# SeedFormer: Patch Seeds based Point Cloud Completion with Upsample Transformer

Haoran Zhou[1], Yun Cao[2], Wenqing Chu[2], Junwei Zhu[2],
Tong Lu[1]*, Ying Tai[2], and Chengjie Wang[2]

[1]State Key Laboratory for Novel Software Technology, Nanjing University
[2]Youtu Lab, Tencent

**Abstract.** Point cloud completion has become increasingly popular among generation tasks of 3D point clouds, as it is a challenging yet indispensable problem to recover the complete shape of a 3D object from its partial observation. In this paper, we propose a novel SeedFormer to improve the ability of detail preservation and recovery in point cloud completion. Unlike previous methods based on a global feature vector, we introduce a new shape representation, namely Patch Seeds, which not only captures general structures from partial inputs but also preserves regional information of local patterns. Then, by integrating seed features into the generation process, we can recover faithful details for complete point clouds in a coarse-to-fine manner. Moreover, we devise an Upsample Transformer by extending the transformer structure into basic operations of point generators, which effectively incorporates spatial and semantic relationships between neighboring points. Qualitative and quantitative evaluations demonstrate that our method outperforms state-of-the-art completion networks on several benchmark datasets. Our code is available at `https://github.com/hrzhou2/seedformer`.

**Keywords:** Point cloud completion, Patch Seeds, Upsample Transformer

## 1 Introduction

As a commonly-used and easily-acquired data format for describing 3D objects, point clouds have boosted wider research in computer vision for understanding 3D scenes and objects. However, raw point clouds, routinely captured by Li-DAR scanners or RGB-D cameras, are inevitably sparse and incomplete due to the limited sensor resolution and self-occlusion. It is an indispensable step to recover complete point clouds from partial/incomplete data in real-world scenes for various downstream applications [24,16,20].

Recent years have witnessed an increasing number of approaches applying deep neural networks on point cloud completion. The dominant architecture [43,1,38] employs a general encoder-decoder structure where a global feature (or called shape code) is extracted from partial inputs and is used to generate a complete point cloud in the decoding phase. However, this global feature
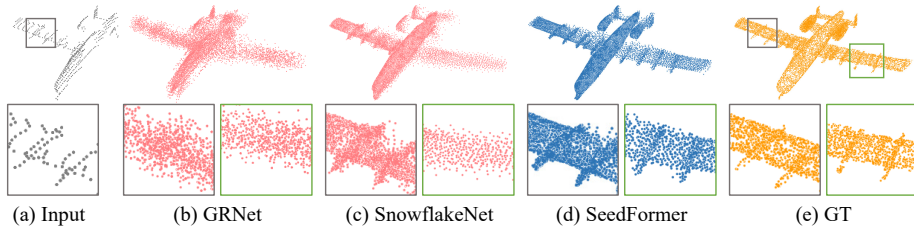
---

* Corresponding author.

Fig. 1: Visual comparison of point cloud completion results. Compared with GRNet [40] and SnowflakeNet [38], SeedFormer is better at preserving existing structures (grey bounding box) and recovering missing details (green bounding box).

structure possesses two intrinsic drawbacks in its representation ability: (i) fine-grained details are easily lost in the pooling operations in the encoding phase and can hardly be recovered from a diluted global feature in the generation, and (ii) such a global feature is captured from a partial point cloud, thus represent-ing only the "incomplete" information of the seen part, and is contrary to the objective of generating the complete shape.

Besides the overall network architecture, another problem is derived from the designs of point cloud generators. Unlike image inpainting which predicts RGB colors of the missing pixels, the 3D point generation is designed differently to predict $(x, y, z)$ coordinates which are unstructured yet continuously distributed in the 3D space. This suggests that the desired points, which can describe the missing parts of the target object, have close spatial and semantic relationships with surrounding points in their local neighborhoods. However, existing methods, either using folding-based generators [41,43,34,39,42] or following a hierarchical structure with MLP/deconvolution implementations [27,11,38], attempt to pro-cess each point independently by splitting one target point into more. These designs neglect the distribution of existing points which results in poor recovery quality of geometric details (see in Fig. 1(b) and Fig. 1(c)).

To resolve the aforementioned problems, we propose a novel point cloud completion network, namely *SeedFormer*, with better detail preservation and recovery ability as shown in Fig. 1. Based on the designed *Patch Seeds* and *Up-sample Transformer*, the decoding phase of our method consists of two main steps of: (i) first generating the complete shape from incomplete features in the seed generator, and (ii) then recovering fine-grained details in a coarse-to-fine manner. In the first stage, unlike previous methods using global features, we introduce Patch Seeds as a latent representation in the point cloud completion architecture. The Patch Seeds characterize the complete shape structure with learned features stored in local seeds. This helps generate more faithful details as it preserves regional information which is highly sparse in a global feature. Secondly, a new point generator is designed for both seed generator and the subsequent layers. Following the idea discussed before, we propose to integrate useful local information into the generation operations by aggregating neighbor-ing points in the proposed Upsample Transformer. In particular, we formulate

the generation of new points as a transformer-style self-attention weighted average of point features in the local field. This leads to a better understanding of local geometric features captured by semantic relationships. Qualitative and quantitative evaluations demonstrate the clear superiority of SeedFormer over the state-of-the-art methods on several widely-used public datasets. Our main contributions can be summarized as follows:

- We propose a novel SeedFormer for point cloud completion, greatly improving the performance of generating complete point clouds in terms of both semantic understanding and detail preservation.
- We introduce Patch Seeds as a new representation in the completion architecture to preserve regional information for recovering fine details.
- We design a new point generator, *i.e.* Upsample Transformer, by extending the transformer structure into basic operations of generating points.

## 2    Related Work

**Voxelization-based shape completion.** The early attempts on 3D shape completion [3,10,25] rely on intermediate representations of voxel grids to describe 3D objects. It is a simple and direct way to apply powerful CNN structures on various 3D applications [18,15,30]. However, this kind of methods inevitably suffers from information loss, and the computational cost increases heavily with regard to voxel resolution. Alternatively, Xie *et al.* [40] propose to use gridding operations and 3D CNNs for coarse completion, followed by refinement steps to generate detailed structures.

**Point cloud completion.** Recently, state-of-the-art deep networks are designed to directly manipulate raw point cloud data, instead of introducing an intermediate representation. The pioneering work PointNet [21] proposes to apply MLPs independently on each point and subsequently aggregate features through pooling operations to achieve permutation invariance. Following this architecture, PCN [43] is the first learning-based method for point cloud completion, which adopts a similar encoder-decoder design with a global feature representing the input shape. It generates a complete point cloud from the global feature using MLPs and folding operations [41]. Focusing on the decoding phase of generating point clouds with more faithful details, several works [27,34,11] extend the generation process into multiple steps in a hierarchical structure. This coarse-to-fine strategy helps produce dense point clouds with details recovered gradually on the missing parts. Furthermore, Wang *et al.* [31] propose a cascaded refinement module to refine predicted points by computing a displacement offset. Similar ideas are also explored in [44,37,39,35]. In order to utilize supervision not only from 3D points but also in the 2D image domain, some other works [45,39] use rendered single-view images to guide the completion task. Among the aforementioned methods, the global feature design is widely used due to its efficiency and simplicity, while it still represents intrinsic drawbacks as we discussed before. PoinTr [42] proposes a set-to-set translation strategy which shares similar ideas

with SeedFormer. However, PoinTr designs local proxies which are used for feature translation in transformer blocks. SeedFormer introduces Patch Seeds that can be propagated to the following upsample layers, focusing on the decoding process of recovering fine details from partial inputs.

