

SVDFormer: Complementing Point Cloud via Self-view Augmentation and Self-structure Dual-generator

Zhe Zhu¹, Honghua Chen¹, Xing He¹, Weiming Wang^{2†}, Jing Qin³, Mingqiang Wei^{1†}

¹Nanjing University of Aeronautics and Astronautics

²Hong Kong Metropolitan University

³The Hong Kong Polytechnic University

zhuzhe0619@nuaa.edu.com; chenhonghuacn@gmail.com; hexing@nuaa.edu.cn;

wmwang@hkmu.edu.hk; harry.qin@polyu.edu.hk; mqwei@nuaa.edu.cn

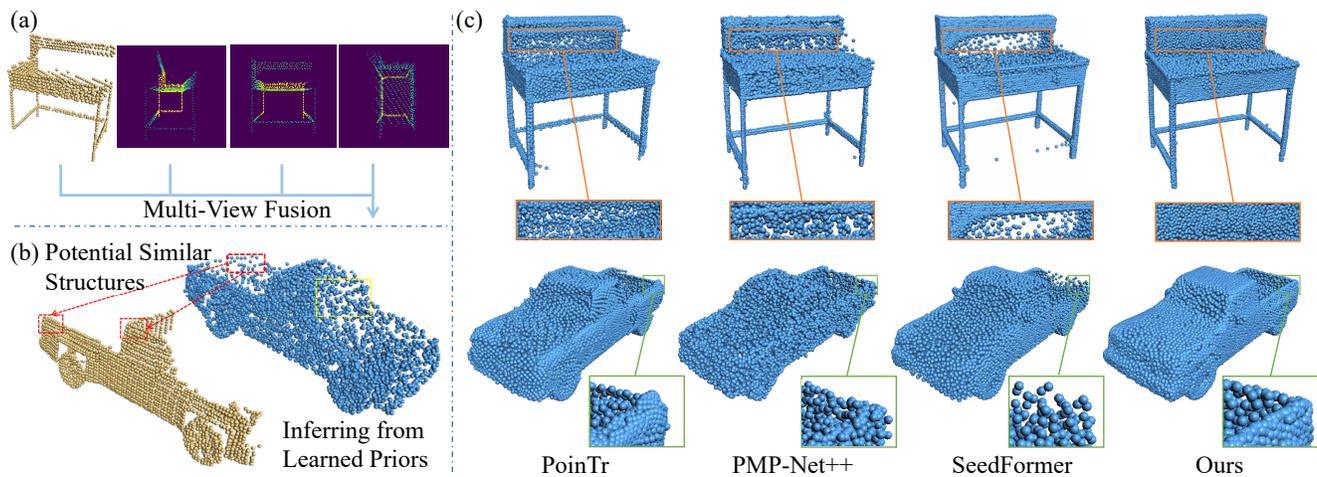


Figure 1. The characteristics of our method and visual comparison of point cloud completion results. (a) SVDFormer understands incomplete shapes from self-projected multiple views. (b) SVDFormer collaborates both geometric similarities (red boxes) and shape priors (yellow boxes) for shape refinement. (c) Qualitative comparison of our SVDFormer with PoinTr [39], PMP-Net++ [29], and SeedFormer [46].

Abstract

In this paper, we propose a novel network, SVDFormer, to tackle two specific challenges in point cloud completion: understanding faithful global shapes from incomplete point clouds and generating high-accuracy local structures. Current methods either perceive shape patterns using only 3D coordinates or import extra images with well-calibrated intrinsic parameters to guide the geometry estimation of the missing parts. However, these approaches do not always fully leverage the cross-modal self-structures available for accurate and high-quality point cloud completion. To this end, we first design a Self-view Fusion Network that leverages multiple-view depth image information to observe incomplete self-shape and generate a compact global shape.

To reveal highly detailed structures, we then introduce a refinement module, called Self-structure Dual-generator, in which we incorporate learned shape priors and geometric self-similarities for producing new points. By perceiving the incompleteness of each point, the dual-path design disentangles refinement strategies conditioned on the structural type of each point. SVDFormer absorbs the wisdom of self-structures, avoiding any additional paired information such as color images with precisely calibrated camera intrinsic parameters. Comprehensive experiments indicate that our method achieves state-of-the-art performance on widely-used benchmarks. Code is available at <https://github.com/czvvvd/SVDFormer>.

†Co-corresponding authors

1. Introduction

Point cloud completion plays an essential role in 3D vision applications and remains an active research topic in recent years. To tackle this task, a variety of learning-based techniques have been proposed, of which many demonstrate encouraging results [40, 13, 43, 39, 31, 34, 41, 20, 46, 42]. However, the sparsity and large structural incompleteness of captured point clouds still limit the ability of these methods to produce satisfactory results.

We observe that there are two primary challenges in this task. The first challenge is that crucial semantic parts may be absent, resulting in a vast solution space for point-based networks [40, 31, 39, 46] to identify plausible global shapes and locate missing regions. Some alternative methods attempt to address this issue by incorporating additional color images [44, 1, 47], but the paired images are hard to obtain, as well as the well-calibrated intrinsic parameters. The second one is how to infer detailed structures. Some recent methods [31, 34] utilize skip-connections between multiple refinement steps, allowing them to better leverage learned shape pattern priors to iteratively recover finer details. Some other methods prioritize preserving the original detail information by no-pooling encoding [46] or structural relational enhancement [16]. However, all above mentioned approaches typically employ a unified refinement strategy for all surface regions, which hinders the generation of geometric details for various missing regions. By observing and analyzing the partial inputs, we find that the missing surface regions can be classified into two types. The first type lacks similar structures in the input shape, and their reconstruction heavily relies on the learned shape prior. The second type is consistent with the local structures that are present in the partial input, and their recovery can be facilitated by appropriate geometric regularity [45]. For instance, LiDAR scans in KITTI [7] are highly sparse and contain limited information for generating fine details. Existing refinement strategies tend to produce and preserve implausible line-like shapes (see from Figure 7).

Based on the above observations, we propose a new neural network for point cloud completion called SVDFormer. Our method makes improvements by fully leveraging self-structure information in a coarse-to-fine paradigm.

First, similar to how a human would perceive and locate the missing areas of a physical object by observing it from different viewpoints, we aim to drive the neural network to absorb this knowledge by augmenting the data representation. To achieve this, we design a Self-View Fusion Network (SVFNet) that learns an effective descriptor, well depicting the global shape from both the point cloud data and depth maps captured from multiple viewpoints (see from Figure 1 (a)). To better exploit such kind of cross-modal information, we specifically introduce a feature fusion module to enhance the inter-view relations and improve the dis-

criminative power of multi-view features.

