

# PointVector: A Vector Representation In Point Cloud Analysis

Xin Deng\* WenYu Zhang\* Qing Ding† XinMing Zhang†  
University of Science and Technology of China

{xin\_deng, wenyuz}@mail.ustc.edu.cn, {dingqing, xinming}@ustc.edu.cn

## Abstract

In point cloud analysis, point-based methods have rapidly developed in recent years. These methods have recently focused on concise MLP structures, such as PointNeXt, which have demonstrated competitiveness with Convolutional and Transformer structures. However, standard MLPs are limited in their ability to extract local features effectively. To address this limitation, we propose a Vector-oriented Point Set Abstraction that can aggregate neighboring features through higher-dimensional vectors. To facilitate network optimization, we construct a transformation from scalar to vector using independent angles based on 3D vector rotations. Finally, we develop a PointVector model that follows the structure of PointNeXt. Our experimental results demonstrate that PointVector achieves state-of-the-art performance **72.3% mIOU** on the S3DIS Area 5 and **78.4% mIOU** on the S3DIS (6-fold cross-validation) with only **58%** model parameters of PointNeXt. We hope our work will help the exploration of concise and effective feature representations. The code will be released soon.

## 1. Introduction

Point cloud analysis is a cornerstone of various downstream tasks. With the introduction of PointNet [25] and PointNet++ [26], the direct processing of unstructured point clouds has become a hot topic. Many point-based networks introduced novel and sophisticated modules to extract local features, e.g., attention-based methods [53] explore attention mechanisms as Fig. 1a with lower consumption, convolution-based methods [36] explore the dynamic convolution kernel as Fig. 1c, and graph-based methods [39] [54] use graph to model relationships of points. The application of these methods to the feature extraction module of PointNet++ brings an improvement in feature quality. However, they are somewhat complicated to design in terms of network structure. PointNeXt [28] adapts the SetAbstrac-

tion (SA) module of PointNet++ [26] and proposes the Inverted Residual MLP (InvResMLP) module. The simple design of MLP network achieves good results. Motivated by this work, we try to further explore the potential of the MLP structure.

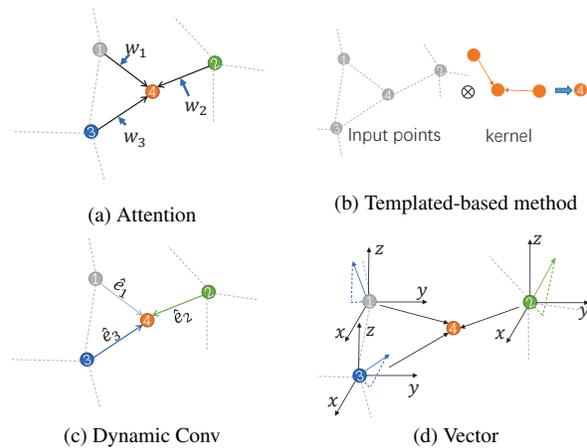


Figure 1. Illustrations of the core operations of the different methods. (a) The features of each point are calculated separately by applying a fixed/isotropic kernel (black arrow) like Linear layer. Then, it imparts anisotropy by weights generated from inputs. (b) The displacement vector is used to filter points that approximate the kernel pattern for features aggregation. (c) It applies unique dynamic kernels with anisotropy for each point feature. (d) Differently, we generate vector representations based on features, and the aggregation methods for vectors are anisotropic due to the direction of the vectors.

PointNeXt uses all standard MLPs, which has insufficient feature extraction capability. In addition to attention and dynamic convolution mechanisms, template-based methods as Fig. 1b such as 3D-GCN [19] employ relative displacement vectors to modulate the association between input points and the convolutional kernel. We introduce a vector representation of features to extend the range of feature variation with the intention of more effectively regulating the connections between local features. Our approach as Fig. 1d differs from template-based methods. Instead of using displacement vectors as a property of the kernel, we

\*Co-first authors with equal contribution to refining the theory and experimental design

†Corresponding authors

generate a vector representation for each neighboring point and aggregate them. Our method introduces less inductive bias, resulting in improved generalization capabilities. Furthermore, we enhance the generation of 3D vector representations by utilizing a vector rotation matrix with two independent angles in 3D space. This method facilitates the network to find the better solution.

Influenced by PointNeXt [28] and PointNet++ [26], we present the VPSA module. This module adheres to the structure of Point Set Abstraction (SA) module of the PointNet series. Vector representations are obtained from input features and aggregated using a reduction function. The vector of each channel is then projected into a scalar to derive local features. By combining VPSA and SA modules, we construct a PointVector model with an architecture akin to that of PointNeXt.

Our model undergoes comprehensive validation on public benchmark datasets. It achieves state-of-the-art performance on the S3DIS [1] semantic segmentation benchmark and competitive results on the ScanObjectNN [48] and ShapeNetPart [49] datasets. By incorporating a priori knowledge of vectors, our model attains superior results with fewer parameters on S3DIS. Detailed ablation experiments further demonstrate the efficacy of our methodology. The contributions are summarized below:

- We propose a novel immediate vector representation with relative features and positions to better guide local feature aggregation.
- We explore the method of obtaining vector representation and propose the generation method of 3D vector by utilizing the vector rotation matrix in 3D space.
- Our proposed PointVector model achieves **72.3%** mean Intersection over Union (mIOU) on S3DIS area5 and **78.4%** mIOU on S3DIS (6-fold cross-validation) with only **58%** model parameters of PointNeXt.

## 2. Related work

**Point-based network.** In contrast to the voxelization [55] [15] [31] and multiview [32] [10] [41] methods, point-based methods deal directly with point clouds. PointNet first proposes using MLP to process point clouds directly. PointNet++ subsequently introduces a hierarchical structure to improve the feature extraction. Subsequent works focused on the design of fine-grained local feature extractors. Graph-based methods [39] [38] rely on a graph neural network and introduce point features and edge features to model local relationships. Conv-based methods [36] [46] [42] [2] [17] propose several dynamic convolution kernels to adaptively aggregate neighborhood features. Many transformer-like networks [11] [51] [50] [9] [14] extract local features with self-attention. Recently, MLP-like

networks are able to obtain good results with simple networks by enhancing the features. PointMLP [24] proposes a geometric affine module to normalize the feature. Rep-Surf [30] fits the surface information through the triangular plane, models umbrella surfaces to provide geometric information. PointNeXt [28] integrates training strategies and model scaling.

**MLP-like Architecture.** The MLP-like structure has recently shown the ability to rival the Transformer with simple architecture. In the image field, MLP-Mixer [37] first use the combination of Spatial MLP and Channel MLP. The subsequent works [3] [18] reduce computational complexity by selecting objects for the spatial MLP while maintaining a large perceptual field to preserve accuracy. Since the point cloud is too large, the MLP-like network determines the perceptual field generally using K-Nearest neighbor sampling or ball sampling methods. The MLP structure in point cloud analysis starts with PointNet [25] and PointNet++ [26], using MLPs to extract features and aggregating them by symmetric functions. Point-Mixer [6] proposes three point-set operators, PointMLP [24] to modify the distribution of features by geometric affine module, and PointNeXt [28] to scale up the PointNet++ model and improve the performance using by training strategies and model scaling.

**Feature Aggregation.** PosPool [21] improves the reduction function defined in PointNet++ by providing a parameter-free position-adaptive pooling operation. AS-SANet [27] introduces a new anisotropic reduction function. Also, the introduction of the attention mechanism [47] provides new dynamic weights for the reduction function. Vectors have direction, and this property is naturally satisfied for anisotropic aggregation functions. GeoCNN [4] projects features based on vectors and angles of neighbor points and centroids in six directions and sums them. WaveMLP [35] represents image patches as waves and describes feature aggregation using wave phase and amplitude. The Vector Neuron [7] constructs a triad of neurons to reconstruct standard neural networks and represent features through vector transformations. The template-based methods represented by 3DGCN [19] uses the cosine value of the relative displacement vectors to filter for aggregation features from neighbors that more conform to the pattern of the kernel. Local displacements [40] use local displacement vectors to update features by combining the weights of fixed kernels. In our method, an intermediate vector representation is generated by modifying the point feature extraction function. The vector direction is determined based on both features and position to fulfill the anisotropic aggregation function.

