

# Point2SpatialCapsule: Aggregating Features and Spatial Relationships of Local Regions on Point Clouds using Spatial-aware Capsules

Xin Wen, Zhizhong Han, Xinhai Liu, Yu-Shen Liu, *Member, IEEE*

**Abstract**—Learning discriminative shape representation directly on point clouds is still challenging in 3D shape analysis and understanding. Recent studies usually involve three steps: first splitting a point cloud into some local regions, then extracting the corresponding feature of each local region, and finally aggregating all individual local region features into a global feature as shape representation using simple max pooling. However, such pooling-based feature aggregation methods do not adequately take the spatial relationships (e.g. the relative locations to other regions) between local regions into account, which greatly limits the ability to learn discriminative shape representation. To address this issue, we propose a novel deep learning network, named Point2SpatialCapsule, for aggregating features and spatial relationships of local regions on point clouds, which aims to learn more discriminative shape representation. Compared with the traditional max-pooling based feature aggregation networks, Point2SpatialCapsule can explicitly learn not only geometric features of local regions but also the spatial relationships among them. Point2SpatialCapsule consists of two main modules. To resolve the disorder problem of local regions, the first module, named *geometric feature aggregation*, is designed to aggregate the local region features into the learnable cluster centers, which explicitly encodes the spatial locations from the original 3D space. The second module, named *spatial relationship aggregation*, is proposed for further aggregating the clustered features and the spatial relationships among them in the feature space using the spatial-aware capsules developed in this paper. Compared to the previous capsule network based methods, the feature routing on the spatial-aware capsules can learn more discriminative spatial relationships among local regions for point clouds, which establishes a direct mapping between log priors and the spatial locations through feature clusters. Experimental results demonstrate that Point2SpatialCapsule outperforms the state-of-the-art methods in the 3D shape classification, retrieval and segmentation tasks under the well-known ModelNet and ShapeNet datasets.

**Index Terms**—point cloud, shape representation, feature aggregation, spatial relationships, capsule network.

## I. INTRODUCTION

3D shape representation learning plays a central role in shape analysis and understanding, which has a wide range of

X. Wen and X. Liu are with the School of Software, Tsinghua University, Beijing 100084, China (e-mail: x-wen16, lxh17@mails.tsinghua.edu.cn).

Z. Han is with the School of Software, Tsinghua University, Beijing 100084, China, and also with the Department of Computer Science, University of Maryland at College Park, College Park, MD 20737 USA (e-mail: h312h@mail.nwpu.edu.cn).

Y.-S. Liu is with the School of Software, Tsinghua University, Beijing 100084, China, and also with the Beijing National Research Center for Information Science and Technology (BNRist), China (e-mail: liuyushen@tsinghua.edu.cn). (Corresponding author: Yu-Shen Liu.)

This work was supported by National Key R&D Program of China (2018YFB0505400)

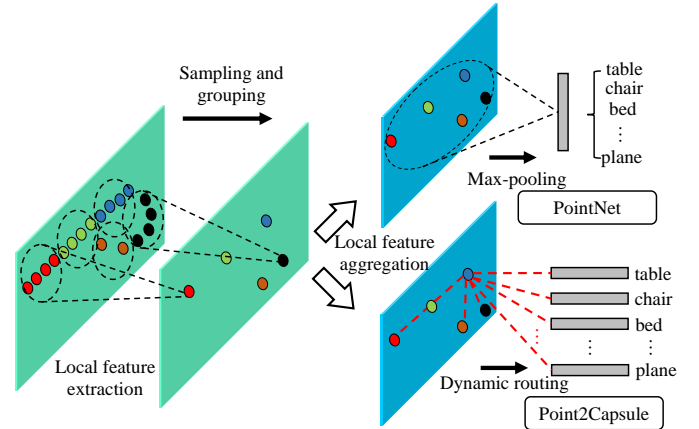


Fig. 1. The illustration of comparison between the max-pooling based PointNet and the dynamic routing based Point2SpatialCapsule in terms of local region feature aggregation. Given an example 2D point cloud like a shape of “P”, the point cloud is split into some local regions (left), from which the corresponding local region features are extracted with sampling and grouping (middle). The comparison of two methods is shown at the right side of this figure. Here, denoted by black dotted line, the max-pooling only keeps the most significant geometric characteristics in local regions (top right), while this causes the spatial relationships between local regions are filtered out. In contrast, the dynamic routing based Point2SpatialCapsule can handle both the geometric characteristics and the spatial relationships of local regions (bottom right), denoted by the red dotted line.

applications such as shape classification [1], [2], [3], retrieval [4], [5], [6], semantic segmentation [7], [8] and instance segmentation [9], [10]. Among the multiple representation forms of 3D shapes, 3D point clouds, benefited from its easy access, have become one of the most popular 3D shape forms in recent years. Specifically, the point clouds consist of a set of unordered points, each of which is composed of 3D coordinates, possibly with some additional attributes such as normal, color and material.

However, learning discriminative shape representation directly on point clouds is still challenging in 3D shape analysis and understanding. Recent studies for learning point cloud representations usually involve the following three steps. Each input point cloud is first split into some local regions. Then, the corresponding features of local regions are extracted using shared Multi-Layer Perceptron (MLP) [1] or kd-trees [11]. Finally, the extracted local region features are aggregated into a global feature vector as the shape representation [12], [13], [7]. Most of the previous methods mainly focus on how to enhance the process of local region feature extraction, while

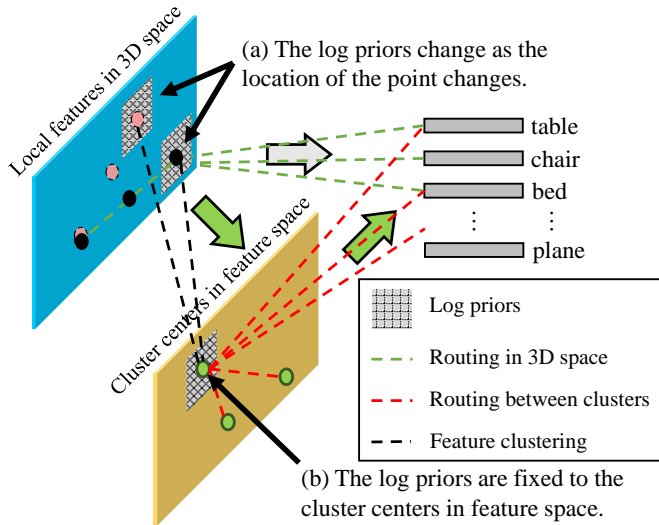


Fig. 2. The illustration of directly applying capsule network to the point cloud (a) and the clustering based Point2SpatialCapsule (b). The shifting and rotation of point cloud change the locations of local regions in 3D space and also change their corresponding log priors. Therefore, routing in 3D space (green dotted line) will cause the shifting of log priors, making the routing algorithm fail to learn the spatial relationships between local regions. In contrast, the geometric feature aggregation aggregates the input local region features into relatively invariant cluster centers. Therefore, routing between clusters (red dotted line) can efficiently learn the log priors for aggregating spatial relationships between local regions.

often employ a simple pooling-based layer [1], [12], [13] to aggregate these extracted features. However, such pooling-based feature aggregation methods do not take adequately the spatial relationships among local regions into account. So far, how to aggregate those learned local region features and their spatial relationships still remains the challenges in existing methods of point cloud representation learning. In this paper, we first argue the importance of learning spatial relationships for aggregating local region features with respect to the following two reasons. (1) For point clouds with similar local regions, the differences in the spatial arrangements of these local regions are important for learning the discriminative features. (2) Considering the permutation invariant nature of point clouds, it is important to learn the intrinsic spatial relationship between each part and the whole, in order to constitute the permutation invariant knowledge for point cloud recognition.

The common strategy for feature aggregation in previous methods [12], [1], [14], [15] is to extract the most significant characteristics (such as the engine of an airplane) in local regions of the point cloud step by step, through a deep neural network with the pooling structure (such as max-pooling). However, the problem is that the pooling-based methods will filter out the spatial relationships of different areas on the feature map [16]. Thus, it only considers the existences of characteristics in local regions, while the spatial arrangements between these regions will not be preserved. As a result, most of the existing methods usually fail to learn the spatial relationships among the local regions, which further limits the ability of the network for learning discriminative 3D shape representation.

To address the aforementioned problem, we propose a novel deep learning network, named Point2SpatialCapsule, for aggregating geometric features and spatial relationships of local regions on point clouds, which aims to learn more discriminative shape representation. Inspired by the recently developed capsule network [16], Point2SpatialCapsule employs the dynamic routing to aggregate the local region features and their spatial relationships. Fig. 1 illustrates the comparison between our dynamic routing based Point2SpatialCapsule and the previous feature aggregation methods like max-pooling used in PointNet [1]. Note that the max-pooling feature aggregation in PointNet [1] only considers the existences of characteristics in local regions, but in contrast our Point2SpatialCapsule can explicitly handle the spatial relationships between local region features through dynamic routing of capsule network. This advantage encourages us to consider adopting capsule network for 3D point clouds representation learning.

However, the problem is that, the original implementation of capsule network is designed for 2D image recognition, where the log priors in capsule network are bounded to the fixed locations on the 2D feature maps [16]. In contrast, for 3D point clouds, the locations of random sampled input points are disordered and their absolute position coordinates may not always keep consistent. As a result, it is difficult to find a direct mapping that can generate the features encoded with fixed spatial locations. What's worse, the previous capsule based methods failed to address such problem, most of which directly generate the capsules from a single global feature vectors. Such practice leads to the loss of spatial relationships between local regions. As a result, the log priors in routing algorithm between capsules can not learn the spatial relationships of local regions, which greatly limits the representation ability of capsules. In this paper, we argue the importance for encoding the fixed spatial locations into capsules, which aims to efficiently utilize the representation ability of log priors for learning the spatial relationships between local regions on point clouds.