Regarding the second challenge, our insight is to disentangle refinement strategies conditioned on the structural type of each point to reveal detailed geometric structures. Therefore, we design a Self-structure Dual-Generator (SDG) with a pair of parallel refinement units, called *Structure Analysis* and *Similarity Alignment*, respectively. The former unit analyzes the generated coarse point clouds by explicitly encoding local incompleteness, which enables it to match learned geometric patterns of training data to infer underlying shapes. The *Similarity Alignment* unit finds the features of similar structures for every point, thus making it easier to refine its local region by mimicking the geometry of input local structures. With the aid of this dual-path design, our method can generate reasonable results for different types of input shapes, including symmetrical synthetic models with various degrees of incompleteness and highly sparse real-world scans.

Extensive experiments demonstrate that SVDFormer achieves state-of-the-art performance on widely-used benchmarks. Our key contributions are listed below:

- We design a novel network called SVDFormer, which significantly improves point cloud completion in terms of global shape understanding and details recovery.
- We propose the novel Self-view Fusion Network (SVFNet) equipped with a feature fusion module to enhance the multi-view and cross-modal feature, which can output a plausible global shape.
- We introduce a Self-structure Dual-Generator (SDG) for refining the coarse completion. It enables our method to handle various kinds of incomplete shapes by jointly learning the local pattern priors and self-similarities of 3D shapes.

2. Related Work

2.1. Learning-based Shape Completion

Early learning-based methods [4, 8, 18, 22] often rely on voxel-based representations for 3D convolutional neural networks. However, these approaches are limited by their high computational cost and limited resolution. Alternatively, GRNet [33] and VE-PCN [24] use 3D grids as an intermediate representation for point-based completion.

In recent years, several methods are proposed to directly process points by end-to-end networks. One pioneering point-based work is PCN [40], which uses a shared multi-layer perceptron (MLP) to extract features and generates additional points using a folding operation [35] in a coarse-to-fine manner. Inspired by it, a lot of point-based methods [25, 14, 28, 31, 46, 39] have been proposed.

Later, to address the issue of limited information available in partial shapes using only point data, several works [44, 1, 47, 12] have explored the use of auxiliary in-

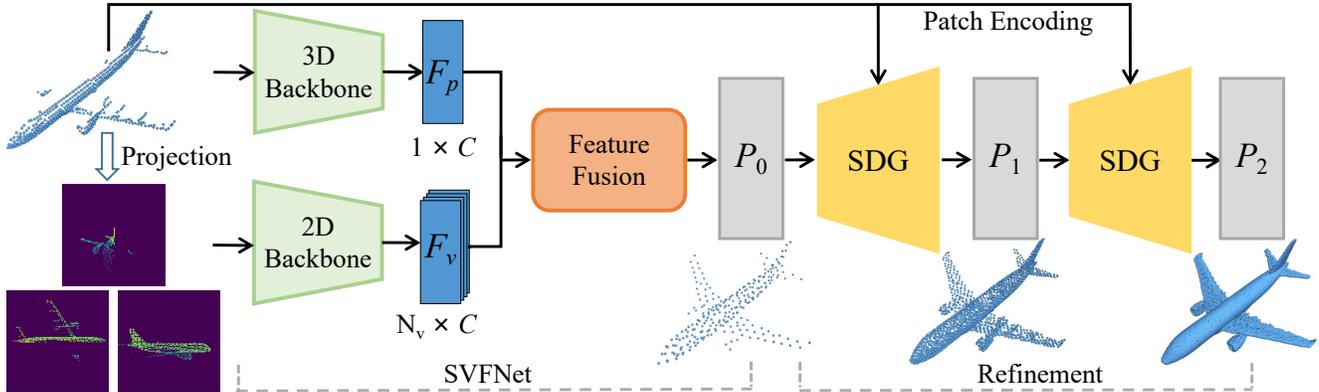


Figure 2. The architecture of SVDFormer. SVFNet first generates a global shape from the cross-modal input. The coarse completion is then upsampled and refined with two SDGs.

put to enhance performance. We call them cross-modal-based methods. These approaches involve the combination of rendered color images and partial point clouds, along with the corresponding camera parameters. Although these methods have shown promising results, they often require additional input that is difficult to obtain in practical settings. Different from these 3D data-driven methods, MVCN [11] operates completion solely in the 2D domain using a conditional GAN. However, it lacks the ability to supervise the results using ground truth with rich space information. Besides, some other methods [32, 1] seek to supervise point cloud completion in the 2D domain. The 2D projections of completed points are used to calculate loss by comparing them to the ground-truth depth. In contrast to these methods, we propose to utilize 2D input by observing self-structures to understand the overall shape. As a result, our method achieves a more comprehensive perception of the overall shape without requiring additional information or differentiable rendering during training.

Considering the high-quality details generation, a variety of strategies have been introduced by learning shape context and local spatial relationships. To achieve this goal, state-of-the-art methods design various refinement modules to learn better shape priors from the training data. SnowflakeNet [31] introduces Snowflake Point Deconvolution (SPD), which leverages skip-transformer to model the relation between parent points and child points. FB-Net [34] adopts the feedback mechanism during refinement and generates points in a recurrent manner. LAKe-Net [20] integrates its surface-skeleton representation into the refinement stage, which makes it easier to learn the missing topology part. Another type of method tends to preserve and exploit the local information in partial input. One direct approach is to predict the missing points by combining the results with partial input data [13, 39]. As the point set can be viewed as a token sequence, PoinTr [39] employs the transformer architecture [23] to predict the miss-

ing point proxies. SeedFormer [46] introduces a shape representation called patch seeds for preventing the loss of local information during pooling operation. Some other approaches [15, 16, 42] propose to enhance the generated shapes by exploiting the structural relations in the refinement stage. However, these strategies employ a unified refinement strategy for all points, which limits their ability to generate pleasing geometric details for different points. Our approach differs from theirs by breaking down the shape refinement task into two sub-goals, and adaptively extracting reliable features for different partial areas.

