

# PCT: Point cloud transformer

Meng-Hao Guo<sup>1</sup>, Jun-Xiong Cai<sup>1</sup>, Zheng-Ning Liu<sup>1</sup>, Tai-Jiang Mu<sup>1</sup>, Ralph R. Martin<sup>2</sup>, and Shi-Min Hu<sup>1</sup> (✉)

© The Author(s) 2021.

**Abstract** The irregular domain and lack of ordering make it challenging to design deep neural networks for point cloud processing. This paper presents a novel framework named *Point Cloud Transformer* (PCT) for point cloud learning. PCT is based on Transformer, which achieves huge success in natural language processing and displays great potential in image processing. It is inherently permutation invariant for processing a sequence of points, making it well-suited for point cloud learning. To better capture local context within the point cloud, we enhance input embedding with the support of farthest point sampling and nearest neighbor search. Extensive experiments demonstrate that the PCT achieves the state-of-the-art performance on shape classification, part segmentation, semantic segmentation, and normal estimation tasks.

**Keywords** 3D computer vision; deep learning; point cloud processing; Transformer

## 1 Introduction

Extracting semantics directly from a point cloud is an urgent requirement in some applications such as robotics, autonomous driving, augmented reality, etc. Unlike 2D images, point clouds are disordered and unstructured, making it challenging to design neural networks to process them. Charles et al. [1] pioneered PointNet for feature learning on point clouds by using multi-layer perceptrons (MLPs), max-pooling,

and rigid transformations to ensure invariance under permutations and rotations. Inspired by strong progress made by convolutional neural networks (CNNs) in the field of image processing, many recent works [2–5] have considered to define convolution operators that can aggregate local features for point clouds. These methods either reorder the input point sequence or voxelize the point cloud to obtain a canonical domain for convolutions.

Recently, *Transformer* [6], the dominant framework in natural language processing, has been applied to image vision tasks, giving better performance than popular convolutional neural networks [7, 8]. Transformer is a decoder–encoder structure that contains three main modules for input (word) embedding, positional (order) encoding, and self-attention. The self-attention module is the core component, generating refined attention feature for its input feature based on global context. First, self-attention takes the sum of input embedding and positional encoding as input, and computes three vectors for each word: *query*, *key*, and *value* through trained linear layers. Then, the attention weight between any two words can be obtained by matching (dot-producting) their query and key vectors. Finally, the attention feature is defined as the weighted sum of all value vectors with the attention weights. Obviously, the output attention feature of each word is related to all input features, making it capable of learning the global context. All operations of Transformer are parallelizable and order-independent. In theory, it can replace the convolution operation in a convolutional neural network and has better versatility. For more detailed introduction of self-attention, please refer to Section 3.2.

Inspired by the Transformer’s success in vision and NLP tasks, we propose a novel framework PCT

1 BNRist, Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China. E-mail: M.-H. Guo, gmh20@mails.tsinghua.edu.cn; J.-X. Cai, junxiong20@mails.tsinghua.edu.cn; Z.-N. Liu, lzhenngning@gmail.com; T.-J. Mu, taijiang@tsinghua.edu.cn; S.-M. Hu, shimin@tsinghua.edu.cn (✉).

2 Cardiff University, Cardiff CF243AA, UK. E-mail: ralph@cs.cf.ac.uk.

Manuscript received: 2021-03-04; accepted: 2021-03-26

for point cloud learning based on the principles of traditional Transformer. The key idea of PCT is using the inherent order invariance of Transformer to avoid the need to define the order of point cloud data and conduct feature learning through the attention mechanism. As shown in Fig. 1, the distribution of attention weights is highly related to part semantics, and it does not seriously attenuate with spatial distance.

Point clouds and natural language are rather different kinds of data, so our PCT framework must make several adjustments for this. These include:

- **Coordinate-based input embedding module.** In Transformer, a positional encoding module is applied to represent the word order in natural language. This can distinguish the same word in different positions and reflect the positional relationships between words. However, point clouds do not have a fixed order. In our PCT framework, we merge the raw positional encoding and the input embedding into a coordinate-based input embedding module. It can generate distinguishable features, since each point has a **unique** coordinate which represents its spatial position.
- **Optimized offset-attention module.** The offset-attention module approach we proposed is an effective upgrade over the original self-attention. It works by replacing the attention feature with the offset between the input of self-attention module and attention feature. This has two advantages. Firstly, the absolute coordinates of the same object can be completely different

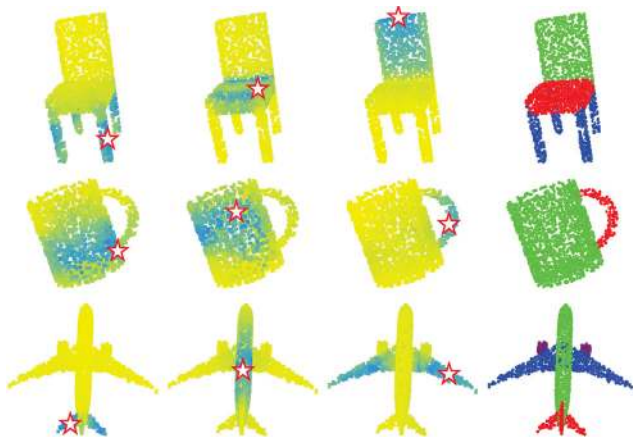
with rigid transformations. Therefore, relative coordinates are generally more robust. Secondly, the Laplacian matrix (the offset between degree matrix and adjacency matrix) has been proven to be very effective in graph convolution learning [9]. From this perspective, we regard the point cloud as a graph with the “float” adjacency matrix as the attention map. Also, the attention map in our work will be scaled with all the sum of each rows to 1. So the degree matrix can be understood as the identity matrix. Therefore, the offset-attention optimization process can be approximately understood as a Laplace process, which will be discussed detailed in Section 3.3. In addition, we have conducted sufficient comparative experiments, introduced in Section 4, on offset-attention and self-attention to prove its effectiveness.

- **Neighbor embedding module.** Obviously, every word in a sentence contains basic semantic information. However, the independent input coordinates of the points are only weakly related to the semantic content. Attention mechanism is effective in capturing global features, but it may ignore local geometric information which is also essential for point cloud learning. To address this problem, we use a neighbor embedding strategy to improve upon point embedding. It also assists the attention module by considering attention between local groups of points containing semantic information instead of individual points.

With the above adjustments, the PCT becomes more suitable for point cloud feature learning and achieves the state-of-the-art performance on shape classification, part segmentation, semantic segmentation, and normal estimation tasks. All experiments are implemented with Jittor [10] deep learning framework. Codes are available at <https://github.com/MenghaoGuo/PCT>.

The main contributions of this paper are summarized as following:

1. We proposed a novel transformer based framework named PCT for point cloud learning, which is exactly suitable for unstructured, disordered point cloud data with irregular domain.
2. We proposed offset-attention with implicit Laplace operator and normalization refinement which is inherently permutation-invariant and more suitable for point cloud learning compared to the original self-attention module in Transformer.



