

Shape-Oriented Convolution Neural Network for Point Cloud Analysis

Chaoyi Zhang,¹ Yang Song,² Lina Yao,² Weidong Cai¹

¹ School of Computer Science, University of Sydney, Australia

² School of Computer Science and Engineering, University of New South Wales, Australia
 {chaoyi.zhang, tom.cai}@sydney.edu.au, {yang.song1, lina.yao}@unsw.edu.au

Abstract

Point cloud is a principal data structure adopted for 3D geometric information encoding. Unlike other conventional visual data, such as images and videos, these irregular points describe the complex shape features of 3D objects, which makes shape feature learning an essential component of point cloud analysis. To this end, a shape-oriented message passing scheme dubbed *ShapeConv* is proposed to focus on the representation learning of the underlying shape formed by each local neighboring point. Despite this *intra-shape relationship* learning, *ShapeConv* is also designed to incorporate the contextual effects from the *inter-shape relationship* through capturing the long-ranged dependencies between local underlying shapes. This shape-oriented operator is stacked into our hierarchical learning architecture, namely Shape-Oriented Convolutional Neural Network (SOCNN), developed for point cloud analysis. Extensive experiments have been performed to evaluate its significance in the tasks of point cloud classification and part segmentation.

Introduction

As a principal data structure adopted for 3D geometric information encoding, point cloud has been widely used in several practical applications, such as self-driving cars and computer graphics. Although geometrical and topological information can be well-described as the raw point coordinates in point cloud data, the further analysis step of these irregular points can be quite challenging, as the local underlying shapes may not be modeled or recognized appropriately. Following the significant success recently achieved by Convolution Neural Networks (CNNs) on regular-formatted visual data, such as images and videos, several attempts have been made to transform raw point cloud data to either 3D volumetric representations or a collection of 2D views, so that it can be handled directly by the CNNs defined on regular grids.

However, these approaches all have their own drawbacks and limitations. For the voxelization-based approaches, the potential information loss occurred in the voxelization stage may irrupt the object shapes by introducing quantization er-

ror. For the view-based methods, although an accurate classification or a descent segmentation can be reached, it requires a large number of view images collected from different angles, to make sure the generated 2D projections contain enough discriminative representations of the 3D objects for further point cloud analysis. Thus, a more shape-oriented geometric learning approach needs to be developed, which is not only expected to be able to manipulate the point cloud data directly, but is also designed to be capable of understanding the discriminative information of each local underlying shape, by modeling the shape-oriented contextual information encoded in the point cloud.

PointNet (Qi et al. 2017a) was the first deep learning based approach to manipulate the point cloud data directly, learning the pointwise features independently and outputting a global shape signature from the symmetric aggregation function applied to these pointwise features. Following the PointNet, the entire community started to pay more attention to local point cloud structure modeling, which was neglected completely by PointNet. DGCNN (Wang et al. 2019) could be seen as the next milestone, as they generally define *EdgeConv* as:

$$\mathbf{x}_i^k = \square_{j:(i,j) \in \mathcal{E}} h_{\Theta}(\mathbf{x}_i^{k-1}, \mathbf{x}_j^{k-1}), \quad (1)$$

where k is the network layer index, X_i and X_j indicate the central point and the neighboring points, \mathcal{E} denotes the local graph dynamically constructed among these points, and \square is the feature aggregation function covering the entire local neighbourhood. As one of their main contributions, the *EdgeConv* proposed helps to reformulate the learning process of the geometric information as the aggregation outcomes from the edge features computed pairwise. It indicates that to update the pointwise feature X_i , the dynamic local graphs are firstly constructed using k-nearest neighbors searching technique, then centralized connections can be built between the centroid point and its neighbours, and the edge features $h_{\Theta}(\cdot, \cdot)$ are computed for the updating of the pointwise feature for the centroid points. The definition of *EdgeConv* reveals the idea that most of the authors claim that the local underlying shapes can be learnt as the aggregated pairwise interactions between the centroid point and its neighbor points.

Unlike their approaches, we propose the shape-oriented

convolution (*ShapeConv*) to link all points forming the local shape to enhance their individual interactions and compute the moment point of this densely-connected local graph. After that, we show that the overall interaction between each point and the local shape can be simplified as the pointwise interaction placed between each point and the moment point. Furthermore, we study the contextual information encoding, via defining and modeling two shape-oriented relationships for point clouds, namely the intra-shape relationship and the inter-shape relationship. By incorporating this contextual information, our *ShapeConv* module proposed is capable of performing several advanced point cloud analyses. To this end, we stack *ShapeConv* into our shape-oriented convolution neural network (*SOCNN*) and evaluate its significance in the tasks of point cloud classification and part segmentation.

Related Work

In this section, we will mainly review the previous works performed toward point cloud analysis, from three perspectives: the voxelization-based approaches, the view-based approaches, and the geometric learning approaches.