2.2. Multi-view fusion for Shape Learning

View-based 3D shape recognition techniques have gained significant attention in recent years. The classic Multi-View Convolutional Neural Network (MVCNN) model was introduced in [19], where color images are fed into a CNN and subsequently combined by a pooling operation. However, this approach has the fundamental drawback of ignoring view relations. Following works [6, 27, 9, 36] propose various strategies to tackle this problem. For example, Yang *et al.* [36] obtains a discriminative 3D object representation by modeling region-to-region relations. LSTM is used to build the inter-view relations [5]. Since the cross-modal data are more available recently, methods are proposed to fuse features of views and point clouds [37, 38]. Inspired by the success of multi-view fusion, our method utilizes point cloud features to enhance relationships between multiple views obtained by self-view augmentation.

3. Method

The input of our SVDFormer consists of three parts: a partial and low-res point cloud $P_{in} \subseteq \mathbb{R}^{N \times 3}$, N_V camera locations $VP \subseteq \mathbb{R}^{N_V \times 3}$ (three orthogonal views in our experiments), and N_V depth maps $D \subseteq \mathbb{R}^{N_V \times 1 \times H \times W}$. Given these inputs, our goal is to estimate a complete point cloud $P_2 \subseteq \mathbb{R}^{N_2 \times 3}$ in a coarse-to-fine manner. The over-

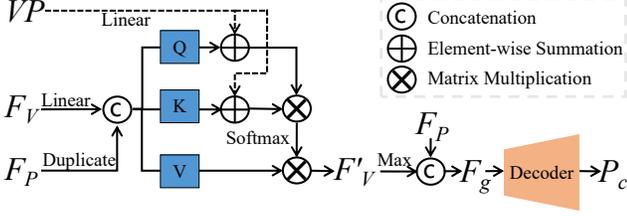


Figure 3. Illustration of the feature fusion module.

all architecture is exhibited in Figure 2, which comprises two parts: an SVFNet and a refiner equipped with two SDG modules. SVFNet first leverages multiple self-projected depth maps to produce a globally completed shape $P_0 \subseteq \mathbb{R}^{N_0 \times 3}$. Subsequently, two SDGs gradually refine and up-sample P_0 to yield the final point cloud P_2 , which exhibits geometric structures with high levels of detail. Note that unlike some recent cross-modal approaches [44, 47, 1, 41], *our method makes full use of self-structures and does not require any additional paired information such as color images with precisely calibrated camera intrinsic parameters* [44, 47]. Depth maps are directly yielded by projecting point clouds themselves from controllable viewpoints during the data-preprocessing stage.

3.1. SVFNet

The SVFNet aims to observe the partial input from different viewpoints and learns an effective descriptor to produce a globally plausible and complete shape. We first extract a global feature F_p from P_{in} using a point-based 3D backbone network and a set of view features F_V from the N_V depth maps using a CNN-based 2D backbone network. We directly adopt well-established backbone networks. In detail, the PointNet++ [17] with three set abstraction layers encodes P_{in} in a hierarchical manner and the ResNet-18 model [10] is employed as the 2D backbone.

However, how to effectively fuse the above cross-modal features is challenging. In our early experiments, we directly concatenate these features, but the produced shape is less pleasing (see the ablation studies in Section 4.5). This may be caused by the domain gap between 2D and 3D representations. To resolve this problem, we propose a new feature fusion module, to fuse F_p and F_V , and output a global shape descriptor F_g , followed by a decoder to generate the global shape P_c . The decoder uses a 1D Conv-Transpose layer to transform F_g to a set of point-wise features and regresses 3D coordinates with a self-attention layer [23]. Finally, we adopt a similar approach to previous studies [31, 46], where we merge P_c and P_{in} and resample the merged output to generate the coarse result P_0 .

Feature Fusion. As shown in Figure 3, F_V is first transformed to query, key, and value tokens via linear projection and the guidance of global shape feature F_P . Then, to

enhance the discriminability of view features, the attention weights are calculated based on the query and key tokens conditioned on the projected viewpoints VP . Detailedly, we map VP into the latent space through a linear transformation and then use them as positional signals for feature fusion. After the elemental-wise product, each feature in F'_V combines the relational information from other views under the guidance of F_P . Finally, the output shape descriptor F_g is derived from F'_V via maximum pooling.

3.2. SDG

The SDG seeks to generate a set of coordinate offsets to fine-tune and upsample the coarse shape, based on the structural type of the missing surface region. To achieve it, SDG is designed as a dual-path architecture as shown in Figure 4, which consists of two parallel units named *Structure Analysis* and *Similarity Alignment*, respectively. Overall, fed with the partial input P_{in} and coarse point cloud P_{l-1} outputted in the last step, we obtain the combined point-wise feature $F_l \subseteq \mathbb{R}^{N \times 2C}$. F_l comprises two kinds of sources of shape information: one is derived from learned shape priors, while the other is learned from the similar geometric patterns found within P_{in} . F_l is then projected to a higher dimensional space and reshaped to produce a set of up-sampled offsets $O_l \subseteq \mathbb{R}^{rN \times 3}$, where r represents the upsampling rate. The predicted offsets are further added back to P_{l-1} to obtain a new completion result. Note that we iterate SDG twice, as shown in Figure 2.

3.2.1 Structure Analysis

Since detailed geometries from missing regions are harder to be recovered, we embed an incompleteness-aware self-attention layer to explicitly encourage the network to focus more on the miss regions. Specifically, P_{l-1} is first concatenated with the shape descriptor F_g , and then embedded into a set of point-wise feature $F_{l-1} = \{f_i\}_{i=1}^{N_{l-1}}$ by a linear layer. Next, F_{l-1} is fed to the incompleteness-aware self-attention layer to obtain a set of features $F_Q = \{q_i\}_{i=1}^{N_{l-1}}$, which encodes the point-wise incompleteness information. q_i is computed by:

$$q_i = \sum_{j=1}^{N_{l-1}} a_{i,j} (f_j W_V) \quad , \quad (1)$$

$$a_{i,j} = \text{Softmax}((f_i W_Q + h_i)(f_j W_K + h_j)^T)$$

where W_Q , W_K , and W_V are learnable matrix with the size of $C \times C$. h_i is a vector that represents the degree of incompleteness for each point x in P_{l-1} . Intuitively, points in missing regions tend to have a larger distance value to the partial input. We thus calculate the incompleteness by:

$$h_i = \text{Sinusoidal}\left(\frac{1}{\gamma} \min_{y \in P_{in}} \|x - y\|\right), \quad (2)$$

