

# SPU-Net: Self-Supervised Point Cloud Upsampling by Coarse-to-Fine Reconstruction with Self-Projection Optimization

Xinhai Liu<sup>1</sup>, Xincheng Liu<sup>1</sup>, Zhizhong Han<sup>2</sup>, and Yu-Shen Liu<sup>1</sup>

<sup>1</sup>School of Software, BNRist, Tsinghua University, Beijing, China

<sup>2</sup>Department of Computer Science, University of Maryland, College Park, USA

{lxh17, lxc19}@mails.tsinghua.edu.cn, h312h@umd.edu, liuyushen@tsinghua.edu.cn

## Abstract

The task of point cloud upsampling aims to acquire dense and uniform point sets from sparse and irregular point sets. Although significant progress has been made with deep learning models, they require ground-truth dense point sets as the supervision information, which can only be trained on synthetic paired training data and are not suitable for training under real-scanned sparse data. However, it is expensive and tedious to obtain large scale paired sparse-dense point sets for training from real scanned sparse data. To address this problem, we propose a self-supervised point cloud upsampling network, named SPU-Net, to capture the inherent upsampling patterns of points lying on the underlying object surface. Specifically, we propose a coarse-to-fine reconstruction framework, which contains two main components: point feature extraction and point feature expansion, respectively. In the point feature extraction, we integrate self-attention module with graph convolution network (GCN) to simultaneously capture context information inside and among local regions. In the point feature expansion, we introduce a hierarchically learnable folding strategy to generate the upsampled point sets with learnable 2D grids. Moreover, to further optimize the noisy points in the generated point sets, we propose a novel self-projection optimization associated with uniform and reconstruction terms, as a joint loss, to facilitate the self-supervised point cloud upsampling. We conduct various experiments on both synthetic and real-scanned datasets, and the results demonstrate that we achieve comparable performance to the state-of-the-art supervised methods.

## 1. Introduction

Point clouds, as one of the most concise 3D representation, have drawn an increasing research attention due to

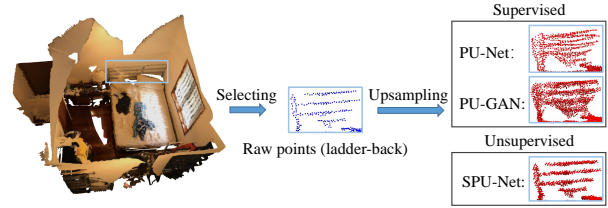


Figure 1. The upsampling results under real-scanned object from ScanNet [6]. The supervised methods (e.g. PU-Net [38] and PU-GAN [17]) destroy the surface distribution of raw points (the ladder-back of the bed). In contrast, our SPU-Net keeps the origin point distribution.

its convenient access from various popular depth sensors, such as LiDARs and RGB-D cameras. However, raw point clouds obtained from devices are usually sparse, noisy and non-uniform, which leads to a huge challenge for deep networks to directly deal with such irregular data. Given a sparse, noisy and non-uniform point cloud, the task of upsampling aims to generate a point set that is dense, complete, and uniform, as a trusted representation of the underlying object surface. Therefore, point cloud upsampling, as an amended operation, is meaningful for many downstream applications like rendering, analysis, reconstruction and other general processing.

Current deep point cloud upsampling methods, such as PU-Net [38], MPU [36] and PU-GAN [17], have achieved outperforming results on some synthetic datasets such as ShapeNet [5] and Visionair repository [1]. However, these methods usually require ground-truth dense point sets as the supervision information in the network training. And the training data is usually constructed by sampling on synthetic CAD models from public datasets [17, 38], which is unavailable for the real-scanned data. Due to the absence of paired ground-truth dense point sets, the aforementioned methods cannot be trained under the real-scanned datasets

such as ScanNet [6] and KITTI [8]. In addition, when data distributions from synthetic data do not match those from real scans, upsampling networks trained on synthetic data do not generalize well to real (sparse) scans. For example, the supervised methods trained on synthetic data, such as PU-Net [38] and PU-GAN [17], easily changed the origin topological structure of underlying object surface on the real-scanned data from ScanNet as shown in Figure 1. Therefore, it is promising to propose a self-supervised point cloud upsampling method, which is not restricted by the supervision information and can keep the original data distributions.

Recently, some unsupervised image super-resolution methods [39, 26] have been proposed and achieved satisfactory performances in generating high-resolution images. However, due to the irregular and unordered nature of point clouds, it is non-trivial to directly apply such image super-resolution methods to the unsupervised point cloud upsampling. Specifically, there are two challenges in the unsupervised point cloud upsampling with deep learning models. (1) *How to establish effective self-supervision information without the supervision of dense point sets?* Previous methods, such as PU-GAN [17] and L2G-AE [20], first generate a dense point set with deep networks and then downsample the dense point set back into a sparse point set, where some supervision information is usually applied to the sparse point sets. However, there is no direct supervision for the dense point sets in the unsupervised upsampling task, which makes it difficult to capture the inherent upsampling patterns. To resolve this issue, we propose a coarse-to-fine reconstruction framework to formulate the self-supervised point cloud upsampling. Specifically, we first downsample the input patch into some coarse patches, and then capture the inherent upsampling patterns by reconstructing the input patch itself from each coarse patch. Next, a dense patch can be obtained by aggregating multiple fine patches, which all follow the distribution of the input patch. (2) *The point clouds upsampled by deep networks should be a faithful representation of the underlying object surface.* Due to the irregular nature of the point cloud and network bias in generation of dense point clouds from sparse ones, it is inevitable to bring some noisy points around the underlying surface when generating the upsampled dense point set, especially for unsupervised upsampling methods. Therefore, some specific loss functions are needed to constrain spatial distribution of generated points, including uniformity and flatness. To resolve these challenge, we propose a novel self-projection optimization to project the noisy points around the underlying surface to the surface itself, and associate it with uniform and reconstruction terms as a joint loss to facilitate the generation of upsampled points. In general, our contributions are summarised as follows.

- We propose a self-supervised point cloud upsampling network (SPU-Net) without the supervision of 3D ground-truth dense point clouds. SPU-Net repeatedly upsamples from downsampled patches, which is not restricted by the paired training data and can preserve the original data distribution.
- We propose a coarse-to-fine reconstruction framework to capture the inherent upsampling patterns inside local patches, which provides a novel self-supervision way to learn to upsample point clouds.
- To constrain the distribution of generated points without the ground-truth dense point sets, we introduce a novel self-projection optimization, which interactively projects the generated points onto the underlying object surface along the projection direction.