Voxelization-based Approaches Voxelization is a particular kind of transformation, which takes irregular point cloud data as input and transforms them into several volumetric objects represented under a regular 3D coordinate system. Benefiting from the new volumetric representations, point cloud data can thus be processed conveniently by the 3D Convolution operators defined on regular 3D grids. Hence, several advanced point cloud analyses, such as the classification and segmentation tasks, could be easily achieved by the CNN, whose discriminative capabilities have been broadly evaluated in the computer vision community (Krizhevsky, Sutskever, and Hinton 2012; He et al. 2016; Huang et al. 2017; Hu, Shen, and Sun 2018). However, the performance of these voxelization-based approaches (Wu et al. 2015; Maturana and Scherer 2015) to point cloud analysis would be largely constrained by the critical shape information loss during the quantization step of the voxelization procedure. Although several subvolume-related works have been proposed to alleviate this spatial information loss (Klokov and Lempitsky 2017; Wang et al. 2017; Riegler, Ulusoy, and Geiger 2017), their resulting subdivision representations can still be suffering from the potential quantization error, compared to other approaches modeling the local underlying shapes directly.

View-based Approaches Rather than representing point cloud data from 3D regular grids, like the voxelization-based approaches mentioned above, view-based approaches are focused on recognizing and analyzing the point cloud data, from collections of 2D views. Due to the promising results achieved by 2D CNNs, view-based approaches can achieve excellent outcomes for point cloud analysis (Su et al. 2015; Qi et al. 2016; Xie et al. 2016). However, their side effects should be taken into consideration as well. That is, to reach an accurate classification or a descent segmentation, it requires a large number of views to be used for model training process. Meanwhile, to ensure the information encoded in

the views are discriminative enough for the shape recognition, these views should be captured from as many different angles as possible (Feng et al. 2018), which results in a long rendering time needed for the data sampling stage.

Geometric Learning Approaches In contrast to the approaches above, geometric learning approaches have been developed recently, with the aim of processing 3D point cloud directly. PointNet (Qi et al. 2017a) and PointNet++ (Qi et al. 2017b) were the first to propose this manner of direct point cloud data manipulation, which could fundamentally solve the quantization error problem caused by the voxelization-based approaches and requires significantly less sampling data than the view-based approaches. Among all geometric learning approaches developed (Qi et al. 2017a; 2017b; Monti et al. 2017; Atzmon, Maron, and Lipman 2018; Liu et al. 2019; Zhao et al. 2019; Hua, Tran, and Yeung 2018), DGCNN (Wang et al. 2019) can be seen as a domain milestone, due to their general definition of *EdgeConv*. As clearly illustrated in their work, most of the geometric learning approaches, including PointNet and PointNet++, can be considered as special instances from the general *EdgeConv* family. However, unlike these instances, we revisit the geometric learning of point cloud data from a shape-oriented level and, based on this, we define two shape-oriented relationships occurring in point cloud data, which encode the global context information and local context information respectively. The pointwise features can therefore be updated by incorporating the contextual effects caused by the intra-shape relationships and inter-shape relationships captured for point cloud data.

Method

Let $\mathbb{P} = \{P_1, P_2, \dots, P_N\}$ and $\mathbb{X} = \{X_1, X_2, \dots, X_N\}$ denote the point cloud to be analyzed and its pointwise features, respectively, where N is the number of points sampled and C is the number of channels of each point feature, such that $X_i \in \mathbb{R}^C$. For a given sampled point P_i , we represent its neighborhood as $\mathcal{N}(P_i)$, and the local shape formed by this neighbourhood area as $\mathcal{S}_{\mathcal{N}(P_i)}$.

ShapeConv: Shape-Orientated Convolution

We argue that the pointwise feature X_i should be capable of encoding the characterization of local shape $\mathcal{S}_{\mathcal{N}(P_i)}$. Meanwhile, considering the global influences between each local shape, X_i should also be designed to capture the long-ranged dependencies between $\mathcal{S}_{\mathcal{N}(P_i)}$ and \mathcal{S}_{others} including all other local shapes.

To this end, *ShapeConv* is proposed to model and aggregate these two shape-oriented relationships contained in point cloud, namely, intra-shape relationship and inter-shape relationship. The output of *ShapeConv* for point P_i is thus described as:

$$X'_i = A_S(L_S(\mathcal{S}_{\mathcal{N}(P_i)}), G_S(\mathcal{S}_{\mathcal{N}(P_i)}, \mathcal{S}_{others})), \quad (2)$$

where $\mathcal{S}_{others} = \{\mathcal{S}_{\mathcal{N}(P_j)}\}$ for $1 \leq j \leq N$ and $i \neq j$, $L_S(\cdot)$ and $G_S(\cdot, \cdot)$ are learning functions of local intra-shape relationship and global inter-shape relationship, and $A_S(\cdot, \cdot)$ is the aggregation function of these shape-oriented