**Point cloud generators.** Generating 3D points is a fundamental step for point cloud processing and can be generalized to wider research areas. Following PointNet, early generative models [1] use fully-connected layers to produce point coordinates directly from a latent representation in auto-encoders [23,14] or GANs [8]. Then, FoldingNet [41] presents a new type of generators by adding variations from a canonical 2D grid in the point deformation. It requires lower computational cost while assuming that 3D object lies on 2D-manifold. Following this design, [39] proposes a style-based folding operator by injecting shape information into point generation using a StyleGAN [12] design. Moreover, Snowflak-eNet [38] proposes to generate new points by splitting parent point features through deconvolution. More recently, with the success of transformers in natural language processing [29,36,5,4], an increasing number of research have developed such architecture for encoding 3D point clouds. In this work, we further extend the application of transformer-based structure into the basic operations of point cloud generation in the decoding phase, which also represents a new pattern of point generators with local aggregation.

## 3   Method

### 3.1   Architecture Overview

The overall architecture of SeedFormer is shown in Fig. 2. We will introduce our method in detail as follows.

**Encoder.** Denote the input point cloud as $\mathcal{P} = \{\mathbf{p}_i | i = 1, 2, ..., N\} \in \mathbb{R}^{N \times 3}$ where $N$ is the total number of points and $\mathbf{p}_i$ possesses the $(x, y, z)$ coordinates. The encoder applies point transformer [46] and set abstraction layers [22] to extract features from the incomplete shape. The number of points is reduced progressively in each layer and then we obtain patch features $\mathcal{F}_p \in \mathbb{R}^{N_p \times C_p}$ and the corresponding patch center coordinates $\mathcal{P}_p \in \mathbb{R}^{N_p \times 3}$, which represent the local structures of the partial point cloud.

**Seed generator.** Seed generator aims to predict the overall structure of the complete shape. It is designed to produce a coarse yet complete point cloud (seed points) as well as seed feature of each point which can capture regional information of local patterns. Given the extracted patch features $\mathcal{F}_p$ and center coordinates $\mathcal{P}_p$, we use Upsample Transformer (Sec. 3.3) to generate a new set of seed features:

$$\mathcal{F} = \mathrm{UpTrans}(\mathcal{F}_p, \mathcal{P}_p). \tag{1}$$

Then, we obtain the corresponding seed points $\mathcal{S} = \{x_i\}_{i=1}^{N_s} \in \mathbb{R}^{N_s \times 3}$ by applying MLPs on the generated seed features $\mathcal{F} = \{f_i\}_{i=1}^{N_s} \in \mathbb{R}^{N_s \times C_s}$. Grouping seed points and features together yields our Patch Seeds representation (more details are discussed in Sec. 3.2).
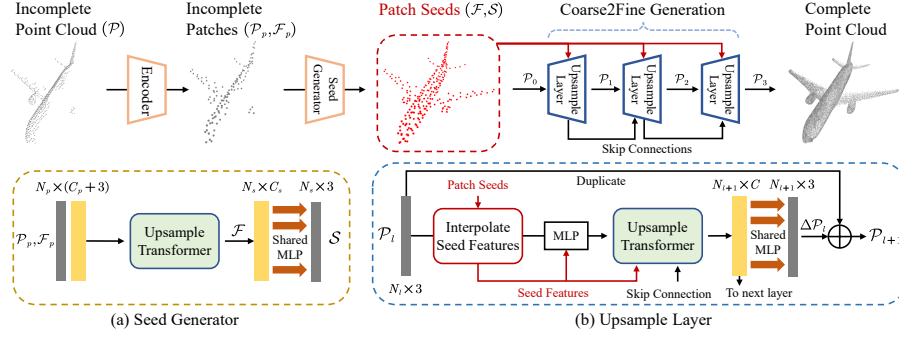
Fig. 2: The overall architecture of SeedFormer is shown in the upper part. (a) The seed generator is applied to obtain Patch Seeds which are subsequently propagated into each of the following layers by interpolating seed features. (b) Several upsample layers are used in the coarse-to-fine generation where an Upsample Transformer is applied for producing new points with skip-connection and seed feature encoding.

**Coarse-to-fine generation.** Afterwards, we follow the coarse-to-fine generation [41,11,38] to progressively recover faithful details in a hierarchical structure. This process consists of several upsample layers, each of which takes the previous point cloud and produces a dense output $\mathcal{P}_l$ ($l = 1, 2, ...$) with an upsampling rate of $r_l$ where our Patch Seeds are closely involved. Specifically, each point in the input point cloud is upsampled into $r_l$ points using the Upsample Transformer. The coarse point cloud $\mathcal{P}_0$ which is fed to the first layer is produced by fusing seeds $\mathcal{S}$ and the input point cloud $\mathcal{P}$ using Farthest Point Sampling (FPS) [22] to preserve partial structure of the original input [31].

### 3.2   Point Cloud Completion with Patch Seeds

**Patch Seeds.** Patch Seeds serve as a new shape representation in point cloud completion. It consists of both seed coordinates $\mathcal{S}$ and features $\mathcal{F}$ where each seed covers a small region around this point with seed feature containing semantic clues about this area. With rich information provided in seed features, SeedFormer can recover ambiguous missing details as shown in Fig. 1(d). Thus, compared with a global feature, the Patch Seeds representation possesses two advantages: (i) it can preserve *regional information of local patterns*, thus tiny details can be recovered in the complete point cloud; (ii) it represents the *complete shape structure* which is recovered from the partial input with explicit supervision (Sec. 3.4).

**Usage.** Throughout the following steps, Patch Seeds are incorporated into each of the upsample layers to provide regional information. Given an input point cloud $\mathcal{P}_l$, we propagate seed features to each point $p_i \in \mathcal{P}_l$ by interpolating feature values in its neighborhood $\mathcal{N}_s(i)$ ($k$ nearest seeds of $p_i$). We use weighted
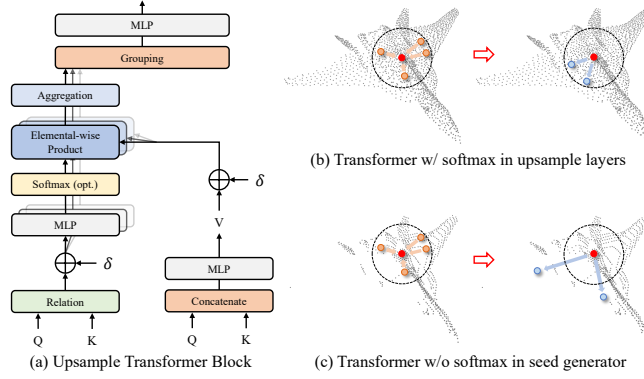
Fig. 3: (a) shows the design of our Upsample Transformer. The generation processes of new points around a target point (red) are shown: (b) an upsample layer can produce dense points within its local neighborhood; (c) the seed generator tends to generate seed points covering unseen parts of the object where the softmax function can be optionally disabled.

average based on inverse distances [22] to compute the interpolated features:

$$s_i^l = \frac{\sum_{j \in \mathcal{N}_s(i)} \hat{d}_{ij} f_j}{\sum_{j \in \mathcal{N}_s(i)} \hat{d}_{ij}} \quad \text{where} \quad \hat{d}_{ij} = \frac{1}{d_{ij}}. \tag{2}$$

Here, $d_{ij}$ denotes the distance between $p_i$ and seed point $x_j$, and $f_j$ is the corresponding seed feature. The interpolated features in this layer are $\{s_i^l\}_{i=1}^{N_l}$.

