# SoftPoolNet: Shape Descriptor for Point Cloud Completion and Classification

Yida Wang[1], David Joseph Tan[2], Nassir Navab[1], and Federico Tombari[1,2]

[1] Technische Universität München
[2] Google Inc.

**Abstract.** Point clouds are often the default choice for many applications as they exhibit more flexibility and efficiency than volumetric data. Nevertheless, their unorganized nature – points are stored in an unordered way – makes them less suited to be processed by deep learning pipelines. In this paper, we propose a method for 3D object completion and classification based on point clouds. We introduce a new way of organizing the extracted features based on their activations, which we name soft pooling. For the decoder stage, we propose regional convolutions, a novel operator aimed at maximizing the global activation entropy. Furthermore, inspired by the local refining procedure in Point Completion Network (PCN), we also propose a patch-deforming operation to simulate deconvolutional operations for point clouds. This paper proves that our regional activation can be incorporated in many point cloud architectures like AtlasNet and PCN, leading to better performance for geometric completion. We evaluate our approach on different 3D tasks such as object completion and classification, achieving state-of-the-art accuracy.

## 1 Introduction

Point clouds are unorganized sparse representations of a 3D point set. Compared to other common representations for 3D data such as 3D meshes and voxel maps, they are simple and flexible, while being able to store fine details of a surface. For this reason, they are frequently employed for many applications within 3D perception and 3D computer vision such as robotic manipulation and navigation, scene understanding, and augmented/virtual reality. Recently, deep learning approaches have been proposed to learn from point clouds for 3D perception tasks such as point cloud classification [4,14,18,19] or point cloud segmentation [11,13,14,17,23]. Among them, one of the key breakthroughs in handling unorganized point clouds was proposed by PointNet [18], introducing the idea of a max pooling in the feature space to yield permutation invariance.

An interesting emerging research trend focusing on 3D data is the so-called 3D completion, where the geometry of a partial scene or object acquired from a single viewpoint, *e.g.* through a depth map, is completed of the missing part due to (self-)occlusion as visualized in Fig. 1. This can be of great use to aid standard 3D perception tasks such as object modeling, scene registration, part-based segmentation and object pose estimation. Most approaches targeting 3D completion
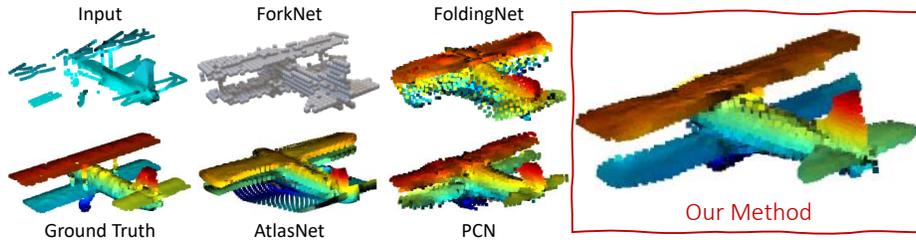
Fig. 1: This paper proposes a method that reconstructs 3D point cloud models with more fine details.

have been proposed for volumetric approaches, since 3D convolutions are naturally suited to this 3D representation. Nevertheless, such approaches bring in the limitations of this representation, including loss of fine details due to discretization and limitations in scaling with the 3D size. Recently, a few approaches have explored the possibility of learning to complete a point cloud [9,25,26].

This paper proposes an encoder-decoder architecture called SoftPoolNet, which can be employed for any task that processes a point cloud as input in order to regress another point cloud as output. One of the tasks and a main focus for this work is 3D object completion from a partial point cloud.

The theoretical contribution of SoftPoolNet is twofold. We first introduce soft pooling, a new module that replaces the max-pooling operator in PointNet by taking into account multiple high-scoring features rather than just one. The intuition is that, by keeping multiple features with high activations rather than just the highest, we can retain more information while keeping the required permutation invariance. A second contribution is the definition of a regional convolution operator that is used within the proposed decoder architecture. This operator is designed specifically for point cloud completion and relies on convolving local features to improve the completion task with fine details.

In addition to evaluating SoftPoolNet for point cloud completion, we also evaluate on the point cloud classification to demonstrate its applicability to general point cloud processing tasks. In both evaluations, SoftPoolNet obtains state of the art results on the standard benchmarks.

## 2    Related work

*Volumetric completion.* Object [7] and scene completion [20,22] are typically carried out by placing all observed elements into a 3D grid with fixed resolution. 3D-EPN [7] completes a single object using 3D convolutions while 3D-RecGAN [24] further improves the completion performance by using discriminative training. As scene completion contains objects in different scales and more random relative position among all of them, SSCNet [20] proposes a 3D volumetric semantic completion architecture using dilated convolutions to recognize

objects with different scales. ForkNet [22] designs a multi-branch architecture to generate realistic training data to supplement the training.

*Point cloud completion.* Object completion based on point cloud data change partial geometries without using a 3D fixed grid. They represent completed shapes as a set of points with 3D coordinates. For instance, FoldingNet [25] deforms a 2D grid from a global feature such as PointNet [18] feature to an output with a desirable shape. AtlasNet [9] generates an object with a set of local patches to simulate mesh data. But overlaps between different local patches makes the reconstruction noisy. MAP-VAE [10] predicts the completed shape by joining the observed part with the estimated counterpart.

*CNNs for point clouds.* Existing works like PointConv [23] and PointCNN [13] index each point with $k$-nearest neighbour search to find local patches, where they then apply the convolution kernels on those local patches. Regarding point cloud deconvolutional operations, FoldingNet [25] uses a 2D grid to help generate a 3D point cloud from a single feature. PCN [26] further uses local FoldingNet to obtain a fine-grained output from a coarse point cloud with low resolution which could be regarded as an alternative to point cloud deconvolution.

## 3   Soft pooling for point features

Given the partial scan of an object, the input to our network is a point cloud with $N_{\text{in}}$ points written in the matrix form as $\mathbf{P}_{\text{in}} = [\mathbf{x}_i]_{i=1}^{N_{\text{in}}}$ where each point is represented as the 3D coordinates $\mathbf{x}_i = [x_i, y_i, z_i]$. We then convert each point into a feature vector $\mathbf{f}_i$ with $N_f$ elements by projecting every point with a point-wise multi-layer perceptron [18] (MLP) $\mathbf{W}_{\text{point}}$ with three layers. Thus, similar to $\mathbf{P}_{\text{in}}$, we define the $N_{\text{in}} \times N_f$ feature matrix as $\mathbf{F} = [\mathbf{f}_i]_{i=1}^{N_{\text{in}}}$. Note that we applied a softmax function to the output neuron of MLP so that the elements in $\mathbf{f}_i$ ranges between 0 and 1.

The main challenge when processing a point cloud is its unstructured arrangement. This implies that changing the order of the points in $\mathbf{P}_{\text{in}}$ describes the same point cloud, but generates a different feature matrix that flows into our architecture with convolutional operators. To solve this problem, we propose to organize the feature vectors in $\mathbf{F}$ so that their $k$-th element are sorted in a descending order, which is denoted as $\mathbf{F}'_k$. Note that $k$ should not be larger than $N_f$. A *toy example* of this process is depicted in Fig. 2(a) where we assume that there are only five points in the point cloud and arrange the five feature vectors from $\mathbf{F} = [\mathbf{f}_i]_{i=1}^{5}$ to $\mathbf{F}'_k = [\mathbf{f}_i]_{i=\{3,5,1,2,4\}}$ by comparing the $k$-th element of each vector. Repeating this process for all the $N_f$ elements in $\mathbf{f}_i$, all $\mathbf{F}'_k$ together result to a 3D tensor $\mathbf{F}' = [\mathbf{F}'_1, \mathbf{F}'_2, \dots \mathbf{F}'_{N_f}]$ with the dimension of $N_{\text{in}} \times N_f \times N_f$. As a result, any permutation of the points in $\mathbf{P}_{\text{in}}$ generate the same $\mathbf{F}'$.