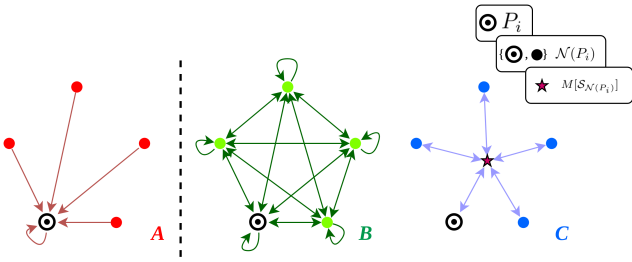


Figure 1: How the pointwise features are aggregated in a local shape $\mathcal{S}_{\mathcal{N}(P_i)}$ for its **intra-shape relationship** learning, by modeling the pairwise interactions (including self-interaction) between points $P_j \in \mathcal{N}(P_i)$ in different manners: A. Only the pairwise interactions targeting at P_i are to be considered, with potential loss of overall geometric information; B. All pairwise interactions will be taken into the account via the dense connections built on $\mathcal{N}(P_i)$; C. One possible simplified version of B, where the moment of shape will be firstly computed as $M[\mathcal{S}_{\mathcal{N}(P_i)}]$ aggregating the overall geometric information, and the interactions between each point and their moment are to be analyzed.

relationships. To keep it simple, elementwise-sum operation is adopted as our $A_{\mathcal{S}}(\cdot)$ here. The other two shape-oriented learning functions will be explained and formulated in the following two sections, and their semantic diagrams are shown as Fig. 1 and Fig. 2. The overall design of *ShapeConv* is demonstrated in Fig. 3.

Intra-Shape Relationship As the local shape $\mathcal{S}_{\mathcal{N}(P_i)}$ is formed by all neighbouring points $P_j \in \mathcal{N}(P_i)$, each neighbouring point P_j is expected to contribute equally to the generation of the complex geometric information locally encoded at $\mathcal{S}_{\mathcal{N}(P_i)}$. This expectation is consistent with the definition of conventional convolution on regular grids, where all pixels within a kernel placed would play the same role for the computation of convolved value, in a “sum of product” manner.

For this reason, we make $P_j \in \mathcal{N}(P_i)$ densely connected to form $\mathcal{S}_{\mathcal{N}(P_i)}$ (shown in Fig. 1 B), rather than treating any point uniquely, such as only building up the connections centralized at P_i among $\mathcal{N}(P_i)$ like in most of the previous works (shown as Fig. 1 A). That is, for any point P_a forming our densely-connected shape $\mathcal{S}_{\mathcal{N}(P_i)}$, we can compute its aggregated features E_a representing the pairwise interactions centralized at P_a as:

$$E_a = \frac{1}{N} \sum_{P_b \in \mathcal{N}(P_i)} g(X_a, X_b), \quad (3)$$

where $g(\cdot, \cdot)$ is a pairwise function and it can be upgraded to represent the directed pairwise interactions targeting at P_a with the subtraction implementation, which has been experimentally proven to be more efficient than other similar implementations, such as sum and concatenation (Zhao et al. 2019). Interestingly, these aggregated features E_a can there-

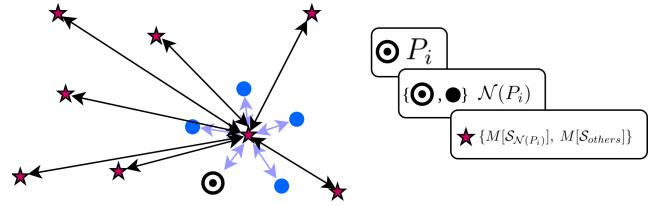


Figure 2: Our *ShapeConv* proposed is also designed to learn the **inter-shape relationship**, via capturing the long-ranged dependencies between different local underlying shapes.

fore be simplified as:

$$E_a = \frac{1}{N} \sum_{P_b \in \mathcal{N}(P_i)} (X_a - X_b) = X_a - X_{M[\mathcal{S}_{\mathcal{N}(P_i)}]}, \text{ and}$$

$$X_{M[\mathcal{S}_{\mathcal{N}(P_i)}]} = \frac{1}{N} \sum_{P_b \in \mathcal{N}(P_i)} X_b, \quad (4)$$

where $M[\cdot]$ geometrically denotes the moment point of local shape and $X_{M[\mathcal{S}_{\mathcal{N}(P_i)}]}$ is the feature of this moment point, which is computed as the averaged feature over $P_j \in \mathcal{N}(P_i)$. Therefore, E_a can be seen as an interaction between P_a and the entire local shape $\mathcal{S}_{\mathcal{N}(P_i)}$, which is demonstrated in Fig. 1 C. This formulation can also be seen as Context Normalization (Moo Yi et al. 2018) performed on each dynamically constructed local shape, with the division step excluded.