### 3.3   Upsample Transformer

**Point generator via local aggregation.** The generation process in point cloud completion aims to produce a set of new points, either maintaining data fidelity by preserving geometric details or well inferring missing parts based on the existing shape structure. The commonly-used folding operations [41,43,34,39] are designed to upsample each point independently with fixed 2D variants which leads to poor detail recovery. Differently, our Upsample Transformer is designed to incorporate closely-related local information into the generation by aggregating features from neighboring points. To this end, based on the point transformer [46] structure, we propose to formulate the generation process of new points as a self-attention weighted average of point features in the local field (Fig. 3(a)). Moreover, following the idea of local aggregation, we can extend other successful encoder designs, *e.g.*, PointNet++ [22], Transformers [29,46] or point cloud convolutions [32,28,47], into point generators in point cloud completion, to produce a new set of features with upsampled points. This is further evaluated in Sec. 4.5.

**Upsample layer.** We first explain how to apply the designed Upsample Transformer in an upsample layer (Fig. 2(b)). Given the input points $\mathcal{P}_l \in \mathbb{R}^{N_l \times 3}$ and

the corresponding interpolated seed features $\{s_i^l\}_{i=1}^{N_l}$, we concatenate them and apply a shared MLP to form the point-wise queries $\{q_i^l\}_{i=1}^{N_l}$. Following [38], we use the output features from the previous layer as keys $\{k_i^l\}_{i=1}^{N_l}$ in the transformer. It is used to preserve learned features of the existing points in the input. Then, the values $\{v_i^l\}_{i=1}^{N_l}$ are obtained by applying a MLP on the concatenated keys and queries. Upsample Transformer applies a channel-wise attention using the subtraction relation in a local neighborhood $\mathcal{N}(i)$ ($k$ nearest neighbors) of each point:

$$\hat{a}_{ijm} = \alpha_m(\beta(q_i^l) - \gamma(k_j^l) + \delta), j \in \mathcal{N}(i). \tag{3}$$

Here, $\alpha_m$, $\beta$ and $\gamma$ are feature mapping functions (*i.e.*, MLPs) to produce attention vectors. $\delta$ is a positional encoding vector to learn spatial relations. For each upsampled point relating to the centered point $p_i$, a specific kernel $\alpha_m$ is defined where $m = 1, 2, ..., r_l$ indicates one of the $r_l$ generation processes in this layer ($r_l N_l = N_{l+1}$). Each kernel learns a certain geometric pattern of local characteristics which outputs a separate group of self-attention weights to form a new point. In addition, to normalize the computed weights into a balanced scale, $\hat{a}_{ijm}$ is applied by a softmax function:

$$a_{ijm} = \frac{\exp(\hat{a}_{ijm})}{\sum_{j \in \mathcal{N}(i)} \exp(\hat{a}_{ijm})}. \tag{4}$$

Then, we compute the generated point features by combining weights with the duplicated values:

$$h_{im} = \sum_{j \in \mathcal{N}(i)} a_{ijm} * (\psi(v_j^l) + \delta), \tag{5}$$

where $\psi$ is also a feature mapping function and $*$ denotes the elemental-wise product. Grouping all $h_{im}$ from each kernel yields our produced upsampled point features $\mathcal{H}_l = \{h_{im} | i = 1, 2, ... N_l; m = 1, 2, ..., r_l\} \in \mathbb{R}^{r_l N_l \times C}$. $\mathcal{H}_l$ also serves as the keys $\{k_i^{l+1}\}_{i=1}^{N_{l+1}}$ of next layer in the skip connection. Finally, we apply shared MLPs on point features to obtain a set of point displacement offsets $\Delta\mathcal{P}_l$. The output point cloud is defined as $\mathcal{P}_{l+1} = \hat{\mathcal{P}}_l + \Delta\mathcal{P}_l$, where $\hat{\mathcal{P}}_l$ is the duplicated point cloud, for both refining existing points and generating new points.

**Positional encoding with seed features.** The basic positional encoding in the self-attention computations is designed to capture spatial relations between 3D points [46]. Besides, since the introduced seeds contain regional features with regard to the seed positions, we encode the interpolated features into *regional* encoding as follows:

$$\delta = \rho(p_i - p_j) + \theta(s_i - s_j). \tag{6}$$

Here, $\delta$ encodes both positional relation and seed feature relation between $p_i$ and $p_j$. $\rho$ and $\theta$ are encoding functions using a two-layer MLP.

**Transformer without softmax.** Similarly, Upsample Transformer is also used in the seed generator but is implemented without skip connection and seed features (Fig. 2(a)). Generating seeds involves a different objective of producing seed points which can predict missing regions and cover the complete shape structure. As shown in Fig. 3(c), when the input point cloud represents an incomplete

shape, it is essential for the seed generator to produce seed points outside the local neighborhood. However, the standard transformer structure represents an intrinsic limitation that the softmax normalization explicitly produces attention weights within a specific range of $(0, 1)$. This may limit the learning ability especially in the seed generator. To solve this, we simply remove the softmax function in the transformer; that is, $a_{ijm} = \hat{a}_{ijm}$ is applied without normalization. The experiments in Sec. 4.5 show that it is easier to generate better seed points by disabling the softmax function or using other alternatives.

### 3.4   Loss Function

We use Chamfer Distance (CD) as our loss function to measure the distance between two unordered point sets [6]. To ensure that the generated seeds $\mathcal{S}$ can cover the complete shape, we down-sample the ground-truth point cloud to the same point number as $\mathcal{S}$ and compute the corresponding CD loss between them. The same loss function is also applied to each output $\mathcal{P}_l$ separately in the upsample layers. Then, we define the sum of all CDs as the completion loss $\mathcal{L}_{comp}$. Besides, following [33], we compute the partial matching loss $\mathcal{L}_{part}$ on the final predicted result to preserve the shape structure of the input point cloud. The total training loss is defined as:

$$\mathcal{L} = \mathcal{L}_{comp} + \mathcal{L}_{part}. \tag{7}$$

## 4   Evaluation

### 4.1   Implementation Details

The SeedFormer encoder applies two layers of set abstraction [22] and obtain $N_p = 128$ patches of the incomplete point cloud. Then, in the seed generator, we produce a set of seed features $\mathcal{F} \in \mathbb{R}^{N_s \times C_s}$ where $N_s = 256$ and $C_s = 128$. The coarse point cloud $\mathcal{P}_0$ contains $N_0 = 512$ points which is obtained by merging $\mathcal{S}$ and $\mathcal{P}$ using FPS. Three upsample layers are used in the following coarse-to-fine generation procedure which output dense point clouds $\{\mathcal{P}_1, \mathcal{P}_2, \mathcal{P}_3\}$ and $\mathcal{P}_3$ corresponds to the final predicted result. Channel number of generated features is set to $C = 128$ for all upsample layers.