### 3. Method

We propose an intermediate vector representation to enhance local feature aggregation in point cloud analysis. This section includes a review of the Point Set Abstraction(SA) operator of the PointNet family in Section 3.1, the presentation of our Vector-oriented Point Set Abstraction module in Section 3.2, a description of our method of extending vectors from scalars in Section 3.3, and the network structure of PointVector in Section 3.4.

#### 3.1. Preliminary

The SA module include a grouping layer (K-NN or Ball-Query) to query each point’s neighbors, shared MLPs, and a reduction layer to aggregate neighbor features. The SA module has an subsample layer to downsample the point cloud in the first layer. We denote  $f_i^{l+1}$  as the extracted feature of point  $i$  after stage  $l+1$ ,  $N_i$  as the neighbors of point  $i$  and  $n$  is the number of incoming points. The content of the SA module can be formulated as follows:

$$f_i^{l+1} = R\{H\{[f_j^l, p_j - p_i] | j \in N_i\}, \quad (1)$$

where  $R$  is the reduction function that aggregates features for point  $i$  from its neighbors  $N_i$  and  $H$  means the shared MLPs.  $f_j^l, p_j, p_i$  denote the input features of point  $j$ , the position of point  $j$  and the position of point  $i$ , respectively.

In the local aggregation operation, the classical method assigns weights to components of  $c$ -dimensional features as shown in Eq.2 and sums the neighboring features in spatial dimensions. We consider the component  $f_i$  of the  $c$ -dimensional feature  $f$  as a base vector with only one non-zero value, and define the vector transformation as follows:

$$f_i * w = w f_i, i = 0 \cdots c, \quad (2)$$

$$[f_i \ 0 \ \cdots \ 0] \begin{bmatrix} w & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix} = [w f_i \ 0 \ \cdots \ 0], \quad (3)$$

where  $w$  is the scalar weight. In Eq.3, the transformation changes one value of the vector. The two equations above are equivalent. Unchanged zeros in the equation do not contribute to subsequent operations and can be disregarded. *In Physics, the degree of freedom of a motion is equal to the number of state quantities that the motion causes the system to change.* A greater number of degrees of freedom in a physical system indicates a larger range of independent variation in the parameters that define its state. Similarly, the degrees of freedom of a vector transformation refer to the number of values in the vector that can change independently. **So, the 3D vector we mentioned means the degrees of freedom of the vector transformation is 3.**

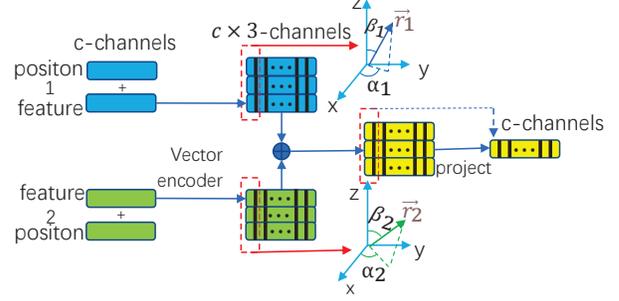


Figure 2. The vector-oriented point set abstraction (VPSA) module of PointVector. It illustrates that VPSA module obtains vector representations from input features, aggregates them, and projects them back to the original feature style. As shown in the figure, each channel of the feature can be considered a 3D vector, with channels being independent of one another.

#### 3.2. Vector-oriented Point Set Abstraction

As discussed in Section 3.1, feature components can be represented as vectors. A higher degree of freedom in vector transformations allows for increased variation and improved representation of connections between neighboring elements. Vectors, with their size and direction properties, are more expressive than scalars for representing features. When aggregated, they exhibit anisotropy due to their directional nature. So, we introduce an intermediate vector representation as Fig.2.

It should be noted that in our assumptions, the component of a  $c$ -dimensional feature represents the projection of the feature vector along the  $c$  coordinate axes. After aggregating the vectors to obtain the  $c \times 3$  centroid feature, where the number of changing values in the component vectors is 3. To merge them into a  $c$ -dimensional feature vector requires aligning the  $c$  components and then summing them. Due to the difficulty in implementing component alignment with this method, we directly project the  $c$  components into scalars and combine them into centroid features. Similar to the intermediate features in a convolutional network, the values on each channel’s feature map represent the response strength to a specific feature at that location.

The input features in our method are transformed into a series of vectors and then aggregated by the reduction function. Note that the element in each channel of the vector representation is vector. We obtain a vector representation that is channel independent. We denote  $f p_j$  as a mixed feature of relative features  $f_j - f_i$  and relative positions  $p_j - p_i$ . The content of the vector-guided aggregation module can be formulated as:

$$f_i^{l+1} = \eta(f_i^l) + H_c\{H_p\{R\{H_v\{f p_j\} | j \in N_i\}\}, \quad (4)$$

where  $H_v$  is the function that generates the vector representation,  $H_p$  denotes the projection Linear transform vector to a scalar, and  $H_c$  is the channel mixing Linear that

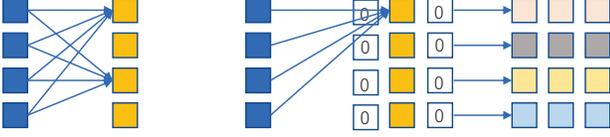


Figure 3. Extension from general feature to vector representation. For simplicity, we tentatively set  $c=4$  and  $m=3$ . The left side represents the process of generating features by standard MLP, and the right side adds 2 components to each scalar of the features to form a vector and then rotates it.

interacts with the information of each channel while transforming dimensions to fit the network. However, the feature representation we introduce is actually represented using a triplet form. We denote  $m$  as the dimension of the vector, and  $c$  is the channel of the feature. In fact, the set of  $c$   $m$ -dimensional vectors is represented in the same form as the  $(m \times c)$ -dimensional feature vectors. The reduction function is followed by a grouped convolution [13] that transforms the vectors to scalars for each channel, which distinguishes the intermediate vector representation from the general feature vector.

When the reduction function  $R$  selects sum, the  $R$  and  $H_p$  functions together constitute a special case of GroupConv [13]. Let  $k$  denote the number of neighbor features. For one group, the convolution kernel of GroupConv is a  $k \times m \times 1$  parameter matrix, while our method can be viewed as  $k$  identical  $m \times 1$  parameter matrices. This is because we treat vectors as wholes and assign equal weight to each element. We will explain in the supplementary material why the original groupconv operation is not suitable for our vector-guided feature aggregation.

### 3.3. Extended Vector From Scalar

The simplest idea for the  $H_v$  function defined in Eq.4 is to obtain  $c$   $m$ -dimensional vectors of point  $j$  directly with MLPs. However, while single-layer MLPs may have limited expressive capability, multi-layer MLPs can be resource-intensive. As discussed in Section 3.1, input features are considered as vectors and we aim to design a transformation with high degrees of freedom. This transformation combines rotation and scaling, represented by a rotation matrix and a learnable parameter respectively. This method achieves better results with lower resource consumption.