Sorting the feature vectors in a descending order highlights the ones with the highest activation values. Thus, by selecting the first $N_r$ feature vectors from all the $\mathbf{F}'_k$ as shown in Fig. 2(b), we assemble $\mathbf{F}^*$ that accumulates the features
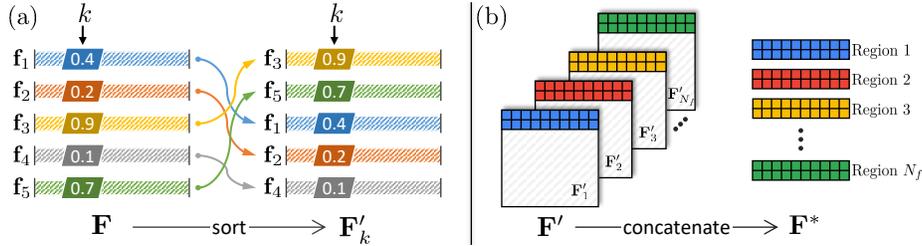
Fig. 2: Toy examples of (a) sorting the the $k$-th element of the vectors in the feature matrix $\mathbf{F}$ to build $\mathbf{F}'_k$ and consequently $\mathbf{F}'$ and (b) concatenation of the first $N_r$ rows of $\mathbf{F}'_k$ to construct the 2D matrix $\mathbf{F}^*$ which corresponds to the regions with high activations.

with the highest activations. Altogether, the output of soft pooling is the $N_f \cdot N_r$ point features. Since each feature vector corresponds to a point in $\mathbf{P}_{\text{in}}$, we can interpret the first $N_r$ feature vectors as a region in the point cloud. The effects of the activations on the 3D reconstruction are illustrated in Fig. 3, where the point cloud is divided into $N_f$ regions. Later in Sec. 6, we discuss on how to learn $\mathbf{W}_{\text{point}}$ by incorporating these regions. That section introduces several loss functions which optimize towards entropy, Chamfer distance and earth-moving distance such that each point is optimized to fall into only one region and to be selected for $\mathbf{F}^*$ by maximizing the $k$-th element of the feature vector associated to the same region.

Similar to PointNet [18], we also rely on MLP to build the feature matrix $\mathbf{F}$. However, PointNet directly applies max-pooling on $\mathbf{F}$ to produce a vector while we try to generalize this approach and sort the feature vectors in order to assemble a matrix $\mathbf{F}^*$ as illustrated in Fig. 2. Considering the distinction between the two approaches, we refer our approach as *soft pooling*. Fundamentally, in addition to the increased amount of information from our feature vectors, the advantage of our method is the ability to apply regional convolutional operations to $\mathbf{F}^*$, as discussed in Sec. 4. The differences are evident in Fig. 4, where the proposed method achieves detailed results on reconstructing all the six legs while
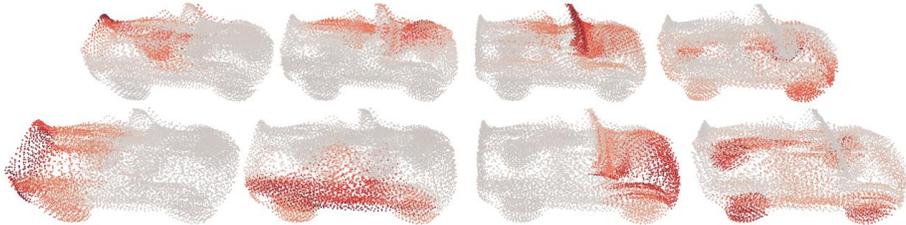


Fig. 3: Deconstructing the learned regions (unsupervised) that correspond to different parts of the car.

(a) Input          (b) Ground Truth          (c) PointNet          (d) Ours
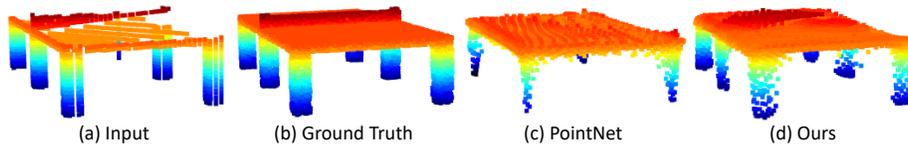
Fig. 4: Comparison of our method and PointNet [18] where PointNet reconstructs the more typical four-leg table instead of six in (c).

PointNet follows the more generic structure of the table with four. This proves that *soft pooling* makes our decoder able to take all observable geometries into account to complete the shape, while the max-pooled PointNet feature cannot reveal the rarely seen geometry.

## 4    Regional convolution

Operating on $\mathbf{F}^*$, we introduce the convolutional kernel $\mathbf{W}_{\text{conv}}$ that transforms $\mathbf{F}^*$ to a new set of points $\mathbf{P}_{\text{conv}}$ by taking several point features into consideration. We structure $\mathbf{W}_{\text{conv}}$ with a dimension of $N_p \times N_f \times 3$ where $N_p$ represent the number of points which are taken into consideration such that

$$\mathbf{P}_{\text{conv}}(i,j) = \sum_{l=1}^{N_f}\sum_{k=1}^{N_p} \mathbf{F}^*(i+k,l)\mathbf{W}_{\text{conv}}(j,k,l) \ . \tag{1}$$

Here, the kernel slides only inside each region of features without taking features from two different regions in one convolutional operation. As the kernel size allows it to cover $N_p$ features, we pad each region with $N_p-1$ duplicated samples at the end of each region in order to keep the output resolution the same as $N_{\text{in}}$. Experimentally, we tried different numbers of $N_p$ ranging from 4 to 64 and evaluated that 32 generates the best results. Learning the values in $\mathbf{W}_{\text{conv}}$ is discussed in Sec. 6.

In addition, we use a convolution stride which is set as a value smaller than $N_p$ to change the output resolution in terms of the number of point features. With a stride of $S$, we then take samples every $S$ point feature in $\mathbf{F}^*$. Notice that, by using a stride which is smaller than 1, we can also upsample $\mathbf{F}^*$ by interpolating $\frac{1}{S} - 1$ new points between two points then apply the convolution kernel again. This is an essential tool in reconstructing the object from a partial scan.

## 5    Network architecture

We build an encoder-decoder architecture which consists of MLP and our regional convolutions, respectively. Serving as the input to our network, we permutate the input scans and resample 1,024 points. If the partial scans have less than 1,024 points, we then duplicate the missing samples.
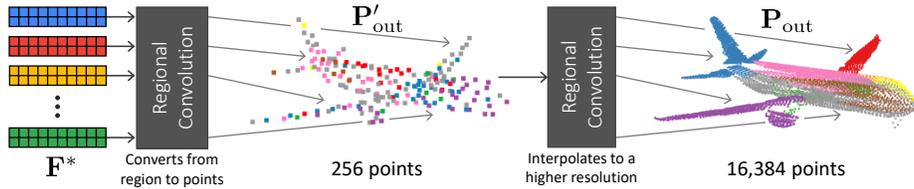
Fig. 5: Decoder architecture of SoftPoolNet with two regional convolution that converts the features from the regions to point clouds and interpolates from the coarse 256 points to a higher resolution with 16,384 points.