**Fig. 1** Attention map and part segmentation generated by PCT. First three columns: point-wise attention map for different query points (indicated by ☆), yellow to blue indicating increasing attention weight. Last column: part segmentation results.

- Extensive experiments demonstrate that the PCT with explicit local context enhancement achieves state-of-the-art performance on shape classification, part segmentation, and normal estimation tasks.

## 2 Related work

### 2.1 Transformer in NLP

Bahdanau et al. [11] proposed a neural machine translation method with an attention mechanism, in which attention weight is computed through the hidden state of an RNN. Self-attention was proposed by Lin et al. [12] to visualize and interpret sentence embeddings. Building on these, Vaswani et al. [6] proposed Transformer for machine translation; it is based solely on self-attention, without any recurrence or convolution operators. Devlin et al. [13] proposed bidirectional transformers (BERT) approach, which is one of the most powerful models in the NLP field. More lately, language learning networks such as XLNet [14], Transformer-XL [15], and BioBERT [16] have further extended the Transformer framework.

However, in natural language processing, the input is in order, and word has basic semantic, whereas point clouds are unordered, and individual points have no semantic meaning in general.

### 2.2 Transformer for vision

Many frameworks have introduced attention into vision tasks. Wang et al. [17] proposed a residual attention approach with stacked attention modules for image classification. Hu et al. [18] presented a novel spatial encoding unit, the SE block, whose idea was derived from the attention mechanism. Zhang et al. [19] designed SAGAN, which uses self-attention for image generation. There has also been an increasing trend to employ Transformer as a module to optimize neural networks. Wu et al. [8] proposed visual transformers that apply Transformer to token-based images from feature maps for vision tasks. Recently, Dosovitskiy [7], proposed an image recognition network, ViT, based on patch encoding and Transformer, showing that with sufficient training data, Transformer provides better performance than a traditional convolutional neural network. Carion et al. [20] presented an end-to-end detection transformer that takes CNN features as input and generates bounding boxes with a Transformer encoder-decoder.

Inspired by the local patch structures used in ViT and basic semantic information in language word, we present a neighbor embedding module that aggregates features from a point's local neighborhood, which can capture the local information and obtain semantic information.

### 2.3 Point-based deep learning

PointNet [1] pioneered point cloud learning. Subsequently, Qi et al. [21] proposed PointNet++, which uses query ball grouping and hierarchical PointNet to capture local structures. Several subsequent works considered how to define convolution operations on point clouds. One main approach is to convert a point cloud into a regular voxel array to allow convolution operations. Tchapmi et al. [2] proposed SEGCloud for pointwise segmentation. It maps convolution features of 3D voxels to point clouds using trilinear interpolation and keeps global consistency through fully connected conditional random fields. Atzmon et al. [4] presented the PCNN framework with extension and restriction operators to map between point-based representation and voxel-based representation. Volumetric convolution is performed on voxels for point feature extraction. MCCNN by Hermosilla et al. [22] allows non-uniformly sampled point clouds; convolution is treated as a Monte Carlo integration problem. Similarly, in PointConv proposed by Wu et al. [5], 3D convolution is performed through Monte Carlo estimation and importance sampling.

A different approach redefines convolution to operation on irregular point cloud data. Li et al. [3] introduced a point cloud convolution network, PointCNN, in which a  $\chi$ -transformation is trained to determine a 1D point order for convolution. Tatarchenko et al. [23] proposed tangent convolution, which can learn surface geometric features from projected virtual tangent images. SPG proposed by Landrieu and Simonovsky [24] divides the scanned scene into similar elements, and establishes a superpoint graph structure to learn contextual relationships between object parts. Yang et al. [25] used a parallel framework to extend CNN from the conventional domain to a curved two-dimensional manifold. However, it requires dense 3D gridded data as input so is unsuitable for 3D point clouds. Wang et al. [26] designed an EdgeConv operator for dynamic graphs, allowing point cloud learning by recovering local topology.

Various other methods also employ attention and Transformer. Yan et al. [27] proposed PointASNL to deal with noise in point cloud processing, using a self-attention mechanism to update features for local groups of points. Hertz et al. [28] proposed PointGMM for shape interpolation with both multi-layer perceptron (MLP) splits and attentional splits.

Unlike the above methods, our PCT is based on Transformer rather than using self-attention as an auxiliary module. While a framework by Wang and Solomon [29] uses Transformer to optimize point cloud registration, our PCT is a more general framework which can be used for various point cloud tasks.

### 3 Transformer for point cloud representation

In this section, we first show how the point cloud representation learned by our PCT can be applied to various tasks of point cloud processing, including point cloud classification, part segmentation, and normal estimation. Thereafter, we detail the design of PCT. We first introduce a naive version of PCT by directly applying the original Transformer [6] to point clouds. We then explain full PCT with its special attention mechanism, and neighbor aggregation to provide enhanced local information.

#### 3.1 Point cloud processing with PCT

**Encoder.** The overall architecture of PCT is presented in Fig. 2. PCT aims to transform (encode) the input points into a new higher dimensional feature space, which can characterize the semantic affinities between points as a basis for various point cloud processing tasks. The encoder of PCT starts by embedding the input coordinates into a new feature space. The embedded features are later fed into 4 stacked attention module to learn a semantically rich and discriminative representation for each point, followed by a linear layer to generate the output feature. Overall, the encoder of PCT shares almost the same philosophy of design as the original Transformer, except that the positional embedding is discarded, since the point's coordinates already contain this information. We refer the reader to Ref. [6] for details of the original NLP Transformer.

Formally, given an input point cloud  $\mathcal{P} \in \mathbb{R}^{N \times d}$  with  $N$  points each having a  $d$ -dimensional feature description, a  $d_e$ -dimensional embedded feature  $\mathbf{F}_e \in$

$\mathbb{R}^{N \times d_e}$  is first learned via the *Input Embedding* module. The point-wise  $d_o$ -dimensional feature representation  $\mathbf{F}_o \in \mathbb{R}^{N \times d_o}$  output by PCT is then formed by concatenating the attention output of each attention layer through the feature dimension, followed by a linear transformation:

$$\begin{aligned} \mathbf{F}_1 &= \text{AT}^1(\mathbf{F}_e) \\ \mathbf{F}_i &= \text{AT}^i(\mathbf{F}_{i-1}), \quad i = 2, 3, 4 \\ \mathbf{F}_o &= \text{concat}(\mathbf{F}_1, \mathbf{F}_2, \mathbf{F}_3, \mathbf{F}_4) \cdot \mathbf{W}_o \end{aligned} \quad (1)$$