We train our network end-to-end using pytorch [19] implementation. We use Adam [13] optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. All the models are trained with a batch size of 48 on two NVIDIA TITAN Xp GPUs. The initial learning rate is set to 0.001 with continuous decay of 0.1 for every 100 epochs.

### 4.2   Experiments on PCN Dataset

**Data.** The PCN dataset [43] is one of the most widely used benchmark datasets for point cloud completion. It is a subset of ShapeNet [2] with shapes from 8

Table 1: Completion results on PCN dataset in terms of per-point L1 Chamfer Distance ×1000 (lower is better).

| Methods | Average | Plane | Cabinet | Car | Chair | Lamp | Couch | Table | Boat |
|---|---|---|---|---|---|---|---|---|---|
| FoldingNet [41] | 14.31 | 9.49 | 15.80 | 12.61 | 15.55 | 16.41 | 15.97 | 13.65 | 14.99 |
| TopNet [27] | 12.15 | 7.61 | 13.31 | 10.90 | 13.82 | 14.44 | 14.78 | 11.22 | 11.12 |
| AtlasNet [9] | 10.85 | 6.37 | 11.94 | 10.10 | 12.06 | 12.37 | 12.99 | 10.33 | 10.61 |
| PCN [43] | 9.64 | 5.50 | 22.70 | 10.63 | 8.70 | 11.00 | 11.34 | 11.68 | 8.59 |
| GRNet [40] | 8.83 | 6.45 | 10.37 | 9.45 | 9.41 | 7.96 | 10.51 | 8.44 | 8.04 |
| CRN [31] | 8.51 | 4.79 | 9.97 | 8.31 | 9.49 | 8.94 | 10.69 | 7.81 | 8.05 |
| NSFA [44] | 8.06 | 4.76 | 10.18 | 8.63 | 8.53 | 7.03 | 10.53 | 7.35 | 7.48 |
| PMP-Net [35] | 8.73 | 5.65 | 11.24 | 9.64 | 9.51 | 6.95 | 10.83 | 8.72 | 7.25 |
| PoinTr [42] | 8.38 | 4.75 | 10.47 | 8.68 | 9.39 | 7.75 | 10.93 | 7.78 | 7.29 |
| SnowflakeNet [38] | 7.21 | 4.29 | 9.16 | 8.08 | 7.89 | 6.07 | 9.23 | 6.55 | 6.40 |
| SeedFormer | **6.74** | **3.85** | **9.05** | **8.06** | **7.06** | **5.21** | **8.85** | **6.05** | **5.85** |



(a) Input          (b) PCN          (c) GRNet    (d) Snowflake    (e) Ours          (f) GT

Fig. 4: Visual comparisons on PCN dataset.

categories. The incomplete point clouds are generated by back-projecting 2.5D depth images from 8 viewpoints in order to simulate real-world sensor data. For each shape, 16,384 points are uniformly sampled from the mesh surfaces as complete ground-truth, and 2,048 points are sampled as partial input. We follow the same experimental setting with PCN for a fair comparison.

**Results.** Following previous methods, we report the Chamfer Distances with L1 norm (×1000) in Tab. 1. Detailed results of each category are also provided. SeedFormer achieves the best scores on all categories of this dataset, outperforming previous state-of-the-art methods by a large amount. Moreover, in Fig. 4, we show visual results of shapes from three categories (Lamp, Chair and Couch), compared with PCN [43], GRNet [40] and SnowflakeNet [38]. It shows that Seed-Former can produce clearly better results with more faithful details. As shown

Table 2: Completion results on ShapeNet-55 dataset evaluated as L2 Chamfer Distance ×1000 (lower is better) and F-Score@1% (higher is better).

| Methods | Table | Chair | Plane | Car | Sofa | CD-S | CD-M | CD-H | CD-Avg | F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| FoldingNet [41] | 2.53 | 2.81 | 1.43 | 1.98 | 2.48 | 2.67 | 2.66 | 4.05 | 3.12 | 0.082 |
| PCN [43] | 2.13 | 2.29 | 1.02 | 1.85 | 2.06 | 1.94 | 1.96 | 4.08 | 2.66 | 0.133 |
| TopNet [27] | 2.21 | 2.53 | 1.14 | 2.18 | 2.36 | 2.26 | 2.16 | 4.3 | 2.91 | 0.126 |
| PFNet [11] | 3.95 | 4.24 | 1.81 | 2.53 | 3.34 | 3.83 | 3.87 | 7.97 | 5.22 | 0.339 |
| GRNet [40] | 1.63 | 1.88 | 1.02 | 1.64 | 1.72 | 1.35 | 1.71 | 2.85 | 1.97 | 0.238 |
| PoinTr [42] | 0.81 | 0.95 | 0.44 | 0.91 | 0.79 | 0.58 | 0.88 | 1.79 | 1.09 | 0.464 |
| SeedFormer | **0.72** | **0.81** | **0.40** | **0.89** | **0.71** | **0.50** | **0.77** | **1.49** | **0.92** | **0.472** |

Table 3: Completion results on ShapeNet-34 dataset evaluated as L2 Chamfer Distance ×1000 (lower is better) and F-Score@1% (higher is better).

| Methods | 34 seen categories | | | | | 21 unseen categories | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | CD-S | CD-M | CD-H | Avg | F1 | CD-S | CD-M | CD-H | Avg | F1 |
| FoldingNet [41] | 1.86 | 1.81 | 3.38 | 2.35 | 0.139 | 2.76 | 2.74 | 5.36 | 3.62 | 0.095 |
| PCN [43] | 1.87 | 1.81 | 2.97 | 2.22 | 0.154 | 3.17 | 3.08 | 5.29 | 3.85 | 0.101 |
| TopNet [27] | 1.77 | 1.61 | 3.54 | 2.31 | 0.171 | 2.62 | 2.43 | 5.44 | 3.50 | 0.121 |
| PFNet [11] | 3.16 | 3.19 | 7.71 | 4.68 | 0.347 | 5.29 | 5.87 | 13.33 | 8.16 | 0.322 |
| GRNet [40] | 1.26 | 1.39 | 2.57 | 1.74 | 0.251 | 1.85 | 2.25 | 4.87 | 2.99 | 0.216 |
| PoinTr [42] | 0.76 | 1.05 | 1.88 | 1.23 | 0.421 | 1.04 | 1.67 | 3.44 | 2.05 | 0.384 |
| SeedFormer | **0.48** | **0.70** | **1.30** | **0.83** | **0.452** | **0.61** | **1.07** | **2.35** | **1.34** | **0.402** |

in the 2-nd and 3-rd rows, our network can preserve complicated details in the regions of chair arms and backs, without introducing undesired components.