To evaluate the performances of SPU-Net, we adopt four widely used metrics to compare with the state-of-the-art methods under a variety of synthetic and real-scanned datasets. Experimental results show that our self-supervised method achieves good performance, which is comparable to the supervised methods in both qualitative and quantitative comparisons.

## 2. Related Work

**Traditional point cloud upsampling methods.** Traditional methods have tried various optimization strategies to generate the upsampled point clouds without using deep learning models. For example, Alexa *et al.* [4] upsampled point set by computing the Voronoi diagrams and adding points at vertices of this diagram. Afterwards, Locally optimal projection (LOP) operator [18, 13] has been proved to be effective for points resampling and surface reconstruction based on  $L_1$  median, especially for point sets with noise and outliers. Furthermore, Huang *et al.* [14] introduced a progressive strategy for edge-aware point set resampling. To fill large holes and complete missing regions, Dpoints [34] developed a new point representation. In general, all these methods are not data-driven, which heavily rely on shape priors, such as normal estimation and smooth surface assumption.

**Deep-learning-based point cloud upsampling methods.** In recent years, deep neural networks have achieved outperforming performances in various point cloud processing tasks, including shape classification [21, 19, 31], object detection [24, 28], semantic scene segmentation [12, 27] and point cloud completion [32, 29]. In the field of point cloud upsampling, Yu *et al.* [38] as a pioneer first proposed a deep neural network PU-Net to upsample point set, which works on patches by learning multi-level per-point features and expanding the point set via multi-branch convolutions. Later, they designed another edge-aware point cloud upsampling

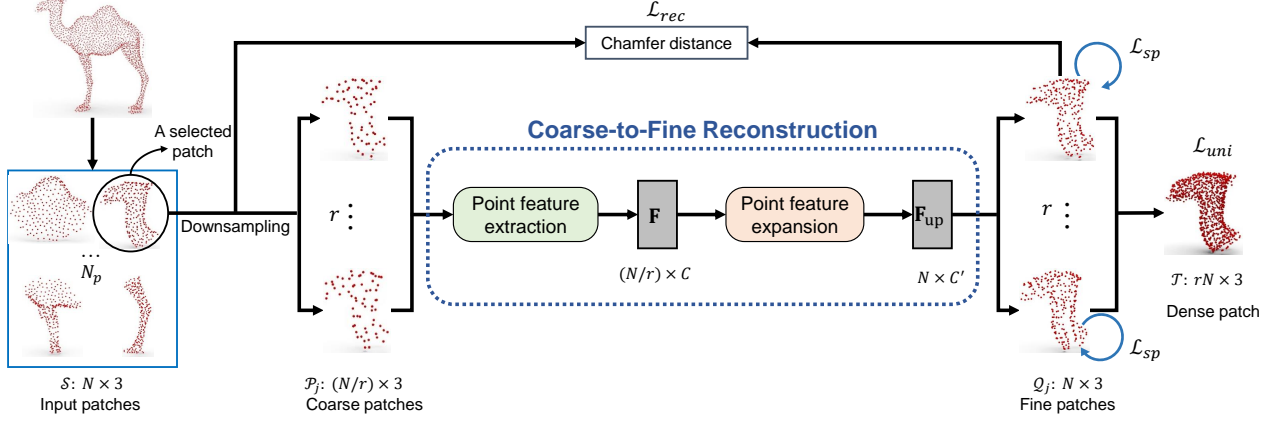


Figure 2. The architecture of SPU-Net. Given an input patch with  $N$  points, we first downsample the input patch into  $r$  coarse patches  $\mathcal{P}_j$ , each with  $N/r$  points, where  $r$  is the upsampling rate. And then we generate a fine patch  $\mathcal{Q}_j$  for each coarse patch  $\mathcal{P}_j$  in a coarse-to-fine reconstruction framework, which consists of point feature extraction and expansion. Finally, we obtain the target dense patch  $\mathcal{T}$  by aggregating all the fine patches. In addition,  $C$  and  $C'$  are the number of feature channels that are 480 and 128, respectively, in the implementation;  $\mathcal{L}_{rec}$ ,  $\mathcal{L}_{uni}$  and  $\mathcal{L}_{sp}$  denote the reconstruction, uniform, and self-projection terms, respectively.

network named EC-Net [37] to achieve point expansion by minimizing the point-to-edge distances. Wang *et al.* [36] presented a multi-step progressive upsampling network to further maintain the patch details. Recently, Li *et al.* [17] proposed a GAN-based framework to generate high quality upsampled point sets. Existing methods can already generate high-quality upsampled point clouds under synthetic datasets. However, these supervised methods require the ground-truth dense point clouds as the supervision information and are not suitable for the real-scanned data. Liu *et al.* [20] proposed a deep neural network named L2G-AE in representation learning, which can be applied for unsupervised point cloud upsampling by reconstructing overlapped local regions. However, L2G-AE concentrated on the capturing of global shape information via local-to-global reconstruction, which limits the network in capturing inherent upsampling patterns and generating high-quality upsampled point sets. In contrast, we present a new method to upsample from downsampled patches in a coarse-to-fine framework, which enables us to generate high-quality upsampled point clouds in a self-supervised manner.

**Unsupervised point cloud processing methods.** Recently, several methods have investigated unsupervised learning strategies for point cloud analysis. The auto-encoder, such as FoldingNet [35] and L2G-AE [20], is a widely used framework for unsupervised point cloud learning, which adopts input point cloud itself as the reconstruction target. During the self-reconstruction process, the corresponding point features is obtained, which can be applied in various applications, such as shape classification [35], semantic segmentation [10, 11] and shape generation [10]. Based on the encoder-decoder architecture, some works conducted self-supervised information as the optimization target, such

as coordinate transformation [7], deformation reconstruction [2] and part partition [40]. In addition, the generative adversarial network [3] was also applied to distinguish the generated point set or the real point set. In this work, we adopt the generalized encoder-decoder framework to build our self-supervised upsampling network based on local patches. In order to capture the inherent upsampling patterns, the key issue is how to reproduce the upsampling process without requiring the supervised dense point sets. To resolve this issue, we propose a coarse-to-fine reconstruction framework to reproduce point cloud upsampling inside local patches. Specifically, we first downsample the input patch into several subsets, and then the coarse-to-fine reconstruction framework is applied to upsample the subsets. By repeating the downsampling and upsampling steps, our method can generate the dense upsampled point sets.