Furthermore, the intra-shape relationship among $\mathcal{S}_{\mathcal{N}(P_i)}$ can be defined as:

$$L_{\mathcal{S}}(\mathcal{S}_{\mathcal{N}(P_i)}) = A_{intra} (f_{intra}(E_a)), \quad (5)$$

where intra-shape aggregation function A_{intra} is supposed to be a symmetry function to achieve the permutation invariance required by unordered point cloud data (Qi et al. 2017a), and $f_{intra} : \mathbb{R}^{C_{in}} \rightarrow \mathbb{R}^{C_{out}}$ is a channel mapping function.

Inter-Shape Relationship To learn the inter-shape relationship in point cloud, as demonstrated in Fig. 2, the long-ranged context between each local underlying shape $\mathcal{S}_{\mathcal{N}(P_i)}$ should be taken into the consideration. Inspired by Non-Local Neural Networks (Wang et al. 2018), several attention-involved modules (Zhang et al. 2019; Fu et al. 2019) have been proposed to learn the long-ranged dependency in the computer vision domain, which are mainly focused on the conventional visual data, such as images and videos. Similar to (Xie et al. 2018), we therefore modify and extend their works to our pointwise version, namely, *pointwise long-ranged attentional context enhancement (PLACE)* module.

As clearly illustrated in Fig. 4, we first obtain the moment features of all local shapes and group them as a global shape feature matrix $M_{Global} \in \mathbb{R}^{C \times N}$. Then we implement g , θ , ϕ , and α as four separate *conv1d* with kernel size = 1. Finally, we achieve the global feature enhancement by using the *PLACE* module to model their long-ranged dependencies, as:

$$M'_{Global} = PLACE_{g,\theta,\phi,\alpha}(M_{Global}). \quad (6)$$

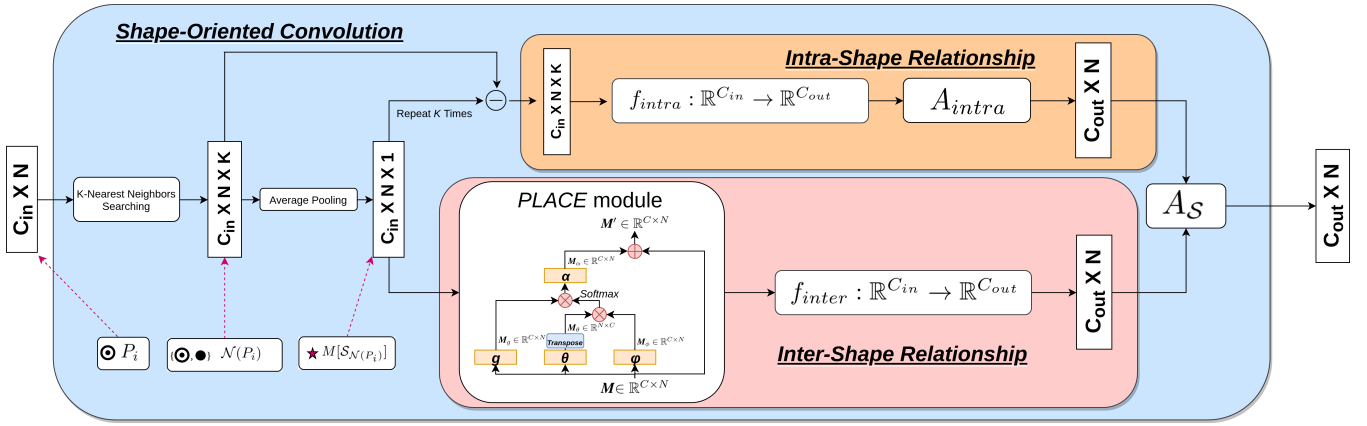


Figure 3: *ShapeConv* operator proposed, modeling two shape-oriented relationships, which are the *intra-shape relationship* and the *inter-shape relationship*. A_S is the shape-oriented aggregation function of these two relationships, while f_{intra} and f_{inter} are the two channel mapping functions for the learning of these two relationships. \ominus denotes the elementwise subtraction between the features of each neighboring point $P_j \in \mathcal{N}(P_i)$ and that of the moment point of the local shape they formed as $\mathcal{S}_{\mathcal{N}(P_i)}$. A_{intra} is the intra-shape aggregation function. The details of the *PLACE* module can be viewed in Fig. 4.

Our inter-shape relationship learning can thus be defined as:

$$G_S(\mathcal{S}_{\mathcal{N}(P_i)}, \mathcal{S}_{others}) = f_{inter}(M'_{Global\ i}), \quad (7)$$

where $M'_{Global\ i}$ is the enhanced global shape feature for $\mathcal{S}_{\mathcal{N}(P_i)}$ and $f_{inter} : \mathbb{R}^{C_{in}} \rightarrow \mathbb{R}^{C_{out}}$ is a channel mapping function.