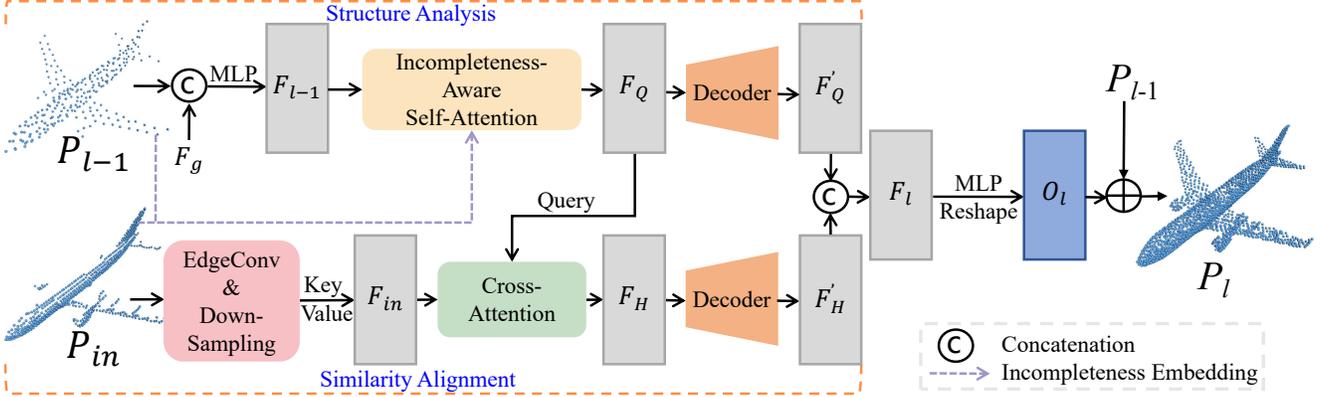


Figure 4. The architecture of SDG. The upper path represents Structure Analysis and the lower path represents Similarity Alignment. Each sub-network generates an offset feature which is then concatenated with each other and used to regress into the coordinate offsets.

where γ is a scaling coefficient. We set it as 0.2 in our experiment. The sinusoidal function [23] is used to ensure that h_i has the same dimension as the embeddings of query, key, and value. We then decode F_Q into F'_Q for further analysis of the coarse shape.

3.2.2 Similarity Alignment

The *Similarity Alignment* unit exploits the potential similar local pattern in P_{in} for each point in P_{l-1} and addresses the feature mismatch problem caused by the unordered nature of point clouds. Inspired by the point proxies in [39], we begin by using three EdgeConv layers [26] to extract a set of downsampled point-wise feature F_{in} . Each vector in F_{in} captures the local context information. Since there could exist long-range similar structures, we then perform feature exchange by cross-attention, which is a classical solution for feature alignment. The calculation process is similar to vanilla self-attention. The only difference lies in that the query matrix is produced by F_Q , while F_{in} serves as the key and value vectors. The cross-attention layer outputs point-wise feature $F_H \subseteq \mathbb{R}^{N_{l-1} \times C}$, which integrates similar local structures in P_{in} for each point in the coarse shape P_{l-1} . In this way, this unit can model the geometric similarity between two point clouds and facilitate the refinement of points with similar structures in the input. Similar with the structure analysis unit, F_H is also decoded into a new feature F'_H . These two decoders have the same architecture, which is implemented with two layers of self-attention [23]. For more details of the used self-attention, cross-attention, and decoders, please refer to the supplemental files.

3.3. Loss Function

In order to measure the differences between the generated point cloud and the ground truth P_{gt} , we use the Chamfer Distance (CD) as our loss function, which is a common choice in recent works. To facilitate the coarse-to-fine gen-

eration process, we regularize the training by computing the loss function as:

$$\mathcal{L} = \mathcal{L}_{CD}(P_c, P_{gt}) + \sum_{i=1,2} \mathcal{L}_{CD}(P_i, P_{gt}) \quad (3)$$

It is worth noting that we downsample the P_{gt} to the same density as P_c, P_1, P_2 in order to compute the losses.

4. Experiment

4.1. Dataset and Evaluation Metric

We first use the PCN [40] and ShapeNet-55/34 [39] dataset for evaluation. To ensure a fair comparison, we follow the same experiment settings as previous methods [39, 46]. PCN contains shapes from 8 categories in ShapeNet [2]. The ground-truth complete point cloud has 16,384 points and the partial input has 2,048 points. ShapeNet-55 [39] is also created based on ShapeNet [2] and contains shapes from 55 categories. The ground-truth point cloud has 8,192 points and the partial input has 2,048 points. ShapeNet-34 contains 34 categories for training and leaves 21 unseen categories for testing, which is used for evaluation of generalization ability on novel categories that are unseen during training. Secondly, to evaluate the generalization ability in real-world scenarios, we test our method on both KITTI [7] and ScanNet [3], which contain partial point clouds extracted from LiDAR scans and RGB-D scans, respectively. Specifically, we test on 2,401 KITTI cars extracted by [40], and 100 chair point clouds from ScanNet. We use CD, Density-aware CD (DCD) [30], and F1-Score as evaluation metrics. Following the recent work [46], we report the ℓ^1 version of CD for PCN and the ℓ^2 version of CD for Shapenet-55, for an easier comparison.

4.2. Results on the PCN Dataset

We compare SVDFormer with state-of-the-art methods [40, 33, 25, 43, 39, 31, 29, 34, 46, 42] in Table 1. CD

Table 1. Quantitative results on the PCN dataset. (ℓ^1 CD $\times 10^3$ and F1-Score@1%)

Methods	Plane	Cabinet	Car	Chair	Lamp	Couch	Table	Boat	CD-Avg↓	DCD↓	F1↑
PCN [40]	5.50	22.70	10.63	8.70	11.00	11.34	11.68	8.59	9.64	-	0.695
GRNet [33]	6.45	10.37	9.45	9.41	7.96	10.51	8.44	8.04	8.83	0.622	0.708
CRN [25]	4.79	9.97	8.31	9.49	8.94	10.69	7.81	8.05	8.51	-	0.652
NSFA [43]	4.76	10.18	8.63	8.53	7.03	10.53	7.35	7.48	8.06	-	0.734
PoinTr [39]	4.75	10.47	8.68	9.39	7.75	10.93	7.78	7.29	8.38	0.611	0.745
SnowflakeNet [31]	4.29	9.16	8.08	7.89	6.07	9.23	6.55	6.40	7.21	0.585	0.801
SDT [42]	4.60	10.05	8.16	9.15	8.12	10.65	7.64	7.66	8.24	-	0.754
PMP-Net++ [29]	4.39	9.96	8.53	8.09	6.06	9.82	7.17	6.52	7.56	0.611	0.781
FBNet [34]	3.99	9.05	7.90	7.38	5.82	8.85	6.35	6.18	6.94	-	-
Seedformer [46]	3.85	9.05	8.06	7.06	5.21	8.85	6.05	5.85	6.74	0.583	0.818
Ours	3.62	8.79	7.46	6.91	5.33	8.49	5.90	5.83	6.54	0.536	0.841