Our encoder is a point-wise MLP that generates the output neuron with a dimension of $[512, 512, 8]$. We then perform soft pooling as described in Sec. 3 that produces $\mathbf{F}^*$ with the size of $[256, 8]$ by setting $N_r$ to 32 and $N_f$ to 8, resulting an output of $N_f \cdot N_r = 256$ features.

Finally, for the decoder, we propose a two-stage point cloud completion architecture which is trained end-to-end. The output of the first is used as the input of the second point cloud completion network. Both of them produces the completed point cloud but with different resolutions. Illustrated in Fig. 5, we construct the decoder with two regional convolutions from Sec. 4. The first output $\mathbf{P}'_{\mathrm{out}}$ is fixed at 256 while the second $\mathbf{P}_{\mathrm{out}}$ produces a maximum resolution of 16,384.

## 6 Loss functions

During learning, we evaluate whether the predicted point feature $\mathbf{P}_{\mathrm{out}}$ matches the given ground truth $\mathbf{P}_{\mathrm{gt}}$ through the Chamfer distance. Similar to [9,25,26], we use the regression loss function for the shape completion from a point cloud

$$\mathcal{L}_{\mathrm{complete}}(\mathbf{W}_{\mathrm{point}}, \mathbf{W}_{\mathrm{conv}}) = \mathrm{Chamfer}(\mathbf{P}_{\mathrm{out}}, \mathbf{P}_{\mathrm{gt}}) . \tag{2}$$

We observed that there are two major drawbacks in using this loss function alone – the reconstructed surface tends to be either curved on the sharp edges such FoldingNet [25] or having noisy points appear on flat surfaces such as AtlasNet [9] and PCN [26]. In this work, we tackle these problems by finding local regions first, then by optimizing the inter- and intra-regional relationships.

Moreover, while FoldingNet [25] sacrifices local details to present the entire model with a single mesh having smooth surface, AtlasNet [9] and PCN [26] use local regions (or patches) to increase the details in the 3D model. However, both of them [9,26] have severe overlapping effects between adjacent regions which makes the generated object noisy and the regions discontinuous. To solve this problem, we aim at reducing the overlaps between two adjacent regions.

### 6.1    Learning activations through regional entropy

Considering that the dimension of a single feature is $N_f$, we can directly define $N_f$ regions for all features. Given the probabilities of regions to which the feature

$\mathbf{f}_i$ belong, we want to optimize the inter- and intra-regional relationships among the features. We directly present the probability of the feature $\mathbf{f}_i$ belonging to all $N_f$ regions by applying the softmax function to $\mathbf{f}_i$ as

$$P(\mathbf{f}_k, i) = \frac{\mathbf{f}_k[i]}{\sum_{j=1}^{N_f} \mathbf{f}_k[j]} \ . \tag{3}$$

Since the information entropy evaluates both the distribution and the confidence of the probabilities of a set of data, we define the feature entropy and the regional entropy based on the regional probability of the feature.

The goal of the inter-regional loss function is to similarly distribute the number of points throughout the regions. We define the regional entropy as

$$\mathcal{E}_r = -\frac{1}{B} \sum_{j=1}^{B} \sum_{i=1}^{R} \left[ \left( \frac{1}{N} \sum_{k=1}^{N} P(\mathbf{f}_k, i) \right) \cdot \log \left( \frac{1}{N} \sum_{k=1}^{N} P(\mathbf{f}_k, i) \right) \right] \tag{4}$$

where $B$ is the batch-size. Here, we want to maximize $\mathcal{E}_r$. Considering that the upper-bound of $\mathcal{E}_r$ is $-\log \frac{1}{R} = \log(R)$, we can then define the inter-regional loss function as

$$\mathcal{L}_{\text{inter}}(\mathbf{W}_{\text{point}}) = \log(R) - \mathcal{E}_r \tag{5}$$

in order to acquire a positive loss function. Once $\mathcal{E}_r$ is close to $\log(R)$, each region would contain similar amount of point features. Interestingly, we can select the number of regions by evaluating how much the regional entropy $\mathcal{E}_r$ differs from its upper-bound. The best number of regions should be the one with a small $\mathcal{L}_{\text{inter}}$. This is evaluated later in Table 6.

On the other hand, the goal of the intra-regional loss function is to boost the confidence of each feature to be in a single region. The intra-regional loss function then minimize the feature entropy

$$\mathcal{L}_{\text{intra}}(\mathbf{W}_{\text{point}}) = -\frac{1}{N} \frac{1}{B} \sum_{k=1}^{N} \sum_{j=1}^{B} \sum_{i=1}^{N_f} P(\mathbf{f}_k, i) \log P(\mathbf{f}_k, i) \ . \tag{6}$$

The optimum case of the feature entropy is for each feature to be a one-hot code, *i.e.* when only one element is 1 while the others are zero.

### 6.2   Reducing the overlapping regions

Although $\mathcal{L}_{\text{intra}}$ tries to make each point feature confident about the region to which it belongs, instances exist where many adjacent points would fall under different regions. For example, we observe in Fig. 6 that patches from different regions are stacked on top of each other, producing noisy reconstructions. Notably, this introduces unexpected results when fitting a mesh to the point cloud. Thus, we want to minimize region overlap by optimizing the network to restrict the connection between adjacent regions to their boundaries.
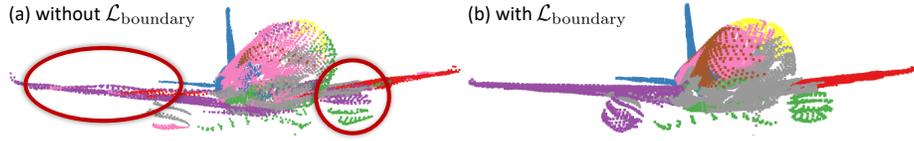
Fig. 6: Effects of without and with $\mathcal{L}_{\text{boundary}}$ where the wings are not planar and the engines are less visible in (a). Note that the colors represent different regions.
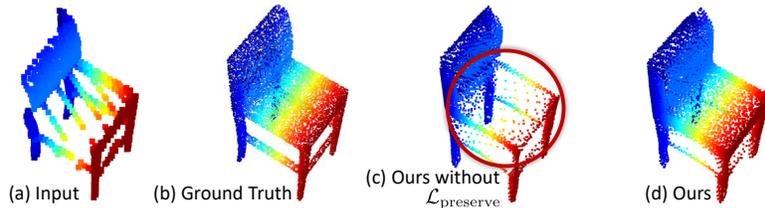


Fig. 7: Effects of without and with $\mathcal{L}_{\text{preserve}}$ where the seat is missing in (c).

First, each point is assigned to a region with the highest activation. All points that belong to region $i$ but has activation for region $j$ larger than a threshold $\tau$ are included in the set $\mathcal{B}_i^j$. Inversely, the points that belong to region $j$ but have activation for region $i$ larger than $\tau$ are added in the set $\mathcal{B}_j^i$. Note that, if both sets $\mathcal{B}_i^j$ and $\mathcal{B}_j^i$ are not empty, the regions $i$ and $j$ are then adjacent. Thus, by minimizing the Chamfer distance between $\mathcal{B}_i^j$ and $\mathcal{B}_j^i$, we can make the overlapping sets of points smaller such that the optimal result is a line. We then define the loss function for the boundary as

$$\mathcal{L}_{\text{boundary}}(\mathbf{W}_{\text{point}}, \mathbf{W}_{\text{conv}}) = \sum_{i=1}^{N_f}\sum_{j=i}^{N_f}\text{Chamfer}(\mathcal{B}_i^j, \mathcal{B}_j^i) \tag{7}$$

where both $\mathbf{W}_{\text{point}}$ and $\mathbf{W}_{\text{conv}}$ are optimized. After experimenting on different values of $\tau$ from 0.1 to 0.9, we set $\tau$ to be 0.3.