Design of *ShapeConv* Module. We dynamically select k -nearest neighbors around a sampled point P_i to form its local underlying shape $\mathcal{S}_{\mathcal{N}(P_i)}$, where $k = 16$ in our implementation. Then an average pooling is applied to compute the averaged feature of the neighboring points $P_j \in \mathcal{N}(P_i)$, which geometrically represents the moment point of local shape $\mathcal{S}_{\mathcal{N}(P_i)}$. The pointwise interaction between each neighboring point and their moment point is calculated as the fea-

ture difference between X_j and $X_{M[\mathcal{S}_{\mathcal{N}(P_i)}]}$ and be further taken as inputs to learn the intra-shape relationship. Meanwhile, the features of all moment points formed by the local underlying shapes are grouped together and fed into our proposed *PLACE* module for the learning of the inter-shape relationship. Within our proposed *ShapeConv* module, f_{intra} and f_{inter} are two channel mapping functions designed, which can be approximated by multi-layer perceptron (MLP) (Hornik 1991). Max-pooling is chosen as the symmetric intra-shape aggregation function for A_{intra} , and A_S is implemented using the elementwise sum operation to incorporate the pointwise features from global context and local context.

SOCNN: Shape-Oriented Convolutional Neural Network

As illustrated in Fig. 5, there are three consecutive *ShapeConv* modules stacked in our *SOCNN* architecture, to capture the intra-shape relationship and inter-shape relationship encoded in point cloud data. Then the multi-scale features from these *ShapeConv* modules are combined through the shortcut connections, while a global shape signature of this point cloud object is further symmetrically aggregated using max-pooling (Qi et al. 2017a). f_{head} and f_{tail} are the two channel-raising mapping functions, which are approximated by MLP.

The classification branch is implemented by another channel mapping function $f_{classification}$ applied on the global shape signature computed, while the segmentation branch would output the per-point classifications, by taking both of the learned multi-level representations and the global shape signature into the consideration. The two extra channel mapping functions $f_{classification}$ and $f_{segmentation}$ are implemented by MLP as well.

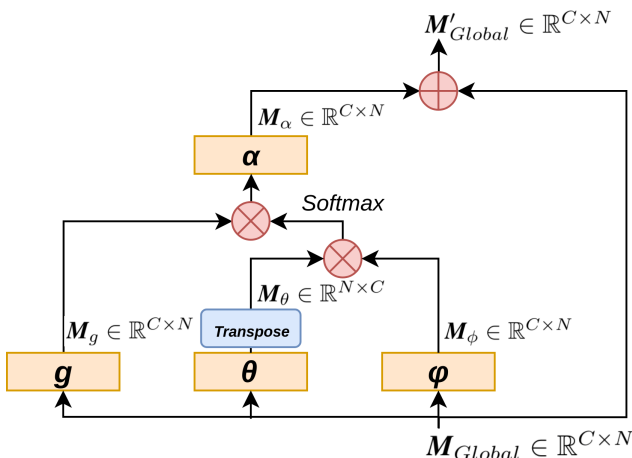


Figure 4: Our *PLACE* module for long-ranged dependency modeling in point cloud, which is simplified from Non-Local Block (Wang et al. 2018).

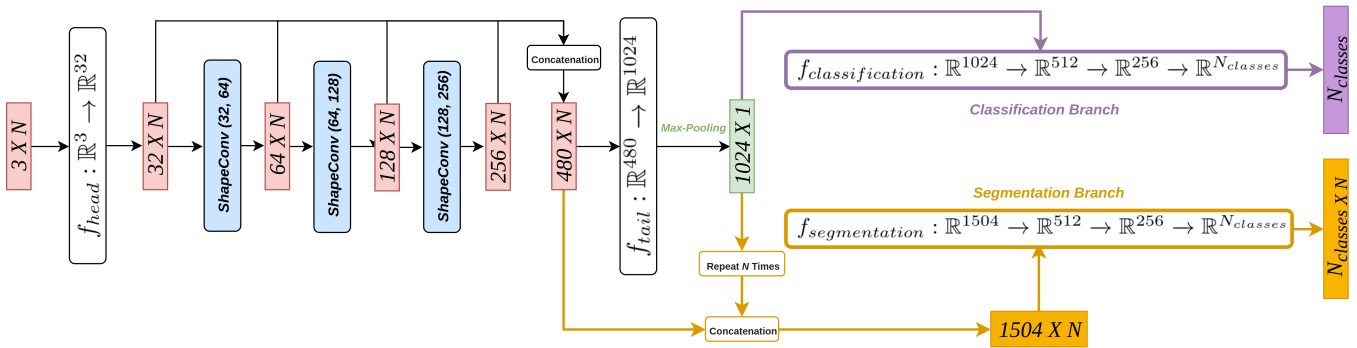


Figure 5: Shape-Oriented Convolutional Neural Network proposed, containing the *classification branch* and the *segmentation branch*. N is the number of sampled points. $ShapeConv(m, n)$ represents the module demonstrated in Fig. 3, with $C_{in} = m$ and $C_{out} = n$. f_{head} , f_{tail} , $f_{classification}$, and $f_{segmentation}$ are the channel mapping functions applied.