### 3. The Architecture of SPU-Net

#### 3.1. Overview

Given a 3D point cloud, we take the same patch-based approach as PU-Net [38] and PU-GAN [17]. We first build  $N_p$  local patches according to the geodesic distance for the 3D point cloud, as shown in Figure 2. For each local patch  $\mathcal{S} = \{\mathbf{s}_i\}_{i=1}^N$  with  $N$  points, our goal is to output a dense and uniform point set  $\mathcal{T} = \{\mathbf{t}_i\}_{i=1}^{rN}$  with  $rN$  points, while keeping the original data distribution, where  $\mathbf{s}_i$ ,  $\mathbf{t}_i$  are the coordinates of 3D points and  $r$  is the upsampling rate. Without the supervision of ground-truth dense point sets, we propose a coarse-to-fine reconstruction framework to reproduce the upsampling process inside each local patch. By upsampling the sparse patch to obtain the fine patch, we are able to capture the inherent upsampling patterns for gen-

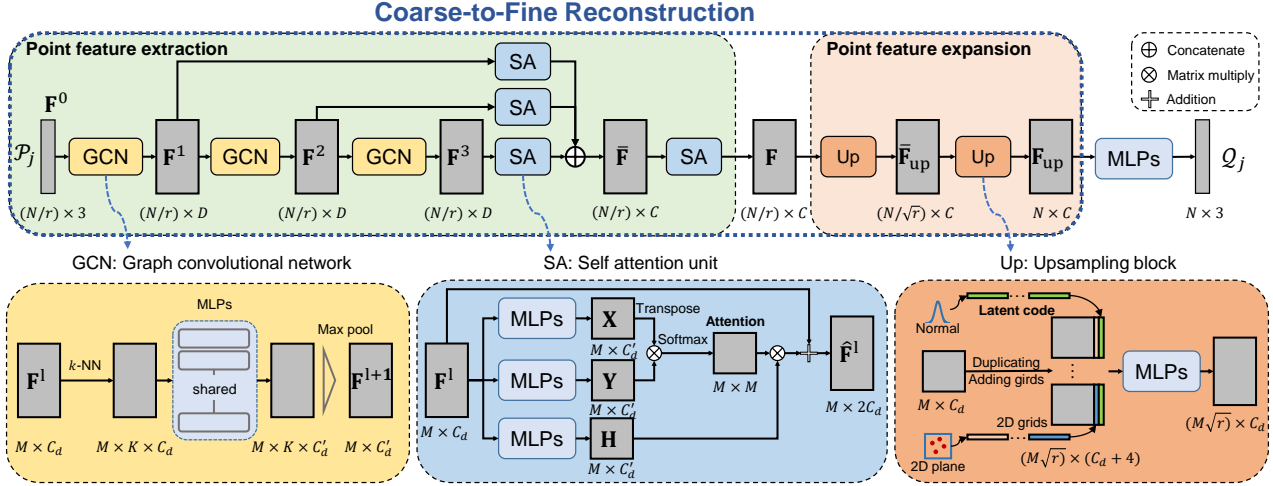


Figure 3. The coarse-to-fine reconstruction framework. Given a coarse patch  $\mathcal{P}_j$  with  $N/r$  points, we aim to generate the corresponding fine patch  $\mathcal{Q}_j$  with  $N$  points. In this framework, there are two main components: point feature extraction and point feature expansion. Here,  $D$  and  $C$  are the number of feature channels that are 64 and 480, respectively, in the implementation.

erating dense patch  $\mathcal{T}$  that is uniformly distributed on the underlying object surface.

Figure 2 illustrates the architecture of our SPU-Net. For an input patch  $\mathcal{S}$ , we first downsample  $\mathcal{S}$  into  $r$  different coarse patches  $\{\mathcal{P}_1, \dots, \mathcal{P}_j, \dots, \mathcal{P}_r\}$ , each with  $N/r$  points (Section 3.2). We then introduce the coarse-to-fine framework to explore the inherent upsampling patterns inside local patches (Section 3.3), which contains two main components: point feature extraction and point feature expansion. Lastly, we present the patch-based training strategy with a joint loss function formed by reconstruction, uniform, self-projection terms (Section 3.4).

### 3.2. Point Set Downsampling

Without the supervision of dense point sets, we have to construct some self-supervision information to support deep networks to capture the inherent upsampling patterns. To take advantage of the input patch without supervision information, we propose a coarse-to-fine reconstruction framework to generate the upsampled dense patch. Specifically, as shown in Figure 2, we first downsample the input patch  $\mathcal{S}$  into some coarse patches  $\mathcal{P}_j$  to formulate the self-supervision with the input patch itself. The downsampling method largely determines the ability of deep networks in capturing the surface distribution of 3D shapes, which should preserve the distribution information of the input patch as much as possible. Here, the farthest point sampling (FPS) algorithm is adopted in our method, which preserves the patch distribution information and generates more uniform coarse patches than other sampling methods such as random sampling. Specifically, we repeat the process of pick out  $N/r$  points from the input local patch  $\mathcal{S}$  for  $r$  times

with the FPS algorithm. Finally, we denote the downsampled coarse patches as  $\{\mathcal{P}_1, \dots, \mathcal{P}_j, \dots, \mathcal{P}_r\}$ . From the downsampled coarse patches, we can reveal the inherent upsampling process by reconstructing the input patch itself, which makes it possible to inference more dense patches.

### 3.3. Coarse-to-Fine Reconstruction

In the coarse-to-fine reconstruction framework, there are two main components: point feature extraction and point feature expansion. In the point feature extraction, we integrate self-attention with graph convolution network (GCN) to simultaneously capture the spatial context of points both inside and among local regions. In the point feature expansion, we propose a hierarchical learnable folding strategy to facilitate the feature propagation from sparse to dense in the feature space.

**Point feature extraction.** To capture the context information from discrete point sets, it is important to extract the spatial correlation of points inside local regions. The graph convolutional network (GCN) [25, 30, 21, 19] has been widely applied to capture the context information inside local regions in existing methods. However, these methods often ignore to capture the correlation among local regions. To simultaneously extract the context information both inside and among local regions, we propose a point feature extraction module, which integrates self-attention units with GCNs, as shown in Figure 3.