### 6.3   Preserving the features from MLP

After sorting and filtering the features to produce $\mathbf{F}^*$, some feature vectors in $\mathbf{F}^*$ are duplicated while some vectors from $\mathbf{F}$ are missing in $\mathbf{F}^*$. To avoid these, we introduce the loss function

$$\mathcal{L}_{\text{preserve}}(\mathbf{W}_{\text{point}}) = \text{Earth-moving}(\mathbf{F}^*, \mathbf{F}) \; . \tag{8}$$

Since the earth moving distance [12] is not efficient when the size of the samples is large, we then randomly select 256 vectors from $\mathbf{F}$ and $\mathbf{F}^*$. Considering that the feature dimensions in $\mathbf{F}$ and $\mathbf{F}^*$ are both $N_f$, the earth moving distance

then takes features with $N_f$ dimension as input. In practice, Fig. 7 visualizes the effects of $\mathcal{L}_{\text{preserve}}$ in the reconstruction, where removing this loss produce a large hole on the seat while incorporating this loss builds a well-distributed point cloud.

## 7    Experiments

For all evaluations, we train our model with an NVIDIA Titan V and parameterize it with a batch size of 8. Moreover, we apply the Leaky ReLU with a negative slope of 0.2 on the output of each regional convolution output.

### 7.1    Object completion on ShapeNet

We evaluate the performance of the geometric completion of a single object on the ShapeNet [5] database where training data are paired point clouds of the partial scanning and the completed shape. To make it comparable to other approaches, we adopt the standard 8 category evaluation [26] for a single object completion. As rotation errors are common in the partial scans, we further evaluate our approach against other works on the ShapeNet database with rotations. We also evaluate the performance on both high and low resolutions which contain 16,384 and 2,048 points, respectively.

We compare against other point cloud completion approaches such as PCN [26], FoldingNet [25], AtlasNet [9] and PointNet++ [19]. To show the advantages over volumetric completion, we also compare against 3D-EPN [24] and ForkNet [22] with an output resolution of $64 \times 64 \times 64$. Notably, we achieve the best results on most objects and in all types of evaluations as presented in Table 1, Table 2 and Table 3.

An interesting hypothesis is the capacity of $\mathcal{L}_{\text{boundary}}$ to be integrated in other existing approaches. Thus, Table 1 and Table 2 also evaluate this hypothesis and prove that this activation helps FoldingNet [25], PCN [26] and AtlasNet [9]

Output Resolution = 16,384

| Method | plane | cabinet | car | chair | lamp | sofa | table | vessel | Avg. |
|---|---|---|---|---|---|---|---|---|---|
| 3D-EPN [7] | 13.16 | 21.80 | 20.31 | 18.81 | 25.75 | 21.09 | 21.72 | 18.54 | 20.15 |
| ForkNet [22] | 9.08 | 14.22 | 11.65 | 12.18 | 17.24 | 14.22 | 11.51 | 12.66 | 12.85 |
| PointNet++ [19] | 10.30 | 14.74 | 12.19 | 15.78 | 17.62 | 16.18 | 11.68 | 13.52 | 14.00 |
| FoldingNet [25] | 5.97 | 10.80 | 9.27 | 11.25 | 12.17 | 11.63 | 9.45 | 10.03 | 10.07 |
| FoldingNet + $\mathcal{L}_{\text{boundary}}$ | 5.79 | 10.61 | 8.62 | 10.33 | 11.56 | 11.05 | 9.41 | 9.79 | 9.65 |
| PCN [26] | 5.50 | 10.63 | 8.70 | 11.00 | 11.34 | 11.68 | 8.59 | 9.67 | 9.64 |
| PCN + $\mathcal{L}_{\text{boundary}}$ | 5.13 | 9.12 | 7.58 | 9.35 | 9.40 | 9.31 | 7.30 | 8.91 | 8.26 |
| **Our Method** | **4.01** | **6.23** | **5.94** | **6.81** | **7.03** | **6.99** | **4.84** | **5.70** | **5.94** |

Table 1: Completion evaluated by means of the Chamfer distance (multiplied by $10^3$) with the output resolution of 16,384.

Output Resolution = 2,048

| Method | plane | cabinet | car | chair | lamp | sofa | table | vessel | *Avg.* |
|---|---|---|---|---|---|---|---|---|---|
| FoldingNet [25] | 11.18 | 20.15 | 13.25 | 21.48 | 18.19 | 19.09 | 17.80 | 10.69 | 16.48 |
| FoldingNet + $\mathcal{L}_{\mathrm{boundary}}$ | 11.09 | 19.95 | 13.11 | 21.27 | 18.22 | 19.06 | 17.62 | 10.10 | 16.30 |
| AtlasNet [9] | 10.37 | 23.40 | 13.41 | 24.16 | 20.24 | 20.82 | 17.52 | 11.62 | 17.69 |
| AtlasNet + $\mathcal{L}_{\mathrm{boundary}}$ | 9.25 | 22.51 | 12.12 | 22.64 | 18.82 | 19.11 | 16.50 | 11.53 | 16.56 |
| PCN [26] | 8.09 | 18.32 | 10.53 | 19.33 | 18.52 | 16.44 | 16.34 | 10.21 | 14.72 |
| PCN + $\mathcal{L}_{\mathrm{boundary}}$ | 6.39 | 16.32 | 9.30 | 18.61 | 16.72 | 16.28 | 15.29 | 9.00 | 13.49 |
| TopNet [21] | 5.50 | 12.02 | 8.90 | 12.56 | 9.54 | 12.20 | 9.57 | 7.51 | 9.72 |
| **Our Method** | **4.76** | **10.29** | **7.63** | **11.23** | **8.97** | **10.08** | **7.13** | **6.38** | **8.31** |
| – *without* $\mathcal{L}_{\mathrm{inter}}$ | 10.82 | 20.45 | 15.21 | 20.19 | 18.05 | 18.58 | 15.65 | 8.81 | 15.97 |
| – *without* $\mathcal{L}_{\mathrm{intra}}$ | 5.23 | 16.10 | 12.49 | 14.62 | 13.90 | 12.37 | 12.96 | 5.72 | 11.67 |
| – *without* $\mathcal{L}_{\mathrm{inter}}$, $\mathcal{L}_{\mathrm{intra}}$ | 10.91 | 20.54 | 15.27 | 20.28 | 18.16 | 18.66 | 15.75 | 8.91 | 16.06 |
| – *without* $\mathcal{L}_{\mathrm{boundary}}$ | 5.46 | 10.98 | 8.27 | 11.95 | 9.51 | 10.92 | 7.78 | 7.40 | 9.03 |
| – *without* $\mathcal{L}_{\mathrm{preserve}}$ | 10.29 | 19.75 | 14.13 | 19.35 | 17.88 | 18.21 | 15.23 | 8.11 | 15.37 |

Table 2: Completion evaluated using the Chamfer distance (multiplied by $10^3$) with the output resolution of 2,048.