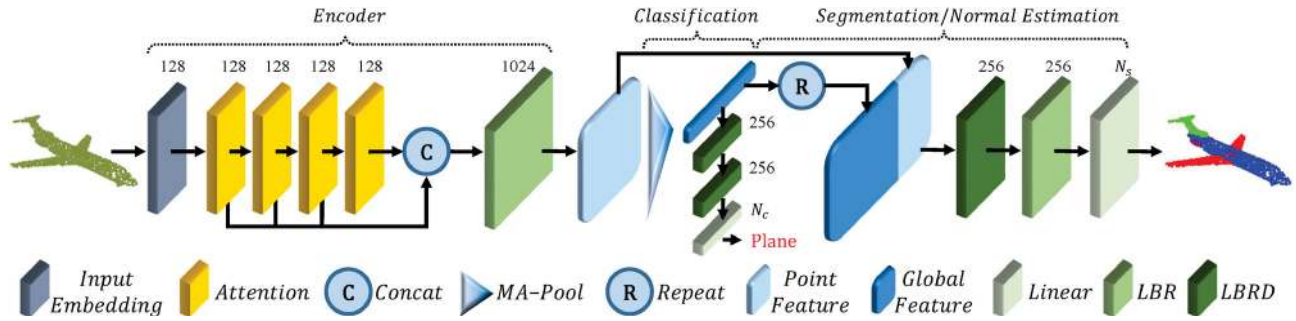
where  $\text{AT}^i$  represents the  $i$ th attention layer, each having the same output dimension as its input, and  $\mathbf{W}_o$  is the weights of the linear layer. Various implementations of input embedding and attention will be explained later.

To extract an effective global feature vector  $\mathbf{F}_g$  representing the point cloud, we choose to concatenate the outputs from two pooling operators: a max-pooling (MP) and an average-pooling (AP) on the learned point-wise feature representation [26].

**Classification.** The details of classification network using PCT is shown in Fig. 2. To classify a point cloud  $\mathcal{P}$  into  $N_c$  object categories (e.g., desk, table, chair), we feed the global feature  $\mathbf{F}_g$  to the classification decoder, which comprises two cascaded feed-forward neural networks LBRs (combining Linear, BatchNorm (BN), and ReLU layers) each with a dropout probability of 0.5, finalized by a Linear layer to predict the final classification scores  $\mathcal{C} \in \mathbb{R}^{N_c}$ . The class label of the point cloud is determined as the class with maximal score.

**Segmentation.** For the task of segmenting the point cloud into  $N_s$  parts (e.g., table top, table legs; a part need not be contiguous), we must predict a part label for each point, we first concatenate the global feature  $\mathbf{F}_g$  with the point-wise features in  $\mathbf{F}_o$ . To learn a common model for various kinds of objects, we also encode the one-hot object category vector as a 64-dimensional feature and concatenate it with the global feature, following most other point cloud segmentation networks [21]. As shown in Fig. 2, the architecture of the segmentation network decoder is almost the same as that for the classification network, except that dropout is only performed on the first LBR. We then predict the final point-wise segmentation scores  $\mathcal{S} \in \mathbb{R}^{N \times N_s}$  for the input point cloud: Finally, the part label of a point is also determined as the one with maximal score.





**Fig. 2** PCT architecture. The encoder mainly comprises an *Input Embedding* module and four stacked *Attention* module. The decoder mainly comprises multiple *Linear* layers. Numbers above each module indicate its output channels. *MA-Pool* concatenates *Max-Pool* and *Average-Pool*. *LBR* combines *Linear*, *BatchNorm*, and *ReLU* layers. *LBRD* means *LBR* followed by a *Dropout* layer.

**Normal estimation.** For the task of normal estimation, we use the same architecture as in segmentation by setting  $N_s = 3$ , without the object category encoding, and regard the output point-wise score as the predict normal.

### 3.2 Naive PCT

The simplest way to modify Transformer [6] for point cloud use is to treat the entire point cloud as a sentence and each point as a word, an approach we now explain. This naive PCT is achieved by implementing a coordinate-based point embedding and instantiating the attention layer with the self-attention introduced in Ref. [6].

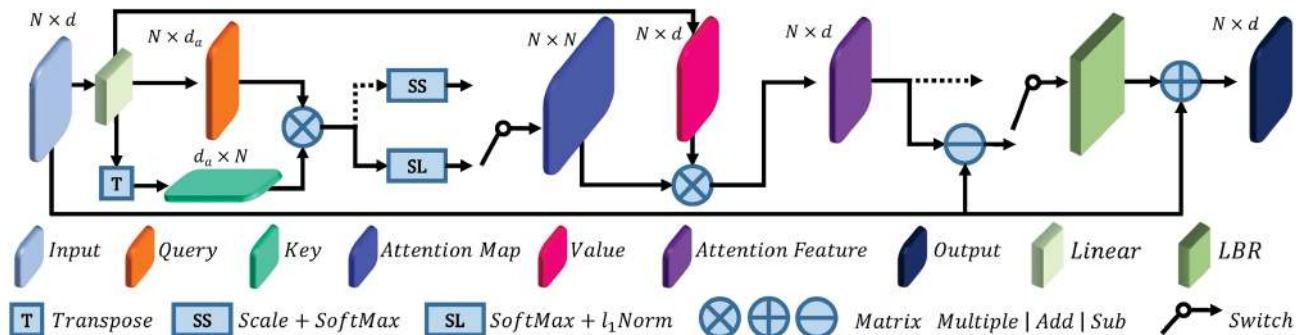
First, we consider a naive point embedding, which ignores interactions between points. Like word embedding in NLP, point embedding aims to place points closer in the embedding space if they are more semantically similar. Specifically, we embed a point cloud  $\mathcal{P}$  into a  $d_e$ -dimensional space  $\mathbf{F}_e \in \mathbb{R}^{N \times d_e}$ , using a shared neural network comprising two cascaded LBRs, each with a  $d_e$ -dimensional output.

We empirically set  $d_e = 128$ , a relatively small value, for computational efficiency. We simply use the point's 3D coordinates as its input feature description (i.e.,  $d_p = 3$ ) (as doing so still outperforms other methods) but additional point-wise input information, such as point normals, could also be used.

For the naive implementation of PCT, we adopt self-attention (SA) as introduced in the original Transformer [6]. Self-attention, also called intra-attention, is a mechanism that calculates semantic affinities between different items within a sequence of data. The architecture of the SA layer is depicted in Fig. 3 by switching to the dotted data flows. Following the terminology in Ref. [6], let  $\mathbf{Q}, \mathbf{K}, \mathbf{V}$  be the *query*, *key*, and *value* matrices, respectively, generated by linear transformations of the input features  $\mathbf{F}_{in} \in \mathbb{R}^{N \times d_e}$  as follows:

$$\begin{aligned} (\mathbf{Q}, \mathbf{K}, \mathbf{V}) &= \mathbf{F}_{in} \cdot (\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v) \\ \mathbf{Q}, \mathbf{K} &\in \mathbb{R}^{N \times d_a}, \quad \mathbf{V} \in \mathbb{R}^{N \times d_e} \\ \mathbf{W}_q, \mathbf{W}_k &\in \mathbb{R}^{d_e \times d_a}, \quad \mathbf{W}_v \in \mathbb{R}^{d_e \times d_e} \end{aligned} \quad (2)$$

where  $\mathbf{W}_q, \mathbf{W}_k$ , and  $\mathbf{W}_v$  are the shared learnable



**Fig. 3** Architecture of Offset-Attention. Numbers above tensors are numbers of dimensions  $N$  and feature channels  $D/D_a$ , with switches showing alternatives of Self-Attention or Offset-Attention: dotted lines indicate Self-Attention branches.

linear transformation, and  $d_a$  is the dimension of the query and key vectors. Note that  $d_a$  may not be equal to  $d_e$ . In this work, we set  $d_a$  to be  $d_e/4$  for computational efficiency.