Given a coarse patch  $\mathcal{P}_j$  with size of  $N/r \times 3$  as input, three GCNs are first employed to capture the local contexts inside local regions by building a local graph around each point  $\mathbf{p}_i^j$ . We introduce a hierarchical feature extraction strategy with multiple semantic levels. Suppose the



input feature map of a GCN at level  $l \in \{0, 1, 2, 3\}$  is  $\mathbf{F}^l = \{\mathbf{f}_i^l\}_{i=1}^M$  with the size of  $M \times C_d$ , where  $\mathbf{f}_i^l$  is the  $i$ -th point feature in  $\mathbf{F}^l$ . In particular, the feature map  $\mathbf{F}^0$  at level 0 is the raw points from  $\mathcal{P}_j$ . To calculate the point feature  $\mathbf{f}_i^{l+1}$ , we first dynamically build a local region  $\mathcal{N}_i^l$  with  $K$  neighbors around each point feature  $\mathbf{f}_i^l$  with  $k$ -NN algorithm. And then, we formulate the propagation of point feature  $\mathbf{f}_i^l$  ( $l \in \{0, 1, 2\}$ ) as

$$\mathbf{f}_i^{l+1} = \max_{\mathbf{f}_j^l \in \mathcal{N}_i^l} \{\sigma(\mathbf{h}_\theta(\mathbf{f}_j^l - \mathbf{f}_i^l))\}. \quad (1)$$

Here,  $(\mathbf{f}_j^l - \mathbf{f}_i^l)$  can be regard as the edge from point  $\mathbf{f}_j^l$  to center point  $\mathbf{f}_i^l$ ,  $\mathbf{h}_\theta$  indicates the learnable parameters in multi-layer-perceptrons (MLPs),  $\sigma$  is a non-linear layer, such as ReLU [23],  $\max$  is a max-pool operation. The max-pool layer is applied to aggregate point features in the local region  $\mathcal{N}_i^l$ .

Following previous works [20, 17], we integrate self-attention units to capture correlation between local regions. As shown in Figure 3, the self-attention unit cooperates with the GCN to capture the detailed spatial context information inside local patches. Suppose that the input feature map is  $\mathbf{F}^l$  with size of  $M \times C_d$ . Three MLPs are used to embed  $\mathbf{F}^l$  into different feature spaces,  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{H}$ , respectively. In particular,  $\mathbf{X}$  and  $\mathbf{Y}$  are applied to calculate the attention values with a simple matrix multiplication, the updated feature map  $\hat{\mathbf{F}}^l$  is calculated as

$$\hat{\mathbf{F}}^l = \mathbf{F}^l + \text{softmax}(\mathbf{Y}\mathbf{X}^\top)\mathbf{H}. \quad (2)$$

After multiple self-attention units, the correlation among local regions is captured in different semantic levels. We then aggregate the multi-level point features by concatenation as

$$\bar{\mathbf{F}} = \text{MLP}(\mathbf{F}^1 \oplus \hat{\mathbf{F}}^1 \oplus \hat{\mathbf{F}}^2 \oplus \hat{\mathbf{F}}^3). \quad (3)$$

To further explore the correlation in the aggregated feature  $\bar{\mathbf{F}}$ , we add another self-attention unit to obtain the final point feature  $\mathbf{F}$ .

**Point feature expansion.** The target of point feature expansion is to construct the mapping from current points to more points, which is widely used in some point-wise applications, such as semantic segmentation [25], point cloud reconstruction [3] and point cloud upsampling [17]. In general, current point feature expansion methods can be roughly categorized into interpolating-based method [25], folding-based method [35] and reshaping-based method [3]. Existing point interpolating-based methods often use the point interpolation relationship to guide the expansion of point features. However, in many scenarios, the interpolation relationship between point sets is usually unknown. In addition, reshaping-based methods usually first expand the feature dimensions with deep networks, such as MLPs

or fully-connected (FC) layers, and then generate the target point features via a simple reshaping operation. In recent years, the folding-based methods have been developed, which first duplicates point features and then concatenate 2D grids to guide point feature expansion.

Compared with other feature expansion methods, the folding-based method is more flexible and has achieved satisfactory performances in various applications [35, 33]. However, previous fixed 2D grids are not adaptive to diverse feature distributions. To resolve this problem, we propose novel learnable 2D grids as the latent code to cooperate with fixed 2D grids in guiding the point feature expansion. As shown in Figure 3, given an input feature map with size of  $M \times C_d$ , we first duplicate the point features and then concatenate two kinds of grids. In particular, the latent code is initialized from standard normal distribution and can be optimized in the network training process. Moreover, in order to smooth the point feature expansion process, we introduce the hierarchical folding strategy. By using two upsampling blocks, we hierarchically obtain the upsampled point features  $\bar{\mathbf{F}}_{up}$  and  $\mathbf{F}_{up}$  with a upsampling rate  $\sqrt{r}$ . With subsequent MLPs, the point features  $\mathbf{F}_{up}$  is applied to reconstruct the fine patches.

### 3.4. Loss Function

**Joint loss.** In our SPU-Net, we optimize the upsampled point set with a joint loss  $\mathcal{L}_{joint}$ , consisting of reconstruction term  $\mathcal{L}_{rec}$ , uniform term  $\mathcal{L}_{uni}$  and self-projection term  $\mathcal{L}_{sp}$ . Overall, we train our SPU-Net by minimizing the joint loss function in an end-to-end manner as

$$\mathcal{L}_{joint} = \alpha\mathcal{L}_{rec} + \beta\mathcal{L}_{uni} + \gamma\mathcal{L}_{sp}, \quad (4)$$

where  $\alpha$ ,  $\beta$  and  $\gamma$  are hyper-parameters in the training.

**Reconstruction term.** For an input local patch  $\mathcal{S}$ , we first downsample the input patch into some coarse patches, and then we repeat upsampling from each one of coarse patches to explore the inherent upsampling patterns in a self-supervised manner. Assume that  $\mathcal{Q}_j = \{\mathbf{q}_i^j\}_{i=1}^N$  is the upsampled fine patch from the  $j$ -th coarse patch  $\mathcal{P}_j$ . Thus, we formulate reconstruction term using chamfer distance (CD) for  $\mathcal{Q}_j$  as

$$\begin{aligned} \mathcal{L}_{rec} = d_{CD}(\mathcal{S}, \mathcal{Q}_j) &= \sum_{i=1}^N \frac{1}{|\mathcal{S}|} \sum_{\mathbf{s}_i \in \mathcal{S}} \min_{\mathbf{q}_i^j \in \mathcal{Q}_j} \|\mathbf{s}_i - \mathbf{q}_i^j\|_2 \\ &+ \sum_{i=1}^N \frac{1}{|\mathcal{Q}_j|} \sum_{\mathbf{q}_i^j \in \mathcal{Q}_j} \min_{\mathbf{s}_i \in \mathcal{S}} \|\mathbf{s}_i - \mathbf{q}_i^j\|_2. \end{aligned} \quad (5)$$

**Uniform term.** The reconstruction term encourages the upsampled fine patch  $\mathcal{Q}_j$  to fit the input patch  $\mathcal{S}$ . However, the target dense patch  $\mathcal{T} = \{\mathcal{Q}_1, \dots, \mathcal{Q}_j, \dots, \mathcal{Q}_r\}$  should be uniformly distributed on the underlying object surface.