### 4.3    Experiments on ShapeNet-55/34

**Data.** We further evaluate our model on ShapeNet-55 and ShapeNet-34 datasets from [42]. These two datasets are also generated from the synthetic ShapeNet [2] dataset while they contain more object categories and incomplete patterns. All 55 categories in ShapeNet are included in ShapeNet-55 with 41,952 shapes for training and 10,518 shapes for testing. ShapeNet-34 uses a subset of 34 categories for training and leaves 21 unseen categories for testing where 46,765 object shapes are used for training, 3,400 for testing on seen categories and 2,305 for testing on novel (unseen) categories. In both datasets, 2,048 points are sampled as input and 8,192 points as ground-truth. Following the same evaluation strategy with [42], 8 fixed viewpoints are selected and the number of points in partial point cloud is set to 2,048, 4,096 or 6,144 (25%, 50% or 75% of the complete point cloud) which corresponds to three difficulty levels of *simple*, *moderate* and *hard* in the test stage.

Table 4: Completion results on KITTI dataset evaluated as Fidelity Distance and Minimal Matching Distance (MMD). Lower is better.

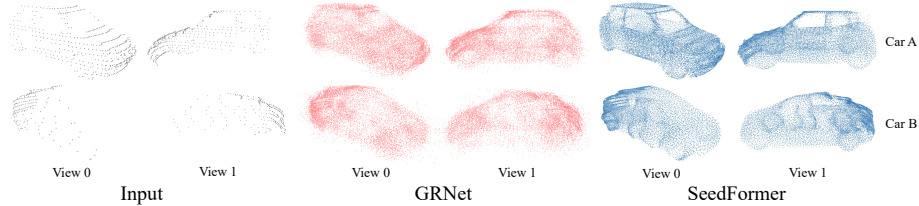|          | PCN [43] | FoldingNet [41] | TopNet [27] | MSN [17] | GRNet [40] | SeedFormer |
|----------|----------|-----------------|-------------|----------|------------|------------|
| Fidelity | 2.235    | 7.467           | 5.354       | 0.434    | 0.816      | **0.151**  |
| MMD      | 1.366    | 0.537           | 0.636       | 2.259    | 0.568      | **0.516**  |



Fig. 5: Visual comparison of point cloud completion results on KITTI dataset. For a clearer comparison, we show two different views of each object.

**Results.** The ShapeNet-55 dataset tests the ability of dealing with more diverse objects and incompleteness levels. Tab. 2 reports the average L2 Chamfer Distances ($\times 1000$) on three difficulty levels and the overall CDs. Additionally, we show results from 5 categories (Table, Chair, Plane, Car and Sofa) with more than 2,500 samples in the training set. Complete results for all 55 categories are available in the supplemental material. We also provide results under the F-Score@1% metric [26]. Compared with previous methods, SeedFormer achieves the best scores on all categories and evaluation metrics. In particular, our method outperforms the SOTA model PoinTr [42] by 15.6% in terms of overall CD and 16.8% in terms of average CD on hard difficulty.

On ShapeNet-34, the networks should handle novel objects from unseen categories which do not appear in the training phase. We show results on two test sets in Tab. 3. Our method again achieves the best scores. Especially, when dealing with unseen objects, SeedFormer shows better generalization ability achieving an average CD of 1.34 which is 34.6% lower than PoinTr.

### 4.4 Experiments on KITTI

**Data.** In order to evaluate the proposed model on real-scanned data, we further conduct experiments on the KITTI [7] dataset for completing sparse point clouds of cars in real-world environments. This dataset consists of a sequence of LiDAR scans from outdoor scenes where car objects are extracted in each frame according to the 3D bounding boxes, resulting in a total of 2,401 partial point clouds. Unlike other datasets which are built from synthetic models, the scanned data in KITTI can be highly sparse and does not have complete point cloud as ground-truth. Thus, we follow the experimental settings of GRNet [40] and evaluate our method using two metrics: (i) Fidelity Distance, which is the

Table 5: Ablation study on shape representations. We also evaluate the density of Patch Seeds.

| Methods | CD-Avg |
|---|---|
| global feature | 6.97 |
| seed number = 128 | 7.00 |
| seed number = 256 | 6.74 |
| seed number = 512 | 6.75 |

Table 6: Ablation study on different generator designs.

| Methods | CD-Avg |
|---|---|
| folding operation | 6.93 |
| deconvolution | 6.90 |
| graph convolution | 6.88 |
| point-wise attention | 6.85 |
| Upsample Transformer | 6.74 |

Table 7: Alternatives to softmax function in seed generator.

| Methods | CD-Avg |
|---|---|
| w/ softmax | 6.83 |
| w/o softmax | 6.74 |
| w/ scaled-softmax | 6.80 |
| w/ log-softmax | 6.80 |

Table 8: Ablation study on different positional encoding in Upsample Transformer.

| Methods | CD-Avg |
|---|---|
| none | 6.88 |
| positional | 6.80 |
| positional + regional | 6.74 |
| combined | 6.78 |

average distance from each point in the input to its nearest neighbour in the output. This measures how well the input is preserved; (ii) Minimal Matching Distance (MMD), which is the Chamfer Distance between the output and the car point cloud from ShapeNet that is closest to the output point cloud in terms of CD. This measures how much the output resembles a typical car.

**Results.** Following GRNet, we fine-tune our model (pretrained on PCN dataset) on ShapeNetCars (the cars from ShapeNet) for a fair comparison. Quantitative evaluation results are shown in Tab. 4 compared with previous methods. We also show some visual comparisons of the predicted point clouds in Fig. 5. Each of the car objects is visualized in two different views for a clearer comparison. We can see that our method performs clearly better on real-scanned data. Even with a very sparse input (see the 2-nd row in Fig. 5), SeedFormer can produce general structures of the desired object.

### 4.5   Ablation Studies

In this section, we demonstrate the effectiveness of several architecture designs in SeedFormer and provide some optional choices that can be applied in the network. All the networks are trained on PCN dataset with identical settings.

**Patch Seeds.** The Patch Seeds representation attempts to preserve geometric features locally, which shows clear superiority of detail recovery over previous global feature designs. We first ablate these two network architectures by replacing Patch Seeds with a global feature in the SeedFormer network. Specifically, in the seed generator, the global feature is obtained through max-pooling from

input patch features, and is used to produce a coarse point cloud by a deconvolution layer [38]. Then, in the subsequent upsample layers, this global feature is concatenated with each input point which is fed to the Upsample Transformer similarly. All other architectures are identical. Results in Tab. 5 show clear improvement of our Patch Seeds design (6.74 vs 6.97). In addition, we ablate the number of seeds in the network and show that 256 seed points are suitable for covering an input 3D object with proper density.

**Upsample Transformer and other generator designs.** We compare Upsample Transformer with other point cloud generators. Among previous methods for point cloud completion, folding-based operation [41] and deconvolution [38] are commonly-used. Unlike our generator which considers semantic relationships between points, these two designs are applied to process each point independently without using contextual information of point clouds.

The information from local areas is vital for point cloud completion, and we provide two optional choices for point generation operators. As discussed in Sec. 3.3, we adopt graph convolution [32] and vanilla Transformer (point-wise attention) to further demonstrate the effectiveness of local aggregation. In this experiment, we replace the Upsample Transformers in both seed generator and upsample layers, and keep other architectures the same. All the generators are designed to produce new features with upsampled points. Tab. 6 shows point generation by local aggregation performs better than previous methods, and the performance of Upsample Transformer stands out in similar designs.

