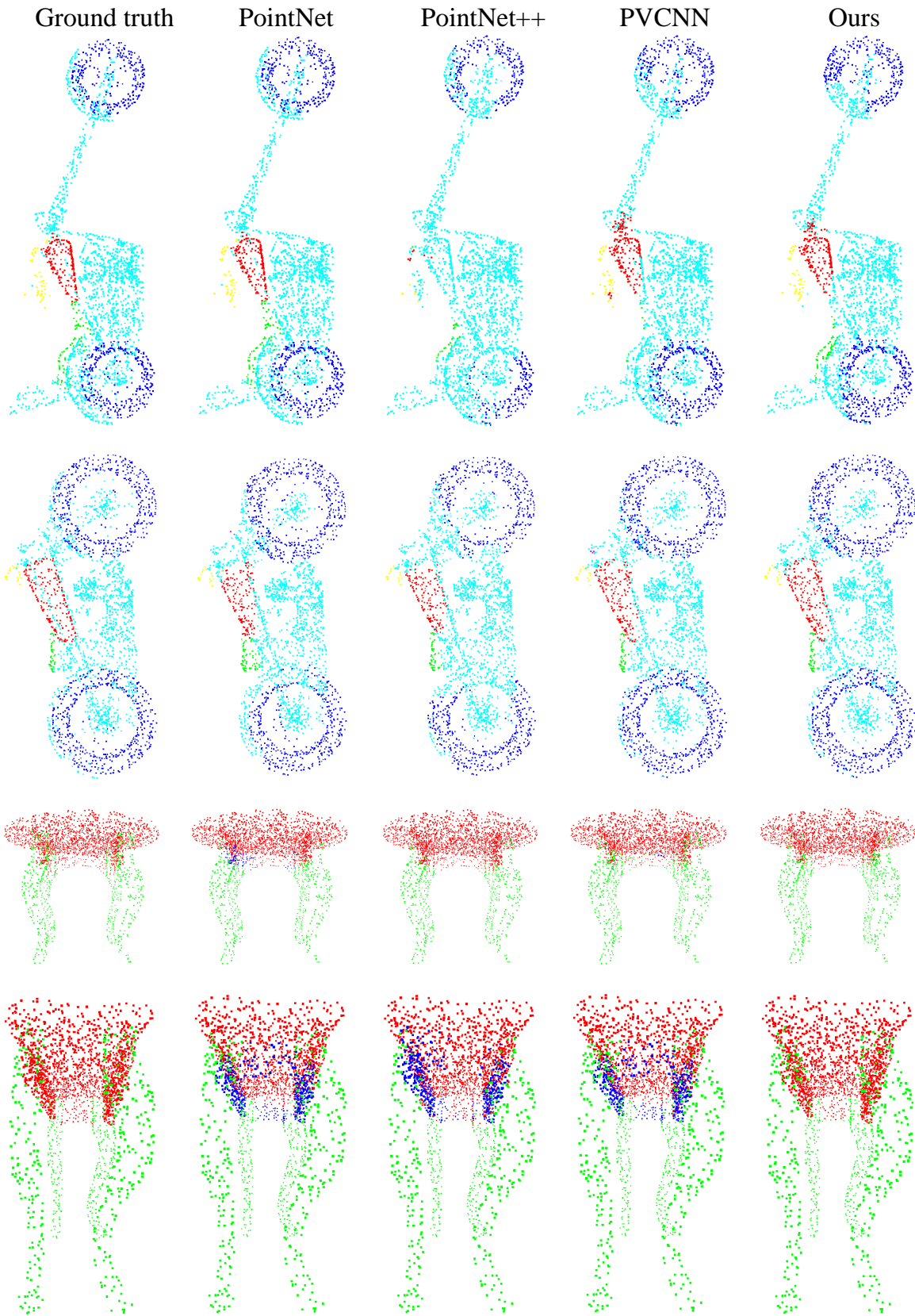


Fig. 6. Part segmentation results of PointNet [6], PointNet [7] PVCNN [9] and our MPVCNN on Shapenet Part.