As shown in Fig.3, a scalar can be directly converted into an  $m$ -dimensional vector by adding  $m-1$  zero-value components. Each channel of the extended vector representation can then be considered as an  $m$ -dimensional vector along a specific coordinate axis direction. Therefore, we can obtain the proper vector direction by additionally training a rotation matrix. Directly predicting the rotation matrix can cause difficulties for nonlinear optimization because the

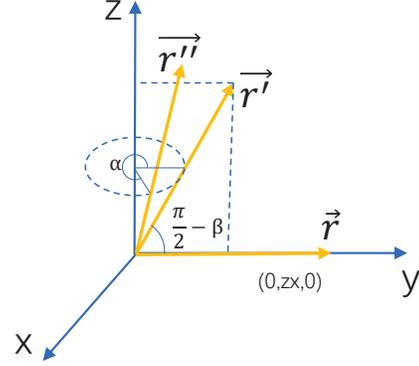


Figure 4. The rotation of a 3D vector. The vector  $\vec{r}$  can be obtained by two rotations to obtain another vector  $\vec{r}'$

matrix elements are interdependent. Instead, we first predict the rotation angle and then derive the rotation matrix based on this angle. The rotation of a 3D vector can be decomposed into rotations around three axes. However, we have not yet determined how to represent the rotation of a 4D vector around a plane. As shown in Fig.4, since the extended 3D vector is on the coordinate axis, one rotation around that axis can be omitted. We keep the default rotation direction as counterclockwise. The vector  $\vec{r}$  is first rotated around the  $x$ -axis by an angle  $\pi/2 - \beta$  and then rotated around the  $z$ -axis by an angle  $\alpha$  to finally obtain the vector  $\vec{r}'$ . The rotation can be formulated as follows:

$$\begin{aligned}
 \vec{r}' &= Rot_z Rot_x \vec{r} \\
 &= \begin{bmatrix} \cos(\alpha) & -\sin(\alpha) & 0 \\ \sin(\alpha) & \cos(\alpha) & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \sin(\beta) & -\cos(\beta) \\ 0 & \cos(\beta) & \sin(\beta) \end{bmatrix} \begin{bmatrix} 0 \\ zx \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} \cos(\alpha) & -\sin(\alpha)\sin(\beta) & \sin(\alpha)\cos(\beta) \\ \sin(\alpha) & \cos(\alpha)\sin(\beta) & -\cos(\alpha)\cos(\beta) \\ 0 & \cos(\beta) & \sin(\beta) \end{bmatrix} \begin{bmatrix} 0 \\ zx \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} -zx \cdot \sin(\alpha)\sin(\beta) \\ zx \cdot \cos(\alpha)\sin(\beta) \\ zx \cdot \cos(\beta) \end{bmatrix},
 \end{aligned} \tag{5}$$

where  $Rot_x, Rot_z$  denote the rotation matrix rotated around the  $x$ -axis and the rotation matrix rotated around the  $z$ -axis, respectively, and  $zx$  is generated by Linear. The independence of  $\alpha$  and  $\beta$  facilitates network optimization.

Therefore, we can expand each scalar value of the features into a 3D vector according to Fig.4 and Eq.5. The feature aggregation in a local area is influenced by the relationship between neighboring points and centroids. Methods such as PointTransformer [53], PACConv [45], and AdaptConv [54] model this relationship using relative position and features. Our approach also extracts rotation angles using MLP on relative positions and features. The acquisition of the vector can be formulated as follows:

$$\begin{aligned}
 zx_j &= Linear(fp_j) \\
 [\alpha_j, \beta_j] &= Relu(BN(Linear([fp_j])),)
 \end{aligned} \tag{6}$$

where  $fp_j$  denotes a mixed feature of relative features  $f_j - f_i$  and relative positions  $p_j - p_i$ , and  $f_j$  means the fea-

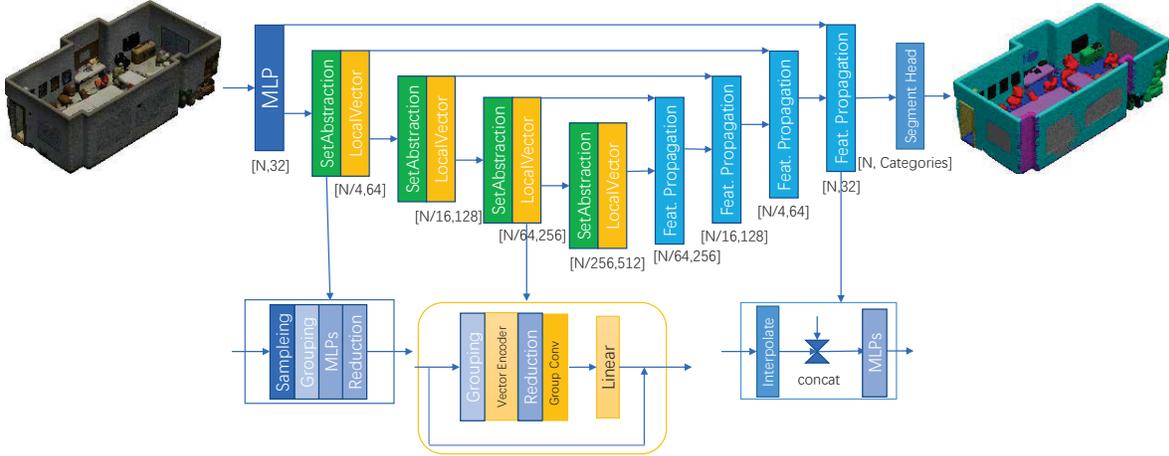


Figure 5. Overall Architecture. We reuse the SA module and Feature Propagation module of PointNet++ and propose the VPSA module to improve the feature extraction of sampled point clouds.

ture of point  $j$ . Therefore, we can obtain the intermediate vector representation from the input features and positions by using Eq.5 and Eq.6.

### 3.4. Architecture

In summary, we propose PointVector, modified from PointNeXt [28] by replacing its InvResMLP module with our proposed VPSA module, we defining its vector dimension  $m = 3$ . The architecture is illustrated in Fig.5. Referring to the classical PointNet++, we use a hierarchical structure containing an encoder and a decoder. For the segmentation task, we use an encoder and a decoder. For the classification task, we only use an encoder. For a fair comparison with PointNeXt, we set up three sizes of models with reference to the parameter settings of PointNeXt. We denote  $C$  as the channel of embedding MLP in the beginning,  $S$  as the numbers of the SA module,  $V$  as the numbers of the VPSA module. The three sizes of models are shown as follows:

- PointVector-S:  $C=32, S=0, V=[1,1,1,1]$
- PointVector-L:  $C=32, S=[1,1,1,1], V=[2,4,2,2]$
- PointVector-XL:  $C=64, S=[1,1,1,1], V=[3,6,3,3]$

Since PointNeXt uses only the PointNeXt-S model for classification, we use our VPSA module instead of the SA module in PointVector-S for a fair comparison. The detailed structure of the classification tasks will appear in the supplementary material. There is a skip connection path in the VPSA module in Fig.5, which is added to the main path and then through a ReLU layer. The reason for using this summation method is that RepSurf [30] indicates how two features with different distributions should be combined. For the segmentation task, finer local information is needed, and

we set reduction function as sum. For the classification task, which favors aggregating global information, we choose the original reduction function such as max.

## 4. Experiments

We evaluate our model on three standard benchmarks: S3DIS [1] for semantic segmentation and ScanObjectNN [48] for real-world object classification and ShapeNetPart [49] for part segmentation. Note that our model is implemented on the basis of PointNeXt. Since we use the training strategy provided by PointNeXt, we refer to the metrics reported by PointNeXt for a fair comparison.

**Experimental setups.** We train PointVector using CrossEntropy loss with label smoothing [33], AdamW optimizer [23], and initial learning rate  $lr=0.002$ , weight\_decay  $10^{-4}$ , with Cosine Decay, and a batch\_size of 32. The above are the base settings for all tasks, and specific parameters will be changed for specific tasks. We follow the train, valid, and test divisions for the dataset. The best model on the validation set will be evaluated on the test set. For S3DIS segmentation task, point clouds are downsampled with a voxel size of 0.4 m following previous methods [36] [27] [53]. The initial learning rate on this task is set to 0.01. For 100 epochs, we use a fixed 24000 points as a batch and set batch\_size to 8. During training, the input points are selected from the nearest neighbors of the random points. Similar to Point Transformer [53], we evaluate our model using the entire scene as input. For ScanObjectNN [48] classification task, we set the weight\_decay to 0.05 for 250 epochs. Following Point-BERT [51], the number of input points is 1024. The training points are randomly sampled from the point cloud, and the testing points are uniformly sampled

during evaluation. The details of data augmentation are the same as those in PointNeXt. For ShapeNetPart part segmentation, we train PointVector-S with a batch size of 32 for 300 epochs. Following PointNet++ [26], 2,048 randomly sampled points with normals are used as input for training and testing.