In order to solve the above limitations, two novel modules are specially designed in Point2SpatialCapsule to achieve local region feature aggregation as follows. (1) The first module, named *geometric feature aggregation*, aims to aggregate the extracted local region features in the feature space. Here, the term “*geometric*” indicates that this module aggregates the geometric information, like the coordinates of central points and the shapes of local regions represented by the feature vectors, into the centers of local feature clusters, which aims to resolve the disorder problem of local regions. (2) The second module, named *spatial relationship aggregation*, is to apply routing algorithm on the learned feature clusters. The term “*spatial-aware*” indicates that the capsules are encoded with the spatial locations, which is to guarantee the direct mapping between the log priors and the fixed locations in the 3D space. Therefore, we call them the *spatial-aware capsules*, which allows the network to efficiently learn the spatial relationships between local regions. Fig. 2 shows the visualized demonstration of the advantage of spatial relationship aggregation. Because of the shifting and rotation of point clouds, the changing locations of local region features in 3D space also change the log priors. To resolve this issue, the geometric

feature aggregation clusters the input local region features into the learnable cluster centers, which are irrelevant to the input points and relatively invariant in the feature space. Therefore, the routing algorithm can efficiently learn the log priors for aggregating the spatial relationships between local regions. Our main contributions are summarized as follows.

- We propose a novel deep network, i.e. Point2SpatialCapsule, for learning more discriminative shape representations of point clouds. Compared with the traditional pooling-based methods, Point2SpatialCapsule can explicitly learn not only geometric features of local regions but also the spatial relationships among them.
- We propose the geometric feature aggregation to resolve the disorder problem of local regions, where the local region features are aggregated into the learnable cluster centers, which are explicitly encoded with the spatial locations from the original 3D space.
- We propose the spatial relationship aggregation to further utilize the spatial locations encoded in the feature clusters. Compared to the previous capsule network based methods, the spatial relationship aggregation can learn more discriminative spatial relationships between local regions by establishing a direct mapping between log priors and the spatial locations through feature clusters.

The remainder of this paper is organized as follows. First, the related work is introduced in Sec.II. Then we detail the proposed Point2SpatialCapsule in Sec.III. The experiments and the ablation studies are given in Sec.IV. Finally, we conclude this paper in Sec.V.

## II. RELATED WORK

In this section, we mainly review the methods related to 3D shape representation learning based on deep learning networks. The existing methods can be roughly divided into four categories according to various 3D shape forms that are learned from, including voxels, point clouds, views and meshes.

### A. Point Cloud Based Methods

Recent studies of point cloud representation learning mainly focus on the local feature extraction and integration. PointNet [1] is the pioneering work of introducing deep learning into point cloud representation learning, which independently learns the features of each point and aggregates the learned features into a global feature with the max-pooling layer. After that, plenty of the follow-up studies [13], [7], [17], [18] focus on how to better integrate the contextual information of local regions on point clouds. For example, PointNet++ [12] designed the hierarchical feature learning architecture based on PointNet to encode multi-scale local areas. Following the convolutional structure of PointNet++, successors such as PointCNN [13] and SpiderCNN [17] investigated some improved convolution operations which aggregate the neighbors of a given point by edge attributes in the local region graph. Different from the idea of using convolution structure, Point2Sequence [7] introduced the sequential model (i.e. RNN) to capture the fine-grained contextual information

of features in local regions. Specifically, Point2Sequence arranges the features into a sequence according to the size of the region scale, and then uses a RNN to capture the contextual information within the local regions. However, the problem is that most of the above methods fail to consider the spatial relationships among different local regions when aggregating the extracted local region features, where the usual practice for these methods is to use the pooling layer to learn the global feature from the local ones.

More recent studies focus on how to improve the local region feature extraction [3], [19], [20]. These methods have shown impressive potentials in the semantic segmentation task on point cloud. For examples, A-CNN [3] was proposed to annularly arrange the neighbor points and apply the convolution network on these arranged points to learn the local region features. RS-CNN [19] designed a shape-aware convolution to learn the local region features from the relation between points.

The proposed Point2SpatialCapsule mainly focus on how to aggregate the feature and relationships of local regions after extracting local features. The usual practice for previous methods is to apply the strategy of bottom-to-top point cloud feature aggregation [11], [21], [14], [8], [22]. For example, Kd-Net [11] performs multiplicative transformations according to the subdivisions of point clouds based on the kd-trees. SO-Net [14] employs a SOM to build the spatial distribution of the input point cloud, which allows hierarchical feature extraction on both individual points and SOM nodes. However, most of the above methods use max-pooling as a feature aggregation method, which inevitably filters out the spatial relationships among local regions. On the other hand, PVNet [23] is also a notable method that considers the local feature aggregation, which focuses on mining the difference in importance between the local features. It employs high-level global features from the multi-view data of input 3D shapes to mine the relative correlations between different local features from the point cloud data. Same as the above-mentioned methods, PVNet only learns the different contributions among local regions, while the spatial relationships among these regions are not considered.

### B. View-based Methods

The dominant performance of multi-view based methods on the task of 3D shape retrieval comes from the research progress of measuring the similarities between 2D image features [24], [25], [26], [27]. As one of the pioneering work, GIFT [25] adopted the Hausdorff distance to measure the similarity between the view sets of two 3D shapes. Another notable research direction is to focus on PANORAMA views of 3D shapes, where a PANORAMA view can be regarded as the seamless aggregation of multiple views captured on a circle. For examples, DeepPano [28] introduced a row-wise max-pooling to relief the effect of rotation about the up-oriented direction, and Sfikas et al. [29] introduced CNN for learning the global features from the PANORAMA views in a consistent order. To explore the potential of attention mechanism, the methods like 3DViewGraph [2] have been proposed to integrate the spatial pattern correlations of unordered

views with attention weights, and Part4Features [5] developed a novel multi-attention mechanism for aggregating the learned local parts.

More recently, SeqViews2SeqLabels [6] was proposed to learn 3D features via aggregating sequential views by RNN, which aims to eliminate the effect of rotation of 3D shapes. Compared with the previous pooling based methods, the RNN-based SeqViews2SeqLabels suffers less from the content and the spatial location loss. Similarly, as an unsupervised approaches, VIP-GAN [4] trains an RNN-based neural network architecture to solve multiple view inter-prediction tasks for each shape.

### C. Voxel-based Methods

Voxel-based methods often rasterize a 3D shape as a function or distribution sampled on voxels [30], [31]. For supervised learning the representation of 3D voxels, 3DShapeNets [32] adopted the convolutional restricted Boltzmann machine to learn the representation of 3D voxels. O-CNN [21] learns the representation of 3D voxel based on a novel octree structure. And Han et al. [33] proposed a novel permutation voxelization strategy to learn high-level and hierarchical 3-D local features from raw 3-D voxels. For unsupervised learning, methods like VConv-DAE [34] use the fully convolutional autoencoder for unsupervised learning the voxel representation by reconstruction. However, the problem is, considering the induced complexity and limitations of directly exploiting the sparsity of voxel grids, it is difficult to introduce the large scale or flexible deep networks for representation learning. Therefore, more recent methods such as OctNet [22] and kd-net [11] consider to utilize the scalable indexing structures for solving this problem, where deep neural networks can be further adopted for achieving more impressive results.

### D. Mesh-based Methods

As for mesh-based methods, to explore the effectiveness of the heat diffusion based descriptor, Xie et al. [35] proposed a shape feature learning scheme based on auto-encoders, where the model can extract the features that are insensitive to the deformations. By fully utilizing the spectral domain, Xie et al. [36] further proposed to learn a novel binary spectral shape descriptor with the deep neural network for 3D shape correspondence. Recently, BoSCC [37] was introduced for a spatially enhanced 3D shape representation based on bag of spatial context correlations. And more recently, Deep Spatiality [38] was also proposed to simultaneously learn 3D global and local features with novel coupled softmax.

### E. Capsule Networks

The ability of capsule network [16] for capturing spatial relationships comes from the dynamic routing algorithm and the log priors, which are bound to the absolute location on the input feature maps. Specifically, the capsule network learns the log priors by considering the relationships between the absolute locations on the feature map and the high-level capsules. Then, through the dynamic routing algorithm, which

is based on the learned log priors, the high-level capsules can integrate the low-level features and their spatial relationships among different locations on the feature maps. This advantage promotes us to consider applying the capsule network to 3D point cloud representation learning.

So far, the capsule network has shown the great potentials in many research areas, such as image processing [39], [40] and natural language processing (NLP) [41], [42], [43]. However, as for the application of capsule network in 3D shape representation learning, there are a few methods proposed in recent years. For example, 3D-CapsNet [44] adopts the capsule network for 3D shape classification tasks based on volumetric data, and 3D-Point-Capsule [45] learns the point cloud representation and part segmentations in an unsupervised way. And for supervised learning, 3DCapsule [46] applies the capsule network as an extension of fully-connected layers for point cloud classification.

An important problem of the above methods is that they all build the capsule layers over the global feature (usually produced by the fully-connected layer or max-pooling) of point clouds, where the spatial relationships between local region features have been filtered out by the network. Therefore, the log priors in routing algorithm cannot learn the spatial distribution among the extracted local features, which limits the biggest advantage of capsule network for aggregation spatial relationships of local regions.