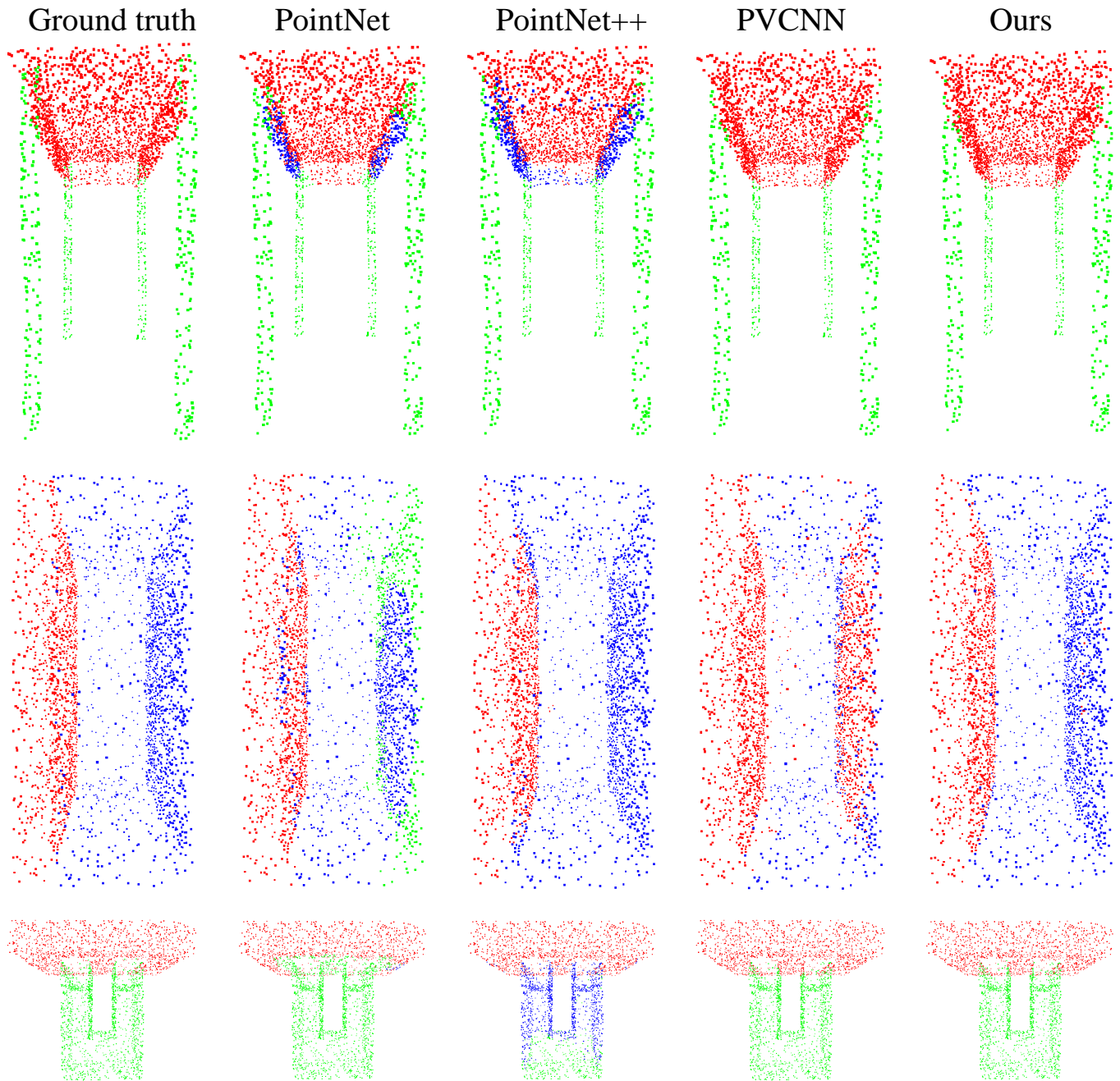


Fig. 7. Part segmentation results of PointNet [6], PointNet [7] PVCNN [9] and our MPVCNN on Shapenet Part.



Fig. 8. Indoor scene segmentation results of PointNet [6], PVCNN [9] and our MPVCNN on S3DIS are5.