For voting strategy [20], we keep it the same as PointNeXt and use it only on part segmentation task. To ensure a fair comparison with standard methods, we do not use any ensemble methods, such as SimpleView [8]. We also provide the model parameters (Params) and GFLOPs. We additionally, similar to PointNeXt, provide throughput (instance per second) as an indicator of inference speed. The input data for the throughput calculation are kept consistent with PointNeXt for fair comparison. The throughput of all methods is measured using  $128 \times 1024$  (batch size 128, number of points 1024) as input on ScanObjectNN and  $64 \times 2048$  on ShapeNetPart. On S3DIS,  $16 \times 15,000$  points are used to measure the throughput following [28] [27]. We evaluate our model using an NVIDIA Tesla V100 32 GB GPU and a 48 core Intel Xeon @ 2.10 Hz CPU.

#### 4.1. 3D Semantic segmentation on S3DIS

S3DIS [1] (Stanford Large-Scale 3D Indoor Spaces) is a challenging benchmark composed of 6 large-scale indoor areas, 271 rooms, and 13 semantic categories in total. For our models in S3DIS, the number of neighbors in SetAbstraction is 32, and the number of neighbors in the Local Vector module is 8. PointTransformer [53] also employs most of the training strategies and data enhancements used by PointNeXt, so it is fair for us to compare with it. For a comprehensive comparison, we report the experimental results of PointVector-L and PointVector-XL on S3DIS with 6-fold cross-validation in Table 1 and S3DIS Area 5 in Table 2, respectively. As shown in table 1&2, we achieve state-of-the-art performance on both validation options. Table 1 shows that our largest mode PointVector-XL outperforms PointNeXt-XL by **1.6%**, **3.1%** and **3.5%** in terms of **overall accuracy (OA)**, **mean accuracy(mAcc)** and **mIOU**, respectively, while has only **58% Params**. At the same time, the computational consumption of ours is only **69%** of PointNeXt-XL in terms of GFLOPs. The reduction in computational consumption because the number of neighbors is reduced to 8. *The limitation is that we make heavy use of GroupConv (groups=channel), which is not well optimized in PyTorch and is slower than standard convolution.* Therefore, **our inference speed is 6 instances/second lower** than PointNeXt-XL. Our model shows better results at all sizes.

On S3DIS Area 5, we selected the best results reported by PointNeXt for comparison and did not repeat the experiment. Our PointVector-XL model outperforms StratifiedFormer [14] and PointNeXt-XL by **0.3%** and **1.8%** in mIOU, respectively. StratifiedFormer expands the scope of

Method	OA	mAcc	mIOU	Params	FLOPs	Throughput
	%	%	%	M	G	(ins./sec.)
PointNet [25]	78.5	66.2	47.6	3.6	35.5	162
PointCNN [17]	88.1	75.6	65.4	0.6	-	-
DGCNN [39]	84.1	-	56.1	1.3	-	8
DeepGCN [16]	85.9	-	60.0	3.6	-	3
KPConv [36]	-	79.1	70.6	15.0	-	30
RandLA-Net [12]	88.0	82.0	70.0	1.3	5.8	159
Point Transformer [53]	90.2	81.9	73.5	7.8	5.6	34
CBL [34]	89.6	79.4	73.1	18.6	-	-
RepSurf [30]	90.9	82.6	74.3	0.976	-	-
PointNet++ [26]	81.0	67.1	54.5	1.0	7.2	186
PointNeXt-L [28]	89.8	82.2	73.9	7.1	15.2	115
PointNeXt-XL [28]	90.3	83.0	74.9	41.6	84.8	46
<b>PointVector-L</b>	<b>91.4</b>	<b>85.5</b>	<b>77.4</b>	4.2	10.7	98
<b>PointVector-XL(Ours)</b>	<b>91.9</b>	<b>86.1</b>	<b>78.4</b>	24.1	58.5	40

Table 1. **Semantic segmentation on S3DIS with 6-fold cross-validation.** Methods are in chronological order. The highest and second scores are marked in bold.

the query by combining high-resolution and low-resolution keys while efficiently extracting contextual information. Even though its *receptive field is much wilder* than our model, we still show a competitive performance. Additionally, there are some differences in the experimental setup between our model and it, in which it has **80k points** of input, much larger than our **24k points** of input. In addition it uses *KPConv [36] instead of Linear* in the first layer. It seems that these measures have significant effects. However, the comparison is not fair enough for us due to the difference of the experimental configurations. We will synchronize its experimental configuration later. Additionally, our models of the same size on Area 5 show better results than PointNeXt. PointVector-L and PointVector-XL perform better than PointNeXt-L and PointNeXt-XL by 1.7% and 1.5% in mIOU, respectively, and we performs better on most of categories.

#### 4.2. 3D Object Classification on ScanObjectNN

ScanObjectNN [48] contains approximately 15,000 real scanned objects that are categorized into 15 classes with 2,902 unique object instances. The dataset has significant challenges due to occlusion and noise. As with PointNeXt, we chose the hardest variant PB.T50\_RS of ScanObjectNN and report the mean±std Overall Accuracy and Mean Accuracy score. For our model in ScanObjectNN, the number of neighbors in SetAbstraction is 32. As shown in table.3, our PointVector-S model achieves a comparable performance on ScanObjectNN in *OA*, while outperforms PointNeXt-S by 0.4% in *mAcc*. This illustrates that our approach is not more biased toward certain categories and is relatively robust. Our approach is at a disadvantage in terms of speed and scale compared to the SA module. Since we introduce high-dimensional vectors, we generate more computations before the reduction compared to the standard SA module. Due to group convolution operations and trigonometric functions, there is a speed bottleneck. Although the inference speed is slower than PointNeXt, we are still faster

Method	OA	mAcc	mIoU	ceiling	floor	wall	beam	column	window	door	table	chair	sofa	bookcase	board	clutter
	%	%	%													
PointNet [25]	-	49.0	41.1	88.8	97.3	69.8	0.1	3.9	46.3	10.8	59.0	52.6	5.9	40.3	26.4	33.2
PointCNN [17]	85.9	63.9	57.3	92.3	98.2	79.4	0.0	17.6	22.8	62.1	74.4	80.6	31.7	66.7	62.1	56.7
DGCNN [39]	83.6	-	47.9	-	-	-	-	-	-	-	-	-	-	-	-	-
DeepGCN [16]	-	-	52.5	-	-	-	-	-	-	-	-	-	-	-	-	-
KPConv [36]	-	72.8	67.1	92.8	97.3	82.4	0.0	23.9	58.0	69.0	81.5	91.0	75.4	75.3	66.7	58.9
PVCNN [22]	87.1	-	59.0	-	-	-	-	-	-	-	-	-	-	-	-	-
PAConv [45]	-	73.0	66.6	94.6	<b>98.6</b>	82.4	0.0	26.4	58.0	60.0	89.7	80.4	74.3	69.8	73.5	57.7
ASSANet-L [27]	-	-	66.8	-	-	-	-	-	-	-	-	-	-	-	-	-
Point Transformer [53]	90.8	76.5	70.4	94.0	98.5	<b>86.3</b>	0.0	38.0	<b>63.4</b>	74.3	89.1	82.4	74.3	<b>80.2</b>	76.0	59.3
PatchFormer [52]	-	-	68.1	-	-	-	-	-	-	-	-	-	-	-	-	-
CBL [34]	90.6	75.2	69.4	93.9	98.4	84.2	0.0	37.0	57.7	71.9	<b>91.7</b>	81.8	77.8	75.6	69.1	62.9
RepSurf-U [30]	90.2	76.0	68.9	-	-	-	-	-	-	-	-	-	-	-	-	-
StratifiedFormer* [14]	<b>91.5</b>	<b>78.1</b>	<b>72.0</b>	<b>96.2</b>	<b>98.7</b>	<b>85.6</b>	0.0	<b>46.1</b>	60.0	<b>76.8</b>	<b>92.6</b>	84.5	77.8	75.2	78.1	<b>64.0</b>
PointNet++ [26]	83.0	-	53.5	-	-	-	-	-	-	-	-	-	-	-	-	-
PointNeXt-L [28]	90.1	76.1	69.5	94.0	98.5	83.5	0.0	30.3	57.3	74.2	82.1	91.2	74.5	75.5	76.7	58.9
PointNeXt-XL [28]	90.7	77.5	70.8	94.2	98.5	84.4	0.0	37.7	59.3	74.0	83.1	91.6	77.4	77.2	78.8	60.6
<b>PointVector-L(Ours)</b>	90.8	77.3	71.2	94.8	98.2	84.1	0.0	31.7	60.0	<b>77.7</b>	83.7	<b>91.9</b>	<b>81.8</b>	<b>78.9</b>	<b>79.9</b>	<b>63.3</b>
<b>PointVector-XL(Ours)</b>	<b>91.0</b>	<b>78.1</b>	<b>72.3</b>	<b>95.1</b>	<b>98.6</b>	85.1	0.0	<b>41.4</b>	<b>60.8</b>	76.7	84.4	<b>92.1</b>	<b>82.0</b>	77.2	<b>85.1</b>	61.4

Table 2. **Semantic segmentation on S3DIS Area5.** \* denotes StratifiedFormer use 80k points as input points. The highest and second scores are marked in bold.