First, we can use the query and key matrices to calculate the attention weights via the matrix dot-product:

$$\tilde{\mathbf{A}} = (\tilde{\alpha})_{i,j} = \mathbf{Q} \cdot \mathbf{K}^T \quad (3)$$

These weights are then normalized (denoted SS in Fig. 3) to give  $\mathbf{A} = (\alpha)_{i,j}$ :

$$\begin{aligned} \bar{\alpha}_{i,j} &= \frac{\tilde{\alpha}_{i,j}}{\sqrt{d_a}} \\ \alpha_{i,j} &= \text{softmax}(\bar{\alpha}_{i,j}) = \frac{\exp(\bar{\alpha}_{i,j})}{\sum_k \exp(\bar{\alpha}_{i,k})} \end{aligned} \quad (4)$$

The self-attention output features  $\mathbf{F}_{sa}$  are the weighted sums of the value vector using the corresponding attention weights:

$$\mathbf{F}_{sa} = \mathbf{A} \cdot \mathbf{V} \quad (5)$$

As the query, key, and value matrices are determined by the shared corresponding linear transformation matrices and the input feature  $\mathbf{F}_{in}$ , they are all order independent. Moreover, softmax and weighted sum are both permutation-independent operators. Therefore, the whole self-attention process is permutation-invariant, making it well-suited to the disordered, irregular domain presented by point clouds.

Finally, the self-attention feature  $\mathbf{F}_{sa}$  and the input feature  $\mathbf{F}_{in}$ , are further used to provide the output feature  $\mathbf{F}_{out}$  for the whole SA layer through an LBR network:

$$\mathbf{F}_{out} = \text{SA}(\mathbf{F}_{in}) = \text{LBR}(\mathbf{F}_{sa}) + \mathbf{F}_{in} \quad (6)$$

### 3.3 Offset-Attention

Graph convolution networks [9] show the benefits of using a Laplacian matrix  $\mathbf{L} = \mathbf{D} - \mathbf{E}$  to replace the adjacency matrix  $\mathbf{E}$ , where  $\mathbf{D}$  is the diagonal degree matrix. Similarly, we find that we can obtain better network performance if, when applying Transformer to point clouds, we replace the original self-attention (SA) module with an offset-attention (OA) module to enhance our PCT. As shown in Fig. 3, the offset-attention layer calculates the offset (difference) between the self-attention (SA) features and the input features by element-wise subtraction. This offset feeds the LBR network in place of the SA feature used in the naive version. Specifically, Eq. (5)

is modified to

$$\mathbf{F}_{out} = \text{OA}(\mathbf{F}_{in}) = \text{LBR}(\mathbf{F}_{in} - \mathbf{F}_{sa}) + \mathbf{F}_{in} \quad (7)$$

$\mathbf{F}_{in} - \mathbf{F}_{sa}$  is analogous to a discrete Laplacian operator, as we now show. First, from Eqs. (2) and (5), the following holds:

$$\begin{aligned} \mathbf{F}_{in} - \mathbf{F}_{sa} &= \mathbf{F}_{in} - \mathbf{A}\mathbf{V} \\ &= \mathbf{F}_{in} - \mathbf{A}\mathbf{F}_{in}\mathbf{W}_v \\ &\approx \mathbf{F}_{in} - \mathbf{A}\mathbf{F}_{in} \\ &= (\mathbf{I} - \mathbf{A})\mathbf{F}_{in} \approx \mathbf{L}\mathbf{F}_{in} \end{aligned} \quad (8)$$

Here,  $\mathbf{W}_v$  is ignored since it is a weight matrix of the *Linear layer*.  $\mathbf{I}$  is an identity matrix comparable to the diagonal degree matrix  $\mathbf{D}$  of the Laplacian matrix and  $\mathbf{A}$  is the attention matrix comparable to the adjacency matrix  $\mathbf{E}$ .

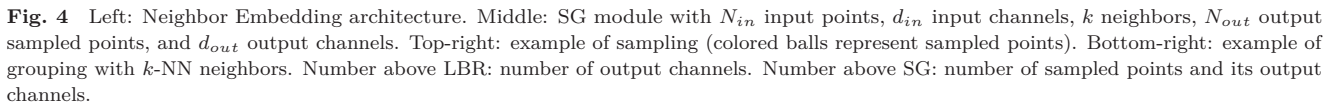
In our enhanced version of PCT, we also refine the normalization by modifying Eq. (4) as follows:

$$\begin{aligned} \bar{\alpha}_{i,j} &= \text{softmax}(\tilde{\alpha}_{i,j}) = \frac{\exp(\tilde{\alpha}_{i,j})}{\sum_k \exp(\tilde{\alpha}_{k,j})} \\ \alpha_{i,j} &= \frac{\bar{\alpha}_{i,j}}{\sum_k \bar{\alpha}_{i,k}} \end{aligned} \quad (9)$$

Here, we use the softmax operator on the first dimension and an  $l_1$ -norm for the second dimension to normalize the attention map. The traditional Transformer scales the first dimension by  $1/\sqrt{d_a}$  and uses softmax to normalize the second dimension. However, our offset-attention sharpens the attention weights and reduces the influence of noise, which is beneficial for downstream tasks. Figure 1 shows example offset attention maps. It can be seen that the attention maps for different query points vary considerably, but are generally semantically meaningful. We refer to this refined PCT, i.e., with point embedding and OA layer, as simple PCT (SPCT) in the experiments.

### 3.4 Neighbor embedding for augmented local feature representation

PCT with point embedding is an effective network for extracting global features. However, it ignores the local neighborhood information which is also essential in point cloud learning. We draw upon the ideas of PointNet++ [21] and DGCNN [26] to design a local neighbor aggregation strategy, *neighbor embedding*, to optimize the point embedding to augment PCT's ability of local feature extraction. As shown in Fig. 4, neighbor embedding module comprises two LBR layers and two SG (sampling and grouping) layers. The LBR layers act as the basis point embedding in Section 3.2. We use two cascaded SG



More specifically, assume that SG layer takes a point cloud  $\mathcal{P}$  with  $N$  points and corresponding features  $\mathbf{F}$  as input and outputs a sampled point cloud  $\mathcal{P}_s$  with  $N_s$  points and its corresponding aggregated features  $\mathbf{F}_s$ . First, We adopt the farthest point sampling (FPS) algorithm [21] to downsample  $\mathcal{P}$  to  $\mathcal{P}_s$ . Then, for each sampled point  $p \in \mathcal{P}_s$ , let  $\text{knn}(p, \mathcal{P})$  be its  $k$ -nearest neighbors in  $\mathcal{P}$ . We then compute the output feature  $\mathbf{F}_s$  as follows:

where  $\mathbf{F}(p)$  is the input feature of point  $p$ ,  $\mathbf{F}_s(p)$  is the output feature of sampled point  $p$ , MP is the max-pooling operator, and  $\text{RP}(\mathbf{x}, k)$  is the operator for repeating a vector  $\mathbf{x}$   $k$  times to form a matrix. The idea of concatenating the feature among sampled point and its neighbors is drawn from *EdgeConv* [26].