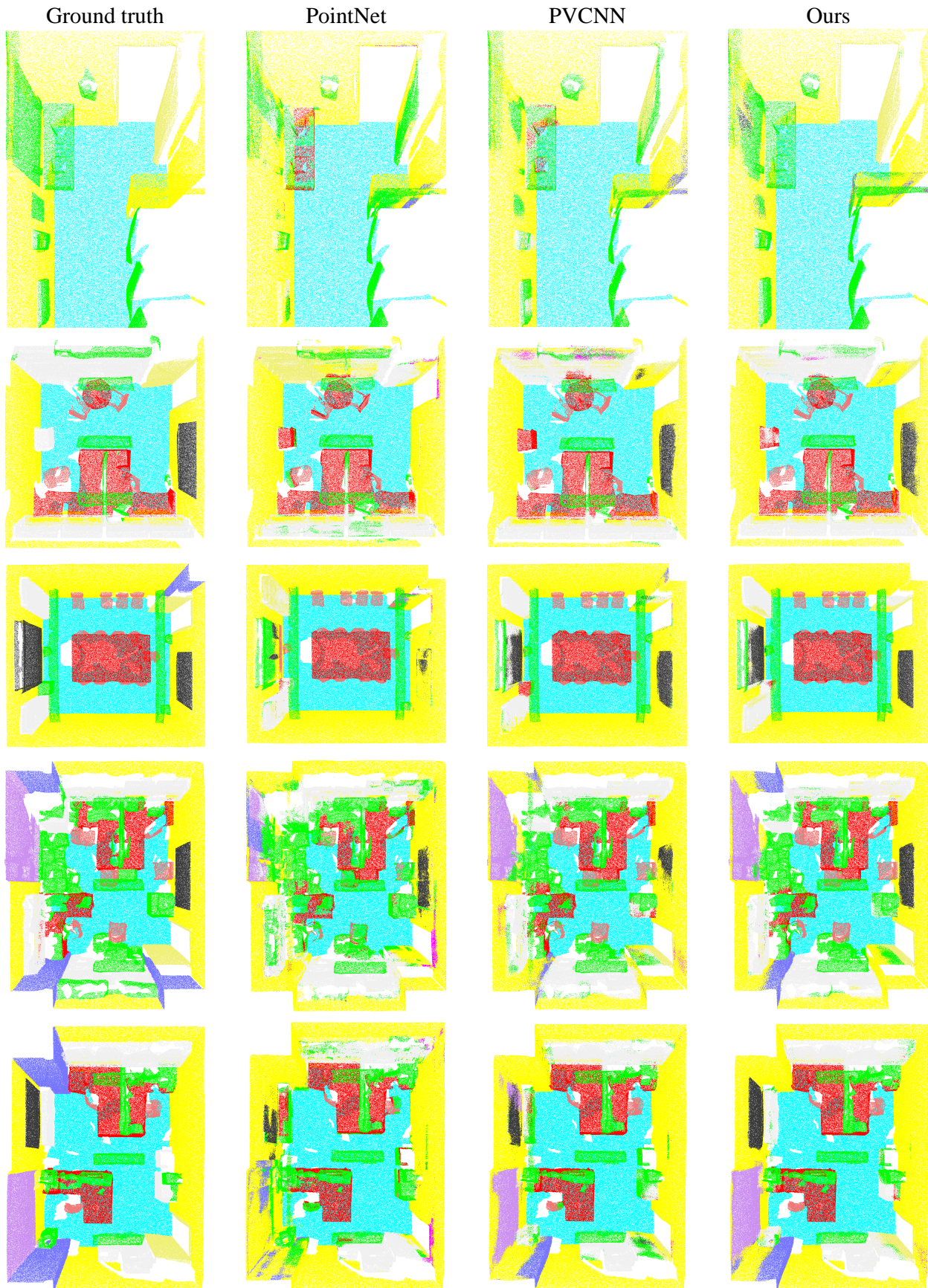


Fig. 9. Indoor scene segmentation results of PointNet [6], PVCNN [9] and our MPVCNN on S3DIS are5.