Method	OA %	mAcc %	Params. M	Throughput ins./sec.
PointNet [25]	68.2	63.4	3.5	<b>4212</b>
PointCNN [17]	78.5	75.1	0.6	44
DGCNN [39]	78.1	73.6	1.8	402
GBNet [29]	80.5	77.8	8.8	194
PRANet [5]	82.1	79.1	2.3	493
PointMLP [24]	85.4 ± 1.3	83.9 ± 1.5	13.2	191
RepSurf-U [30]	86.0	83.1	6.8	-
PointNet++ [26]	77.9	75.4	1.5	1872
PointNeXt-S [28]	<b>87.7 ± 0.4</b>	<b>85.8 ± 0.6</b>	1.4	2040
<b>PointVector-S(Ours)</b>	<b>87.8 ± 0.4</b>	<b>86.2 ± 0.5</b>	1.55	901

Table 3. **Object classification on ScanObjectNN.** The highest and second scores are marked in bold.

than other methods [24] [39] [29]. Our method does not perform well on the classification task, where the downsampling phase of the classification task requires a max reduction function to retain salient contour information.

### 4.3. 3D Object Part Segmentation on ShapeNetPart

ShapeNetPart [49] is an object-level dataset for part segmentation. It consists of 16,880 models from 16 different shape categories, 2-6 parts for each category, and 50 part labels in total. As shown in Tab.4, our PointVector-S and PointVector-S\_C64 models both achieve results that are comparable to PointNeXt. For the PointNeXt-S model with C=160, the number of parameters is large, and we do not give a corresponding version of the model.

### 4.4. Ablation Study

We perform ablation experiments at S3DIS to verify the effectiveness of the module, and because PointVector-XL is too large, we make changes to PointVector-L. To make the comparison fair, we did not change the training parameters.

Method	Ins.mIoU	Throughput
PointNet [25]	83.7	<b>1184</b>
DGCNN [39]	85.2	147
KPConv [36]	86.2	44
3D-GCN [19]	85.1	-
CurveNet [44]	86.8	97
ASSANet-L [27]	86.1	640
Point Transformer [53]	86.6	297
PointMLP [24]	86.1	270
Stratifiedformer [14]	86.6	398
PointNet++ [26]	85.1	560
PointNeXt-S* [28]	86.5	776
PointNeXt-S* (C=64)	<b>86.9±0.1</b>	330
PointNeXt-S* (C=160)	<b>87.2</b>	75
<b>PointVector-S(Ours)</b>	86.5	446
<b>PointVector-S(C=64)</b>	<b>86.9</b>	211

Table 4. **Object Part Segmentation on ShapeNetPart.** \*Our evaluation results on this task alone are not consistent with the throughput results derived from that paper. Other works we did not test one by one.

**Vector-oriented Point Set Abstraction.** We abstract the module into two key operations: sum and GroupConv(groups=Channel), which shows that this part of the module is channel independent, so we add a FC to mix the channel information. Considering that the channel information is already mixed using non-GroupConv operations, the channel mixing Linear will be removed. The convolution and grouped convolution parts have a convolution kernel size of  $1 \times k$  and a stride size of 1. As shown in Tab.5, the direct use of fixed convolution brings a large number of parameters and fits very poorly with the irregular structure of the point cloud. max+FC shows better performance

Method	OA %	mAcc %	mIOU M	Params
max+FC*	90.6	76.4	70.6	6.35
Conv	90.4	75.7	69.4	24.56
GroupConv	90.6	76.5	70.8	4.76
sum+FC	90.7	76.6	71.0	6.35
max+GroupConv	90.6	76.2	70.6	4.71
sum+GroupConv	90.8	77.3	71.2	4.71

Table 5. Core operation of VPSA. We abstract the module into sum and GroupConv operations, and replace this part. FC means Channel-FC as Linear. \* means it acts as a baseline.

because intuitively aggregating features with higher dimensionality retains more information. GroupConv obtains a lower mIOU because it assigns independent weights to each element of the group; however, the three elements of a 3D vector of a channel should be given the same weight when summing. Furthermore, sum+FC is not very different from sum+GroupConv because GroupConv and channel mixing Linear can be combined into a specific layer of FC. In contrast, sum+GroupConv has the smallest number of parameters and best performance, so we chose it.

Method	OA %	mAcc %	mIOU %	Params M
MLP	91.0	76.5	70.8	5.55
Linear+direction	90.8	76.6	70.8	5.55
Linear+rotation	90.8	77.3	71.2	4.71

Table 6. Methods for obtaining vector representations.

**Extended Vector From Scalar.** To verify the effectiveness of our vector rotation-based method, we compare it with two other methods. As shown in Tab.6, MLP is represented by two Linear layers and a ReLU activation and BatchNorm layers. Linear+direction means that Linear predicts the vector modulus length, then uses MLP to obtain the unit vector as direction, and the final modulus length is multiplied by the unit vector. The rotation-based vector expansion method proposed in Section 3.3 is ahead of other methods and has fewer parameters. This shows that the rotation-based approach can use fewer parameters to obtain a vector representation more suitable for neighbor features.

Method	OA %	mAcc %	mIOU %	Params M
Scalar	90.4	76.1	69.8	3.87
2D vector	90.4	77.2	70.9	3.9
3D vector	90.8	77.3	71.2	4.7

Table 7. Different dimensional vector.

**Vector dimension.** We need to explore the connection between the effect of vector representation and dimension. Intuitively, higher dimensional vectors will be more expressive of features than lower dimensional vectors. Tab.7 shows that the 3D vector has a better ability to express the features and that the increase in the number of parameters is not very large. The mIOU without our vector representation is still higher than the results of PointNeXt. We will discuss the validity of the other parts of our network in the supplementary material.

**Robustness.** Table.8 shows that our method is extremely robust to various perturbations as Stratified Transformer. The ball query we use cannot get the same neighbors in the scaled point cloud. If the query radius is scaled together, then mIOU is invariant. It indicates that our method also has scale invariance.

Method	None	$\pi/2$	$\pi$	$3\pi/2$	+0.2	-0.2	$\times 0.8$	$\times 1.2$	jitter
PointNet++ [25]	59.75	58.15	58.18	58.19	22.33	29.85	56.24	59.74	59.05
PointTr [51]	70.36	65.94	67.78	65.72	70.44	70.43	65.73	66.15	59.67
Stratified	71.96	72.59	72.37	71.86	71.99	71.93	70.42	71.21	72.02
Ours	72.29	72.27	72.30	72.32	72.29	72.29	69.34	69.26	72.16

Table 8. Robustness study on S3DIS (mIOU %). We apply z-axis rotation ( $\pi/2$ ,  $\pi$ ,  $3\pi/2$ ), shifting ( $\pm 0.2$ ), scaling ( $\times 0.8$ ,  $\times 1.2$ ) and jitter in testing. PointTr: Point Transformer. Stratified: Stratified Transformer.