For point cloud segmentation or normal estimation, we need to determine point-wise part labels or normal, so the process above is only used for local feature extraction without reducing the point cloud size.

We now evaluate the performance of naive PCT (NPCT, with point embedding and self-attention), simple PCT (SPCT, with point embedding and offset-attention), and full PCT (with neighbor embedding and offset-attention) on two public datasets, ModelNet40 [30] and ShapeNet [31], giving a comprehensive comparison with other methods. The same soft cross-entropy loss function as Ref. [26] and the stochastic gradient descent (SGD) optimizer with momentum 0.9 were adopted for training in each case. Other training parameters, including the learning rate, batch size, and input format, were particular to each specific dataset and are given later.

ModelNet40 [30] contains 12,311 CAD models in 40 object categories; it is widely used in point cloud shape classification and surface normal estimation benchmarking. For a fair comparison, we used the official split with 9843 objects for training and 2468 for evaluation. The same sampling strategy as used in PointNet [1] was adopted to uniformly sample each object to 1024 points. During training, a random translation in  $[-0.2, 0.2]$ , a random anisotropic scaling in  $[0.67, 1.5]$ , and a random input dropout were applied to augment the input data. During testing, no data augmentation or voting methods were used. For all the three models, the mini-batch sizes were 32. 250 training epochs were used and the initial learning

rates were 0.0001, with a cosine annealing schedule to adjust the learning rate at every epoch.

Experimental results are shown in Table 1. Compared to PointNet and NPCT, SPCT makes a 2.8% and 1.0% improvement respectively. PCT achieves the best result of 93.2% overall accuracy. Note that our network currently does not consider normals as inputs which could in principle further improve network performance.

#### 4.2 Normal estimation on ModelNet40 dataset

The surface normal estimation is to determine the normal direction at each point. Estimating surface normal has wide applications in, e.g., rendering. The task is challenging because it requires the approach to understand the shapes completely for dense regression. We again used ModelNet40 as a benchmark, and used average cosine distance to measure the difference between ground truth and predicted normals. For all the three models, a batch size of 32,200 training epochs were used. The initial learning rates were also set as 0.01, with a cosine annealing schedule used to adjust learning rate every epoch. As indicated in Table 2, both our NPCT and SPCT make a

**Table 1** Comparison with state-of-the-art methods on the ModelNet40 classification dataset. Accuracy means overall accuracy. All results quoted are taken from the cited papers. P = points, N = normals

Method	Input	#Points	Accuracy
PointNet [1]	P	1k	89.2%
A-SCN [32]	P	1k	89.8%
SO-Net [33]	P, N	2k	90.9%
Kd-Net [34]	P	32k	91.8%
PointNet++ [21]	P	1k	90.7%
PointNet++ [21]	P, N	5k	91.9%
PointGrid [35]	P	1k	92.0%
PCNN [4]	P	1k	92.3%
PointWeb [36]	P	1k	92.3%
PointCNN [3]	P	1k	92.5%
PointConv [5]	P, N	1k	92.5%
A-CNN [37]	P, N	1k	92.6%
P2Sequence [38]	P	1k	92.6%
KPCConv [39]	P	7k	92.9%
DGCNN [26]	P	1k	92.9%
RS-CNN [40]	P	1k	92.9%
PointASNL [27]	P	1k	92.9%
NPCT	P	1k	91.0%
SPCT	P	1k	92.0%
PCT	P	1k	<b>93.2%</b>

**Table 2** Normal estimation average cosine-distance error on ModelNet40 dataset

Method	#Points	Error
PointNet [1]	1k	0.47
PointNet++ [21]	1k	0.29
PCNN [4]	1k	0.19
RS-CNN [40]	1k	0.15
NPCT	1k	0.24
SPCT	1k	0.23
PCT	1k	<b>0.13</b>

significant improvement compared with PointNet and PCT achieves the lowest average cosine distance.

#### 4.3 Part segmentation task on ShapeNet dataset

Point cloud part segmentation is a challenging task which aims to divide a 3D model into multiple meaningful parts. We performed an experimental evaluation on the ShapeNet Parts dataset [31], which contains 16,880 3D models with a training to testing split of 14,006 to 2,874. It has 16 object categories and 50 part labels; each instance contains no fewer than two parts. Following PointNet [1], all models were downsampled to 2,048 points, retaining point-wise part annotation. During training, random translation in  $[-0.2, 0.2]$ , and random anisotropic scaling in  $[0.67, 1.5]$  were applied to augment the input data. During testing, we used a multi-scale testing strategy, where the scales are set in  $[0.7, 1.5]$  with a step of 0.1. For all the three models, the batch size, training epochs, and the learning rates were set the same as the training of normal estimation task.

Table 3 shows the class-wise segmentation results. The evaluation metric used is part-average Intersection-over-Union, and is given both overall and for each object category. The results show that our SPCT makes an improvement of 2.1% and 0.6% over PointNet and NPCT respectively. PCT achieves the best results with 86.4% part-average Intersection-over-Union. Figure 5 shows further segmentation examples provided by PointNet, NPCT, SPCT, and PCT.