**Alternatives to softmax function.** The softmax function can provide balanced weights in the self-attention mechanism. However, in the case of seed generators, it also presents intrinsic limitations on generating new points. In Tab. 7, we show that the standard transformer structure (w/ softmax) is not the best choice (6.83) in the seed generator. Our improvement (w/o softmax) by applying self-attention without softmax function aims to release the points from limited weights within a specific range of $(0,1)$. Similar results can be achieved by using a log-softmax or scaled-softmax (multiplied by a scale parameter $\lambda$) function, which are also alternatives to the original softmax function.

**Positional encoding.** Upsample Transformer applies both positional encoding from spatial relations and regional encoding from seed feature relations. This extends the original positional encoding in transformers by utilizing information from interpolated seed features. Tab. 8 shows that this design performs better which also demonstrates the effectiveness of incorporating information from Patch Seeds into the generation process of Upsample Transformer. Moreover, we design another ablation which combines both features in one encoding (6.78). We concatenate the two inputs and obtain a joint version through MLPs.

## 5   Visualization of Patch Seeds

In this section, we give some insights to achieve a deeper understanding of the generation process of Patch Seeds. The seed generator takes $N_p = 128$ patches (Fig. 6(b)) extracted from the input point cloud and produces a new set of

(a) Input          (b) Patch          (c) Patch Seeds    (d) Prediction          (e) GT
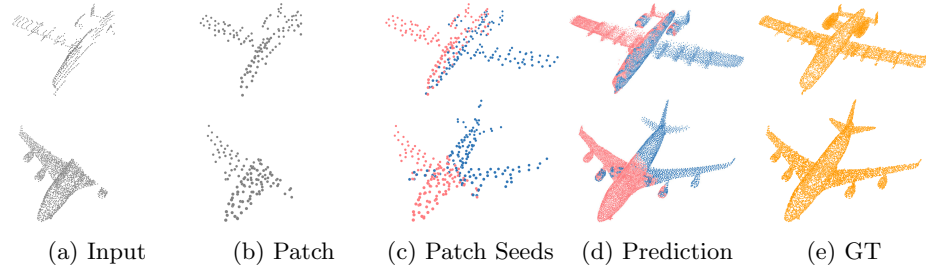
Fig. 6: Illustration of Patch Seeds and other intermediate results. (a) Input point cloud. (b) Extracted patch centers from partial input. (c) Generated seeds (red: seed0, blue: seed1) according to the seed generator process. (d) Our predicted results with each point colored according its nearest seed. (e) Ground truth.

$N_s = 256$ seed points (Fig. 6(c)). Instead of using a global feature, adopting a patch-to-patch translation allows us to track the paths in the seed generation process. Specifically, the generated seeds can be divided into two groups, denoted as seed0 and seed1, since each patch center is split into two seeds according to different groups of self-attention weights in Upsample Transformer (see in Eq. 3). Thus, we visualize the obtained Patch Seeds in Fig. 6(c) with specific colors (red: seed0, blue: seed1). We can see that the red points (seed0) are close to the partial patches where the network tries to preserve the input structures. As for another group (blue points), if the input point cloud is symmetric to the ground-truth (the 1-st row in Fig. 6(a)), SeedFormer learns to duplicate the existing points as well as features, thus it can recover the missiles which are inferred from the seen parts. If the input is asymmetric (in the 2-nd row), our network can also predict missing parts according to the similar observations in the seen point cloud. For a clearer illustration, we also visualize the predicted complete point clouds in Fig. 6(d) where each point is colored according to its nearest seed.

## 6   Conclusion

In this paper, we propose a novel point cloud completion network, termed Seed-Former. As opposed to previous methods, SeedFormer introduce a new shape representation, namely Patch Seeds, in point cloud completion. The idea of Patch Seeds is to capture both global shape structure and fine-grained local details by learning regional features which are stored in several local seeds. This leads to a better performance of shape recovery and detail preservation in the generated point clouds. Furthermore, by extending the transformer structure into point generation, we propose a novel Upsample Transformer to capture useful neighborhood information. Comprehensive experiments show that our method achieves clear improvements on several challenging benchmark dataset compared to state-of-the-art competitors.

# References

1. Achlioptas, P., Diamanti, O., Mitliagkas, I., Guibas, L.: Learning representations and generative models for 3d point clouds. In: International conference on machine learning. pp. 40–49. PMLR (2018)
2. Chang, A.X., Funkhouser, T., Guibas, L., Hanrahan, P., Huang, Q., Li, Z., Savarese, S., Savva, M., Song, S., Su, H., et al.: Shapenet: An information-rich 3d model repository. arXiv preprint arXiv:1512.03012 (2015)
3. Dai, A., Ruizhongtai Qi, C., Nießner, M.: Shape completion using 3d-encoder-predictor cnns and shape synthesis. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 5868–5877 (2017)
4. Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q.V., Salakhutdinov, R.: Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint arXiv:1901.02860 (2019)
5. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)
6. Fan, H., Su, H., Guibas, L.J.: A point set generation network for 3d object reconstruction from a single image. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 605–613 (2017)
7. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The kitti dataset. The International Journal of Robotics Research $32$(11), 1231–1237 (2013)
8. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. Advances in neural information processing systems $27$ (2014)
9. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: A papier-mâché approach to learning 3d surface generation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 216–224 (2018)
10. Han, X., Li, Z., Huang, H., Kalogerakis, E., Yu, Y.: 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. pp. 85–93 (2017)
11. Huang, Z., Yu, Y., Xu, J., Ni, F., Le, X.: Pf-net: Point fractal network for 3d point cloud completion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7662–7670 (2020)
12. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 4401–4410 (2019)
13. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
14. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 (2013)
15. Le, T., Duan, Y.: Pointgrid: A deep network for 3d shape understanding. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 9204–9214 (2018)
16. Liang, M., Yang, B., Wang, S., Urtasun, R.: Deep continuous fusion for multi-sensor 3d object detection. In: Proceedings of the European conference on computer vision (ECCV). pp. 641–656 (2018)
17. Liu, M., Sheng, L., Yang, S., Shao, J., Hu, S.M.: Morphing and sampling network for dense point cloud completion. In: Proceedings of the AAAI conference on artificial intelligence. vol. 34, pp. 11596–11603 (2020)

18. Maturana, D., Scherer, S.: Voxnet: A 3d convolutional neural network for real-time object recognition. In: 2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS). pp. 922–928. IEEE (2015)

19. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems **32** (2019)

20. Qi, C.R., Liu, W., Wu, C., Su, H., Guibas, L.J.: Frustum pointnets for 3d object detection from rgb-d data. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 918–927 (2018)

21. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 652–660 (2017)

22. Qi, C.R., Yi, L., Su, H., Guibas, L.J.: Pointnet++: Deep hierarchical feature learning on point sets in a metric space. Advances in neural information processing systems **30** (2017)

23. Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning representations by back-propagating errors. nature **323**(6088), 533–536 (1986)

24. Rusu, R.B., Marton, Z.C., Blodow, N., Dolha, M., Beetz, M.: Towards 3d point cloud based object maps for household environments. Robotics and Autonomous Systems **56**(11), 927–941 (2008)