## 5. Conclusion and Limitation.

We introduce PointVector, which achieves state-of-the-art results on the S3DIS semantic segmentation task. Our vector-oriented point set abstraction improves local feature aggregation with fewer parameters. The rotation-based vector expansion method bridges the gap between vector representation and standard feature forms. By optimizing two independent perspectives, it achieves better results. Additionally, our method exhibits robustness to various perturbations. It is noteworthy that further exploration of vector representation’s meaning may reveal additional applications, i.e. dominant neighbor selection.

The speed of our approach is constrained by the grouped convolution implementation. An interesting avenue for future work includes exploring rotations above three dimensions and decomposing four-dimensional rotations into combinations of plane rotations. Additionally, summing after component alignment aligns with our assumptions better than scalar projection.

## Acknowledgement

This work was supported in part by the National Key Research and Development Program of China under Grant 2020YFB2103803.

## A. Preliminary

### A.1. Problem of WaveMLP.

WaveMLP [35] views the patch of each picture as a wave representation, and considers that the feature of that patch should have two attributes, phase and amplitude, with amplitude representing the actual property of the feature and phase modulating the amplitude that this wave exhibits at a moment. It thus considers that the feature extraction of the patches can be viewed as a superposition of waves. However, there is an important problem, WaveMLP gets an absolute representation of a patch, i.e. the patch is the same when participating in aggregation in any local region. The representation of a patch should be different in different local regions, so we focus on modulating the feature aggregation in local regions. That is, we use a vector representation to better express the relative relationship between neighbor points and centroids in the local region.

In addition, WaveMLP use GroupConv [13] to implement the aggregation and projection process with kernel sizes of  $1 \times 7$  and  $7 \times 1$ . In this paper we take the form of a combination of the reduction function and GroupConv for aggregation. We give an example of why the original GroupConv is not suitable for this representation of vectors. We take two-dimensional vectors  $(x_1, y_1)$  and  $(x_2, y_2)$  as an example. The vectors are represented in coordinate form, and then the original vector aggregation method can be formulated as:

$$\begin{aligned} f_{12} &= (w_1(x_1, y_1) + w_2(x_2, y_2)) \cdot (w_3, w_4)^T \\ &= w_3w_1x_1 + w_3w_2x_2 + w_4w_1y_1 + w_4w_2y_2 \quad (7) \\ &= a_1x_1 + a_2x_2 + a_3y_1 + a_4y_2, \end{aligned}$$

where  $f_{12}$  denotes the result of aggregating two vectors,  $w_1$  and  $w_2$  are the weights of two vectors in summation,  $\{w_3, w_4\}$  is the projection matrix, and  $a_i$  is the weight of each component. We can obtain the equation that should be satisfied between the coefficients of each component:  $a_1 * a_4 = a_2 * a_3$ . That is, the final trained weights need to satisfy this equation for the weighted summation formula of the vectors to hold. However, the network does not impose this restriction on these parameters. So the original groupconv does not preserve the totality of the vector.

### A.2. Methodology Review.

The point-based approach was first introduced by PointNet [25]. We denote  $f_i^{l+1}$  as the extracted feature of point  $i$  after stage  $l+1$ ,  $N_i$  as the neighbors of point  $i$  and  $n$  is the number of incoming points. The simplest point-set operator can be expressed as follows:

$$f_i^{l+1} = R\{H\{[f_j^l, p_j - p_i]\}|j \in N_i\}, \quad (8)$$

where  $R$  is the reduction function that aggregates features for point  $i$  from its neighbors  $N_i$  and  $H$  means the shared MLPs.

The subsequent dynamic convolution-based network [36] [45] can be similarly represented as PointNet-like point set operators:

$$f_i^{l+1} = Sum\{\phi\{f_j^l, p_j - p_i\} \cdot f_j^l | j \in N_i\}, \quad (9)$$

where  $\phi()$  means the dynamic weight generation function that generates dynamic weights for each point based on the input feature and location information. Eq.9 shows that the reduction function of dynamic convolution chooses sum and uses dynamic weights to generate a new  $f_j$ .

Similarly, the attention network [53] can be expressed as a similar point set operator. The core operation can be formulated as follows:

$$f_i^{l+1} = Sum\{att\{f_j^l, f_i^l, pos\} \cdot \sigma\{f_j^l, pos\} | j \in N_i\}, \quad (10)$$

where  $att()$  means the attention function that generates attention weights for each point,  $pos$  denotes the position information, and  $\sigma()$  means the linear transform function without anisotropy. Eq.10 shows that it uses the attention mechanism to update the features of each point  $j$  and then uses sum as the reduction function.

Furthermore, template-based methods such as 3D-GCN make use of kernels with relative displacement vectors and weights. These weights are influenced by the cosine similarity between the relative displacement vector of the input features and the relative displacement vector of the kernel. The core operation can be formulated as follows:

$$\begin{aligned} f_i^{l+1} &= f_i \cdot kernel_c + \sum_{m=1}^k max\{sim\{kernel_m, f_j\} | j \in N_i\}, \\ sim\{kernel_m, f_j\} &= cos\{dk_m, dp_j\} \cdot kernel_m \cdot f_j, \end{aligned} \quad (11)$$

where  $k$  means the kernel size,  $N_i$  means the neighbors of point  $i$ ,  $kernel_c$  means the center element of kernel,  $cos\{dk_m, dp_j\}$  means Cosine similarity of  $m$ -th kernel element and  $j$ -th point feature,  $kernel_m$  means  $m$ -th kernel element,  $dk_m, dp_j$  means displacement vector of  $m$ -th kernel element and  $j$ -th point feature respectively.

We propose a unique method for generating new features  $f_j$  by introducing a vector representation, where the direction of the vector guides the aggregation method.

## B. Architecture

### B.1. Vector encoder

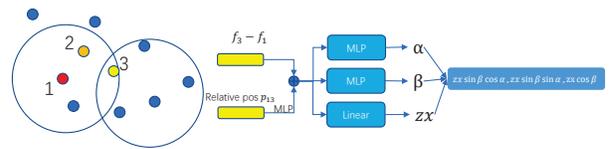


Figure 6. **The Vector encoder module.** Two angles are predicted by MLP and  $z_x$  is transformed by linear.

We provide detailed definitions in the manuscript, and we provide illustrations to illustrate the exact process. As shown in Fig. 6, the local information is obtained by a combination of relative features and relative positions. Note that the sum symbol in the figure means sum and ReLU operations. We use the simplest method to predict the angles using MLP, and by default the two angles are independent of each other. For  $zx$ , a simple transformation is performed with linear, and then a vector representation is obtained by rotation. The vector representation  $v \in R^{B,C \times 3,N}$ , where  $B$  is the batch size,  $C$  is the channel of module and  $N$  is the spatial size of the input feature of the module.

## B.2. Classification architecture.



Figure 7. **The Classification architecture PointVector-S.** For comparison with PointNext [28], we replaced the SetAbstract module with the LocalVector module, keeping the other parameters the same.

As shown in Fig. 7, we use the LocalVector module to replace the 4 SetAbstract modules and keep the downsampling parameters unchanged. The last SetAbstract was originally used to aggregate all the remaining points, so we leave it as it is. In the classification task, the max reduction function has a greater advantage by retaining the most intense part of the variation.