Table 2. Quantitative results on ShapeNet-55. CD-S, CD-M, and CD-H stand for CD values under the easy, moderate, and hard difficulty levels, respectively. (ℓ^2 CD $\times 10^3$ and F1-Score@1%)

Methods	CD-S	CD-M	CD-H	CD-Avg↓	DCD-Avg↓	F1↑
FoldingNet [35]	2.67	2.66	4.05	3.12	-	0.082
PCN [40]	1.94	1.96	4.08	2.66	0.618	0.133
TopNet [21]	2.26	2.16	4.3	2.91	-	0.126
PFNet [13]	3.83	3.87	7.97	5.22	-	0.339
GRNet [33]	1.35	1.71	2.85	1.97	0.592	0.238
PoinTr [39]	0.58	0.88	1.79	1.09	0.575	0.464
SeedFormer [46]	0.50	0.77	1.49	0.92	0.558	0.472
Ours	0.48	0.70	1.30	0.83	0.541	0.451

values are courtesy of [46, 39], while F1-Score and DCD values are computed using their pre-trained models. The quantitative results demonstrate that SVDFormer achieves almost the best performance across all metrics. Especially, our method outperforms SeedFormer by 8.06% in DCD.

Figure 5 provides a visual comparison of the results produced by the different methods. In the case of the car and plane models, all methods are successful in generating the overall shapes. However, SVDFormer outperforms the other methods by producing sharper and more complete edges for detailed structures, such as plane wings and car spoilers. This is due to the generation ability of SDG. In the case of chair and couch models, SVDFormer can accurately locate missing regions and generate points along the holes in the models, leading to more faithful results.

4.3. Results on the ShapeNet-55/34 Dataset

The test set of ShapeNet-55 can be classified into three levels of difficulty: simple (S), moderate (M), and hard (H), which correspond to different numbers of missing points (2,048, 4,096, and 6,144). The quantitative results are pre-

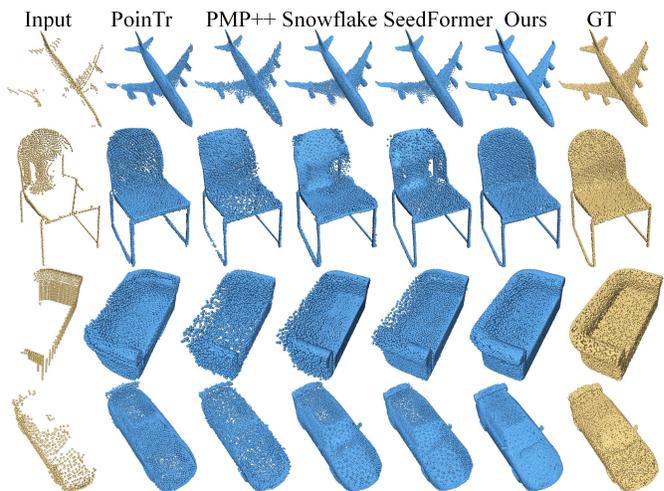


Figure 5. Visual comparisons with recent methods [39, 29, 31, 46] on the PCN dataset. Our method produces the most complete and detailed structures compared to its competitors.

sented in Table 2, consisting of CD values for three difficulty levels and the average value of two additional metrics. Our method achieves the best result in CD across all difficulty settings. Notably, SVDFormer outperforms the state-of-the-art method SeedFormer, achieving a 12.8% improvement in CD for the hard difficulty level. The results under different difficulty levels are shown in Figure 6. Compared to PoinTr and SeedFormer, our method produces smoother surfaces. The visual results clearly demonstrate that SVDFormer is capable of efficiently recovering geometries from shapes with varying degrees of incompleteness.

We further evaluate SVDFormer on the ShapeNet-34 dataset. The results on both seen and unseen categories are detailed in Table 3, which shows that SVDFormer achieves the best performance in terms of all three metrics.

Table 3. Quantitative results on ShapeNet-34.

Methods	34 seen categories						21 unseen categories					
	CD-S	CD-M	CD-H	CD-Avg↓	DCD-Avg↓	F1↑	CD-S	CD-M	CD-H	CD-Avg↓	DCD-Avg↓	F1↑
FoldingNet [35]	1.86	1.81	3.38	2.35	-	0.139	2.76	2.74	5.36	3.62	-	0.095
PCN [40]	1.87	1.81	2.97	2.22	0.624	0.150	3.17	3.08	5.29	3.85	0.644	0.101
TopNet [40]	1.77	1.61	3.54	2.31	-	0.171	2.62	2.43	5.44	3.50	-	0.121
PFNet [13]	3.16	3.19	7.71	4.68	-	0.347	5.29	5.87	13.33	8.16	-	0.322
GRNet [33]	1.26	1.39	2.57	1.74	0.600	0.251	1.85	2.25	4.87	2.99	0.625	0.216
PoinTr [39]	0.76	1.05	1.88	1.23	0.575	0.421	1.04	1.67	3.44	2.05	0.604	0.384
SeedFormer [46]	0.48	0.70	1.30	0.83	0.561	0.452	0.61	1.07	2.35	1.34	0.586	0.402
Ours	0.46	0.65	1.13	0.75	0.538	0.457	0.61	1.05	2.19	1.28	0.554	0.427

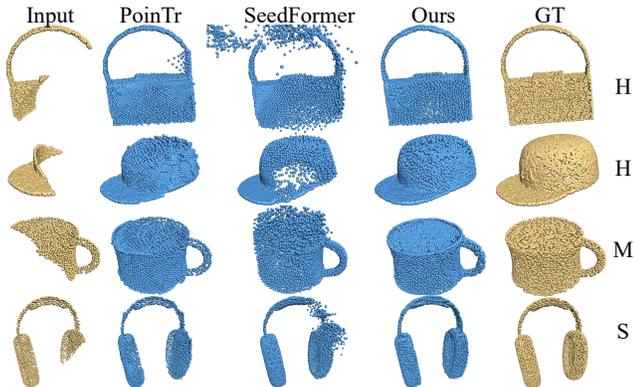


Figure 6. Visual comparison with two representative approaches [39, 46] on ShapeNet-55. H (Hard), M (Moderate), and S (Simple) stand for the three difficulty levels.