Experiment

Implementation Details

We select Adam as the optimizer, with learning rate 0.001 and cosine annealing applied (Loshchilov and Hutter 2017). Batch size is set to 32, while the momentum of batch normalization is initially set as 0.9 and decays with a rate of 0.5 for every 30 epochs. BatchNorm and LeakyRelu are used in all layers and omitted in figures above for simplification purpose. Dropout layers (with dropout rate = 0.5) are adopted within $f_{classification}$. The overall training framework is implemented on Pytorch with two NVIDIA GTX 1080Ti GPUs, using a distributed training scheme with Synchronized BatchNorm proposed (Zhang et al. 2018).

Shape Classification Task

We firstly evaluate our model on the ModelNet40 dataset (Wu et al. 2015) for point cloud classification task. This dataset consists of 9843 training 3D models and 2468 testing models, which are collected for 40 shape categories. We follow the same experimental setting used by (Qi et al. 2017a). For each raw 3D model from ModelNet40, we discard the mesh data after generating their corresponding point cloud data, by uniformly sampling 1024 points with (x, y, z) coordinates as their initial pointwise features. During the point sampling processing, the meshes data are discarded and their (x, y, z) coordinates are normalized to re-scale the 3D objects into unit spheres. The *classification branch* of our *SOCNN* is used for this shape classification task. Simple point cloud data augmentation techniques are adopted on the raw point coordinates, which are random scaling, translation, and perturbing. Similar to (Qi et al. 2017a; 2017b; Liu et al. 2019), ten voting tests are applied for each testing instance and their averaged results are computed as the final predictions.

Compared with other state-of-the-art approaches, our *SOCNN* achieves comparably significant results for the task of point cloud classification, which is demonstrated in Table 1. To the best of our knowledge, among all the methods manipulating point cloud data directly, *SOCNN* is the first method proposed to incorporate the global context and lo-

Table 1: Shape classification results (%) on ModelNet40 dataset. * denotes additional points sampled for the classification task.

Method	#points	Acc.
ECC (2017)	1k	87.4
PointNet (2017a)	1k	89.2
PointNet++ (2017b)	1k	90.7
PointNet++* (2017b)	5k	91.9
KD-Net (2017)	1k	90.6
KD-Net* (2017)	5k	91.8
SpiderCNN* (2018)	5k	92.4
SO-Net* (2018)	2k	90.9
PCNN (2018)	1k	92.3
DGCNN (2019)	1k	92.2
DGCNN* (2019)	2k	93.5
PointWeb (2019)	1k	92.3
RS-CNN (2019)	1k	93.6
RS-CNN* (2019)	2k	93.6
Proposed	1k	93.1
Proposed*	2k	93.6

cal context, by modeling the two shape-oriented relationships independently, i.e., the *intra-shape relationship* and the *inter-shape relationship*. As one of the other top-ranked methods, RS-CNN did not consider the global inference, which is captured and processed by our inter-shape module. Compared to PointWeb, which requires normal vectors calculated from the object meshes and extensively modeled pointwise interactions via a learning-based approach, our design achieves a similar effect in encouraging the information exchange within each local neighborhood, but in a more efficient manner.

Shape Part Segmentation Task

We then further evaluate our model on the ShapeNet-Part dataset (Yi et al. 2016) for the point cloud segmentation task. This dataset consists of 16881 3D objects, covering 16 shape categories. Most of the point cloud instances are annotated with less than six part labels, and there exist 50 parts labels

Table 2: Part segmentation results on ShapeNet-Part dataset. Metric is mIoU(%) on points.

	mean	aero	bag	cap	car	chair	ear phone	guitar	knife	lamp	laptop	motor	mog	pistol	rocket	skate board	table
# Shapes		2690	76	55	898	3758	69	787	392	1547	451	202	184	283	66	152	5271
PointNet (2017a)	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
PointNet++ (2017b)	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
PointCNN (2018)	86.1	84.1	86.5	86.0	80.8	90.6	79.7	92.3	88.4	85.3	96.1	77.2	95.3	84.2	64.2	80.0	83.0
SpiderCNN (2018)	85.3	83.5	81.0	87.2	77.5	90.7	76.8	91.1	87.3	83.3	95.8	70.2	93.5	82.7	59.7	75.8	82.8
PCNN (2018)	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
KCNet (2018)	84.7	82.8	81.5	86.4	77.6	90.3	76.8	91.0	87.2	84.5	95.5	69.2	94.4	81.6	60.1	75.2	81.3
DGCNN (2019)	85.1	84.2	83.7	84.4	77.1	90.9	78.5	91.5	87.3	82.9	96.0	67.8	93.3	82.6	59.7	75.5	82.0
RS-CNN (2019)	86.2	83.5	84.8	88.8	79.6	91.2	81.1	91.6	88.4	86.0	96.0	73.7	94.1	83.4	60.5	77.7	83.6
Proposed	85.7	83.9	84.1	85.0	77.4	91.3	78.3	91.7	87.4	83.8	96.4	69.7	93.5	83.1	58.9	76.2	82.9

in total. We split dataset into 12137 training objects, 1870 validation objects, and 2874 testing objects, following the official split policy announced by (Chang et al. 2015). For each 3D shape object, its corresponding point cloud data is generated by 2048 points sampled uniformly with (x, y, z) coordinates as their initial pointwise features. The *segmentation branch* of *SOCNN* is used for this point cloud segmentation task.