25. Stutz, D., Geiger, A.: Learning 3d shape completion from laser scan data with weak supervision. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1955–1964 (2018)

26. Tatarchenko, M., Richter, S.R., Ranftl, R., Li, Z., Koltun, V., Brox, T.: What do single-view 3d reconstruction networks learn? In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3405–3414 (2019)

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

28. Thomas, H., Qi, C.R., Deschaud, J.E., Marcotegui, B., Goulette, F., Guibas, L.J.: Kpconv: Flexible and deformable convolution for point clouds. In: Proceedings of the IEEE/CVF international conference on computer vision. pp. 6411–6420 (2019)

29. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. Advances in neural information processing systems **30** (2017)

30. Wang, P.S., Liu, Y., Guo, Y.X., Sun, C.Y., Tong, X.: O-cnn: Octree-based convolutional neural networks for 3d shape analysis. ACM Transactions On Graphics (TOG) **36**(4), 1–11 (2017)

31. Wang, X., Ang Jr, M.H., Lee, G.H.: Cascaded refinement network for point cloud completion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 790–799 (2020)

32. Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. Acm Transactions On Graphics (tog) **38**(5), 1–12 (2019)

33. Wen, X., Han, Z., Cao, Y.P., Wan, P., Zheng, W., Liu, Y.S.: Cycle4completion: Unpaired point cloud completion using cycle transformation with missing region coding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13080–13089 (2021)

34. Wen, X., Li, T., Han, Z., Liu, Y.S.: Point cloud completion by skip-attention network with hierarchical folding. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 1939–1948 (2020)
35. Wen, X., Xiang, P., Han, Z., Cao, Y.P., Wan, P., Zheng, W., Liu, Y.S.: Pmp-net: Point cloud completion by learning multi-step point moving paths. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7443–7452 (2021)
36. Wu, F., Fan, A., Baevski, A., Dauphin, Y.N., Auli, M.: Pay less attention with lightweight and dynamic convolutions. arXiv preprint arXiv:1901.10430 (2019)
37. Xia, Y., Xia, Y., Li, W., Song, R., Cao, K., Stilla, U.: Asfm-net: Asymmetrical siamese feature matching network for point completion. In: Proceedings of the 29th ACM International Conference on Multimedia. pp. 1938–1947 (2021)
38. Xiang, P., Wen, X., Liu, Y.S., Cao, Y.P., Wan, P., Zheng, W., Han, Z.: Snowflakenet: Point cloud completion by snowflake point deconvolution with skip-transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 5499–5509 (2021)
39. Xie, C., Wang, C., Zhang, B., Yang, H., Chen, D., Wen, F.: Style-based point generator with adversarial rendering for point cloud completion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 4619–4628 (2021)
40. Xie, H., Yao, H., Zhou, S., Mao, J., Zhang, S., Sun, W.: Grnet: Gridding residual network for dense point cloud completion. In: European Conference on Computer Vision. pp. 365–381. Springer (2020)
41. Yang, Y., Feng, C., Shen, Y., Tian, D.: Foldingnet: Point cloud auto-encoder via deep grid deformation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 206–215 (2018)
42. Yu, X., Rao, Y., Wang, Z., Liu, Z., Lu, J., Zhou, J.: Pointr: Diverse point cloud completion with geometry-aware transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12498–12507 (2021)
43. Yuan, W., Khot, T., Held, D., Mertz, C., Hebert, M.: Pcn: Point completion network. In: 2018 International Conference on 3D Vision (3DV). pp. 728–737. IEEE (2018)
44. Zhang, W., Yan, Q., Xiao, C.: Detail preserved point cloud completion via separated feature aggregation. In: European Conference on Computer Vision. pp. 512–528. Springer (2020)
45. Zhang, X., Feng, Y., Li, S., Zou, C., Wan, H., Zhao, X., Guo, Y., Gao, Y.: View-guided point cloud completion. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 15890–15899 (2021)
46. Zhao, H., Jiang, L., Jia, J., Torr, P.H., Koltun, V.: Point transformer. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 16259–16268 (2021)
47. Zhou, H., Feng, Y., Fang, M., Wei, M., Qin, J., Lu, T.: Adaptive graph convolution for point cloud analysis. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4965–4974 (2021)

In this supplementary material, we provide additional information to complement the manuscript. First, we present additional implementation details and experimental settings of SeedFormer (Sec. A). Second, we provide more details about the ablation studies (variants of SeedFormer) in Sec. B. Third, we show complexity analysis of our method compared with previous state-of-the-art methods (Sec. C). At last, we present more qualitative results compared with other methods (Sec. D).

## A    Implementation Details

In this section, we provide additional implementation details of the proposed SeedFormer.

**Encoder.** We use point transformer [46] ($SA$) and set abstraction [22] ($PT$) layers to extract features from an input point cloud. The point cloud (with 2,048 points) is downsampled using FPS in each layer of set abstraction. Detailed network architecture is as follows: $SA(C = 128, N = 512) \rightarrow PT(C = 128) \rightarrow SA(C = 256, N = N_p) \rightarrow PT(C = C_p)$. The outputs are $N_p = 128$ patch centers $\mathcal{P}_p$ as well as the corresponding patch features $\mathcal{F}_p$ with $C_p = 256$ channels.

**Seed Generator.** The inputs to seed generator are the obtained patch centers $\mathcal{P}_p$ and correponding patch features $\mathcal{F}_p$. The implementation of Upsample Transformer (Eq. 1 in the main paper) is slightly different from that in an upsample layer. Since there are no skipped features, we use $\mathcal{F}_p$, applied by linear projections, for both keys and queries in the transformer. Also, the positional encoding of each point is calculated without incorporating seed features. Then, after we obtain seed features $\mathcal{F}$, we concatenate them with a pooled feature from $\mathcal{F}_p$ which is used for encoding global contexts, and apply shared MLPs to produce the seed coordinates $\mathcal{S}$.

**Experimental settings.** For PCN dataset (with 16,384 points in ground-truth), we set the upsampling rate as $r_1 = 1, r_2 = 4, r_3 = 8$ where the first upsample layer is used to refine the input coarse point cloud $\mathcal{P}_0$. For ShapeNet-34/55 (8,192 points), we set $r_1 = 1, r_2 = 4, r_3 = 4$. Each input point cloud from all datasets contains 2,048 points. For those point clouds with less than 2,048 points, we randomly duplicate input points; for those with more than 2,048 points, we pick a subset as input.

## B    Ablation Studies

We give more details about ablation networks used in Sec. 4.5 in the main paper. Following the same idea of point generation by local aggregation, we propose two optional generator designs which are both more effective than previous methods. Compared with the default Upsample Transformer, these options can also achieve decent results (Tab. 6 in the main paper) while being more efficient.