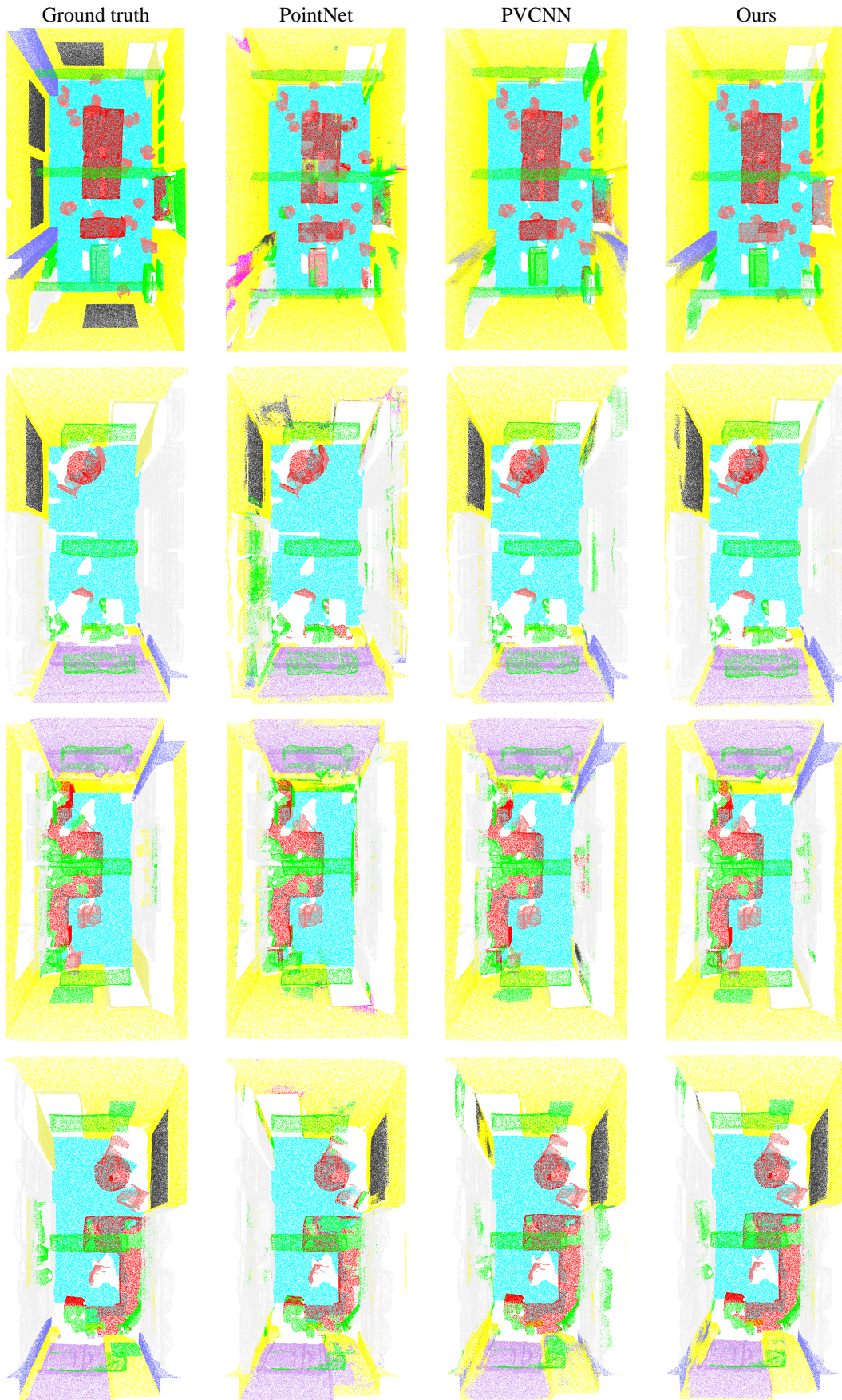


Fig. 10. Indoor scene segmentation results of PointNet [6], PVCNN [9] and our MPVCNN on S3DIS are5.