4.4. Results on Real-world Scans

For the real-world scans, as there is no ground truth available for real-world partial point clouds, we evaluate the models pre-trained on PCN without fine-tuning or re-training. We report the Minimal Matching Distance (MMD) [40] (see from Table 4) as a quantitative evaluation metric to assess the similarity of the output to a typical car/chair for the real-world scans. Also, we demonstrate the visual comparisons in Figure 7. Our method can produce cleaner shapes with detailed structures and sharp edges. It can be concluded that our method can generate more detailed results even when the input is extremely sparse and has a different distribution from the training data.

4.5. Ablation Studies and Discussions

To ablate SVDFormer, we remove and modify the main components. All ablation variants are trained and tested on the PCN dataset. The ablation variants can be categorized as ablations on SVFNet and SDG.

Ablation on SVFNet. To investigate the impact of shape descriptor extraction methods, we compare two variants of SVFNet, and the results are presented in Table 5. In the variant A, we remove the input depth maps, and the completion

Table 4. Quantitative results on Real-world Scans. All the results are produced by models pre-trained on PCN ($MMD \times 10^3$).

Methods	GRNet [33]	PoinTr [39]	SeedFormer [46]	Ours
KITTI [7]	5.350	32.854	1.179	0.967
ScanNet [3]	2.672	2.516	2.231	1.926

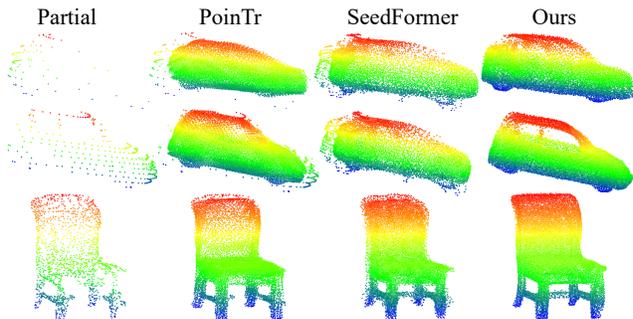


Figure 7. Visual comparison on real-world scans. All results are produced by models pre-trained on PCN.

Table 5. Effect of SVFNet. (ℓ^1 CD $\times 10^3$ and F1-Score@1%)

Methods	CD↓	DCD↓	F1↑
A : w/o Projection	6.63	0.547	0.831
B : w/o Feature Fusion	6.68	0.551	0.827
Ours	6.54	0.536	0.841

Table 6. Effect of SDG. (ℓ^1 CD $\times 10^3$ and F1-Score@1%)

Methods	Analysis	Alignment	Embedding	CD↓	DCD↓	F1↑
C	✓	✓		6.69	0.549	0.829
D	✓		✓	6.76	0.552	0.825
E		✓		6.78	0.556	0.823
F	✓			6.88	0.561	0.819
Ours	✓	✓	✓	6.54	0.536	0.841
SeedFormer [46]				6.74	0.583	0.818
SnowflakeNet-baseline [31]				7.21	0.585	0.801
SnowflakeNet + SDG				6.73	0.553	0.828

performance is limited by relying only on 3D coordinates to understand shapes. In the variant B, we evaluate the impor-

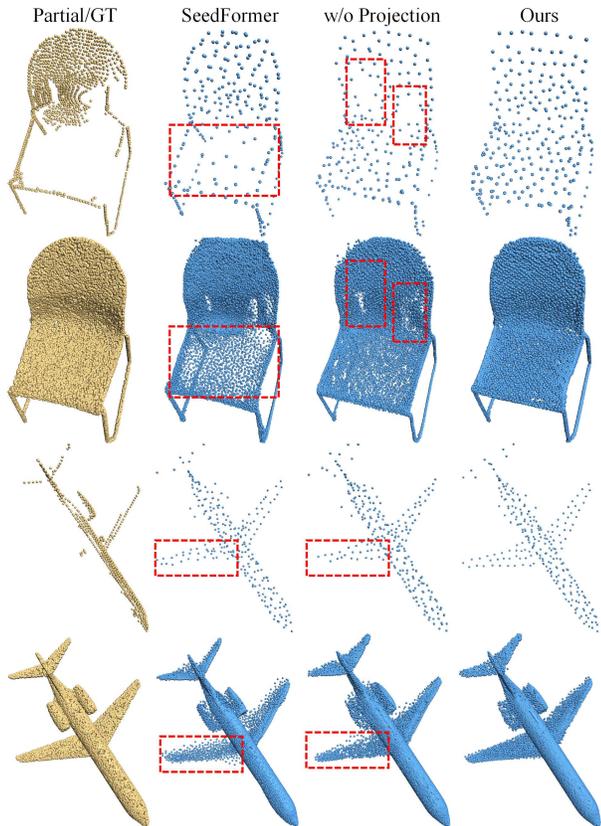


Figure 8. Visual comparison of the representative coarse-to-fine method [46] and our variant A (w/o projection) on two partial models. The upper results are the generated coarse point clouds.

tance of our Feature Fusion module by replacing the fusion of different inputs with late fusion, which directly concatenates F_p and F_v . We observe an evident drop in performance, indicating that the proposed SVFNet can effectively fuse cross-modal features.

Furthermore, to conduct a more thorough analysis of the effectiveness of our SVFNet, we generate visualizations of the results produced by our approach, our variant A, and SeedFormer, which also employ the coarse-to-fine paradigm. In Figure 8, we present the results alongside the coarse point cloud generated directly by SVFNet (patch seeds of [46]). Our analysis reveals that during the initial coarse completion stage, both SeedFormer and the variant A produce suboptimal results, such as generating too few points in missing areas. This presents a challenge for the refinement stage, making it difficult to produce satisfactory final results. Our SVFNet overcomes this challenge by leveraging multiple viewpoints to observe partial shapes. By doing so, our method is able to locate missing areas and generate a compact global shape, leading to the production of fine details in the final results. See the supplementary for additional ablation experiments for the number of views.