Following the same evaluation metrics set by PointNet (Qi et al. 2017a), we calculate the Intersection-over-Union (IoU) of our point cloud part segmentation results. Specifically, the comparisons are made in terms of per-object-category IoUs and the mean IoU (mIoU). To make a fair comparison, we evaluate our model with other state-of-the-arts approaches, which were proposed to manipulate point cloud data directly and would sample 2048 points for each object for the part segmentation task. The visual outputs generated by our *SOCNN* proposed for the part segmentation task are shown in Fig. 6. As presented in Table 2, the quantitative comparisons demonstrate that our model achieves state-of-the-art performance on the part segmentation task of the *CHAIR* objects and the *LAPTOP* objects.

Ablation Studies

We design and perform extensive ablation studies on ModelNet40 dataset to analyze the significance of different components proposed for the shape-oriented relationship modeling. The results of the ablation studies can be viewed in

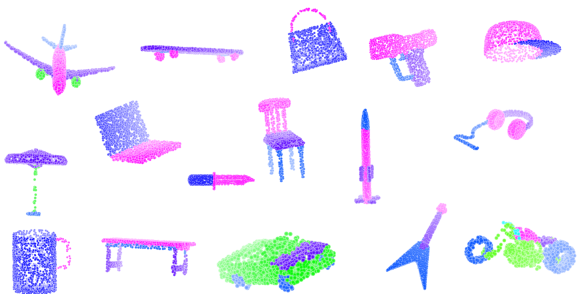


Figure 6: Part segmentation examples on the ShapeNet-Part dataset.

Table 3: Ablation studies designed for *SOCNN* (%). MP denotes whether the model calculates moment points and uses them for the relationship learning. INTRA and INTER indicate whether the model contains intra-shape relationship modeling and inter-shape relationship modeling, respectively.

model	#points	MP	INTRA	INTER	Acc.
A	1k		✓		89.5
B	1k			✓	88.2
C	1k	✓	✓		90.9
D	1k	✓		✓	88.7
E	1k	✓	✓	✓	93.1
F	2k	✓	✓	✓	93.6

Table 3. Models A and B are implemented by manipulating their pointwise features directly, rather than computing the feature difference between each point and the moment point of local shapes. To retain their number of parameters and thus make a fair comparison with other models, the neighbouring features and source feature from each neighbourhood $\mathcal{N}(P_i)$ are selected for the intra- and inter-shape relationship modeling for the formed $\mathcal{S}_{\mathcal{N}(P_i)}$, respectively. It can be seen from the results that the modeling of both the intra-shape relationship and inter-shape relationship has positive influences towards the final classification. Compared to the inter-shape relationship, the intra-shape relationship may contribute more to the final recognition results. Notably, the calculation of moment points itself can be understood as a quick and efficient operation to aggregate the local context and global context, which enhances the performance of models from C to E significantly.

Conclusion

In this paper, we firstly define two shape-oriented relationships existing in point cloud data and reformulate their geometric representation learning as two modeling processes for the global context information and the local context information. Unlike previous geometric information modeling performed for point clouds, the shape-oriented convolution (*ShapeConv*) module is proposed to incorporate the contextual effects caused by the intra-shape relationship and inter-shape relationship and aggregate these effects to update the pointwise features. Notably, we experimentally observe that the computation of the moment point from a local underlying

ing shape can be seen as a simple but efficient way to combine the contextual information captured at both the global level and local level. Finally, we propose the shape-oriented convolutional neural network (SOCNN) for point cloud analysis and evaluate its significance in the point cloud tasks of shape classification and shape part segmentation.