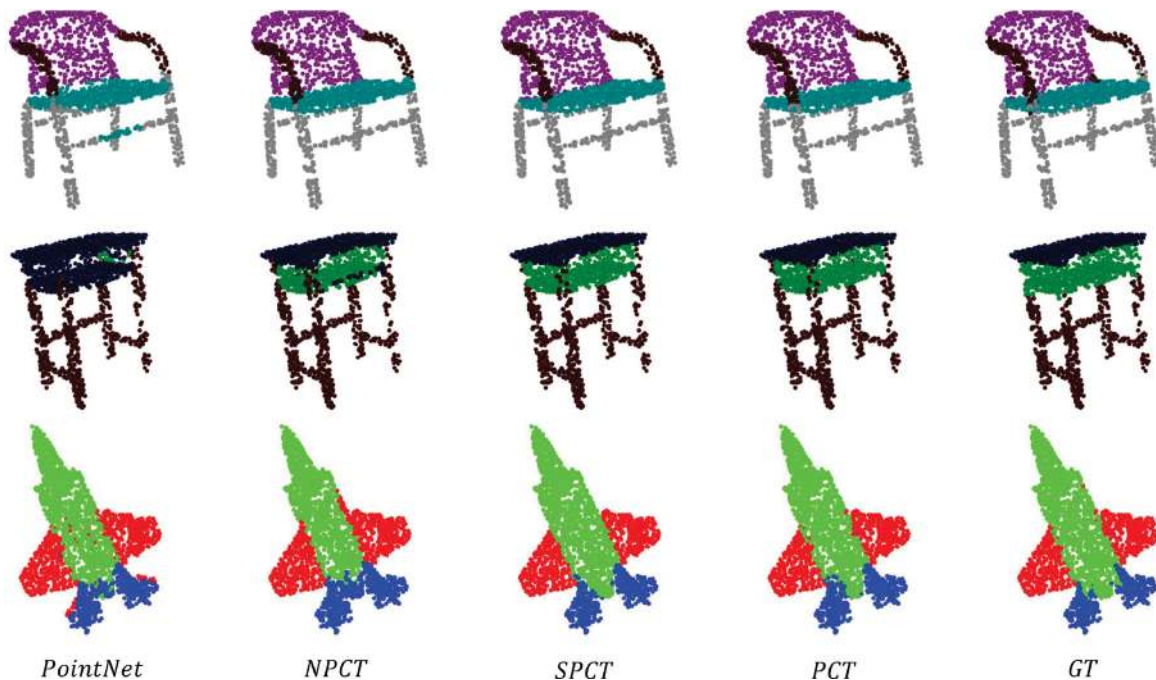
#### 4.4 Semantic segmentation task on S3DIS dataset

The S3DIS is a indoor scene dataset for point cloud semantic segmentation. It contains 6 areas and 271 rooms. Each point in the dataset is divided into 13



**Table 3** Comparison on the ShaperNet part segmentation dataset. pIoU means part-average Intersection-over-Union. All results quoted are taken from the cited papers

Method	pIoU	air-plane	bag	cap	car	chair	ear-phone	guitar	knife	lamp	laptop	motor-bike	mug	pistol	rocket	skate-board	table
PointNet [1]	83.7	83.4	78.7	82.5	74.9	89.6	73.0	91.5	85.9	80.8	95.3	65.2	93.0	81.2	57.9	72.8	80.6
Kd-Net [34]	82.3	80.1	74.6	74.3	70.3	88.6	73.5	90.2	87.2	81.0	94.9	57.4	86.7	78.1	51.8	69.9	80.3
SO-Net [33]	84.9	82.8	77.8	88.0	77.3	90.6	73.5	90.7	83.9	82.8	94.8	69.1	94.2	80.9	53.1	72.9	83.0
PointNet++ [21]	85.1	82.4	79.0	87.7	77.3	90.8	71.8	91.0	85.9	83.7	95.3	71.6	94.1	81.3	58.7	76.4	82.6
PCNN [4]	85.1	82.4	80.1	85.5	79.5	90.8	73.2	91.3	86.0	85.0	95.7	73.2	94.8	83.3	51.0	75.0	81.8
DGCNN [26]	85.2	84.0	83.4	86.7	77.8	90.6	74.7	91.2	87.5	82.8	95.7	66.3	94.9	81.1	63.5	74.5	82.6
P2Sequence [38]	85.2	82.6	81.8	87.5	77.3	90.8	77.1	91.1	86.9	83.9	95.7	70.8	94.6	79.3	58.1	75.2	82.8
PointConv [5]	85.7	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
PointCNN [3]	86.1	84.1	<b>86.5</b>	86.0	80.8	90.6	79.7	<b>92.3</b>	<b>88.4</b>	85.3	<b>96.1</b>	<b>77.2</b>	95.2	<b>84.2</b>	<b>64.2</b>	<b>80.0</b>	83.0
PointASNL [27]	86.1	84.1	84.7	87.9	79.7	<b>92.2</b>	73.7	91.0	87.2	84.2	95.8	74.4	95.2	81.0	63.0	76.3	83.2
RS-CNN [40]	86.2	83.5	84.8	88.8	79.6	91.2	<b>81.1</b>	91.6	<b>88.4</b>	86.0	96.0	73.7	94.1	83.4	60.5	77.7	83.6
<b>NPCT</b>	85.2	83.2	74.5	86.7	76.8	90.7	75.4	91.1	87.3	84.5	95.7	65.2	93.7	82.7	56.9	73.8	83.0
<b>SPCT</b>	85.8	84.5	83.5	85.9	78.7	90.9	75.1	92.1	87.0	85.0	95.9	69.6	94.5	82.2	61.4	76.0	83.0
<b>PCT</b>	<b>86.4</b>	<b>85.0</b>	82.4	<b>89.0</b>	<b>81.2</b>	91.9	71.5	91.3	88.1	<b>86.3</b>	95.8	64.6	<b>95.8</b>	83.6	62.2	77.6	<b>83.7</b>

**Fig. 5** Segmentations from PointNet, NPCT, SPCT, PCT, and ground truth (GT).**Table 4** Comparison on the S3DIS semantic segmentation dataset tested on Area5

Method	mAcc	mIoU	ceiling	floor	wall	beam	column	window	door	chair	table	book-case	sofa	board	clutter
PointNet [1]	48.98	41.09	88.80	97.33	69.80	0.05	3.92	46.26	10.76	58.93	52.61	5.85	40.28	26.38	33.22
SEGCloud [2]	57.35	48.92	90.06	96.05	69.86	0.00	18.37	38.35	23.12	70.40	75.89	40.88	58.42	12.96	41.60
DGCNN [26]	84.10	56.10	—	—	—	—	—	—	—	—	—	—	—	—	—
PointCNN [3]	63.86	57.26	92.31	98.24	79.41	0.00	17.60	22.77	62.09	74.39	80.59	31.67	66.67	62.05	56.74
SPG [24]	66.50	58.04	89.35	96.87	78.12	0.00	42.81	48.93	61.58	84.66	75.41	69.84	52.60	2.10	52.22
PCNN [4]	67.01	58.27	92.26	96.20	75.89	0.27	5.98	69.49	63.45	66.87	65.63	47.28	68.91	59.10	46.22
PointWeb [36]	66.64	60.28	91.95	98.48	79.39	0.00	21.11	59.72	34.81	76.33	88.27	46.89	69.30	64.91	52.46
<b>PCT</b>	<b>67.65</b>	<b>61.33</b>	92.54	98.42	80.62	0.00	19.37	61.64	48.00	76.58	85.20	46.22	67.71	67.93	52.29

categories. For fair comparison, we use the same data processing method as Ref. [1]. Table 4 shows that our PCT achieves superior performance compared to the previous methods.

#### 4.5 Computational requirements analysis

We now consider the computational requirements of NPCT, SPCT, PCT, and several other methods by comparing the floating point operations required (FLOPs) and number of parameters (Params) in Table 5. SPCT has the lowest memory requirements with only 1.36M parameters and also puts a low load on the processor of only 1.82G FLOPs, yet delivers highly accurate results. These characteristics make it suitable for deployment on a mobile device. PCT has best performance, yet modest computational and memory requirements. If we pursue higher performance and ignore the amount of calculation and parameters, we can add a neighbor embedding layer in the input embedding module. The results of 3-layer embedding PCT are shown in Tables 6 and 7.