## C. Experiments

### C.1. Classification on ModelNet40

ModelNet40 [43] is a commonly used dataset for object classification, which is generated by 3D graphic CAD models. It has 40 object categories, each of which contains 100 unique CAD models. Recent works [24] [30] [10] show an increasing interest in the real-world scanned dataset ScanObjectNN [48] than this synthesized 3D dataset ModelNet40. Therefore, we choose to report the results on ScanObjectNN in the manuscript. Furthermore, we report the results of our PointVector-S model on ModelNet40. We use the same parameters as PointNext: CrossEntropy loss with label smoothing, AdamW optimizer, a learning rate of  $1e-3$ , a weight decay of 0.05, cosine learning rate decay, and a batch size of 32 for 600 epochs, while using random scaling and translation as data augmentations. As shown in table 9, the relatively poor performance of our model on

Method	mAcc %	OA %
PointNet [25]	86.2	89.2
PointCNN [17]	88.1	92.2
PointConv [42]	-	92.5
KPConv [36]	-	92.9
DGCNN [39]	90.2	92.9
DeepGCN [16]	90.9	93.6
ASSANet-L [27]	-	92.9
Point Cloud Transformer [9]	-	93.2
Point Transformer [53]	90.6	93.7
CurveNet [44]	-	93.8
PointMLP [24]	90.9±0.4	93.7±0.2
PointNet++	-	91.9
PointNet++(PointNext)	89.9 ± 0.8	92.8 ± 0.1
PointNext(C=32)	90.8 ± 0.2	93.2 ± 0.1
PointNext(C=64)	90.9 ± 0.5	93.7 ± 0.3
PointVector-S(C=32)	90.3 ± 0.2	93.2 ± 0.2
PointVector-S(C=64)	91.0 ± 0.5	93.5 ± 0.2

Table 9. Object Classification on ModelNet40.

the ModelNet40 dataset indicates the limitation of the proposed local vector representation in aggregating global information. We used hyperparameters consistent with PointNext and a training strategy that may not be suitable for our model, which may also account for the relatively poor performance. Note that our network structure on the classification task directly takes vector feature aggregation for downsampling, but max-pooling is probably the simplest and most effective method for downsampling.

### C.2. Ablation study

There is a slight problem with the experimental setup in the manuscript, in the 6-fold cross-validation experiment we report the PointVector-L as the standard setup mentioned in the manuscript, but in the S3DIS Area5 and ablation experiments we report the setup of PointVector-L as  $V=[2, 2, 4, 2]$ . But, the max+groupconv in the manuscript is reported as  $V=[2,4,2,2]$ .

Method	size	OA %	mAcc %	mIOU %
PointVector-L (max+groupconv)	$V=[2,4,2,2]$	90.6	76.2	70.6
	$V=[2,2,4,2]$	90.6	77.1	71.1
PointVector-L (sum+groupconv)	$V=[2,4,2,2]$	90.3	77.21	70.8
	$V=[2,2,4,2]$	90.8	77.3	71.2
PointVector-XL	$V=[3,5,3,3]$	90.8	78.3	72.3
	$V=[3,3,5,3]$	91.0	76.7	71.1

Table 10. Results for models with different number of stagess on S3DIS Area5.

**Number of stages.** Since the PointVector-L with max+groupconv is reported by another configuration in the manuscript, we compare the two configurations here. As the tab.10 shows, the two reduction functions, max and sum, obtain very similar results, but sum has a higher mAcc and OA. This is consistent with our assumption that better results can be obtained by simply using groupconv to process vectors of each channel independently. Small and large models do not behave consistently in terms of the number of stages. This is an interesting phenomenon, but not the main point of our statement, so it will not be discussed for now.

The following experiments are reported by default as PointVector-XL [3,5,3,3], PointVector-L [2,2,4,2] if no special instructions are given.

Method	OA %	mAcc %	mIOU %	Params M
PointNeXt-XL	90.7	77.5	70.8	41.6
PointVector-base	90.9	77.0	71.4	37.2
PointVector-XL	90.8	78.3	72.3	24.1

Table 11. Baseline. The same experimental configuration was used for all three models.

**Baseline.** Our model has some gaps in channel variations and inputs with PointNeXt. To really evaluate whether our model has a greater advantage, we reset a baseline. We take our core operations i.e. Vector encoder and reduction+groupconv+channel mixing Linear was removed and replaced with PointNeXt’s MLP+max+MLPs, where the channel of first MLP was transformed from  $c$  to  $3c$ . The new model is named PointVector-base. The tab.11 shows that our model has a large improvement in each metric compared to baseline. Also this shows that the other parts of our model are superior compared to the original PointNeXt.

type	Method	OA %	mAcc %	mIOU %
feature	$f_j - f_i + \text{pos}$	90.8	77.3	71.2
	$[f_j - f_i, \text{pos}]$	88.8	70.5	64.9
	$f_j + \text{pos}$	90.9	76.6	70.5
residual	linear	90.8	77.3	71.2
	identity	90.3	75.8	69.3

Table 12. Other Components. + means that the two are added together and then passed through the relu layer. [,] means to directly concatenate two elements.

**Other Components.** The manuscript mentions that other operations of our model have a larger role, so we conducted ablation experiments on PointVector-L to explore the effect

of both input features and residuals on the S3DIS segmentation task. Tab.12 shows that the two parts of the features are added together and then relu can better fuse their information. In addition relative features are more robust than absolute features. The key is that residual uses linear compared to identity, which is a huge improvement.

## D. Visualization

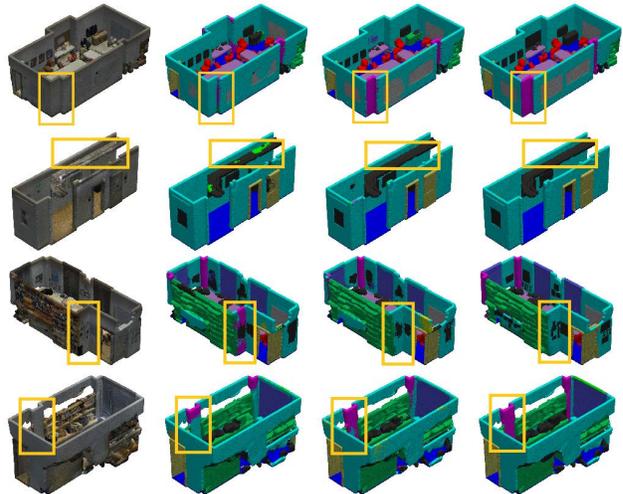


Figure 8. **Qualitative comparisons of PointNext ( $2^{nd}$  column), PointVector++ ( $3^{rd}$  column), and Ground Truth ( $4^{th}$  column) on S3DIS semantic segmentation.** The input point cloud is visualized with original colors in the  $1^{st}$  column. We have circled the different places with a paintbrush.

As shown in the Fig.8, it can be found that our model performs a little better in complex areas. This shows that we are able to extract more detail in such intensely varied areas than the max-pooling operation of PointNeXt. But we are also prone to miscalculation in flat areas, which is our disadvantage.

## E. Code release

Since our model is based on PointNext, we used their code and added a PointVector model. Since our classification and part segmentation and semantic segmentation tasks use different model compositions, the model code is also different, and the corresponding PointVector.py needs to be replaced at runtime. We have not organized the code yet, where PATM represents the core part of our LocalVector module. In the classification and part segmentation tasks, it replaces the convs+max pooling operation in SetAbstraction. See the official instructions for PointNext for related running instructions. And on s3dis our gravity\_dim is set to 1. The code of each task is a little different, on ScanObjectNN classification task we insert leakyrelu in the two

linear after the reduction function, and the relative features of the input after BN, encoder’s activation function all use leakyrelu can reach 88.4% OA, but this is not the main point of our statement, so we do not discuss it for now. The code takes time to organize and we will make it public later.