Therefore, to address this problem of previous methods, Point2SpatialCapsule aggregates the features into clusters in feature space, and applies the routing algorithm between these aggregated clusters. In the research of point cloud representation learning, methods like PointNetVLAD [47] have adopted the similar clustering strategy, i.e. NetVLAD [48], for feature aggregation. However, different from the previous methods that only cluster features for aggregating regions with similar geometric characteristics (e.g. shapes), our method takes one step further to not only considering geometric characteristics, but also explore the potentials for aggregating spatial relationships between these regions. Specifically, Point2SpatialCapsule produces the clusters for both the features and their coordinates, in order to explicitly preserve the features and their spatial location.

## III. SHAPE REPRESENTATION LEARNING WITH POINT2SPATIALCAPSULE

An overview of shape representation learning network with Point2SpatialCapsule is shown in Fig. 3. The whole network consists of three main parts as follows. (1) The first part is the multi-scale local feature extraction, which is a PointNet++ based network for extracting the features from multi-scale local regions on point clouds (see Sec. III-A). (2) The second part is Point2SpatialCapsule, which is composed of two main modules for aggregating the learned features into the global shape representation. Here, the first module, i.e. *geometric feature aggregation*, is to aggregate local region features into clusters (see Sec. III-B). The second module, i.e. *spatial relationship aggregation*, is to aggregate the feature clusters and their spatial relationships into global feature representation

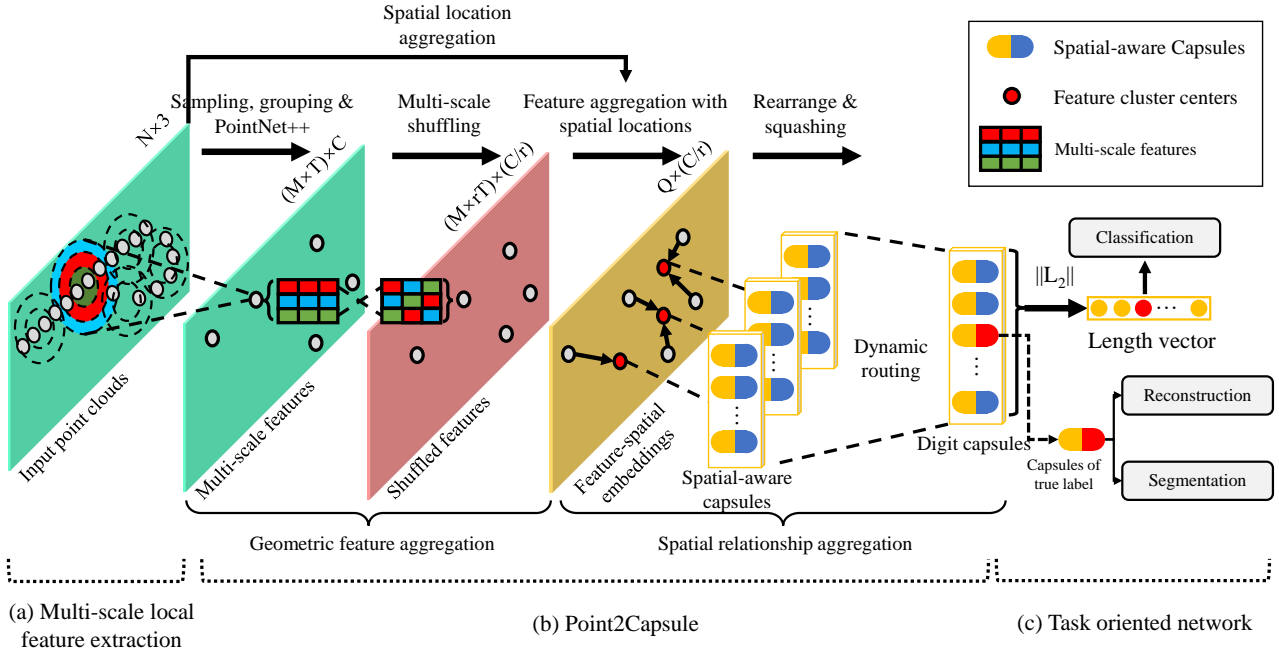


Fig. 3. The architecture of our proposed Point2SpatialCapsule. For input point clouds, (a) the multi-scale local feature extraction first extracts features from multi-scale areas, (b) then the geometric feature aggregation encodes the extracted multi-scale local region features and their locations into the learnable clustering centers to produce the feature-spatial embeddings. The spatial relationship aggregation aggregates the feature-spatial embeddings by considering both the embeddings and their spatial relationships. (c) The task oriented network is adopted for performing on different tasks.

(see Sec. III-C). In this section, we will also detail the training procedure of Point2SpatialCapsule (see Sec. III-D). (3) The third part is the task oriented network used for various tasks such as shape segmentation (see Sec. III-E).

#### A. Multi-scale Local Feature Extraction

The first part of our network is the multi-scale local feature extraction, as shown in Fig. 3(a). Given a set of input points  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ , by following the practice of PointNet++ [12] and ShapeContextNet [49], we iteratively produce a sub-sampling  $\{\mathbf{x}_{k_1}, \mathbf{x}_{k_2}, \dots, \mathbf{x}_{k_M}\}$  with  $M$  points as the centroids of the local regions using farthest point sampling (FPS), such that the newly added point  $\mathbf{x}_{k_j}$  is the farthest point (in metric distance) from the rest sampled points  $\{\mathbf{x}_{k_1}, \mathbf{x}_{k_2}, \dots, \mathbf{x}_{k_{j-1}}\}$ . Then, for each sampled point, the  $K$  nearest neighbor (kNN) searching is employed to find  $\{K_i\}_{i=1, \dots, T}$  neighbors for this point, under  $T$  different scale areas. Followed by a grouping layer, the sampled point and its neighbors are grouped as a  $K_i \times 3$  tensor for scale  $K_i$ . After that, a simple but effective MLP layer is employed to extract the features of all neighbor points, producing a tensor with shape  $K_i \times C$ . Finally, a max-pooling layer is applied to integrate the point features in each scale to produce the scale feature of dimension  $C$  for scale  $K_i$ . For  $M$  points in total and  $T$  scales for each point, the multi-scale local feature extraction layer produces  $M \times T$  multi-scale features, forming a tensor of shape  $M \times T \times C$  as its output.

In the implementation, we apply two layers of multi-scale local feature extraction for hierarchically extracting features from point clouds.

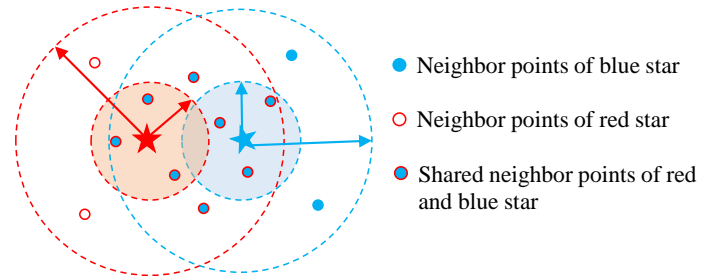


Fig. 4. Illustration of feature similarity between adjacent points. Features of small scales are not very similar because of the small overlap in sampled points. On the contrary, the features of large scales are more similar because of a larger overlap.

#### B. Point2SpatialCapsule: Geometric Feature Aggregation

In this subsection, we detail the first module of Point2SpatialCapsule, which aims to aggregate the extracted features into clusters and encodes these features with spatial locations.

As shown in Fig. 3(b), before clustering features, the module of geometric feature aggregation first applies the multi-scale shuffling to enhance the diversity of features. Then the features are aggregated into clusters and encoded with the spatial locations (e.g. the absolute locations in the 3D space) from the original 3D space.

1) *Multi-scale Shuffling*: Different from the previous methods that apply the pooling-based strategy for integrating the features extracted from multi-scale regions, we propose the multi-scale shuffling layer to build the shuffled features. The

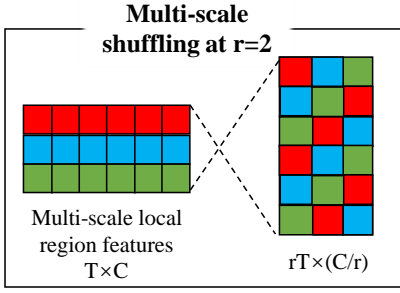


Fig. 5. Illustration of the multi-scale shuffling, which is the solution to the problem of feature similarity.

reason for adding this layer is demonstrated in Fig. 4, as explained below. When searching the neighbor points in a large scale, the searching areas of two adjacent centroids will overlap with each other and output the same neighbor points. As a result, the adjacent points will tend to have the similar features for large scale, which can reduce the diversity of features and introduce an initial clustering center for the subsequent clustering layer. On the other hand, the features of small scales between two centroids are dissimilar because of small overlaps. Therefore, the multi-scale shuffling is introduced to smooth the perceived range of features between different scales and enhance the feature diversity, by mixing the dissimilar features of small scales with the similar features of large scales. As a result, the multi-scale shuffling can promote the network to consider all input features equally and alleviate the problem of similar features.

The effect of multi-scale shuffling is shown in Fig. 5. Specifically, given a point with  $T$  scale features of  $C$  dimension, which forms a tensor with the shape  $T \times C$ , the multi-scale shuffling periodically rearranges the elements in the  $T \times C$  tensor into a tensor of shape  $rT \times (C/r)$ , where  $r$  is an integer. Thus, for  $M$  points in total, the multi-scale shuffling will produce  $M \times rT$  shuffled features of dimension  $C/r$ , resulting in a tensor of size  $(M \times rT) \times (C/r)$ .