Output Resolution = 1,024

| Method | plane | cabinet | car | chair | lamp | sofa | table | vessel | *Avg.* |
|---|---|---|---|---|---|---|---|---|---|
| 3D-EPN [7] | 6.20 | 7.76 | 8.70 | 7.68 | 10.73 | 8.08 | 8.10 | 8.17 | 8.18 |
| PointNet++ [19] | 5.96 | 11.62 | 6.69 | 11.06 | 18.58 | 10.26 | 8.61 | 8.38 | 10.14 |
| FoldingNet [25] | 15.64 | 22.13 | 17.46 | 29.74 | 32.00 | 24.57 | 18.99 | 21.88 | 22.80 |
| PCN [26] | 3.88 | 7.07 | 5.50 | 6.81 | 8.46 | 7.24 | 6.01 | 6.27 | 6.40 |
| DeepSDF [16] | 3.88 | - | - | 5.63 | - | **4.68** | - | - | - |
| LGAN [3] | 3.32 | - | - | 5.59 | - | - | - | - | - |
| MAP-VAE [10] | 3.23 | - | - | 5.57 | - | - | - | - | - |
| **Our Method** | **2.52** | **5.49** | **4.08** | **5.20** | **6.17** | 5.25 | **4.61** | **5.80** | **4.89** |

Table 3: Completion results using the Earth-Moving distance (multiplied by $10^2$) with the output resolution of 1,024. We report the values of DeepSDF [16] from their original paper by rescaling according to the difference of point density.

perform better. Nevertheless, even with such improvements, the complete version of the proposed method still outperforms them.

## 7.2   Car completion on KITTI

The KITTI [8] dataset present partial scans of real-world cars using Velodyne 3D laser scanner. We adopt the same training and validating procedure for car completion as proposed by PCN [26]. We train a car completion model based on the training data generated from ShapeNet [5] and test our completion method on sparse point clouds generated from the real-world LiDAR scans. For each sample, the points within the bounding boxes are extracted with 2,483 partial

| Method | Fidelity | MMD | Consistency |
|---|---|---|---|
| FoldingNet [25] | 0.03155 | 0.02080 | 0.01326 |
| AtlasNet [9] | 0.03461 | 0.02205 | 0.01646 |
| PCN [26] | 0.02800 | 0.01850 | 0.01163 |
| **Our Method** | **0.02171** | **0.01465** | **0.00922** |
| PCN [26] (rotate) | 0.03352 | 0.02370 | 0.01639 |
| **Our Method** (rotate) | **0.02392** | **0.01732** | **0.01175** |

Table 4: Car completion on LiDAR scans from KITTI.

point clouds. Each point cloud is then transformed to the box's coordinates to be completed by our model then transformed back to the world frame. PCN [26] proposed three metrics to evaluate the performance of our model: (1) *Fidelity*, *i.e.* the average distance from each point in the input to its nearest neighbour in the output (*i.e.* measures how well the input is preserved); (2) *Minimal Matching Distance* (MMD), *i.e.* the Chamfer distance between the output and the car's point cloud nearest neighbor from ShapeNet (*i.e.* measures how much the output resembles a typical car); and, (3) *Consistency*, the average Chamfer distance between the completion outputs of the same instance in consecutive frames (*i.e.* measures how consistent the networks outputs are against variations in the inputs).

Table 4 shows that we achieve state of the art on the metrics compared to FoldingNet [25], AtlasNet [9] and PCN [26]. When we introduce random rotations on the bounding box in order to simulate errors in the initial stages, we still acquire the lowest errors.

## 7.3   Classification on ModelNet and PartNet

We evaluate the performance of the features in term of classification on ModelNet10 [27], ModelNet40 [27] and PartNet [15] datasets. ModelNet40 contains 12,311 CAD models in 40 categories. Here, the training data contains 9,843 samples and the testing data contains 2,468 samples. Following RS-DGCNN [2], a linear Support Vector Machine [6] (SVM) is trained on the representations learned in an unsupervised manner on the ShapeNet dataset. RS-DGCNN [2] divides the point cloud of the objects into several regions by positioning the object in a pre-defined voxel grid, then use the regional information to help train latent feature. In Table 5, the proposed method outperforms RS-DGCNN [2] by 1.64% accurracy on ModelNet40 dataset, which shows that our feature contains better categorical information. Notably, similar results are also acquired from ModelNet10 [27] and PartNet [15] with their respective evaluation strategy.

## 7.4   Ablation study

*Loss functions.* In the reconstruction and classification experiments, Table 2 and Table 5 also include the ablation study that investigates the effects of the

| Method | ModelNet40 [27] | ModelNet10 [27] | PartNet [15] |
|---|---|---|---|
| VConv-DAE | 75.50% | 80.50% | - |
| 3D-GAN | 83.30% | 91.00% | 74.23% |
| Latent-GAN | 85.70% | 95.30% | - |
| FoldingNet | 88.40% | 94.40% | - |
| VIP-GAN | 90.19% | 92.18% | - |
| RS-PointNet [2] | 87.31% | 91.61% | 76.95% |
| RS-DGCNN [2] | 90.64% | 94.52% | - |
| KCNet [1] | 91.0% | 94.4% | - |
| **Our Method** | **92.28%** | **96.14%** | **84.32%** |
| – *without* $\mathcal{L}_{\text{inter}}$ | 89.40% | 95.75% | 81.13% |
| – *without* $\mathcal{L}_{\text{intra}}$ | 83.70% | 90.21% | 79.28% |
| – *without* $\mathcal{L}_{\text{inter}}, \mathcal{L}_{\text{intra}}$ | 82.97% | 90.02% | 78.41% |
| – *without* $\mathcal{L}_{\text{boundary}}$ | 88.26% | 95.01% | 80.86% |
| – *without* $\mathcal{L}_{\text{preserve}}$ | 86.09% | 92.27% | 79.05% |

Table 5: Object classification on ModelNet40 [27], ModelNet40 [27] and Part-Net [15] datasets in terms of accuracy.

loss functions from Sec. 6. For both experiments, we notice all loss functions are critical to achieve good results since each of them focuses on different aspects.

*Activations.* Since the number of regions is one of the hyper-parameters in our approach, we evaluate on the performance with different number of regions quantitatively in Table 6. These results demonstrate that the accuracy for the shape completion is increasing as the number of regions increases from 2 to 8, then the performance gradually drops as the number of regions continues to increase from 8 to 32. By observing $\mathcal{L}_{\text{inter}}$ at the same time, we find that it achieves the minimum value of 0.20 when there are 8 regions as well. This proves that $\mathcal{L}_{\text{inter}}$ can be used as an indicator for whether the expected number of regions could be used or not.