## References

- [1] Iro Armeni, Ozan Sener, Amir R. Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. [2](#), [5](#), [6](#)
- [2] Matan Atzmon, Haggai Maron, and Yaron Lipman. Point convolutional neural networks by extension operators. *ACM Transactions on Graphics*, 37(4), 2018. [2](#)
- [3] Shoufa Chen, Enze Xie, Chongjian GE, Runjian Chen, Ding Liang, and Ping Luo. CycleMLP: A MLP-like architecture for dense prediction. In *International Conference on Learning Representations*, 2022. [2](#)
- [4] Yuedong Chen, Guoxian Song, Zhiwen Shao, Jianfei Cai, Tat-Jen Cham, and Jianmin Zheng. Geoconv: Geodesic guided convolution for facial action unit recognition. *Pattern Recognition*, 122, 2022. [2](#)
- [5] Silin Cheng, Xiwu Chen, Xinwei He, Zhe Liu, and Xiang Bai. Pra-net: Point relation-aware network for 3d point cloud analysis. *IEEE Transactions on Image Processing*, 30:4436 – 4448, 2021. [7](#)
- [6] Jaesung Choe, Chunghyun Park, Francois Rameau, Jaesik Park, and In So Kweon. Pointmixer: Mlp-mixer for point cloud understanding. 2021. [2](#)
- [7] Congyue Deng, Or Litany, Yueqi Duan, Adrien Poulénard, Andrea Tagliasacchi, and Leonidas J. Guibas. Vector neurons: A general framework for so(3)-equivariant networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 12200–12209, October 2021. [2](#)
- [8] Ankit Goyal, Hei Law, Bowei Liu, Alejandro Newell, and Jia Deng. Revisiting point cloud shape classification with a simple and effective baseline. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages 3809–3820. PMLR, 2021. [6](#)
- [9] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R. Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187 – 199, 2021. [2](#), [10](#)
- [10] Abdullah Hamdi, Silvio Giancola, and Bernard Ghanem. Mvtn: Multi-view transformation network for 3d shape recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1 – 11, Virtual, Online, Canada, 2021. [2](#), [10](#)
- [11] Qingdong He, Zhengning Wang, Hao Zeng, Yi Zeng, and Yijun Liu. Svga-net: Sparse voxel-graph attention network for 3d object detection from point clouds. 2020. [2](#)
- [12] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randla-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 11105 – 11114, Virtual, Online, United states, 2020. [6](#)
- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc., 2012. [4](#), [9](#)
- [14] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8500–8509, June 2022. [2](#), [6](#), [7](#)
- [15] Alex H. Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2019-June, pages 12689 – 12697, Long Beach, CA, United states, 2019. [2](#)
- [16] Guohao Li, Matthias Mueller, Guocheng Qian, Itzel Carolina Delgadillo Perez, Abdullellah Abualshour, Ali Kassem Thabet, and Bernard Ghanem. Deepgcns: Making gcns go as deep as cnns. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. [6](#), [7](#), [10](#)
- [17] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed systems. In *Advances in Neural Information Processing Systems*, volume 2018-December, pages 820 – 830, Montreal, QC, Canada, 2018. [2](#), [6](#), [7](#), [10](#)
- [18] Dongze Lian, Zehao Yu, Xing Sun, and Shenghua Gao. As-mlp: An axial shifted mlp architecture for vision. In *International Conference on Learning Representations (ICLR)*, 2022. [2](#)
- [19] Zhi-Hao Lin, Sheng-Yu Huang, and Yu-Chiang Frank Wang. Convolution in the cloud: Learning deformable kernels in 3d graph convolution networks for point cloud analysis. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1797–1806, 2020. [1](#), [2](#), [7](#)
- [20] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2019-June, pages 8887 – 8896, Long Beach, CA, United states, 2019. [6](#)
- [21] Ze Liu, Han Hu, Yue Cao, Zheng Zhang, and Xin Tong. A closer look at local aggregation operators in point cloud analysis. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 326–342. Springer, 2020. [2](#)
- [22] Zhijian Liu, Haotian Tang, Yujun Lin, and Song Han. Pointvoxel cnn for efficient 3d deep learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. [7](#)
- [23] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *7th International Conference on Learn-*

- ing Representations, *ICLR 2019*, New Orleans, LA, United states, 2019. 5
- [24] Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local geometry in point cloud: A simple residual mlp framework, 2022. 2, 7, 10
- [25] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017*, volume 2017-January, pages 77 – 85, Honolulu, HI, United states, 2017. 1, 2, 6, 7, 9, 10
- [26] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 1, 2, 6, 7
- [27] Guocheng Qian, Hasan Hammoud, Guohao Li, Ali Thabet, and Bernard Ghanem. Assanet: An anisotropic separable set abstraction for efficient point cloud representation learning. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 28119–28130. Curran Associates, Inc., 2021. 2, 5, 6, 7, 10
- [28] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. 1, 2, 5, 6, 7, 10
- [29] Shi Qiu, Saeed Anwar, and Nick Barnes. Geometric back-projection network for point cloud classification. *IEEE Transactions on Multimedia*, 24:1943 – 1955, 2022. 7
- [30] Haoxi Ran, Jun Liu, and Chengjie Wang. Surface representation for point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18942–18952, June 2022. 2, 5, 6, 7, 10
- [31] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 10526 – 10535, Virtual, Online, United states, 2020. 2
- [32] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015. 2
- [33] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 5
- [34] Liyao Tang, Yibing Zhan, Zhe Chen, Baosheng Yu, and Dacheng Tao. Contrastive boundary learning for point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8489–8499, June 2022. 6, 7
- [35] Yehui Tang, Kai Han, Jianyuan Guo, Chang Xu, Yanxi Li, Chao Xu, and Yunhe Wang. An image patch is a wave: Phase-aware vision mlp. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10935–10944, June 2022. 2, 9
- [36] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, Francois Goulette, and Leonidas Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, volume 2019-October, pages 6410 – 6419, Seoul, Korea, Republic of, 2019. 1, 2, 5, 6, 7, 9, 10
- [37] Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. Mlp-mixer: An all-mlp architecture for vision. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 24261–24272. Curran Associates, Inc., 2021. 2
- [38] Lei Wang, Yuchun Huang, Yaolin Hou, Shenman Zhang, and Jie Shan. Graph attention convolution for point cloud semantic segmentation. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2019-June, pages 10288 – 10297, Long Beach, CA, United states, 2019. 2
- [39] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics*, 38(5), 2019. 1, 2, 6, 7, 10
- [40] Yida Wang, David Joseph Tan, Nassir Navab, and Federico Tombari. Learning local displacements for point cloud completion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [41] Yan Wang, Wanxia Zhong, Hang Su, Fujiang Zheng, Yiran Pang, Hongchuan Wen, and Kun Cai. An improved mvccn for 3d shape recognition. In *Proceedings of 2021 IEEE International Conference on Emergency Science and Information Technology, ICESIT 2021*, pages 469 – 472, Chongqing, China, 2021. 2
- [42] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 2019-June, pages 9613 – 9622, Long Beach, CA, United states, 2019. 2, 10
- [43] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, volume 07-12-June-2015, pages 1912 – 1920, Boston, MA, United states, 2015. 10
- [44] Tiange Xiang, Chaoyi Zhang, Yang Song, Jianhui Yu, and Weidong Cai. Walk in the cloud: Learning curves for point clouds shape analysis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 895 – 904, Virtual, Online, Canada, 2021. 7, 10

- [45] Mutian Xu, Runyu Ding, Hengshuang Zhao, and Xiaojuan Qi. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3172 – 3181, Virtual, Online, United states, 2021. [4](#), [7](#), [9](#)
- [46] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, volume 11212 LNCS, pages 90 – 105, Munich, Germany, 2018. [2](#)
- [47] Jiancheng Yang, Qiang Zhang, Bingbing Ni, Linguo Li, Jinxian Liu, Mengdie Zhou, and Qi Tian. Modeling point clouds with self-attention and gumbel subset sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. [2](#)
- [48] Li Yi, Vladimir G. Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics*, 35(6), 2016. [2](#), [5](#), [6](#), [10](#)
- [49] Li Yi, Vladimir G. Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics*, 35(6), 2016. [2](#), [5](#), [7](#)
- [50] Xumin Yu, Yongming Rao, Ziyi Wang, Zuyan Liu, Jiwen Lu, and Jie Zhou. Pointtr: Diverse point cloud completion with geometry-aware transformers. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 12478 – 12487, Virtual, Online, Canada, 2021. [2](#)
- [51] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19313–19322, June 2022. [2](#), [5](#)
- [52] Cheng Zhang, Haocheng Wan, Xinyi Shen, and Zizhao Wu. Patchformer: An efficient point transformer with patch attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11799–11808, June 2022. [7](#)
- [53] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip H.S. Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 16259–16268, October 2021. [1](#), [4](#), [5](#), [6](#), [7](#), [9](#), [10](#)
- [54] Haoran Zhou, Yidan Feng, Mingsheng Fang, Mingqiang Wei, Jing Qin, and Tong Lu. Adaptive graph convolution for point cloud analysis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4945 – 4954, Virtual, Online, Canada, 2021. [1](#), [4](#)
- [55] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 4490 – 4499, Salt Lake City, UT, United states, 2018. [2](#)