**Table 5** Computational resource requirements

Method	#Params	#FLOPs	Accuracy
PointNet [1]	3.47M	<b>0.45G</b>	89.2%
PointNet++(SSG) [21]	1.48M	1.68G	90.7%
PointNet++(MSG) [21]	1.74M	4.09G	91.9%
DGCNN [26]	1.81M	2.43G	92.9%
NPCT	<b>1.36M</b>	1.80G	91.0%
SPCT	<b>1.36M</b>	1.82G	92.0%
PCT	2.88M	2.32G	<b>93.2%</b>

**Table 6** Comparison on the ModelNet40 classification dataset. PCT-2L means PCT with 2 layer neighbor embedding and PCT-3L means PCT with 3 layer neighbor embedding. Accuracy means overall accuracy. P = points

Method	Input	#Points	Accuracy
PCT-2L	P	1k	93.2%
PCT-3L	P	1k	<b>93.4%</b>

**Table 7** Comparison on the ShaperNet part segmentation dataset. pIoU means part-average Intersection-over-Union. PCT-2L means PCT with 2 layer neighbor embedding and PCT-3L means PCT with 3 layer neighbor embedding

Method	pIoU	air-plane	bag	cap	car	chair	ear-phone	guitar	knife	lamp	laptop	motor-bike	mug	pistol	rocket	skate-board	table
<b>PCT-2L</b>	86.4	85.0	82.4	89.0	<b>81.2</b>	<b>91.9</b>	71.5	91.3	<b>88.1</b>	<b>86.3</b>	95.8	64.6	<b>95.8</b>	<b>83.6</b>	<b>62.2</b>	<b>77.6</b>	<b>83.7</b>
<b>PCT-3L</b>	<b>86.6</b>	<b>85.3</b>	<b>84.5</b>	<b>89.4</b>	81.0	91.7	<b>78.6</b>	<b>91.5</b>	87.5	85.8	<b>96.0</b>	<b>70.6</b>	95.6	82.8	60.9	76.6	<b>83.7</b>

## 5 Conclusions

In this paper, we propose a permutation-invariant point cloud transformer, which is suitable for learning on unstructured point clouds with irregular domain. The proposed offset-attention and normalization mechanisms help to make our PCT effective. Experiments show that PCT has good semantic feature learning capability, and achieves state-of-the-art performance on several tasks, particularly shape classification, part segmentation, and normal estimation.

Transformer has already revealed powerful capabilities given large amounts of training data. At present, the available point cloud datasets are very limited compared to image. In future, we will train it on larger datasets and study its advantages and disadvantages with respect to other popular frameworks. The encoder-decoder structure of Transformer supports more complex tasks, such as point cloud generation and completion. We will extend the PCT to further applications. Besides, we will attempt more precise methods to approximate Laplacian operation and complete offset-attention.

## Acknowledgements

This work was supported by the National Natural Science Foundation of China (Project Number 61521002) and the Joint NSFC-DFG Research Program (Project Number 61761136018).

## References

- [1] Charles, R. Q.; Hao, S.; Mo, K. C.; Guibas, L. J. PointNet: Deep learning on point sets for 3D classification and segmentation. IN: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 77–85, 2017.
- [2] Tchapmi, L. P.; Choy, C. B.; Armeni, I.; Gwak, J.; Savarese, S. SEGCloud: Semantic segmentation of 3D point clouds. In: Proceedings of the International Conference on 3D Vision, 537–547, 2017.

- [3] Li, Y.; Bu, R.; Sun, M.; Wu, W.; Di, X.; Chen, B. PointCNN: Convolution on x-transformed points. In: Proceedings of the 32nd International Conference on Neural Information Processing Systems, 828–838, 2018.
- [4] Atzmon, M.; Maron, H.; Lipman, Y. Point convolutional neural networks by extension operators. *ACM Transactions on Graphics* Vol. 37, No. 4, Article No. 71, 2018.
- [5] Wu, W. X.; Qi, Z.; Fuxin, L. PointConv: Deep convolutional networks on 3D point clouds. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 9613–9622, 2019.
- [6] Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; Polosukhin, I. Attention is all you need. In: Proceedings of the 31st International Conference on Neural Information Processing, 6000–6010, 2017.
- [7] Dosovitskiy, A.; Beyer, L.; Kolesnikov, A.; Weissenborn, D.; Houlsby, N. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [8] Wu, B.; Xu, C.; Dai, X.; Wan, A.; Zhang, P.; Tomizuka, M.; Keutzer, K.; Vajda, P. Visual transformers: Token-based image representation and processing for computer vision. *arXiv preprint arXiv:2006.03677*, 2020.
- [9] Bruna, J.; Zaremba, W.; Szlam, A.; LeCun, Y. Spectral networks and locally connected networks on graphs. In: Proceedings of the International Conference on Learning Representations, 2014.
- [10] Hu, S.-M.; Liang, D.; Yang, G.-Y.; Yang, G.-W.; Zhou, W.-Y. Jittor: A novel deep learning framework with meta-operators and unified graph execution. *Science China Information Sciences* Vol. 63, No. 12, Article No. 222103, 2020.
- [11] Bahdanau, D.; Cho, K. H.; Bengio, Y. Neural machine translation by jointly learning to align and translate. In: Proceedings of the 3rd International Conference on Learning Representations, 2015.
- [12] Lin, Z.; Feng, M.; dos Santos, C. N.; Yu, M.; Xiang, B.; Zhou, B.; Bengio, Y. A structured self-attentive sentence embedding. In: Proceedings of the International Conference on Learning Representations, 2017.
- [13] Devlin, J.; Chang, M.; Lee, K.; Toutanova, K. BERT: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Vol. 1, 4171–4186, 2019.
- [14] Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J. G.; Salakhutdinov, R.; Le, Q. V. XLNet: Generalized autoregressive pretraining for language understanding. In: Proceedings of the 33rd Conference on Neural Information Processing Systems, 5754–5764, 2019.
- [15] Dai, Z. H.; Yang, Z. L.; Yang, Y. M.; Carbonell, J.; Le, Q.; Salakhutdinov, R. Transformer-XL: Attentive language models beyond a fixed-length context. In: Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 2978–2988, 2019.
- [16] Lee, J.; Yoon, W.; Kim, S.; Kim, D.; Kim, S.; So, C. H.; Kang, J. BioBERT: A pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* Vol. 36, No. 4, 1234–1240, 2020.
- [17] Wang, F.; Jiang, M. Q.; Qian, C.; Yang, S.; Li, C.; Zhang, H. G.; Wang, X.; Tang, X. Residual attention network for image classification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 6450–6458, 2017.
- [18] Hu, J.; Shen, L.; Sun, G. Squeeze-and-excitation networks. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 7132–7141, 2018.
- [19] Zhang, H.; Goodfellow, I. J.; Metaxas, D. N.; Odena, A. Self-attention generative adversarial networks. In: Proceedings of the International Conference on Machine Learning, 7354–7363, 2019.
- [20] Carion, N.; Massa, F.; Synnaeve, G.; Usunier, N.; Kirillov, A.; Zagoruyko, S. End-to-end object detection with transformers. In: *Computer Vision – ECCV 2020. Lecture Notes in Computer Science, Vol. 12346*. Vedaldi, A.; Bischof, H.; Brox, T.; Frahm, J. M. Eds. Springer Cham, 213–229, 2020.
- [21] Qi, C. R.; Yi, L.; Su, H.; Guibas, L. J. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In: Proceedings of the 31st Conference on Neural Information Processing Systems, 5099–5108, 2017.
- [22] Hermosilla, P.; Ritschel, T.; Vázquez, P. P.; Vinacua, À.; Ropinski, T. Monte Carlo convolution for learning on non-uniformly sampled point clouds. *ACM Transactions on Graphics* Vol. 37, No. 6, Article No. 235, 2018.
- [23] Tatarchenko, M.; Park, J.; Koltun, V.; Zhou, Q. Y. Tangent convolutions for dense prediction in 3D. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 3887–3896, 2018.