Ablation on SDG. Table 6 compares different variants of

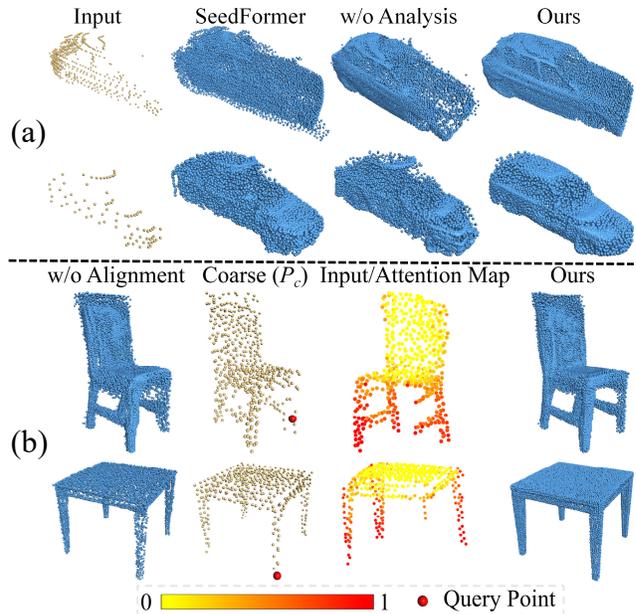


Figure 9. (a): Visual comparison of [46] and variant E (w/o Structure Analysis) on LiDAR scans. (b): Visual comparison of variant E (w/o Similarity Alignment) and generated coarse point cloud P_c on RGB-D scans. We select a query point (marked with red) in P_c and visualize the attention map in the cross-attention layer. The redder the color, the higher the similarity.

SVDFormer on the SDG module. In the variant C, we remove the Incompleteness Embedding of SDG, which results in a higher CD value and a lower F1-Score, indicating that the ability to perceive the incompleteness degree of each part is crucial for the model’s performance. In the variants D and E, we completely remove the Similarity Alignment and Structure Analysis path from SDG, respectively. The results show that the performance of the model drops when any one of these paths is removed.

To better understand and analyze SDG, we show more visual results in Figure 9. Specifically, we investigate the effectiveness of the Structure Analysis path and the Similarity Alignment unit by comparing the performance of different variants of the model on real-world scans. In Figure 9 (a), our method can generate plausible shapes, while the variant E and SeedFormer produce undesired structures due to over-preserving the partial input. This result proves the importance of the Structure Analysis path, particularly when the input contains limited information. In Figure 9 (b), we compare the results of our method with the variant D. We show the generated coarse shape and select a query point (missed in input) in P_c . We then visualize the attention map in the cross-attention layer to demonstrate the effectiveness of the Similarity Alignment unit. The results show that, for shapes with highly-similar regions, the Similarity Alignment unit can locate similar geometries of short or long-range distances, leading to finer details.

Table 7. Complexity analysis. We compare the inference time (ms) and the number of parameters (Params) of our method and three classical methods on ShapeNet-55. Our method achieves a balance between computation cost and performance.

Methods	Time	Params	CD↓	DCD↓
GRNet [33]	10.67ms	73.15M	1.97	0.592
PoinTr [39]	12.36ms	30.09M	1.07	0.575
SeedFormer [46]	40.63ms	3.24M	0.92	0.558
Ours	23.11ms	19.62M	0.83	0.546

Extending SDG to other methods. In addition, we evaluate the generation ability of SDG by replacing the SPD in SnowflakeNet [31] with SDG. As presented in Table 6, SnowflakeNet achieves better performance in terms of all metrics, when paired with our SDG module. This indicates that our disentangled refiner has better generation ability.

Complexity Analysis. We show the complexity analysis in Table 7, where the inference time on a single NVIDIA 3090 GPU, number of parameters, and the results on ShapeNet-55 are shown. The comparison indicates that our method achieves a trade-off between cost and performance.

5. Conclusion

We propose SVDFormer for point cloud completion. We start by identifying the main challenges in the completion and developing new solutions for each of them. SVDFormer leverages self-projected multi-view analysis to comprehend the overall shape and effectively perceive missing regions. Furthermore, we introduce a decoder called Self-structure Dual-generator that breaks down the shape refinement process into two sub-goals, resulting in a disentangled but improved generation ability. Experiments on various shape types demonstrate that SVDFormer achieves the state-of-the-art performance on point cloud completion.

Acknowledgements

The work described in this paper was supported by the National Natural Science Foundation of China (No. 62172218), the Research Grants Council of the Hong Kong Special Administrative Region, China (No. UGC/FDS16/E14/21), the Shenzhen Science and Technology Program (No. JCYJ20220818103401003, No. JCYJ20220530172403007), the Natural Science Foundation of Guangdong Province (No. 2022A1515010170), and Hong Kong RGC General Research Fund (No. 15218521).

References

[1] Emanuele Aiello, Diego Valsesia, and Enrico Magli. Cross-modal learning for image-guided point cloud shape completion. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and

Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 2, 3, 4

[2] 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. 5

[3] 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 conference on computer vision and pattern recognition*, pages 5828–5839, 2017. 5, 7

[4] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5868–5877, 2017. 2

[5] Guoxian Dai, Jin Xie, and Yi Fang. Siamese cnn-bilstm architecture for 3d shape representation learning. In *IJCAI*, pages 670–676, 2018. 3

[6] Yifan Feng, Zizhao Zhang, Xibin Zhao, Rongrong Ji, and Yue Gao. Gvcnn: Group-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 264–272, 2018. 3

[7] 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, 5, 7

[8] Xiaoguang Han, Zhen Li, Haibin Huang, Evangelos Kalogerakis, and Yizhou Yu. High-resolution shape completion using deep neural networks for global structure and local geometry inference. In *Proceedings of the IEEE international conference on computer vision*, pages 85–93, 2017. 2

[9] Zhizhong Han, Honglei Lu, Zhenbao Liu, Chi-Man Vong, Yu-Shen Liu, Matthias Zwicker, Junwei Han, and CL Philip Chen. 3d2seqviews: Aggregating sequential views for 3d global feature learning by cnn with hierarchical attention aggregation. *IEEE Transactions on Image Processing*, 28(8):3986–3999, 2019. 3

[10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 4

[11] Tao Hu, Zhizhong Han, Abhinav Shrivastava, and Matthias Zwicker. Render4completion: Synthesizing multi-view depth maps for 3d shape completion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019. 3

[12] Sheng Yu Huang, Hao-Yu Hsu, and Yu-Chiang Frank Wang. SPoVT: Semantic-prototype variational transformer for dense point cloud semantic completion. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022. 2

[13] Zitian Huang, Yikuan Yu, Jiawen Xu, Feng Ni, and Xinyi Le. Pf-net: Point fractal network for 3d point cloud completion.