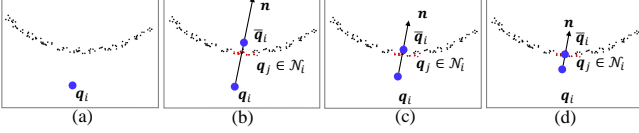


Figure 4. The illustration of self-projection optimization. During training, the noisy point  $q_i$  is gradually pulled onto the underlying object surface, where the projection direction  $n$  is calculated by measuring the local distribution of  $q_i$  and local center  $\bar{q}_i$  in the local region  $\mathcal{N}_i$ .

Similar to [17, 38], we introduce a uniform term to improve the uniformity of the final target patch  $\mathcal{T}$  with size of  $rN \times 3$ . Specifically, we first use the farthest sampling to pick  $M$  seed points in  $\mathcal{T}$  and using a ball query of radius  $r_d$  to build a local region (denoted as  $\mathcal{T}_j, j = 1, \dots, M$ ) in  $\mathcal{T}$  around each seed. Thus, the number of points roughly relies on the small disk of area  $\pi r_d^2$  on the underlying surface. In addition, the patch  $\mathcal{T}$  is normalized in a unit sphere with an area of  $\pi 1^2$ . So, the expected number of points  $\hat{n}$  in  $\mathcal{T}_j$  is  $rNr_d^2$ . Suppose the local region  $\mathcal{T}_j$  satisfies the regular hexagonal distribution and the distance  $d_{j,k}$  is for  $k$ -th point in  $\mathcal{T}_j$ . So, the expected point-to-neighbor distance  $\hat{d}$  should be roughly  $\sqrt{2\pi r_d^2 / (|\mathcal{T}_j| \sqrt{3})}$ . To measure the deviation of  $|\mathcal{T}_j|$  from  $\hat{n}$  and  $d_{j,k}$  from  $\hat{d}$ , we calculate the uniform term with chi-square as

$$\begin{aligned} \mathcal{L}_{uni}(\mathcal{T}_j) &= \sum_{j=1}^M U_{number}(\mathcal{T}_j) \cdot U_{distance}(\mathcal{T}_j) \\ &= \sum_{j=1}^M \frac{(|\mathcal{T}_j| - \hat{n})^2}{\hat{n}} \cdot \sum_{k=1}^{|\mathcal{T}_j|} \frac{(d_{j,k} - \hat{d})^2}{\hat{d}}. \end{aligned} \quad (6)$$

**Self-projection term.** The reconstruction term and uniform term optimize the geometric shape and surface distribution of the generated upsampled point sets, respectively. However, due to the irregular nature of the point cloud and network bias in generation of dense point clouds from sparse ones, it is inevitable to bring some noisy points in the generated upsampled point sets. To resolve this problem, we propose a novel self-projection function to project the generated points onto the underlying object surface. Inspired by traditional directed projection (DP) algorithm [22], each point in the generated set can be projected to the underlying object surface by directed projection. So, we should first find a projection direction for each point in the generated point set. One option is to project each point to the distribution of the input patch. However, due to the distribution differences of between input patches and generated patches, we formulate the function inside local regions with a self-projection manner. Suppose the generated fine patch is  $\mathcal{Q} = \{q_i\}_{i=1}^N$  with a size of  $N \times 3$ . We first build a local region  $\mathcal{N}_i$  around the  $i$ -th point of  $\mathcal{Q}$  with the  $k$ -nearest

Table 1. The effects of the point set downsampling methods.

Metric	FPS	RS	MS
CD ( $10^{-3}$ )	<b>0.38</b>	0.52	0.55
HD ( $10^{-3}$ )	<b>2.24</b>	4.19	4.27
P2F ( $10^{-3}$ )	<b>5.87</b>	7.57	7.22
UNI ( $10^{-3}$ )	<b>8.94</b>	14.89	9.9

neighbors algorithm. Then, we calculate the local region center  $\bar{q}_i = \frac{1}{|\mathcal{N}_i|} \sum_{q_j \in \mathcal{N}_i} (q_j)$  of local region  $\mathcal{N}_i$  to be the projection target of point  $q_j$ . Therefore, the self-projection optimization is formulated with chi-square as

$$\mathcal{L}_{sp} = \frac{1}{|\mathcal{Q}| \cdot |\mathcal{N}_j|} \sum_{i=1}^{|\mathcal{Q}|} \sum_{q_j \in \mathcal{N}_i} \frac{|(q_i - q_j)^2 - (\bar{q}_i - q_j)^2|}{1 + (q_i - q_j)^2}. \quad (7)$$

Here,  $\frac{1}{1 + (q_i - q_j)^2}$  is the weight of each point in the local region  $\mathcal{N}_j$ . As shown in Figure 4,  $(q_i - q_j)^2$  and  $(\bar{q}_i - q_j)^2$  indicate the real distribution and the target distribution of local region, respectively.

## 4. Experiments

### 4.1. Datasets and Network Configurations

For fair comparison, we do quantitative and qualitative comparisons with state-of-the-art methods under the dataset from PU-GAN [17], which contains 147 3D models that is formed with 120 objects in training dataset and 27 objects in test set. And we follow the same patch-based training strategy as in PU-Net [38] and PU-GAN [17]. By default, we crop 24 patches for each training model and set the number of patch points  $N = 256$ , the upsampling rate  $r = 4$ . Moreover, as for the uniform term, the number of seed points  $M$  is 50, and we cropped the same set of  $\mathcal{T}_j$  with radius  $r_d = \sqrt{p}$  for each  $p \in \{0.4\%, 0.6\%, 0.8\%, 1.0\%, 1.2\%\}$  as in [17]. We train the network for 200 epochs with the Adam algorithm [16]. And we set the learning rate as 0.0001, we gradually reduce both rates by a decay rate of 0.7 per 50k iterations until  $10^{-6}$ . The batch size is 24, and  $\alpha$ ,  $\beta$ , and  $\gamma$  are empirically set as 100, 10, and 0.01 in Eq. (4), respectively. We implemented our network using TensorFlow and trained it on Nvidia 2,080 Ti GPU.

To evaluate the performance of point cloud upsampling, we employ the same four evaluation metrics as PU-GAN [17], including uniformity (UNI), point-to-surface (P2F), Chamfer distance (CD), and Hausdorff distance (HD). For quantitative evaluation, we use Poisson disk sampling to sample 2,048 point as training data for each 3D object, and 8,192 points as the ground truth upsampled point cloud in the testing phase.