Moreover, Fig. 8 shows the regional activations when we shuffle the sequence of points in the partial scan. We can see that both the reconstructed geometry relative sub-regions are identical. So, it illustrates that, by using the proposed regional activations, our model is permutation invariant, which indicates that the reordered point cloud is suitable to perform convolutions.

| $(N_f, N_r)$ | $(2, 128)$ | $(4, 64)$ | $(8, 32)$ | $(16, 16)$ | $(32, 8)$ |
|---|---|---|---|---|---|
| Chamfer Distance | 7.80 | 6.31 | **5.94** | 6.27 | 6.75 |
| $\mathcal{L}_{\text{inter}}$ | 0.41 | 0.67 | **0.20** | 0.49 | 1.33 |

Table 6: Influence of $N_f$ and $N_r$ on the Chamfer distance (multiplied by $10^3$) and $\mathcal{L}_{\text{inter}}$.
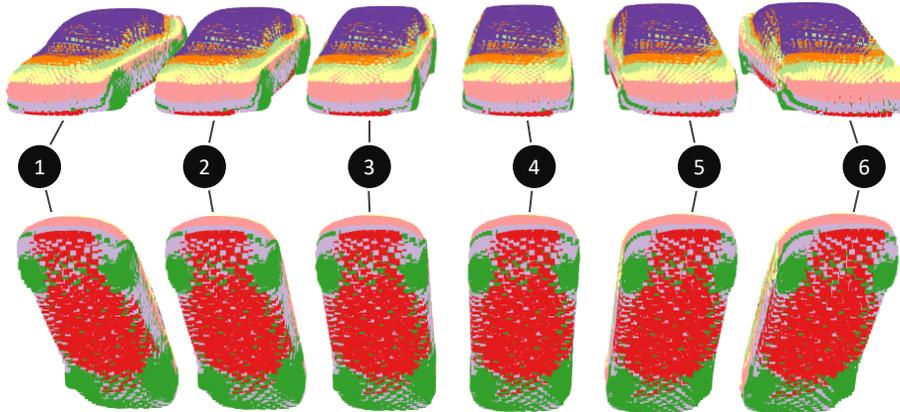
Fig. 8: With identical results, this evaluation shows the robustness of the reconstruction when we randomly shuffle the input point cloud.

*Point cloud versus volumetric data.* In addition to achieving worse numerical results in Sec. 7.1, volumetric approaches have smaller resolutions than the point cloud approaches due to the memory constraints. The difference becomes more evident in Fig. 9, where ForkNet [22] is limited by a $64\times64\times64$ grid. Nevertheless, both the volumetric and point cloud approaches have difficulty in reconstructing thin structures. For instance, the volumetric approach tends to ignore the joints between the wheels and car chassis in Fig. 9 while FoldingNet [25] and AtlasNet [9] only use large surface to cover the area of wheels. In contrast, our approach is capable of reconstructing the thin structures quite well. Moreover, in Table 7, we also achieve the lowest inference time compared to all point cloud and volumetric approaches.

| Method | Size (MB) | Inference Time (s) | Closed Surface | Type of Data |
|---|---|---|---|---|
| 3D-EPN [7] | 420 | - | Yes | Volumetric |
| ForkNet [22] | 362 | - | Yes | Volumetric |
| FoldingNet [25] | 19.2 | 0.05 | Yes | Points |
| AtlasNet [9] | 2 | 0.32 | No | Points |
| PCN [26] | 54.8 | 0.11 | No | Points |
| DeepSDF [16] | 7.4 | 9.72 | Yes | SDF |
| **Our Method** | 37.2 | 0.04 | Yes | Points |

Table 7: Overview of the object completion methods. The inference time is the amount of time to conduct inference on a single sample.
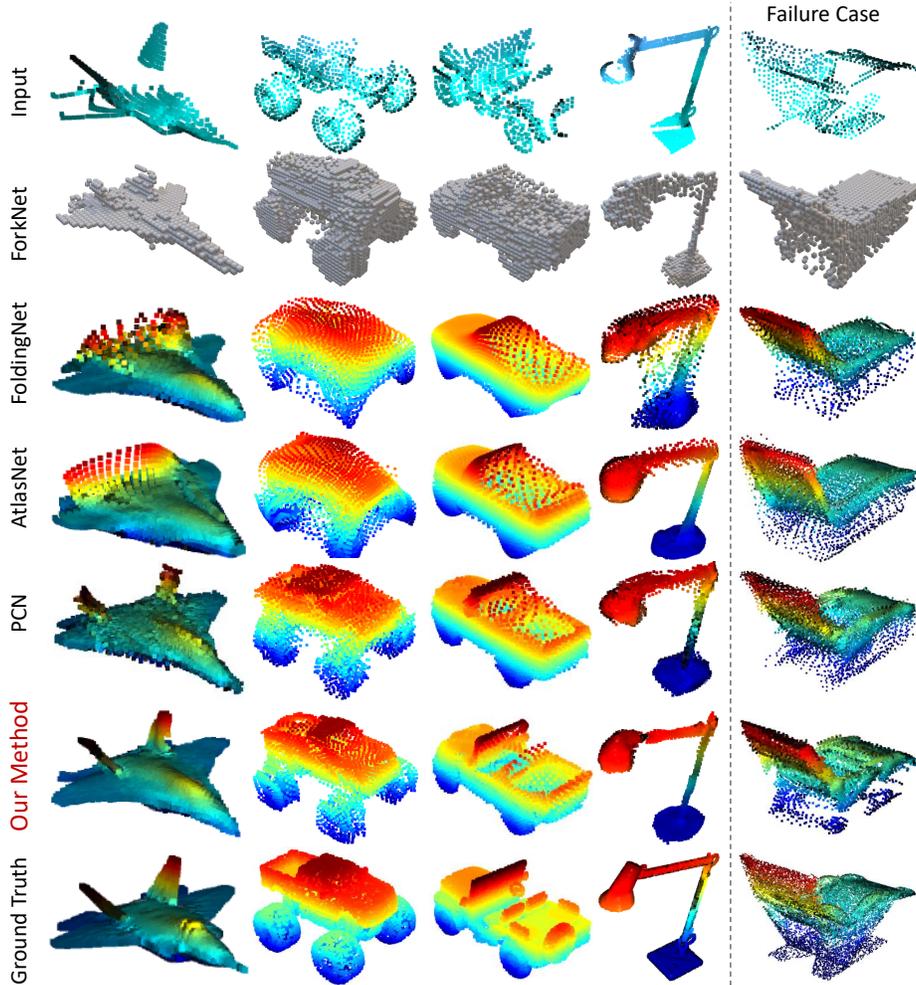
Fig. 9: Evaluated on ShapeNet [5], comparison of shape completion based on ForkNet [22], FoldingNet [25], AtlasNet [9] and PCN [26] against our method.

## 8   Conclusion

This paper introduced the SoftPool idea as a novel and general way to extract rich deep features from unordered point sets such as 3D point clouds. Also, it proposed a state-of-the-art point cloud completion approach by designing a regional convolution network for the decoding stage. Our numerical evaluation reflects that our approach achieves the best results on different 3D tasks, while our quantitative results illustrate the reconstruction and completion ability of our method with respect to ground truth.