- In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7662–7670, 2020. [2](#), [3](#), [6](#), [7](#)
- [14] Minghua Liu, Lu Sheng, Sheng Yang, Jing Shao, and Shi-Min Hu. Morphing and sampling network for dense point cloud completion. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 11596–11603, 2020. [2](#)
- [15] Liang Pan. Ecg: Edge-aware point cloud completion with graph convolution. *IEEE Robotics and Automation Letters*, 5(3):4392–4398, 2020. [3](#)
- [16] Liang Pan, Xinyi Chen, Zhongang Cai, Junzhe Zhang, Haiyu Zhao, Shuai Yi, and Ziwei Liu. Variational relational point completion network. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8524–8533, 2021. [2](#), [3](#)
- [17] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. [4](#)
- [18] David Stutz and Andreas Geiger. Learning 3d shape completion from laser scan data with weak supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1955–1964, 2018. [2](#)
- [19] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015. [3](#)
- [20] Junshu Tang, Zhijun Gong, Ran Yi, Yuan Xie, and Lizhuang Ma. Lake-net: Topology-aware point cloud completion by localizing aligned keypoints. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1726–1735, 2022. [2](#), [3](#)
- [21] Lyne P Tchammi, Vineet Kosaraju, Hamid Rezatofighi, Ian Reid, and Silvio Savarese. Topnet: Structural point cloud decoder. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 383–392, 2019. [6](#)
- [22] Jacob Varley, Chad DeChant, Adam Richardson, Joaquín Ruales, and Peter Allen. Shape completion enabled robotic grasping. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 2442–2447. IEEE, 2017. [2](#)
- [23] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [3](#), [4](#), [5](#)
- [24] Xiaogang Wang, Marcelo H Ang, and Gim Hee Lee. Voxel-based network for shape completion by leveraging edge generation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 13189–13198, 2021. [2](#)
- [25] 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](#), [5](#), [6](#)
- [26] 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)*, 38(5):1–12, 2019. [5](#)
- [27] Xin Wei, Ruixuan Yu, and Jian Sun. Learning view-based graph convolutional network for multi-view 3d shape analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. [3](#)
- [28] 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](#)
- [29] Xin Wen, Peng Xiang, Zhizhong Han, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Yu-Shen Liu. Pmp-net++: Point cloud completion by transformer-enhanced multi-step point moving paths. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. [1](#), [5](#), [6](#)
- [30] Tong Wu, Liang Pan, Junzhe Zhang, Tai Wang, Ziwei Liu, and Dahua Lin. Balanced chamfer distance as a comprehensive metric for point cloud completion. *Advances in Neural Information Processing Systems*, 34:29088–29100, 2021. [5](#)
- [31] Peng Xiang, Xin Wen, Yu-Shen Liu, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Zhizhong Han. Snowflakenet: Point cloud completion by snowflake point deconvolution with skip-transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5499–5509, 2021. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [9](#)
- [32] Chulin Xie, Chuxin Wang, Bo Zhang, Hao Yang, Dong Chen, and Fang Wen. Style-based point generator with adversarial rendering for point cloud completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4619–4628, 2021. [3](#)
- [33] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Jiageng Mao, Shengping Zhang, and Wenxiu Sun. Grnet: Gridding residual network for dense point cloud completion. In *European Conference on Computer Vision*, pages 365–381. Springer, 2020. [2](#), [5](#), [6](#), [7](#), [9](#)
- [34] Xuejun Yan, Hongyu Yan, Jingjing Wang, Hang Du, Zhihong Wu, Di Xie, Shiliang Pu, and Li Lu. Fbnet: Feedback network for point cloud completion. In *European Conference on Computer Vision*, pages 676–693. Springer, 2022. [2](#), [3](#), [5](#), [6](#)
- [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. [2](#), [6](#), [7](#)
- [36] Ze Yang and Liwei Wang. Learning relationships for multi-view 3d object recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7505–7514, 2019. [3](#)
- [37] Haoxuan You, Yifan Feng, Rongrong Ji, and Yue Gao. Pvnet: A joint convolutional network of point cloud and multi-view for 3d shape recognition. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1310–1318, 2018. [3](#)
- [38] Haoxuan You, Yifan Feng, Xibin Zhao, Changqing Zou, Rongrong Ji, and Yue Gao. Pvrnet: Point-view relation neu-

- ral network for 3d shape recognition. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 9119–9126, 2019. [3](#)
- [39] Xumin Yu, Yongming Rao, Ziyi Wang, Zuyan Liu, Jiwen Lu, and Jie Zhou. PointR: Diverse point cloud completion with geometry-aware transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12498–12507, 2021. [1](#), [2](#), [3](#), [5](#), [6](#), [7](#), [9](#)
- [40] Wentao Yuan, Tejas Khot, David Held, Christoph Mertz, and Martial Hebert. Pcn: Point completion network. In *2018 International Conference on 3D Vision (3DV)*, pages 728–737. IEEE, 2018. [2](#), [5](#), [6](#), [7](#)
- [41] Bowen Zhang, Xi Zhao, He Wang, and Ruizhen Hu. Shape completion with points in the shadow. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. [2](#), [4](#)
- [42] Wenxiao Zhang, Zhen Dong, Jun Liu, Qingan Yan, Chunxia Xiao, et al. Point cloud completion via skeleton-detail transformer. *IEEE Transactions on Visualization and Computer Graphics*, 2022. [2](#), [3](#), [5](#), [6](#)
- [43] Wenxiao Zhang, Qingan Yan, and Chunxia Xiao. Detail preserved point cloud completion via separated feature aggregation. In *European Conference on Computer Vision*, pages 512–528. Springer, 2020. [2](#), [5](#), [6](#)
- [44] Xuancheng Zhang, Yutong Feng, Siqi Li, Changqing Zou, Hai Wan, Xibin Zhao, Yandong Guo, and Yue Gao. View-guided point cloud completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15890–15899, 2021. [2](#), [4](#)
- [45] Wenbin Zhao, Jiabao Lei, Yuxin Wen, Jianguo Zhang, and Kui Jia. Sign-agnostic implicit learning of surface self-similarities for shape modeling and reconstruction from raw point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10256–10265, 2021. [2](#)
- [46] Haoran Zhou, Yun Cao, Wenqing Chu, Junwei Zhu, Tong Lu, Ying Tai, and Chengjie Wang. Seedformer: Patch seeds based point cloud completion with upsampler transformer. In *European Conference on Computer Vision*, pages 416–432. Springer, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#), [9](#)
- [47] Zhe Zhu, Liangliang Nan, Haoran Xie, Honghua Chen, Jun Wang, Mingqiang Wei, and Jing Qin. Csdn: Cross-modal shape-transfer dual-refinement network for point cloud completion. *IEEE Transactions on Visualization and Computer Graphics*, 2023. [2](#), [4](#)