### 4.2. Parameters

All the experiments in parameter comparisons are evaluated under the dataset from PU-GAN, where Chamfer

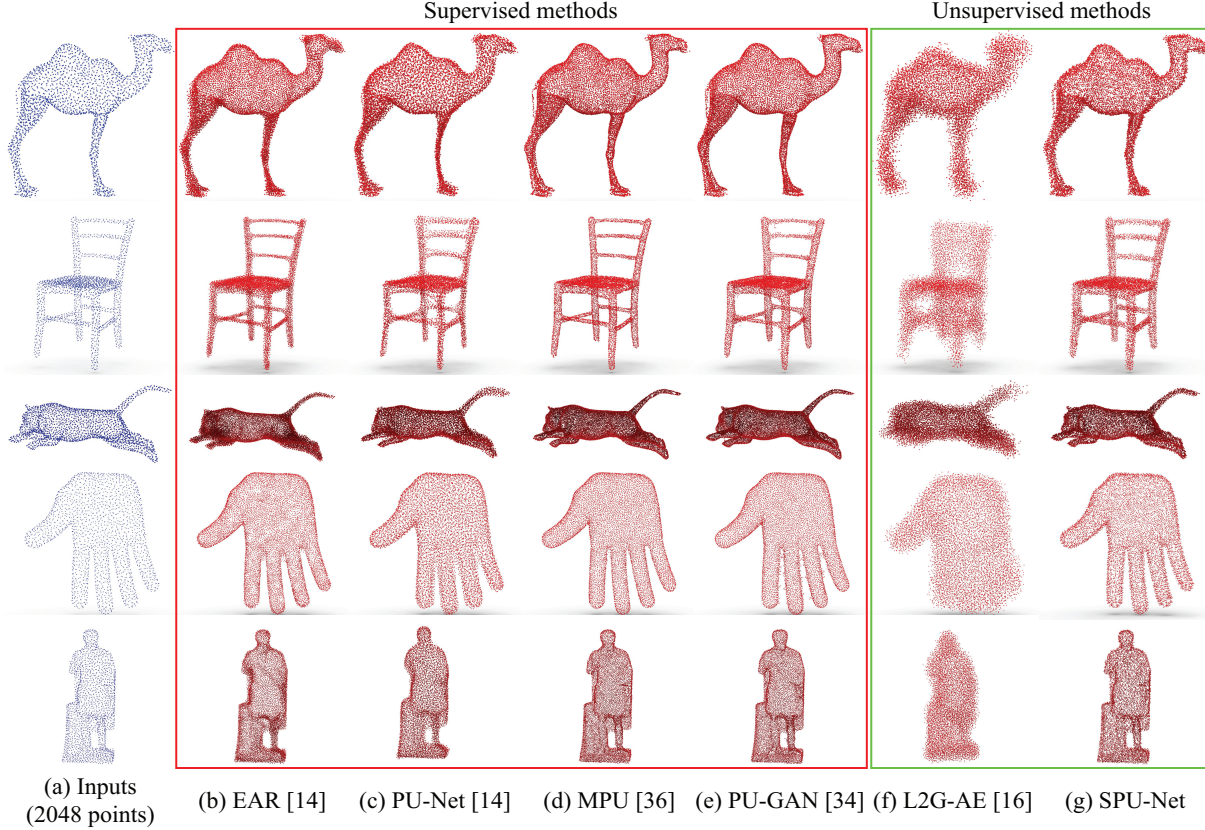


Figure 5. Visualization of upsampling results from 2,048 points to 8,192 points.

Table 2. The effects of the number of local region points  $K$ .

Metric	$K = 5$	10	15	20
CD ( $10^{-3}$ )	0.45	<b>0.38</b>	0.39	0.41
HD ( $10^{-3}$ )	3.03	<b>2.24</b>	2.51	2.48
P2F ( $10^{-3}$ )	5.91	<b>5.87</b>	6.04	6.12
UNI ( $10^{-3}$ )	12.47	<b>8.94</b>	9.65	9.90

distance (CD), Hausdorff distance (HD), point-to-surface (P2F) and uniformity (UNI,  $p = 1.2\%$ ) are adopted as evaluation metrics. For each point cloud with 2,048 points, we sample  $N_p = 24$  seed points with FPS and build local patches according to geodesic distance. We initialize the network hyper-parameters, as depicted in the network configurations. Given an input patch, we first downsample this patch into some coarse patches, which should preserve the patch point distribution in the downsampling process. Thus, the sampling methods influence the neural network in the capturing of inherent upsampling patterns. In Table 1, we report the results of three different sampling strategies, including farthest point sampling (FPS), random sampling (RS) and mixing of these two sampling methods (MS). The results show that FPS can better cover the distribution of input patches and facilitate the point upsampling process.

Then, we investigate the impact of the number of local points  $K$  in GCNs of the point feature extraction. Specif-

ically, we range the number of local region points  $K$  from 5 to 20. Table 2 illustrates the results of different  $K$ . The best performance achieves at  $K = 10$ , which can effectively capture the local region context inside local regions. Moreover, we also evaluate the robustness of our SPU-Net under different sparsity input point clouds. Figure 7 shows the upsampling point sets with varying sizes of input points from 256 to 4,096. Our method is stable even for input with only 256 points. More hyper-parameter comparisons are given in the **Supplemental Material**.

### 4.3. Ablation Study

To evaluate the components in SPU-Net, including the coarse-to-fine framework (i.e., removing the self-attention unit) and loss function terms, we remove each of them and generate upsampling results for testing models. Specifically, we remove the self-attention unit (No self-attention, NSA), the learnable grids (No learnable grid, NLG), the hierarchical grids (No hierarchical grid, NHG), the reconstruction term (No reconstruction term, NRT), the uniform loss (No uniform term, NUT) and the self-projection term (No self-projection term, NST), respectively. All above ablation studies are compared with the full network pipeline (FULL). In addition, we also replace the whole coarse-to-



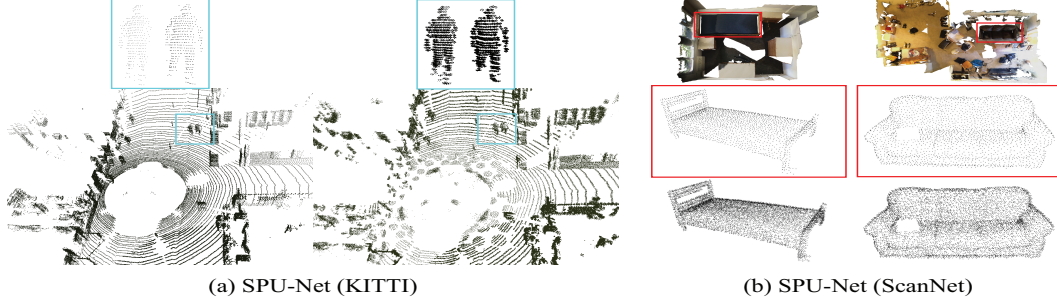


Figure 6. The upsampling results of SPU-Net under real-scanned data from KITTI [9] and ScanNet [6].