The multi-scale shuffling is inspired by the subpixel convolution [50] for image upsampling, where the number of area scales  $T$  can be considered as the size of image, and the feature dimension  $C$  can be regarded as the channels of feature maps. However, different from subpixel convolution which is designed for speeding up the calculations and reducing the amount of parameters in the network, the multi-scale shuffling used in our method aims to enhance the diversity of scale features. We will quantitatively explore the importance of the multi-scale shuffling in ablation studies in Sec. IV-E.

2) *Feature Aggregation with Spatial Encodings*: The purpose of this layer is to aggregate the shuffled features into the learnable feature cluster centers, which can be regarded as the latent embeddings describing the semantic patterns of the local regions features. To achieve this purpose, we propose to cluster the features in the feature space and their coordinates in the original 3D space. After that, the cluster centers in both the feature space and the 3D space are fused to produce the *feature-spatial embeddings*, as illustrated in Fig. 6(a).

Although the traditional clustering methods like k-means

can be adopted to produce the feature cluster centers, their computational cost may be very high because of the huge number of features to be clustered. Therefore, inspired by the recent development of NetVLAD [48], we adopt the soft-assignment for learning the clustering centers for the input shuffled local features. Specifically, the network learns  $Q$  cluster centers for input features, denoted as  $\{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_Q | \mathbf{q}_k \in \mathbb{R}^{C/r}\}$ , as colored by yellow in Fig. 6(a). For each cluster center  $\mathbf{q}_k$ , the layer produces a *feature embedding*  $C(\mathbf{q}_k) \in \mathbb{R}^{C/r}$ , which is an aggregated representation over the whole input shuffled features  $\{\hat{\mathbf{p}}_i\}$ , denoted by

$$C(\mathbf{q}_k) = \sum_{i=1}^n \frac{e^{\mathbf{w}_k^T \hat{\mathbf{p}}_i + b_k}}{\sum_{k'} e^{\mathbf{w}_{k'}^T \hat{\mathbf{p}}_i + b_{k'}}} (\hat{\mathbf{p}}_i - \mathbf{q}_k), \quad (1)$$

where  $\{\mathbf{w}_k\}$  and  $\{b_k\}$  are the weights and biases, respectively, that determine the contribution of each local feature to the cluster center  $\mathbf{q}_k$ . During training, all the weights, biases and the cluster centers are updated through back-propagation algorithm.

To explicitly encode the spatial locations of local features into their cluster centers, we first cluster the coordinates  $\{\mathbf{x}_i\}$  of input points into the coordinates cluster centers  $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_Q | \mathbf{y}_k \in \mathbb{R}^{C/r}\}$ , which is the same process as described above and colored by green in Fig. 6(a). The *spatial embeddings*  $C(\mathbf{y}_k) \in \mathbb{R}^{C/r}$  for coordinates is given as

$$C(\mathbf{y}_k) = \sum_{i=1}^n \frac{e^{\mathbf{w}_k^T \mathbf{x}_i + b'_k}}{\sum_{k'} e^{\mathbf{w}_{k'}^T \mathbf{x}_i + b'_{k'}}} (\mathbf{x}_i - \mathbf{y}_k). \quad (2)$$

Then, the produced local feature embedding and its corresponding spatial embedding are concatenated to form an explicit *feature-spatial embedding*  $C(\mathbf{s}_k) = [C(\mathbf{y}_k) : C(\mathbf{x}_k)]$ .

### C. Point2SpatialCapsule: Spatial Relationship Aggregation

In Fig. 6(b), we show the overall architecture of previous methods [45], [46] for building the capsules, and compare it with our proposed Point2SpatialCapsule shown in Fig. 6(a). The main difference is that Point2SpatialCapsule builds the spatial-aware capsules based on cluster centers with spatial encodings, while the previous studies simply build the capsules based on the single representation vector generated by fully-connection or pooling based local feature aggregator. As a result, the previous methods fail to preserve the spatial relationships between local regions, which further limits the representation learning ability of dynamic routing.

In this subsection, in order to efficiently learn the prior logs, we first independently generate the spatial-aware capsules from the feature-spatial embeddings using *rearrange* and *squashing*. Then, we propose to apply routing algorithm between the spatial-aware capsules.

1) *Rearrange and Squashing*: To build the spatial-aware capsules from the feature-spatial embeddings produced by the geometric feature aggregation module, we deviate from the 2D practice of the original capsule network [16]. In the original capsule network, the spatial-aware capsule aggregates its representation vector by collecting the output logits across

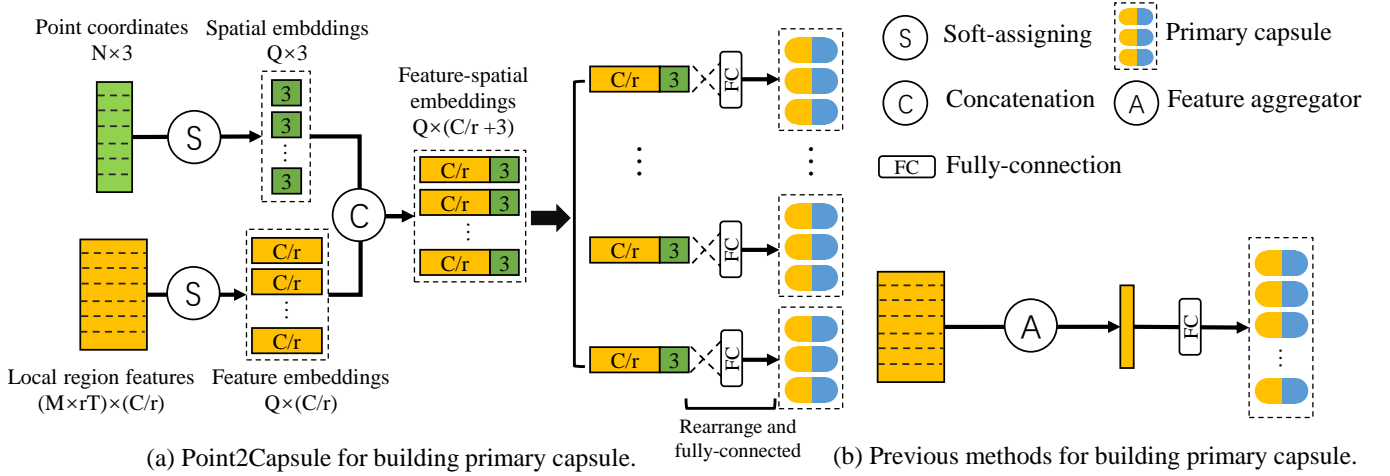


Fig. 6. Comparison of the strategies for applying dynamic routing in local region features between (a) Point2SpatialCapsule and (b) the previous methods [45], [46]. The previous methods directly use a feature aggregator (e.g. max-pooling) to aggregates all local region features, and generate the capsules based on the aggregated global feature. In contrast, Point2SpatialCapsule proposes to cluster the local region features and the point coordinates, and then combines them as spatial-aware cluster centers. The spatial-aware capsules of Point2SpatialCapsule are independently generated according to each cluster centers.

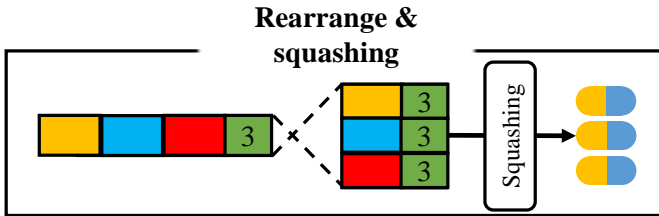


Fig. 7. Illustration of rearrange and squashing layer. The green block is the 3D spatial embedding, while the yellow, blue and red block together constitute a feature embedding.

different channels at the same location on the feature maps. In Point2SpatialCapsule, since we have built the feature-spatial embeddings encoded with the spatial locations, we can consider that each embedding corresponds to a fixed location, which is the learnable cluster center. Therefore, we can directly rearrange the output and use the fully-connected layer with a squashing activation to produce the spatial-aware capsules. The rearrange layer is to split the feature-spatial embeddings  $C(s_k)$  into several short vectors  $\{\mathbf{u}_i\}$ . As shown in Fig. 7, the input feature-spatial embedding is split into  $K = 3$  vectors, each of which is combined with the spatial embedding. Then, we follow a squashing layer, as denote by

$$\text{squashing}(\mathbf{u}_i) = \frac{\|\mathbf{u}_i\|^2}{1 + \|\mathbf{u}_i\|^2} \frac{\mathbf{u}_i}{\|\mathbf{u}_i\|}, \quad (3)$$

where the spatial-aware capsules  $\{\mathbf{u}_i\}$  are generated as the final output of this layer.

2) *Routing Algorithm*: Given the input spatial-aware capsules, we follow [16] to apply dynamic routing algorithm to obtain the digit capsule. Specifically, the digit capsule  $\mathbf{v}_j$  is the output of weighted sum of the prediction vector  $\hat{\mathbf{u}}_{ij}$  followed by the squashing layer, which can be formulated as

$$\hat{\mathbf{u}}_{ij} = W_{ij} \mathbf{u}_i, \quad (4)$$

$$\mathbf{v}_j = \text{squashing}\left(\sum_i c_{ij} \hat{\mathbf{u}}_{ij}\right), \quad (5)$$

where  $\mathbf{u}_i$  is the  $i$ th spatial-aware capsule and  $W_{ij}$  is a learnable matrix. The *coupling coefficients* [16]  $c_{ij}$  is determined by the iterative dynamic routing process, denote by

$$c_{ij} = \frac{e^{b_{ij}}}{\sum_k e^{b_{ik}}}. \quad (6)$$

In 2D capsule network, the  $\{b_{ij}\}$  are log priors that only depend on the fixed locations and the type of two capsules. In our network, because the disordered input features are clustered as the feature-spatial embeddings by soft-assignment, these features are bounded to the fixed locations (which are the cluster centers) in the feature space. Therefore, the dynamic routing can learn the log priors between these centers and the digit capsules.