# References

1. Mining point cloud local structures by kernel correlation and graph pooling. In: CVPR (2018) 12

2. Self-supervised deep learning on point clouds by reconstructing space. In: NIPS (2019) 11, 12

3. Achlioptas, P., Diamanti, O., Mitliagkas, I., Guibas, L.: Learning representations and generative models for 3D point clouds. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 40–49. PMLR, Stockholmsmssan, Stockholm Sweden (10–15 Jul 2018) 10

4. Arief, H.A., Arief, M.M., Bhat, M., Indahl, U.G., Tveite, H., Zhao, D.: Density-adaptive sampling for heterogeneous point cloud object segmentation in autonomous vehicle applications. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops. pp. 26–33 (2019) 1

5. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015) 9, 10, 14

6. Cortes, C., Vapnik, V.: Support-vector networks. Machine learning **20**(3), 273–297 (1995) 11

7. Dai, A., Qi, C.R., Nießner, M.: Shape completion using 3d-encoder-predictor cnns and shape synthesis. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). vol. 3 (2017) 2, 9, 10, 13

8. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The kitti dataset. The International Journal of Robotics Research **32**(11), 1231–1237 (2013) 10

9. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: A papier-mch approach to learning 3d surface generation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018) 2, 3, 6, 9, 10, 11, 13, 14

10. Han, Z., Wang, X., Liu, Y.S., Zwicker, M.: 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: The IEEE International Conference on Computer Vision (ICCV) (October 2019) 3, 10

11. Landrieu, L., Simonovsky, M.: Large-scale point cloud semantic segmentation with superpoint graphs. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4558–4567 (2018) 1

12. Li, P., Wang, Q., Zhang, L.: A novel earth mover's distance methodology for image matching with gaussian mixture models. In: The IEEE International Conference on Computer Vision (ICCV) (December 2013) 8

13. Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B.: Pointcnn: Convolution on x-transformed points. In: Advances in Neural Information Processing Systems. pp. 820–830 (2018) 1, 3

14. Liu, Y., Fan, B., Xiang, S., Pan, C.: Relation-shape convolutional neural network for point cloud analysis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8895–8904 (2019) 1

15. Mo, K., Zhu, S., Chang, A.X., Yi, L., Tripathi, S., Guibas, L.J., Su, H.: PartNet: A large-scale benchmark for fine-grained and hierarchical part-level 3D object understanding. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) 11, 12

16. Park, J.J., Florence, P., Straub, J., Newcombe, R., Lovegrove, S.: Deepsdf: Learning continuous signed distance functions for shape representation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 165–174 (2019) 10, 13

17. Pham, Q.H., Nguyen, T., Hua, B.S., Roig, G., Yeung, S.K.: Jsis3d: Joint semantic-instance segmentation of 3d point clouds with multi-task pointwise networks and multi-value conditional random fields. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 8827–8836 (2019) 1

18. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 652–660 (2017) 1, 3, 4, 5

19. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In: Advances in Neural Information Processing Systems (NIPS) (2017) 1, 9, 10

20. Song, S., Yu, F., Zeng, A., Chang, A.X., Savva, M., Funkhouser, T.: Semantic scene completion from a single depth image. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR). IEEE (2017) 2

21. Tchapmi, L.P., Kosaraju, V., Rezatofighi, H., Reid, I., Savarese, S.: Topnet: Structural point cloud decoder. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 383–392 (2019) 10

22. Wang, Y., Tan, D.J., Navab, N., Tombari, F.: Forknet: Multi-branch volumetric semantic completion from a single depth image. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 8608–8617 (2019) 2, 3, 9, 13, 14

23. Wu, W., Qi, Z., Fuxin, L.: Pointconv: Deep convolutional networks on 3d point clouds. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9621–9630 (2019) 1, 3

24. Yang, B., Rosa, S., Markham, A., Trigoni, N., Wen, H.: Dense 3d object reconstruction from a single depth view. IEEE transactions on pattern analysis and machine intelligence (2018) 2, 9

25. Yang, Y., Feng, C., Shen, Y., Tian, D.: Foldingnet: Point cloud auto-encoder via deep grid deformation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 206–215 (2018) 2, 3, 6, 9, 10, 11, 13, 14

26. Yuan, W., Khot, T., Held, D., Mertz, C., Hebert, M.: Pcn: Point completion network. In: 2018 International Conference on 3D Vision (3DV). pp. 728–737. IEEE (2018) 2, 3, 6, 9, 10, 11, 13, 14

27. Zhirong Wu, Song, S., Khosla, A., Fisher Yu, Linguang Zhang, Xiaoou Tang, Xiao, J.: 3d shapenets: A deep representation for volumetric shapes. In: 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). pp. 1912–1920 (2015) 11, 12

# SoftPoolNet: Shape Descriptor for Point Cloud Completion and Classification

Yida Wang[1], David Joseph Tan[2], Nassir Navab[1], and Federico Tombari[1,2]

[1] Technische Universität München
[2] Google Inc.

## 1 Comparison of SoftPoolNet to PointNet and PCN

Our architecture is composed by two parts: *encoder* and *decoder*. The encoder takes the partial scan as input. We process the partial scans with our novel soft pooling to produce the ordered feature $F^*$. Then, the decoder takes the feature $F^*$ as input and apply our regional convolution twice to produce the point clouds with resolutions of 256 and 16,384 successively.

Notably, there are some similar components between our encoder and Point-Net [**?**], as well as our decoder and PCN [**?**]. The following sections discuss the distinction in more detail.

### 1.1 Distinction of our encoder from PointNet

Each point on the cloud goes through the multi-layer perceptron (MLP) to accumulate the feature vectors into the matrix $\mathbf{F}$. Then, we sort the feature vectors in a descending order based on the $k$-th element of each vector. The sorted matrix is denoted as $\mathbf{F}'_i$. After independently sorting all $N_f$ elements, we collect the matrices to form the tensor $\mathbf{F}'$ as shown in Fig. 1(a). We then build our softpool feature by taking the first $N_r$ elements of each matrix and concatenate them to $\mathbf{F}^*$.



(a) *Soft*-pooling from the proposed     (b) *Max*-pooling from PointNet
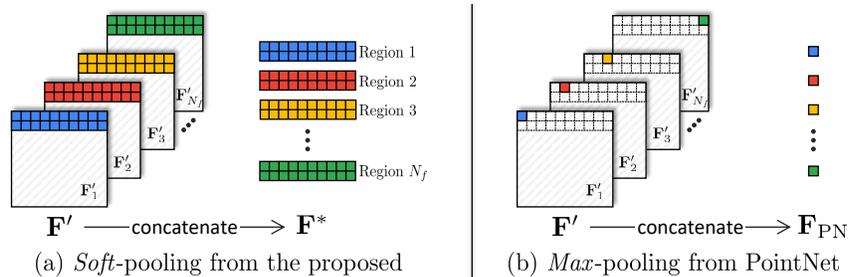
Fig. 1: Comparison between (a) our soft-pool operation and (b) max-pooling from PointNet, where the feature from PointNet is only a subset of our feature.

When comparing our softpool feature $\mathbf{F}^*$ with the feature from PointNet [?] denoted as $\mathbf{F}_{\mathrm{PN}}$, PointNet executes a max-pooling operation on $\mathbf{F}'$ as illustrated in Fig. 1(b). Assuming that both features are derived from the same $\mathbf{F}'$ produced by an MLP, we can conclude that $\mathbf{F}_{\mathrm{PN}}$ is a subset of our feature where

$$\mathbf{F}_{\mathrm{PN}} = \Big[\mathbf{F}'_1[1],\ \mathbf{F}'_2[2],\ \mathbf{F}'_3[3],\ \ldots\ \mathbf{F}'_{N_f}[N_f]\Big] \tag{1}$$

only takes the one value of each matrix while our method takes the first $N_r$ rows. Due to this, the dimensionality of the feature are then distinct. PointNet takes a vector with 1,024 values while we take $N_r \times N_f \times N_f$.

Notably, both our softpool feature and the PointNet feature are permutation invariant, which means that $\mathbf{F}^*$ and $\mathbf{F}_{\mathrm{PN}}$ are the same irrelevant of the order of the input points. This is one of the most important aspect when handling point clouds since this data is unordered.

### 1.2   Distinction of our decoder from PCN

Based on our decoder architecture in the paper, the resulting feature from the encoder undergoes two successive regional convolution operations. The first converts the features to a course point cloud $\mathbf{P}'_{\mathrm{out}}$ with 256 points. From there, the second regional convolution interpolates from the coarse to a fully-packed point cloud with 16,384 points which is denoted as $\mathbf{P}_{\mathrm{out}}$.