- [24] Landrieu, L.; Simonovsky, M. Large-scale point cloud semantic segmentation with superpoint graphs. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 4558–4567, 2018.
- [25] Yang, Y. Q.; Liu, S. L.; Pan, H.; Liu, Y.; Tong, X. PFCNN: Convolutional neural networks on 3D surfaces using parallel frames. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 13575–13584, 2020.
- [26] Wang, Y.; Sun, Y.; Liu, Z.; Sarma, S. E.; Bronstein, M. M.; Solomon, J. M. Dynamic graph CNN for learning on point clouds. *ACM Transactions on Graphics* Vol. 38, No. 5, Article No. 146, 2019.
- [27] Yan, X.; Zheng, C. D.; Li, Z.; Wang, S.; Cui, S. G. PointASNL: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 5588–5597, 2020.
- [28] Hertz, A.; Hanocka, R.; Giryas, R.; Cohen-Or, D. PointGMM: A neural GMM network for point clouds. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 12051–12060, 2020.
- [29] Wang, Y.; Solomon, J. Deep closest point: Learning representations for point cloud registration. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 3522–3531, 2019.
- [30] Wu, Z.; Song, S.; Khosla, A.; Yu, F.; Zhang, L.; Tang, X.; Xiao, J. 3D ShapeNets: A deep representation for volumetric shapes. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 1912–1920, 2015.
- [31] Yi, L.; Kim, V. G.; Ceylan, D.; Shen, I. C.; Yan, M. Y.; Su, H.; Lu, C.; Huang, Q.; Sheffer, A.; Guibas, L. A scalable active framework for region annotation in 3D shape collections. *ACM Transactions on Graphics* Vol. 35, No. 6, Article No. 210, 2016.
- [32] Xie, S. N.; Liu, S. N.; Chen, Z. Y.; Tu, Z. W. Attentional ShapeContextNet for point cloud recognition. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 4606–4615, 2018.
- [33] Li, J. X.; Chen, B. M.; Lee, G. H. SO-net: Self-organizing network for point cloud analysis. In: Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 9397–9406, 2018.
- [34] Klovov, R.; Lempitsky, V. Escape from cells: Deep kd-networks for the recognition of 3D point cloud models. In: Proceeding of the IEEE International Conference on Computer Vision, 863–872, 2017.
- [35] Le, T.; Duan, Y. PointGrid: A deep network for 3D shape understanding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 9204–9214, 2018.
- [36] Zhao, H.; Jiang, L.; Fu, C.; Jia, J. PointWeb: Enhancing local neighborhood features for point cloud processing. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 5560–5568, 2019.
- [37] Komarichev, A.; Zhong, Z. C.; Hua, J. A-CNN: Annularly convolutional neural networks on point clouds. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 7413–7422, 2019.
- [38] Liu, X. H.; Han, Z. Z.; Liu, Y. S.; Zwicker, M. Point2Sequence: Learning the shape representation of 3D point clouds with an attention-based sequence to sequence network. In: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 8778–8785, 2019.
- [39] Thomas, H.; Qi, C. R.; Deschaud, J. E.; Marcotegui, B.; Goulette, F.; Guibas, L. KPConv: Flexible and deformable convolution for point clouds. In: Proceedings of the IEEE/CVF International Conference on Computer Vision, 6410–6419, 2019.
- [40] Liu, Y. C.; Fan, B.; Xiang, S. M.; Pan, C. H. Relation-shape convolutional neural network for point cloud analysis. In: Proceeding of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 8887–8896, 2019.



**Meng-Hao Guo** received his bachelor degree in Xidian University. Now he is a Ph.D. candidate in the Department of Computer Science and Technology, Tsinghua University. His research interests include computer graphics, computer vision, and machine learning.



**Jun-Xiong Cai** is currently a postdoctoral researcher at Tsinghua University, where he received Ph.D. degree in computer science and technology in 2020. His research interests include computer graphics, computer vision, and 3D geometry processing.





**Zheng-Ning Liu** received his bachelor degree in computer science from Tsinghua University in 2017. He is currently a Ph.D. candidate in the Department of Computer Science and Technology, Tsinghua University. His research interests include 3D computer vision, 3D reconstruction, and computer

graphics.



**Tai-Jiang Mu** is currently an assistant researcher at Tsinghua University, where he received his B.S. and Ph.D. degrees in computer science and technology in 2011 and 2016, respectively. His research interests include computer vision, robotics, and computer graphics.



**Ralph R. Martin** received his Ph.D. degree from Cambridge University in 1983. He is currently an emeritus professor with Cardiff University. He has authored over 250 papers and 14 books, covering such topics as solid and surface modeling, intelligent sketch input, geometric reasoning, reverse engineering,

and various aspects of computer graphics. He is a Fellow of the Learned Society of Wales, the Institute of Mathematics and its Applications, and the British Computer Society. He is currently the Associate Editor-in-Chief of *Computational Visual Media*.



**Shi-Min Hu** is currently a professor in the Department of Computer Science and Technology, Tsinghua University, Beijing, China. He received his Ph.D. degree from Zhejiang University in 1996. His research interests include digital geometry processing, video processing, rendering, computer animation, and computer-aided geometric design. He has published more than 100 papers in journals and refereed conferences. He is the Editor-in-Chief of *Computational Visual Media*, and on editorial boards of several journals, including *Computer Aided Design* and *Computer & Graphics*. He is a senior member of IEEE and ACM, and Fellow of CCF and SMA.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made.

The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

Other papers from this open access journal are available free of charge from <http://www.springer.com/journal/41095>. To submit a manuscript, please go to <https://www.editorialmanager.com/cvmj>.