Before training, all of the log priors  $\{b_{ij}\}$  are initialized to zero. During training,  $\{b_{ij}\}$  are learned discriminatively at the same time with other parameters in the network, by adding the scalar product of  $\mathbf{v}_{ij}$  and  $\hat{\mathbf{u}}_{ij}$ , i.e.

$$b_{ij} \leftarrow b_{ij} + \mathbf{v}_{ij} \cdot \hat{\mathbf{u}}_{ij}. \quad (7)$$

#### D. Point2SpatialCapsule: Training

Following the practice of [16], Point2SpatialCapsule uses the reconstruction loss and the classification loss for supervised point cloud representation learning.

The length of each digit capsule indicates the probability that the characteristic represented by this capsule exists in the input point clouds [16]. During training, the margin loss  $\mathcal{L}_{cls}$  is adopted for shape classification defined as

$$\mathcal{L}_{cls} = \sum_j T_j \max(0, m^+ - \|\mathbf{v}_j\|)^2 + \sum_j T_j \lambda (1 - T_j) \max(0, \|\mathbf{v}_j\| - m^-)^2, \quad (8)$$

where  $T_j = 1$  if class  $j$  is the true label; otherwise,  $T_j = 0$ .  $m^+$ ,  $m^-$  and  $\lambda$  are the hyper parameters.

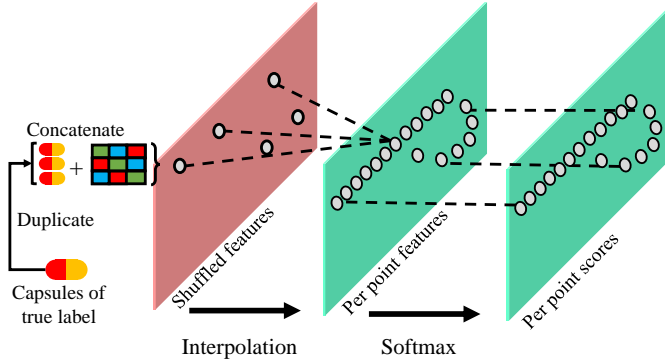


Fig. 8. Illustration of the segmentation network in our Point2SpatialCapsule.

We further reconstruct the input point clouds using four fully-connected layers, with each layer followed by a *relu* activation and batch normalization except for the last layer. The digit capsule corresponding to the true label is used as the input representation vector to the reconstruction network. The chamfer loss between the original point cloud  $\mathbf{X}$  and the reconstructed point cloud  $\hat{\mathbf{X}} = \{\hat{x}_i\}$  is adopted as the reconstruction loss  $\mathcal{L}_{rec}$ , as denoted by

$$\mathcal{L}_{rec} = \frac{1}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} \min_{\hat{\mathbf{x}} \in \hat{\mathbf{X}}} \|\mathbf{x} - \hat{\mathbf{x}}\| + \frac{1}{|\hat{\mathbf{X}}|} \sum_{\hat{\mathbf{x}} \in \hat{\mathbf{X}}} \min_{\mathbf{x} \in \mathbf{X}} \|\hat{\mathbf{x}} - \mathbf{x}\|. \quad (9)$$

The total loss for training is the weighted sum of margin loss and the reconstruction loss, as denote by

$$\mathcal{L} = \mathcal{L}_{cls} + \alpha \mathcal{L}_{rec}, \quad (10)$$

where  $\alpha = 0.0001$  for all the experiments in this paper.

### E. Model Adjustments for Part Segmentation

The goal of part segmentation is to predict a semantic label for each point in the point cloud. There are two alternative ways for acquiring the per-point feature for each point from the global feature: duplicating the global feature with  $N$  times [1], [18], or performing upsampling by interpolation [12], [14]. In this paper, we follow the second way to duplicate the vectors in digit capsules belonging to the true label. Then we concatenate the duplicated vectors with the shuffled features. The interpolation layers are used for propagating the features from shape level to point level by upsampling, as shown in Fig. 8.

## IV. EXPERIMENTS

### A. Experimental Setup

1) *Datasets*: The 3D shape classification and retrieval experiments are conducted on two subsets of the Princeton ModelNet dataset [32], i.e. ModelNet40 and ModelNet10. The ModelNet40 dataset contains 12,311 shapes which belong to 40 categories. We follow the same training and split settings as [14], which contains 9,843 shapes for training and 2,468 shapes for testing, respectively. The ModelNet10 dataset is a relatively small dataset which contains the 10 common categories of ModelNet40. Following [14], we split the ModelNet10 into 2,468 training samples and 909 testing

samples. Since the original ModelNet provides CAD models represented by vertices and faces, we use the prepared ModelNet10/40 data from [12] for fair comparison. The part segmentation task is conducted on the ShapeNet part dataset [51], which contains 16,881 models from 16 categories and is split into training, validation and testing following PointNet++. There are 2048 points sampled for each 3D shape, where each point in a point cloud object belongs to certain one of 50 part classes and each point cloud contains 2 to 5 parts.

2) *Classification and Retrieval Settings*: Because the length of representation vector in digit capsule indicates the probability that certain characteristic exists in the input point clouds. In the case of Point2SpatialCapsule, the characteristic of digit capsule is the class label. Thus, we choose the digit capsule  $v_j$  with the biggest length  $\|v_j\|$  as the predicted label for shape classification. For the shape retrieval task, we use the Euclidean distances between the length vectors  $V = [\|v_1\|, \|v_2\|, \dots, \|v_m\|]$  of point clouds for similarity measurement. Such similarity measurement is in accordance with the way how capsule stores information. What's more, a direct comparison between the length vectors requires less computational cost than comparing representation vectors in capsules.

3) *Implementation Details*: In this paper, we use two multi-scale local feature extraction layer for hierarchically extracting features from point clouds. For the first feature extractor, the input is 1024 points associated with their x, y and z coordinates, from which 512 points is sampled using farthest point sampling. For each sampled point, we select [8, 16, 32, 64] nearest neighbor points of four scales. The MLPs used in the first block have [32, 32, 64] units for each layer. The second feature extractor samples 256 points out of the 512 points. The number of points for kNN search is the same as the first block. The MLPs for the second block have the units of [64, 64, 128] for each layer. The parameter  $r$  for multi-scale shuffling is 2. The number of the cluster centers is  $Q = 64$  and the dimension is  $C = 256$ . In the rearrange and squashing layer, we split each embedding into 16 16-dimensional short vectors, which form 1024 16-dimensional spatial-aware capsules in total.

### B. 3D Shape Classification

Table I compares Point2SpatialCapsule with the existing state-of-the-art methods of point cloud representation learning in terms of shape classification accuracy under ModelNet10 and ModelNet40, respectively. For fair comparison, all the results in Table I are obtained under the same input, which handles with raw point sets. Point2Capsule achieves a superior result (93.4%) under ModelNet40, which is higher than the baseline method PointNet++ by 2.7%. Specially, Point2Capsule with additional normal vectors achieves the best results (95.9% and 93.7%), compared with the best additional-input method SO-Net [14] (95.7% and 93.4%), under ModelNet10 and ModelNet40, respectively.

We note that both PointNet++ and Point2SpatialCapsule use a multi-scale local feature extraction strategy, where the difference lies in the method used for aggregating local features. The PointNet++ applies max-pooling for aggregating the



local features, while Point2SpatialCapsule uses the geometric feature aggregation with spatial relationship aggregation for learning the global representation. Therefore, the improvement in classification accuracy of Point2SpatialCapsule proves the effectiveness of the proposed network for local feature aggregations.

3DCapsule [46] is the work most related to our Point2SpatialCapsule in Table I. As already discussed in Sec.II, 3DCapsule simply applies the capsule network on the global features produced by a pooling/full-connected layer, which falls into the scenario of information loss of the spatial locations.

In contrast, our Point2SpatialCapsule applies dynamic routing on the feature-spatial embeddings generated by the geometric feature aggregation module, which can aggregate both the features and their spatial location. The experimental results in Table I shows the implementation of capsule network in our Point2SpatialCapsule is more effective than the implementation of 3DCapsule.

Compared with PointCNN [13] and DGCNN [18], Point2SpatialCapsule still achieves the best results. We note that, PointCNN and DGCNN are also CNN-based neural network, which aims to preserve the spatial locations and spatial relationships of local regions. However, both of them use the max-pooling for aggregating the local region features, which filters out the spatial locations and relationships, especially when aggregating the local features into the global features. As shown Table I, the proposed Point2SpatialCapsule yields better performance than the PointCNN and DGCNN, which demonstrates the superior advantages of Point2SpatialCapsule for preserving the spatial locations and relationships.

As seen in Table I, our Point2SpatialCapsule outperforms most of the xyz-input methods on point clouds. Specifically, our result is ranked the first place under ModelNet10 (95.8%), and ranked the second place under ModelNet40 (93.4%) which is slightly lower than RS-CNN [19] by 0.2%. As claimed in [19], RS-CNN performed “ten voting tests with random scaling and averages the predictions” during testing. In contrast, we only apply the single model prediction for fair comparison with most of the existing methods [14], [46], [11]. Moreover, when using additional normal vectors as the input, the proposed Point2SpatialCapsule can achieve the best performance among all reported results under ModelNet10 (95.9%) and ModelNet40 (93.7%), respectively. This convincingly verifies the effectiveness of Point2SpatialCapsule.