Compared to PCN [?], both approaches execute a coarse-to-fine approach which is performed by our second regional convolution. However, the architecture and the method are different.

Given $\mathbf{P}'_{\mathrm{out}}$, PCN [?] duplicates $\mathbf{P}'_{\mathrm{out}}$ 64 times and appends a 2D coordinates of an $8 \times 8$ grid. Then, they use MLP to produce $\mathbf{P}_{\mathrm{out}}$ that locally deforms the 2D grids around each point similar FoldingNet [?]. In contrast, we interpolate 63 samples between every 2 points of $\mathbf{P}'_{\mathrm{out}}$ and use the proposed regional convolution to produce $\mathbf{P}_{\mathrm{out}}$. Compared to MLP in PCN, our regional convolution takes more local samples into account to produce a point in the higher resolution.

## 2   Ablation study on the softpool feature F*

Using our architecture trained with $N_r = 32$, we present the qualitative results when only a subset of the rows is selected. The objective is to investigate which parts of the object each region reconstructs first. In Fig. 2, we start by limiting with the first two rows of the feature matrix then increasing $N_r$ to reach 32. By selecting the first 2 features, we observe that the softpool feature focuses on a skeleton of the object without large surfaces. Although the regions reconstruct different parts of the object, they tend to cover the important components like the wings of plane and the wheels of car. As we increase $N_r$ from 2 to 32, the object is slowly completed without huge overlaps between different regions.

In addition to the first 32 rows when setting $N_r$, we also looked into the rows beyond 32. The lamp in Fig. 3 focuses on the following ranges: $[33 : 64]$,
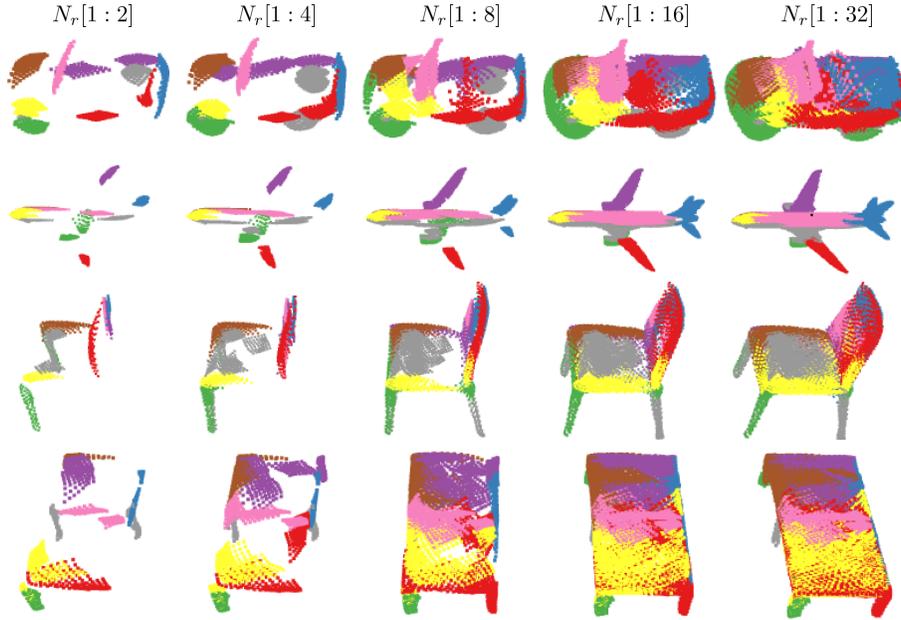
$N_r[1:2]$          $N_r[1:4]$          $N_r[1:8]$          $N_r[1:16]$          $N_r[1:32]$



Fig. 2: Results when choosing the first subsets of $N_r$ with the following ranges: $[1:2]$, $[1:4]$, $[1:8]$, $[1:16]$ and $[1:32]$ when the architecture is trained with $N_r = 32$.

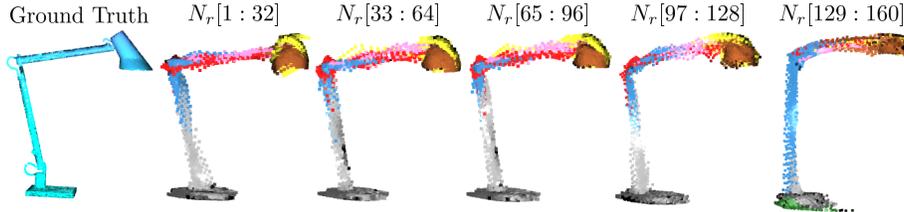Ground Truth    $N_r[1:32]$     $N_r[33:64]$    $N_r[65:96]$    $N_r[97:128]$   $N_r[129:160]$



Fig. 3: Results when choosing different ranges of rows from $\mathbf{F}'$ to form $\mathbf{F}^*$ instead of selecting the first $N_r = 32$ rows.

$[65:96]$, $[97:128]$ and $[129:169]$. Although the shape of the lamp starts to deform as we go beyond 32, our reconstruction results still captures its overall shape even when we select the range $[129:169]$. Therefore, this proves that our feature is not constrained to the first 32 rows when sorting and demonstrates the robustness of our softpool feature.

## 3   Ablation study on $\tau$

When computing for $\mathcal{L}_{\text{boundary}}$, we introduced the threshold $\tau$ to compute the sets. In Table 1, we then evaluate different values of $\tau$ and investigate its behavior with respect to the Chamfer distance. The table demonstrates that the results are not sensitive to the $\tau$, where the thresholds between 0.2-0.9 generate a small difference in the Chamfer distance (with less than 1) from the chosen threshold of 0.3. Notably, compared to the related work, any threshold between 0.1 to 0.9 outperforms the other methods.

| $\tau$ | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 |
|---|---|---|---|---|---|---|---|---|---|
| Chamfer Distance | 7.08 | 5.99 | **5.94** | 6.12 | 6.19 | 6.18 | 6.21 | 6.25 | 6.71 |

Table 1: Sensitivity of the average Chamfer distance (multiplied by $10^3$) to the threshold $\tau$.

## 4   Ablation study on $N_r$, $N_r$ and $\mathbf{L_{boundary}}$

We investigate the influence of increasing the weight of $\mathcal{L}_{\text{boundary}}$ on the reconstruction as we change the number of regions $N_f$ and the number selected rows $N_r$. While we chose the best option with $N_r$ set to 8 and $N_r$ set to 32, Table 2 also shows that a larger weight on $\mathcal{L}_{\text{boundary}}$ improves the performance when the number of regions is larger, e.g. when $N_f$ is 32.

| $(N_f, N_r)$ | (2, 128) | (4, 64) | (8, 32) | (16, 16) | (32, 8) |
|---|---|---|---|---|---|
| $1 \times \mathcal{L}_{\text{boundary}}$ | 7.80 | 6.31 | 5.94 | 6.27 | 6.75 |
| $2 \times \mathcal{L}_{\text{boundary}}$ | 7.80 | 6.31 | **5.91** | 6.25 | 6.72 |
| $10 \times \mathcal{L}_{\text{boundary}}$ | 7.82 | 6.29 | 5.95 | 6.01 | 6.19 |

Table 2: Influence of $N_f$, $N_r$ and the weight of $\mathcal{L}_{\text{boundary}}$ for object completion on the average Chamfer distance (multiplied by $10^3$).