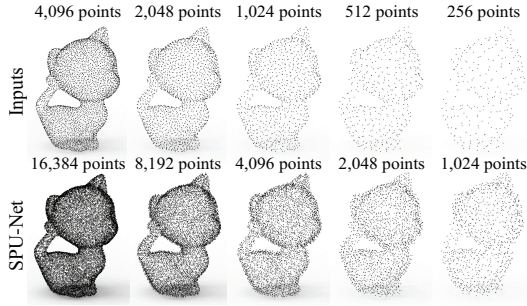


Figure 7. Upsampling point sets with varying sizes.

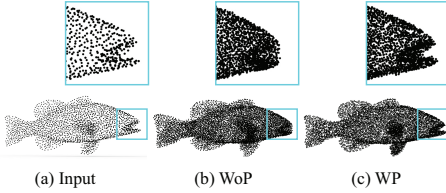


Figure 8. The upsampling results with self-projection optimization (WP) and without self-projection optimization (WoP).

Table 3. The effects of some key components in the SPU-Net.

Metric	FULL	NSA	NLG	NHG	NRT	NUT	NST	GP
CD ( $10^{-3}$ )	<b>0.38</b>	0.44	1.16	0.47	4.68	<b>0.38</b>	0.65	0.41
HD ( $10^{-3}$ )	<b>2.24</b>	2.99	29.18	3.18	34.58	2.55	6.24	2.53
P2F ( $10^{-3}$ )	5.87	6.45	5.87	6.20	40.44	<b>4.49</b>	9.01	5.93
UNI ( $10^{-3}$ )	<b>8.94</b>	12.23	47.08	9.28	802.64	13.49	19.38	10.35

fine framework with the generator of PU-GAN as a new baseline (Generator of PU-GAN, GP) to show the effectiveness of our coarse-to-fine framework. The results in Table 3 suggest that all key components play an important role in improving the performance of our SPU-Net, where the network components and the loss function work together to capture the inherent upsampling patterns. In particular, we show the visualization results in Figure 8, which demonstrate the effectiveness of self-projection optimization in projecting noisy points to the underlying object surface itself.

Table 4. Quantitative comparisons with the state-of-the-arts.

Methods	Uniformity for different $p$ ( $10^{-3}$ )					P2F	CD	HD
	0.4%	0.6%	0.8%	1.0%	1.2%	( $10^{-3}$ )	( $10^{-3}$ )	( $10^{-3}$ )
EAR [15]	16.84	20.27	23.98	26.15	29.18	5.82	0.52	7.37
PU-Net [38]	29.74	31.33	33.86	36.94	40.43	6.84	0.72	8.94
MPU [36]	7.51	7.41	8.35	9.62	11.13	3.96	0.49	6.11
PU-GAN [17]	<b>3.38</b>	<b>3.49</b>	<b>3.44</b>	<b>3.91</b>	<b>4.64</b>	<b>2.33</b>	<b>0.28</b>	4.64
L2G-AE [20]	24.61	34.61	44.86	55.31	64.94	39.37	6.31	63.23
Ours (Train2Test)	4.53	4.82	5.68	6.69	7.95	5.97	0.38	2.24
Ours (All2Test)	4.71	5.02	65.91	7.03	8.50	5.79	0.37	2.55
Ours (Test2Test)	4.82	5.14	5.86	6.88	8.13	6.85	0.41	<b>2.18</b>

#### 4.4. Point Cloud Upsampling

We compare our SPU-Net with several state-of-the-art upsampling methods both quantitatively and qualitatively, including EAR [15], PU-Net [38], MPU [36], PU-GAN [17] and L2G-AE [20]. For EAR, we use its demo code and generate the best results by fine-tuning the associated parameters. In Table 4, our SPU-Net achieves comparable results with existing supervised upsampling networks, such as PU-Net [38], and outperforms another unsupervised method L2G-AE [20]. All evaluation metrics are the same as the ones employed in PU-GAN, where the uniformity is evaluated with varying scales  $p$ . And we report multiple testing results of our SPU-Net under different training dataset, including only training dataset (Train2Test), only test dataset (Test2Test) and the entire dataset (All2Test). In addition, we also show the qualitative comparisons in Figure 5, which shows that our SPU-Net can generate upsampled point sets with more details.

#### 4.5. Upsampling Real-scanned Data

Figure 6 shows that our upsampling results on real-scanned point sets from KITTI [8] and ScanNet [6]. From the results, our SPU-Net can generate dense and uniform upsampled point sets from sparse ones. In particular, all results listed in Figure 6 are trained under the real-scanned data without ground-truth dense point sets. Therefore, our SPU-Net can expand the training data from synthetic data to the real-scanned data.



## 5. Conclusion

In this paper, we propose a novel self-supervised point cloud upsampling method to generate dense and uniform point set from sparse inputs without the supervision of ground-truth dense point clouds. Our coarse-to-fine reconstruction framework effectively facilitates the point upsampling by point feature extraction and point feature expansion. In addition, our self-projection optimization successfully projects noisy points onto the underlying object surface itself, which greatly improves the quality of point cloud upsampling in an unsupervised manner. Our experimental results demonstrate that our method can achieve good performance on both synthetic and real-scanned datasets, even comparable results to the state-of-the-art supervised method.