### C. 3D Shape Retrieval

In Table II, we compare the proposed Point2Capsules with counterpart methods in 3D shape retrieval task, in terms of *mean average precisions* (mAPs). Since most of the methods focusing on 3D shape retrieval are based on multi-views of 3D models, in this subsection, we also quote the experimental results of the multi-view based methods to verify the effectiveness of Point2SpatialCapsule. Note that, the results of PointNet and PointNet++ are obtained by following the same training procedure as described in their original papers, which are denoted by \* in this table.

TABLE I  
THE SHAPE CLASSIFICATION ACCURACY (%) COMPARISON ON MODELNET10 AND MODELNET40.

Method	Input	ModelNet10	ModelNet40
PointNet [1]	1024 × 3	-	89.2
PointNet++ [12]	1024 × 3	-	90.7
PointNet++ [12]	1024 × 3 + norm	-	91.9
SCN [49]	1024 × 3	-	90.0
Kd-Net [11]	2 <sup>15</sup> × 3	94.0	91.8
KC-Net [15]	1024 × 3	94.4	91.0
PointCNN [13]	1024 × 3	-	91.7
DGCNN [18]	1024 × 3	-	92.2
SO-Net [14]	2048 × 3	94.1	90.9
SO-Net [14]	5000 × 3 + norm	95.7	93.4
Point2Sequence [7]	2048 × 3	95.3	92.6
Ψ-tree [8]	2048 × 3	94.6	92.0
PATs [52]	1024 × 3	-	92.2
PointWeb [20]	1024 × 3	-	92.3
A-CNN [3]	1024 × 3	95.5	92.6
RS-CNN [19]	1024 × 3	-	93.6
PointConv [53]	1024 × 3	-	92.5
3DCapsule [46]	1024 × 3	-	91.5
Point2SpatialCapsule(Ours)	1024 × 3	<b>95.8</b>	93.4
Point2SpatialCapsule(Ours)	1024 × 3 + norm	<b>95.9</b>	<b>93.7</b>

TABLE II  
THE SHAPE RETRIEVAL ACCURACY IN TERMS OF MAPS ON MODELNET10 AND MODELNET40.

Method	Input	ModelNet10	ModelNet40
PointNet* [1]	1024 × 3	67.98	62.41
PointNet++* [12]	1024 × 3	72.52	68.97
3DShapeNets [51]	Multi-View	68.26	49.23
DeepPano [28]	Multi-View	84.18	76.81
MVCNN [54]	Multi-View	-	83.0
PANORAMA [29]	Multi-View	87.39	83.45
GIFT [25]	Multi-View	91.12	81.94
SequenceViews [6]	Multi-View	89.55	89.0
SPNet [55]	Multi-View	<b>94.20</b>	85.21
VIPGAN [4]	Multi-View	90.69	89.23
Point2SpatialCapsule(Ours)	1024 × 3	<b>93.43</b>	<b>89.43</b>

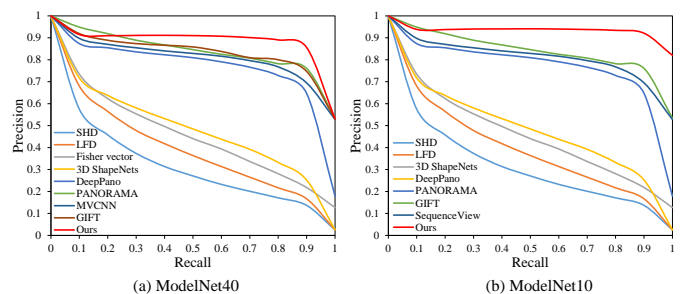


Fig. 9. The comparison of precision and recall curves obtained by different methods under (a) ModelNet40 and (b) ModelNet10.

As shown in Table II, our method has achieved a comparable retrieval accuracy compared with multi-view based methods on both ModelNet10 and ModelNet40. Specifically, Point2SpatialCapsule achieves the best retrieval accuracy 89.43% on ModelNet40, among all reported retrieval results. Point2Capsules achieves the second place result (93.43%) on ModelNet10, which is slight lower than SPNet [55] by 0.77%. However, Point2SpatialCapsule still beats SPNet by 4.22% on ModelNet40 in terms of mAPs, which shows a more balanced

TABLE III  
THE ACCURACIES (%) OF PART SEGMENTATION ON SHAPENET PART SEGMENTATION DATASET.

	mean	Intersection over Union (IoU)															
		air.	bag	cap	car	cha.	ear.	gui.	kni.	lam.	lap.	mot.	mug	pis.	roc.	ska.	tab.
# SHAPES		2690	76	55	898	3758	69	787	392	1547	451	202	184	283	66	152	5271
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
PointNet++ [12]	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
SCN [49]	84.6	83.8	80.8	83.5	79.3	90.5	69.8	91.7	86.5	82.9	96.0	69.2	93.8	82.5	62.9	74.4	80.8
Point2Seq [7]	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
Kd-Net [11]	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
O-CNN [21]	85.2	84.2	86.9	84.6	74.1	90.8	81.4	91.3	87.0	82.5	94.9	59.0	94.9	79.7	55.2	69.4	84.2
KCNet [15]	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 [18]	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
SO-Net [14]	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
RS-CNN [19]	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
PointCNN [13]	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	84.21	64.23	80.0	83.0
Ours	85.3	83.5	83.4	88.5	77.6	90.8	79.4	90.9	86.9	84.3	95.4	71.7	95.3	82.6	60.6	75.3	82.5

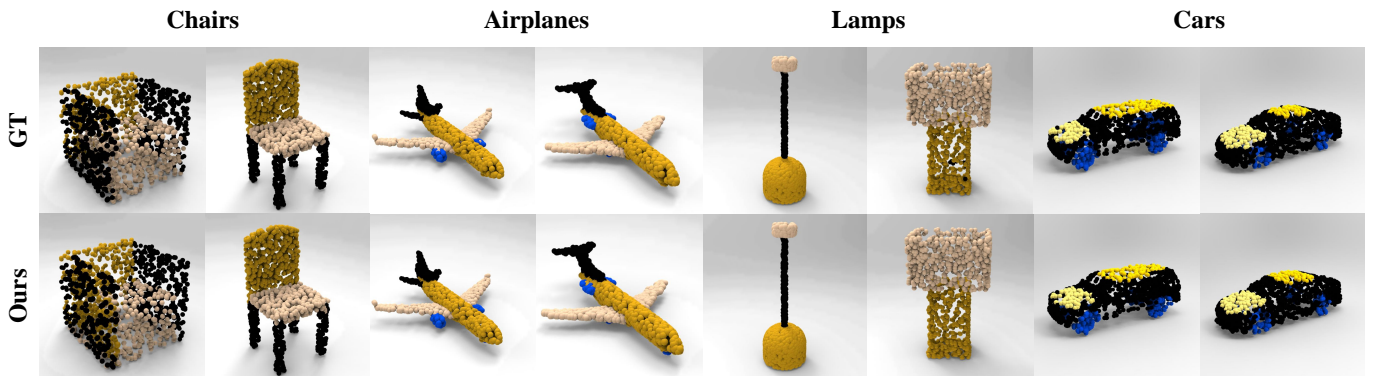


Fig. 10. Visualization of part segmentation results. In each shape pair, the first row is the ground truth (GT), and the second row is our predicted result. Parts with the same color are in the same part class.

performance of Point2SpatialCapsule over different scales of datasets. The comparison of precision-recall (PR) curves under ModelNet40 and ModelNet10 are shown in Fig. 9, where the results of Point2SpatialCapsule show the high performance for the 3D shape retrieval task.

The better performance of Point2SpatialCapsule can be dedicated to the following two reasons. First, the Point2SpatialCapsule is able to learn to encode the spatial locations of local features, which can produce a more discriminative representation for point clouds. Second, the digit capsules provide a more interpretable features for representing the point clouds, which is the length vector. Compared with the traditional single vector representations, in which the high-level characteristics are implicitly encoded in the latent feature space, the length of digit capsule explicitly indicates the probability that the characteristics appear in the point clouds. Therefore, using the distance between length vectors of digit capsules is more effective and interpretable for 3D shape retrieval.

#### D. 3D Shape Part Segmentation

In Table III, we also report the performance of Point2SpatialCapsule on the part segmentation task in terms of the Intersection over Union (IoU) [1]. As shown in Table III, our Point2SpatialCapsule achieves the mean instance IoU of 85.3%, which outperforms the baseline method Point-

Net++ on 13 categories out of total 16 categories. Note that, same as PointNet++, Point2SpatialCapsule also employs the multi-scale sampling and grouping strategy for local feature extraction. Therefore, the experimental results prove that Point2Sequence improves the quality of local feature extraction, and leads to the better performance on the segmentation task. Fig. 10 visualizes some examples of our segmentation results, where our results are highly consistent with the ground truth.

Note that, segmentation application needs discriminative features of local regions. Although Point2Capsule is proposed for global shape features by encoding the information of spatial locations in local regions, rather than producing more discriminative features of local regions like RS-CNN [19] and PointCNN [13], we still achieve comparable results in segmentation results.

#### E. Ablation Studies

In this section, we keep the settings of the network the same as described in Sec.III, except for the specified part for ablation study. We first investigate the influence of each part to our model, and then we analyze three important hyper-parameters in terms of classification accuracy on ModelNet40.