References

- Atzmon, M.; Maron, H.; and Lipman, Y. 2018. Point convolutional neural networks by extension operators. *ACM Transactions on Graphics (TOG)* 37(4):71:1–71:12.
- Chang, A. X.; Funkhouser, T.; Guibas, L.; Hanrahan, P.; Huang, Q.; Li, Z.; Savarese, S.; Savva, M.; Song, S.; Su, H.; Xiao, J.; Yi, L.; and Yu, F. 2015. Shapenet: An information-rich 3d model repository. cite arxiv:1512.03012.
- Feng, Y.; Zhang, Z.; Zhao, X.; Ji, R.; and Gao, Y. 2018. Gvcnn: Group-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 264–272.
- Fu, J.; Liu, J.; Tian, H.; Li, Y.; Bao, Y.; Fang, Z.; and Lu, H. 2019. Dual attention network for scene segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3146–3154.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.
- Hornik, K. 1991. Approximation capabilities of multilayer feedforward networks. *Neural networks* 4(2):251–257.
- Hu, J.; Shen, L.; and Sun, G. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7132–7141.
- Hua, B.-S.; Tran, M.-K.; and Yeung, S.-K. 2018. Pointwise convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 984–993.
- Huang, G.; Liu, Z.; Van Der Maaten, L.; and Weinberger, K. Q. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4700–4708.
- Klokov, R., and Lempitsky, V. S. 2017. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 863–872.
- Krizhevsky, A.; Sutskever, I.; and Hinton, G. E. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 1097–1105.
- Li, Y.; Bu, R.; Sun, M.; Wu, W.; Di, X.; and Chen, B. 2018. Pointcnn: Convolution on x-transformed points. In *Advances in Neural Information Processing Systems (NeurIPS)*, 820–830.
- Li, J.; Chen, B. M.; and Lee, G. H. 2018. So-net: Self-organizing network for point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 9397–9406.
- Liu, Y.; Fan, B.; Xiang, S.; and Pan, C. 2019. Relation-shape convolutional neural network for point cloud analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 8895–8904.
- Loshchilov, I., and Hutter, F. 2017. Sgdr: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR) 2017 Conference Track*.
- Maturana, D., and Scherer, S. 2015. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 922–928.
- Monti, F.; Boscaini, D.; Masci, J.; Rodola, E.; Svoboda, J.; and Bronstein, M. M. 2017. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5115–5124.
- Moo Yi, K.; Trulls, E.; Ono, Y.; Lepetit, V.; Salzmann, M.; and Fua, P. 2018. Learning to find good correspondences. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2666–2674.
- Qi, C. R.; Su, H.; Nießner, M.; Dai, A.; Yan, M.; and Guibas, L. J. 2016. Volumetric and multi-view cnns for object classification on 3d data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5648–5656.
- Qi, C. R.; Su, H.; Mo, K.; and Guibas, L. J. 2017a. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 652–660.
- Qi, C. R.; Yi, L.; Su, H.; and Guibas, L. J. 2017b. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems (NeurIPS)*, 5099–5108.
- Riegler, G.; Ulusoy, A. O.; and Geiger, A. 2017. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3577–3586.
- Shen, Y.; Feng, C.; Yang, Y.; and Tian, D. 2018. Mining point cloud local structures by kernel correlation and graph pooling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4548–4557.
- Simonovsky, M., and Komodakis, N. 2017. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3693–3702.
- Su, H.; Maji, S.; Kalogerakis, E.; and Learned-Miller, E. G. 2015. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 945–953.
- Wang, P.-S.; Liu, Y.; Guo, Y.-X.; Sun, C.-Y.; and Tong, X. 2017. O-CNN: Octree-based Convolutional Neural Networks for 3D Shape Analysis. *ACM Transactions on Graphics (TOG)* 36(4).

- Wang, X.; Girshick, R.; Gupta, A.; and He, K. 2018. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7794–7803.
- Wang, Y.; Sun, Y.; Liu, Z.; Sarma, S. E.; Bronstein, M. M.; and Solomon, J. M. 2019. Dynamic graph cnn for learning on point clouds. *ACM Transactions on Graphics (TOG)*.
- Wu, Z.; Song, S.; Khosla, A.; Yu, F.; Zhang, L.; Tang, X.; and Xiao, J. 2015. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1912–1920.
- Xie, J.; Dai, G.; Zhu, F.; Wong, E.; and Fang, Y. 2016. Deepshape: Deep-learned shape descriptor for 3d shape retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)* 39:1–1.
- Xie, S.; Liu, S.; Chen, Z.; and Tu, Z. 2018. Attentional shapecontextnet for point cloud recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4606–4615.
- Xu, Y.; Fan, T.; Xu, M.; Zeng, L.; and Qiao, Y. 2018. Spider-cnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 87–102.
- Yi, L.; Kim, V. G.; Ceylan, D.; Shen, I.-C.; Yan, M.; Su, H.; Lu, C.; Huang, Q.; Sheffer, A.; and Guibas, L. 2016. A scalable active framework for region annotation in 3d shape collections. *SIGGRAPH Asia*.
- Zhang, H.; Dana, K.; Shi, J.; Zhang, Z.; Wang, X.; Tyagi, A.; and Agrawal, A. 2018. Context encoding for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 7151–7160.
- Zhang, H.; Goodfellow, I.; Metaxas, D.; and Odena, A. 2019. Self-attention generative adversarial networks. In Chaudhuri, K., and Salakhutdinov, R., eds., *Proceedings of the 36th International Conference on Machine Learning (ICML)*, volume 97 of *Proceedings of Machine Learning Research*, 7354–7363. Long Beach, California, USA: PMLR.
- Zhao, H.; Jiang, L.; Fu, C.-W.; and Jia, J. 2019. Pointweb: Enhancing local neighborhood features for point cloud processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5565–5573.