Fig. 11. Indoor scene segmentation results of PointNet [6], PVCNN [9] and our MPVCNN on S3DIS are5.

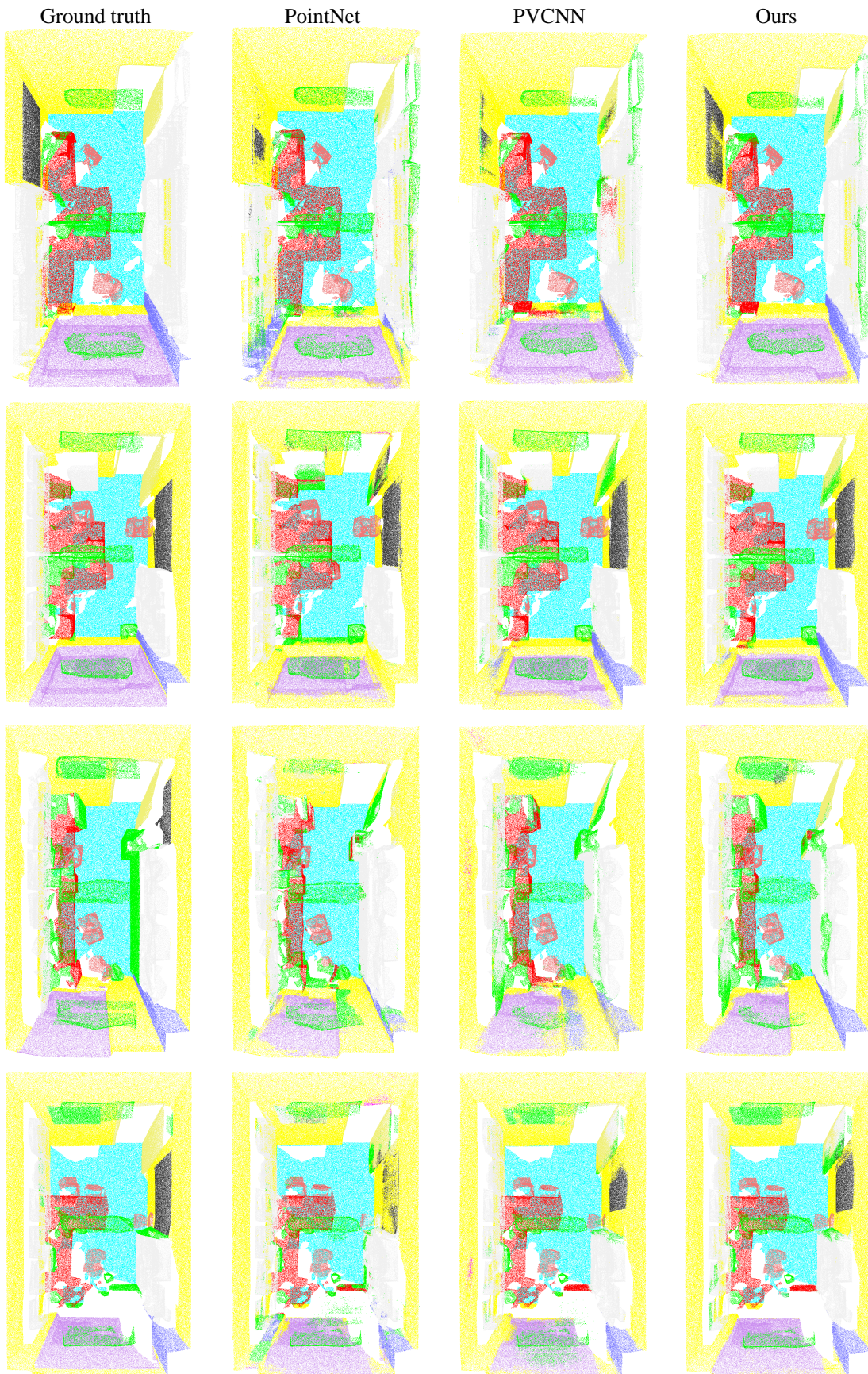


Fig. 12. Indoor scene segmentation results of PointNet [6], PVCNN [9] and our MPVCNN on S3DIS are5.