##### 1) The Influence of Each Part to Point2SpatialCapsule:

In order to investigate the effect of each part in Point2SpatialCapsule, we develop and evaluate three different

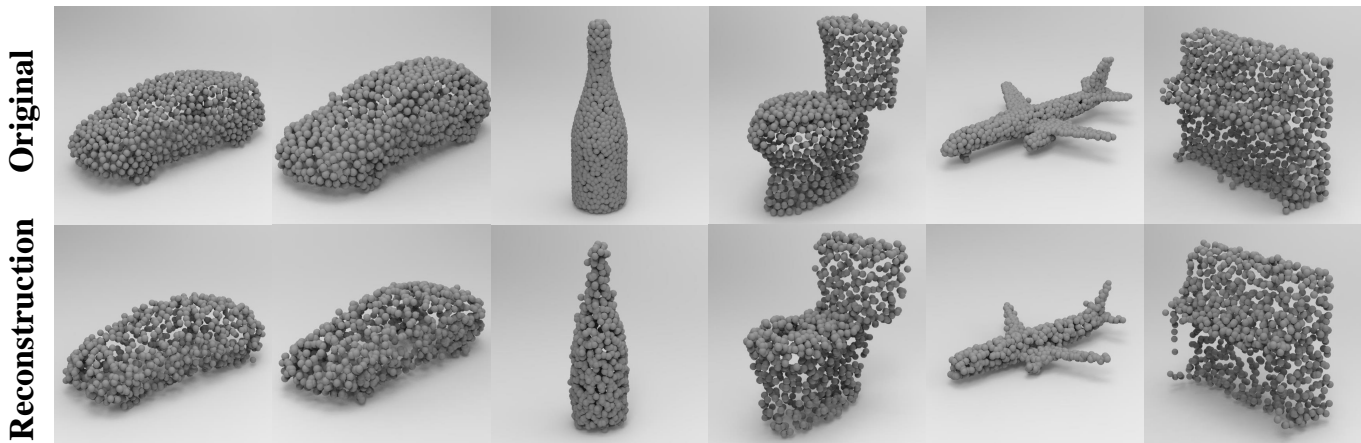


Fig. 11. The visualization of reconstruction results on the test set of ModelNet40. The top roll is the input original point cloud, and the bottom roll is the reconstructed point cloud from the Point2SpatialCapsule’s reconstruction network.

TABLE IV  
THE EFFECT OF EACH PART OF POINT2SPATIALCAPSULE ON MODELNET40.

Model	No-Multi	No-VLAD	No-Caps	Full-Model
Acc (%)	92.5	91.4	92.1	<b>93.44</b>

variations of our model as follows. (1) ‘No-Multi’ is the model without multi-scale shuffling, where the output of the multi-scale feature extractor is the direct input to the soft-assignment layer. (2) ‘No-VLAD’ is the model without the geometric feature aggregation, where the output of multi-scale feature extraction layer is directly reshaped as spatial-aware capsules and input to the dynamic routing layer. (3) ‘No-Caps’ is the model without the capsule net, where the output of geometric feature aggregation module is concatenated as a single vector and directly fed into the fully-connected layer for shape classification.

The experimental results are shown in Table IV. From the results we can find that each part of Point2SpatialCapsule contributes to the model performance. We note that, the No-VLAD model achieves the worst performance among the four models, which means that directly applying the capsule network on the point cloud impairs the model’s representational ability. The result of No-VLAD model supports our point of view that dynamic routing cannot learn the log priors directly from the disordered point clouds, and verifies the effectiveness of the proposed geometric feature aggregation. The results of No-Multi proves the importance of applying multi-scale shuffling for smoothing the perceived range between features of different scales. The significant improvement of Full-Model compared to the No-Caps verifies the superior advantage of capsule network for aggregating local features in point cloud recognition.

TABLE V  
THE EFFECT OF THE ITERATIONS OF DYNAMIC ROUTING ON MODELNET40.

Iterations	1	3	5
Acc (%)	<b>93.44</b>	92.22	91.98

2) *The Analysis of Capsules Net*: Following the common practice of [16], we investigate the influence of iterations in dynamic routing. As shown in Table V, we report the model performance with 1, 3 and 5 iterations of dynamic routing. According to [16], multiple iterations will increase the model’s learning ability but may also cause the problem of overfitting. As for Point2SpatialCapsule, we find that dynamic routing with 1 iteration is already enough for learning the point cloud features.

TABLE VI  
THE INFLUENCE OF THE NUMBER OF CLUSTER CENTERS IN GEOMETRIC FEATURE AGGREGATION MODULE ON MODELNET40.

Number	16	32	64	128
Acc (%)	92.70	92.93	<b>93.44</b>	93.21

3) *The Analysis of Geometric Feature Aggregation*: We also analysis the influence of cluster centers in NetVLAD. As shown in Table VI, the model achieves the best result with 64 cluster centers. The explanations are two-fold: (1) the small number of cluster centers could reduce the representational ability of feature embeddings; (2) the slight reduce in performance of the large number of cluster centers is the result of producing similar feature embeddings, which leads to the information redundancy and hinders the model learning more discriminative local features.

TABLE VII  
THE INFLUENCE OF RECONSTRUCTION LOSS ON MODELNET40.

$\alpha$	$10^{-3}$	$10^{-4}$	$10^{-5}$	0
Acc (%)	91.69	<b>93.44</b>	92.79	92.46

4) *The Analysis of Reconstruction Loss*: In Table VII, we discuss the influence of reconstruction loss, where  $\alpha$  is the weight factor as specified in Eq. (10). From the results, we find that a large  $\alpha$  leads to the decreasing of model performance, which in our opinion is the result of a slower learning process cause by the large reconstruction loss weight, especially during the early stage of training. On the other hand, the experimental results also prove the reconstruction loss useful. Compared

with a small weight ( $\alpha = 10^{-5}$ ) and the model without reconstruction ( $\alpha = 0$ ), the model with  $\alpha = 10^{-4}$  outperforms them by 0.65% and 0.98%, respectively. In Fig. 11, we visualize the reconstruction results on the test set of ModelNet40, from which we can find that Point2SpatialCapsule can learn to produce a relatively satisfactory result, despite of the simple reconstruction network employed in the model.

## V. CONCLUSIONS

In this paper, we propose a spatial-aware network, named Point2SpatialCapsule, to jointly aggregate geometric feature and spatial relationships of local regions on point cloud. The proposed Point2SpatialCapsule has a wide range of potential applications, which can be combined with other local feature extraction methods of multi-scale regions for learning the global shape representation of 3D point clouds. Compared with the previous feature aggregation methods, Point2SpatialCapsule has the ability to integrate both the geometric features of local regions and the spatial relationships among them. The features of local regions are aggregated by spatial-aware capsules with dynamic routing, which can preserve the spatial relationships between the extracted features. Experiments show that our network can achieve superior performance on point cloud classification, retrieval and part segmentation tasks under different datasets.