## References

- [1] Visionair. <http://www.infra-visionair.eu/>. Accessed: 2020-11-09. 1
- [2] Idan Achituve, Haggai Maron, and Gal Chechik. Self-Supervised Learning for Domain Adaptation on Point-Clouds. *arXiv preprint arXiv:2003.12641*, 2020. 3
- [3] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning Representations and Generative Models for 3D Point Clouds. In *International Conference on Machine Learning*, pages 40–49. PMLR, 2018. 3, 5
- [4] Marc Alexa, Johannes Behr, Daniel Cohen-Or, Shachar Fleishman, David Levin, and Claudio T. Silva. Computing and Rendering Point Set Surfaces. *IEEE Transactions on Visualization and Computer Graphics*, 9(1):3–15, 2003. 2
- [5] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. ShapeNet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015. 1
- [6] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-annotated 3D Reconstructions of Indoor Scenes. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017. 1, 2, 8
- [7] Xiang Gao, Wei Hu, and Guo-Jun Qi. GraphTER: Unsupervised Learning of Graph Transformation Equivariant Representations via Auto-Encoding Node-wise Transformations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7163–7172, 2020. 3
- [8] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. *The International Journal of Robotics Research*, 32(11):1231–1237, 2013. 2, 8
- [9] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision Meets Robotics: The KITTI Dataset. *International Journal of Robotics Research (IJRR)*, 2013. 8
- [10] Zhizhong Han, Xiyang Wang, Yu-Shen Liu, and Matthias 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*, pages 10441–10450. IEEE, 2019. 3
- [11] Kaveh Hassani and Mike Haley. Unsupervised Multi-Task Feature Learning on Point Clouds. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8160–8171, 2019. 3
- [12] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. RandLA-Net: Efficient Semantic Segmentation of Large-Scale Point Clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11108–11117, 2020. 2
- [13] Hui Huang, Dan Li, Hao Zhang, Uri Ascher, and Daniel Cohen-Or. Consolidation of Unorganized Point Clouds for Surface Reconstruction. *ACM Transactions on Graphics (TOG)*, 28(5):1–7, 2009. 2
- [14] Hui Huang, Shihao Wu, Minglun Gong, Daniel Cohen-Or, Uri Ascher, and Hao Zhang. Edge-Aware Point Set Resampling. *ACM Transactions on Graphics (TOG)*, 32(1):1–12, 2013. 2
- [15] H. Huang, S. Wu, M. Gong, D. Cohen-Or, U. Ascher, and H. Zhang. Edge-Aware Point Set Resampling. *ACM Transactions on Graphics*, 32:9:1–9:12, 2013. 8
- [16] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015. 6
- [17] Ruihui Li, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. PU-GAN: A Point Cloud Upsampling Adversarial Network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7203–7212, 2019. 1, 2, 3, 5, 6, 8
- [18] Yaron Lipman, Daniel Cohen-Or, David Levin, and Hillel Tal-Ezer. Parameterization-Free Projection for Geometry Reconstruction. *ACM Transactions on Graphics (TOG)*, 26(3):22–es, 2007. 2
- [19] Xinhai Liu, Zhizhong Han, Yu-Shen Liu, and Matthias Zwicker. Point2Sequence: Learning the Shape Representation of 3D Point Clouds with an Attention-Based Sequence to Sequence Network. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8778–8785, 2019. 2, 4
- [20] Xinhai Liu, Zhizhong Han, Xin Wen, Yu-Shen Liu, and Matthias Zwicker. L2G Auto-Encoder: Understanding Point Clouds by Local-to-Global Reconstruction with Hierarchical Self-Attention. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 989–997, 2019. 2, 3, 5, 8
- [21] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-Shape Convolutional Neural Network for Point Cloud Analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8895–8904, 2019. 2, 4
- [22] Yu-Shen Liu, Jean-Claude Paul, Jun-Hai Yong, Pi-Qiang Yu, Hui Zhang, Jia-Guang Sun, and Karthik Ramani. Automatic least-squares projection of points onto point clouds with applications in reverse engineering. *Computer-Aided Design*, 38(12):1251–1263, 2006. 6

- [23] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann Machines. In *International Conference on Machine Learning*, 2010. 5
- [24] Charles R. Qi, Xinlei Chen, Or Litany, and Leonidas J. Guibas. ImVoteNet: Boosting 3D Object Detection in Point Clouds With Image Votes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2020. 2
- [25] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017. 4, 5
- [26] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. SinGAN: Learning a Generative Model from a Single Natural Image. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4570–4580, 2019. 2
- [27] Hanyu Shi, Guosheng Lin, Hao Wang, Tzu-Yi Hung, and Zhenhua Wang. SpSequenceNet: Semantic Segmentation Network on 4D Point Clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4574–4583, 2020. 2
- [28] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From Points to Parts: 3D Object Detection from Point Cloud with Part-Aware and Part-Aggregation Network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [29] Xiaogang Wang, Marcelo H Ang Jr, and Gim Hee Lee. Cascaded Refinement Network for Point Cloud Completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 790–799, 2020. 2
- [30] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic Graph CNN for Learning on Point Clouds. *ACM Transactions on Graphics (TOG)*, 2019. 4
- [31] Xin Wen, Zhizhong Han, Xinhai Liu, and Yu-Shen Liu. Point2SpatialCapsule: Aggregating Features and Spatial Relationships of Local Regions on Point Clouds Using Spatial-Aware Capsules. *arXiv preprint arXiv:1908.11026*, 2019. 2
- [32] Xin Wen, Tianyang Li, Zhizhong Han, and Yu-Shen Liu. Point Cloud Completion by Skip-Attention Network with Hierarchical Folding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1939–1948, 2020. 2
- [33] Xin Wen, Tianyang Li, Zhizhong Han, and Yu-Shen Liu. Point Cloud Completion by Skip-Attention Network With Hierarchical Folding. In *IEEE Conference on Computer Vision and Pattern Recognition*, June 2020. 5
- [34] Shihao Wu, Hui Huang, Minglun Gong, Matthias Zwicker, and Daniel Cohen-Or. Deep Points Consolidation. *ACM Transactions on Graphics (TOG)*, 34(6):1–13, 2015. 2
- [35] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. FoldingNet: Point Cloud Auto-Encoder via Deep Grid Deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 206–215, 2018. 3, 5
- [36] Wang Yifan, Shihao Wu, Hui Huang, Daniel Cohen-Or, and Olga Sorkine-Hornung. Patch-Based Progressive 3D Point Set Upsampling. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5958–5967, 2019. 1, 3, 8
- [37] Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. EC-Net: An Edge-Aware Point Set Consolidation Network. In *Proceedings of the European Conference on Computer Vision*, pages 386–402, 2018. 3
- [38] Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. PU-Net: Point Cloud Upsampling Network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2790–2799, 2018. 1, 2, 3, 6, 8
- [39] Yuan Yuan, Siyuan Liu, Jiawei Zhang, Yongbing Zhang, Chao Dong, and Liang Lin. Unsupervised image super-resolution using cycle-in-cycle generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 701–710, 2018. 2
- [40] Ling Zhang and Zhigang Zhu. Unsupervised Feature Learning for Point Cloud Understanding by Contrasting and Clustering Using Graph Convolutional Neural Networks. In *International Conference on 3D Vision*, pages 395–404. IEEE, 2019. 3