## REFERENCES

- [1] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep learning on point sets for 3D classification and segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [2] Z. Han, X. Wang, C.-M. Vong, Y.-S. Liu, M. Zwicker, and C. Chen, "3DViewGraph: Learning global features for 3D shapes from a graph of unordered views with attention," in *International Joint Conference on Artificial Intelligence*, 2019.
- [3] A. Komarichev, Z. Zhong, and J. Hua, "A-CNN: Annularly convolutional neural networks on point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7421–7430.
- [4] Z. Han, M. Shang, Y.-S. Liu, and M. Zwicker, "View inter-prediction GAN: Unsupervised representation learning for 3D shapes by learning global shape memories to support local view predictions," in *33rd AAAI Conference on Artificial Intelligence*, 2019.
- [5] Z. Han, X. Liu, Y.-S. Liu, and M. Zwicker, "Parts4Feature: Learning 3D global features from generally semantic parts in multiple views," in *International Joint Conference on Artificial Intelligence*, 2019.
- [6] Z. Han, M. Shang, Z. Liu, C.-M. Vong, Y.-S. Liu, J. Han, M. Zwicker, and C. P. Chen, "SeqViews2SeqLabels: Learning 3D global features via aggregating sequential views by RNN with attention," *IEEE Transactions on Image Processing*, vol. 28, no. 2, pp. 658–672, 2019.
- [7] X. Liu, Z. Han, Y.-S. Liu, and M. Zwicker, "Point2Sequence: Learning the shape representation of 3D point clouds with an attention-based sequence to sequence network," in *33rd AAAI Conference on Artificial Intelligence*, 2019.
- [8] H. Lei, N. Akhtar, and A. Mian, "Octree guided CNN with spherical kernels for 3D point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [9] X. Wang, S. Liu, X. Shen, C. Shen, and J. Jia, "Associatively segmenting instances and semantics in point clouds," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4096–4105.
- [10] J. Hou, A. Dai, and M. Nießner, "3D-SIS: 3D semantic instance segmentation of RGB-D scans," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4421–4430.
- [11] R. Klokov and V. Lempitsky, "Escape from cells: Deep kd-networks for the recognition of 3D point cloud models," in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 863–872.
- [12] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep hierarchical feature learning on point sets in a metric space," in *Advances in Neural Information Processing Systems*, 2017, pp. 5099–5108.
- [13] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "PointCNN: Convolution on x-transformed points," in *Advances in Neural Information Processing Systems*, 2018, pp. 820–830.
- [14] J. Li, B. M. Chen, and G. H. Lee, "SO-Net: Self-organizing network for point cloud analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 9397–9406.
- [15] Y. Shen, C. Feng, Y. Yang, and D. Tian, "Mining point cloud local structures by kernel correlation and graph pooling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [16] S. Sabour, N. Frosst, and G. E. Hinton, "Dynamic routing between capsules," in *Advances in Neural Information Processing Systems*, 2017, pp. 3856–3866.
- [17] Y. Xu, T. Fan, M. Xu, L. Zeng, and Y. Qiao, "SpiderCNN: Deep learning on point sets with parameterized convolutional filters," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 87–102.
- [18] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph CNN for learning on point clouds," *arXiv:1801.07829*, 2018.
- [19] Y. Liu, B. Fan, S. Xiang, and C. Pan, "Relation-shape convolutional neural network for point cloud analysis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8895–8904.
- [20] H. Zhao, L. Jiang, C.-W. Fu, and J. Jia, "PointWeb: Enhancing local neighborhood features for point cloud processing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5565–5573.
- [21] P.-S. Wang, Y. Liu, Y.-X. Guo, C.-Y. Sun, and X. Tong, "O-CNN: Octree-based convolutional neural networks for 3D shape analysis," *ACM Transactions on Graphics*, vol. 36, no. 4, p. 72, 2017.
- [22] G. Riegler, A. O. Ulusoy, and A. Geiger, "OctNet: Learning deep 3D representations at high resolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [23] H. You, Y. Feng, R. Ji, and Y. Gao, "PVNet: A joint convolutional network of point cloud and multi-view for 3D shape Recognition," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1310–1318.
- [24] Z. Han, H. Lu, Z. Liu, C.-M. Vong, Y.-S. Liu, M. Zwicker, J. Han, and C. P. Chen, "3D2SeqViews: Aggregating sequential views for 3D global feature learning by CNN with hierarchical attention aggregation," *IEEE Transactions on Image Processing*, vol. 28, no. 8, pp. 3986–3999, 2019.
- [25] S. Bai, X. Bai, Z. Zhou, Z. Zhang, Q. Tian, and L. J. Latecki, "GIFT: Towards scalable 3D shape retrieval," *IEEE Transactions on Multimedia*, vol. 19, no. 6, pp. 1257–1271, 2017.
- [26] Z. Han, X. Wang, Y.-S. Liu, and M. Zwicker, "Multi-Angle Point cloud-VAE: Unsupervised feature learning for 3D point clouds from multiple angles by joint self-reconstruction and half-to-half prediction," in *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- [27] Z. Han, M. Shang, X. Wang, Y.-S. Liu, and M. Zwicker, "Y2Seq2Seq: Cross-modal representation learning for 3D shape and text by joint reconstruction and prediction of view and word sequences," *The Thirty-Third AAAI Conference on Artificial Intelligence (AAAI)*, 2019.
- [28] B. Shi, S. Bai, Z. Zhou, and X. Bai, "DeepPano: Deep panoramic representation for 3-d shape recognition," *IEEE Signal Processing Letters*, vol. 22, no. 12, pp. 2339–2343, 2015.
- [29] S. K. T. T. and I. Pratikakis, "Exploiting the PANORAMA representation for convolutional neural network classification and retrieval," in *3DOR*, 2017.
- [30] Y.-S. Liu, K. Ramani, and M. Liu, "Computing the inner distances of volumetric models for articulated shape description with a visibility graph," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 12, pp. 2538–2544, 2011.
- [31] C. Wang, Y.-S. Liu, M. Liu, J.-H. Yong, and J.-C. Paul, "Robust shape normalization of 3D articulated volumetric models," *Computer-Aided Design*, vol. 44, no. 12, pp. 1253–1268, 2012.
- [32] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3D ShapeNets: A deep representation for volumetric shapes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1912–1920.
- [33] Z. Han, Z. Liu, J. Han, C. Vong, S. Bu, and C. Chen, "Unsupervised learning of 3D local features from raw voxels based on a novel permutation voxelization strategy," *IEEE Transactions on Cybernetics*, vol. 49, no. 2, pp. 481–494, 2019.

- [34] A. Sharma, O. Grau, and M. Fritz, “VConv-DAE: Deep volumetric shape learning without object labels,” in *European Conference on Computer Vision*. Springer, 2016, pp. 236–250.
- [35] J. Xie, Y. Fang, F. Zhu, and E. Wong, “DeepShape: Deep learned shape descriptor for 3D shape matching and retrieval,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1275–1283.
- [36] J. Xie, M. Wang, and Y. Fang, “Learned binary spectral shape descriptor for 3D shape correspondence,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3309–3317.
- [37] Z. Han, Z. Liu, C.-M. Vong, Y.-S. Liu, S. Bu, J. Han, and C. P. Chen, “BoSCC: Bag of spatial context correlations for spatially enhanced 3D shape representation,” *IEEE Transactions on Image Processing*, vol. 26, no. 8, pp. 3707–3720, 2017.
- [38] —, “Deep Spatiality: Unsupervised learning of spatially-enhanced global and local 3D features by deep neural network with coupled softmax,” *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 3049–3063, 2018.
- [39] A. Jaiswal, W. Abdalimageed, Y. Wu, and P. Natarajan, “CapsuleGAN: Generative adversarial capsule network,” in *Proceedings of the European Conference on Computer Vision*, 2018.
- [40] R. Lalonde and U. Bagci, “Capsules for object segmentation,” *arXiv:1804.04241*, 2018.
- [41] M. Yang, W. Zhao, J. Ye, Z. Lei, Z. Zhao, and S. Zhang, “Investigating capsule networks with dynamic routing for text classification,” in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 3110–3119.
- [42] L. Xiao, H. Zhang, W. Chen, Y. Wang, and Y. Jin, “MCapsNet: Capsule network for text with multi-task learning,” in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 4565–4574.
- [43] Z. Ningyu, D. Shumin, S. Zhanlin, C. Xi, Z. Wei, and C. Huajun, “Attention-based capsule networks with dynamic routing for relation extraction,” in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2018.
- [44] K. Burak, A. Ayesha, and V. Senem, “Object classification from 3D volumetric data with 3D capsule networks,” in *IEEE Global Conference on Signal and Information Processing*, 2018.
- [45] Y. Zhao, T. Birdal, H. Deng, and F. Tombari, “3D Point Capsule Networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [46] C. Ali and P. Lars, “3DCapsule: Extending the capsule architecture to classify 3D point clouds,” in *IEEE Winter Conference on Applications of Computer Vision*, 2019.
- [47] M. A. Uy and G. H. Lee, “PointNetVLAD: Deep point cloud based retrieval for large-scale place recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4470–4479.
- [48] R. Arandjelovic, P. Gronat, A. Torii, T. Pajdla, and J. Sivic, “NetVLAD: CNN architecture for weakly supervised place recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5297–5307.
- [49] S. Xie, S. Liu, Z. Chen, and Z. Tu, “Attentional ShapeContextNet for point cloud recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4606–4615.
- [50] W. Shi, J. Caballero, F. Huszár, J. Totz, A. P. Aitken, R. Bishop, D. Rueckert, and Z. Wang, “Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 1874–1883.
- [51] M. Savva, F. Yu, H. Su, M. Aono, B. Chen, D. Cohen-Or, W. Deng, H. Su, S. Bai, X. Bai *et al.*, “SHREC’16 track large-scale 3D shape retrieval from ShapeNet core55,” in *Proceedings of the Eurographics Workshop on 3D Object Retrieval*, 2016.
- [52] J. Yang, Q. Zhang, B. Ni, L. Li, J. Liu, M. Zhou, and Q. Tian, “Modeling point clouds with self-attention and gumbel subset sampling,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3323–3332.
- [53] W. Wu, Z. Qi, and L. Fuxin, “PointConv: Deep convolutional networks on 3D point clouds,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9621–9630.
- [54] C. R. Qi, H. Su, M. Niebner, A. Dai, M. Yan, and L. J. Guibas, “Volumetric and multi-view CNNs for object classification on 3D data,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5648–5656.
- [55] M. Yavartanoo, E. Y. Kim, and K. M. Lee, “SPNet: Deep 3D object classification and retrieval using stereographic projection,” in *Proceedings*

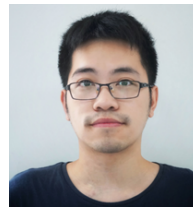
*of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.



**Xin Wen** received the B.S. degree in engineering management from the Tsinghua University, China, in 2012. He is currently the PhD student with the School of Software, Tsinghua University. His research interests include deep learning, shape analysis and pattern recognition, and NLP.



**Zhizhong Han** received the Ph.D. degree from Northwestern Polytechnical University, China, 2017. He is currently a Post-Doctoral Researcher with the Department of Computer Science, University of Maryland at College Park, College Park, USA. He is also a Research Member of the BIM Group, Tsinghua University, China. His research interests include machine learning, pattern recognition, feature learning, and digital geometry processing.



**Xinhai Liu** received the B.S. degree in computer science and technology from the Huazhong University of Science and Technology, China, in 2017. He is currently the PhD student with the School of Software, Tsinghua University. His research interests include deep learning, 3D shape analysis and 3D pattern recognition.



**Yu-Shen Liu** (M’18) received the B.S. degree in mathematics from Jilin University, China, in 2000, and the Ph.D. degree from the Department of Computer Science and Technology, Tsinghua University, Beijing, China, in 2006. From 2006 to 2009, he was a Post-Doctoral Researcher with Purdue University. He is currently an Associate Professor with the School of Software, Tsinghua University. His research interests include shape analysis, pattern recognition, machine learning, and semantic search.