Table 9: Complexity analysis on PCN dataset evaluated as the number of parameters (Params) and theoretical computational cost (FLOPs). We also report the average CDs of all categories as references.

| Methods | Params | FLOPs | CD-Avg |
|---|---|---|---|
| FoldingNet [41] | **2.41M** | 27.65G | 14.31 |
| PCN [43] | 6.84M | 14.69G | 9.64 |
| GRNet [40] | 76.71M | 25.88G | 8.83 |
| SnowflakeNet [38] | 19.32M | 10.32G | 7.21 |
| point-wise attention | 3.12M | **7.87G** | 6.85 |
| SeedFormer | 3.20M | 29.61G | **6.74** |

Similarly, given the features $\{f_i^l\}_{i=1}^{N_l}$ of the input point cloud $\mathcal{P}_l$, we first apply graph convolutions [32] to generate new point features by aggregating local neighborhood information:

$$h_{im} = \max_{j \in \mathcal{N}(i)} \alpha_m(f_j^l). \tag{8}$$

Here, $\alpha_m$ is a feature mapping function (MLPs) and $m = 1, 2, ..., r_l$ corresponds to each of the $r_l$ generation processes in this layer. The output features are $\mathcal{H}_l = \{h_{im} | i = 1, 2, ..., N_l; m = 1, 2, ..., r_l\} \in \mathbb{R}^{r_l N_l \times C}$ with a upsampling rate of $r_l$. Then, we obtain the new point cloud by learning displacement offsets from the produced features using shared MLPs.

Another option applies transformer structures in point generation. Different from an Upsample Transformer, we adopt point-wise attention which is more efficient in the computations. Similar to Eq. 3 in the main paper, we compute self-attention weights by:

$$\hat{a}_{ijm} = \alpha_m(\beta(q_i^l) - \gamma(k_j^l) + \delta), j \in \mathcal{N}(i). \tag{9}$$

The difference is that $\alpha_m : \mathbb{R}^C \to \mathbb{R}$ outputs one-dimensional weights which are assigned to different points in the neighborhood. Then, we obtain the generated point features by combining weights with duplicated values (Eq. 5 in the main paper). To sum up, the point-wise attention is faster (Sec. C) while it can still outperform state-of-the-art methods.

## C  Complexity Analysis

We provide a detailed complexity analysis of our method. Tab. 9 reports the theoretical computational cost (FLOPs) and number of parameters for different models. The scores are measured on PCN dataset (16,384 points) with a batch size of 1. We also provide the overall Chamfer Distances as references. We can see that the model size of SeedFormer is very favorable compared to previous methods. This is owing to our designs of seed generator and upsample layers

Table 10: Detailed results for novel 21 categories on ShapeNet-34 dataset. *S.*, *M.* and *H.* stand for the simple, moderate and hard difficulty levels.

| | PCN [43] | | | GRNet [40] | | | PoinTr [42] | | | SeedFormer | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | S. | M. | H. | S. | M. | H. | S. | M. | H. | S. | M. | H. |
| bag | 2.48 | 2.46 | 3.94 | 1.47 | 1.88 | 3.45 | 0.96 | 1.34 | 2.08 | **0.49** | **0.82** | **1.45** |
| basket | 2.79 | 2.51 | 4.78 | 1.78 | 1.94 | 4.18 | 1.04 | 1.40 | 2.90 | **0.60** | **0.85** | **1.98** |
| birdhouse | 3.53 | 3.47 | 5.31 | 1.89 | 2.34 | 5.16 | 1.22 | 1.79 | 3.45 | **0.72** | **1.19** | **2.31** |
| bowl | 2.66 | 2.35 | 3.97 | 1.77 | 1.97 | 3.90 | 1.05 | 1.32 | 2.40 | **0.60** | **0.77** | **1.50** |
| camera | 4.84 | 5.30 | 8.03 | 2.31 | 3.38 | 7.20 | 1.63 | 2.67 | 4.97 | **0.89** | **1.77** | **3.75** |
| can | 1.95 | 1.89 | 5.21 | 1.53 | 1.80 | 3.08 | 0.80 | 1.17 | 2.85 | **0.56** | **0.89** | **1.57** |
| cap | 7.21 | 7.14 | 10.94 | 3.29 | 4.87 | 13.02 | 1.40 | 2.74 | 8.35 | **0.50** | **1.34** | **5.19** |
| keyboard | 1.07 | 1.00 | 1.23 | 0.73 | 0.77 | 1.11 | 0.43 | 0.45 | 0.63 | **0.32** | **0.41** | **0.60** |
| dishwasher | 2.45 | 2.09 | 3.53 | 1.79 | 1.70 | 3.27 | 0.93 | 1.05 | 2.04 | **0.63** | **0.78** | **1.44** |
| earphone | 7.88 | 6.59 | 16.53 | 4.29 | 4.16 | 10.30 | 2.03 | 5.10 | 10.69 | **1.18** | **2.78** | **6.71** |
| helmet | 6.15 | 6.41 | 9.16 | 3.06 | 4.38 | 10.27 | 1.86 | 3.30 | 6.96 | **1.10** | **2.27** | **4.78** |
| mailbox | 2.74 | 2.68 | 4.31 | 1.52 | 1.90 | 4.33 | 1.03 | 1.47 | 3.34 | **0.56** | **0.99** | **2.06** |
| microphone | 4.36 | 4.65 | 8.46 | 2.29 | 3.23 | 8.41 | 1.25 | 2.27 | 5.47 | **0.80** | **1.61** | **4.21** |
| microwaves | 2.59 | 2.35 | 4.47 | 1.74 | 1.81 | 3.82 | 1.01 | 1.18 | 2.14 | **0.64** | **0.83** | **1.69** |
| pillow | 2.09 | 2.16 | 3.54 | 1.43 | 1.69 | 3.43 | 0.92 | 1.24 | 2.39 | **0.43** | **0.66** | **1.45** |
| printer | 3.28 | 3.60 | 5.56 | 1.82 | 2.41 | 5.09 | 1.18 | 1.76 | 3.10 | **0.69** | **1.25** | **2.33** |
| remote | 0.95 | 1.08 | 1.58 | 0.82 | 1.02 | 1.29 | 0.44 | 0.58 | 0.78 | **0.27** | **0.42** | **0.61** |
| rocket | 1.39 | 1.22 | 2.01 | 0.97 | 0.79 | 1.60 | 0.39 | 0.72 | 1.39 | **0.28** | **0.51** | **1.02** |
| skateboard | 1.97 | 1.78 | 2.45 | 0.93 | 1.07 | 1.83 | 0.52 | 0.80 | 1.31 | **0.35** | **0.56** | **0.92** |
| tower | 2.37 | 2.40 | 4.35 | 1.35 | 1.80 | 3.85 | 0.82 | 1.35 | 2.48 | **0.51** | **0.92** | **1.87** |
| washer | 2.77 | 2.52 | 4.64 | 1.83 | 1.97 | 5.28 | 1.04 | 1.39 | 2.73 | **0.61** | **0.87** | **1.94** |
| mean | 3.22 | 3.13 | 5.43 | 1.84 | 2.23 | 4.95 | 1.05 | 1.67 | 3.45 | **0.61** | **1.07** | **2.35** |

which process each point with a shared generator. On the other hand, since Upsample Transformer is required to aggregate local points in each generation, the computational cost of SeedFormer is relatively high but it is comparable with GRNet [40] and FoldingNet [41]. We also provide a more efficient version of SeedFormer using a faster transformer structure (point-wise attention) as discussed in Sec. 4.5 of the main manuscript. This model achieves a decent result (6.85) with lower computational cost. In all, our model can offer a pleasant trade-off between cost and performance.

## D    Additional Experimental Results

**Qualitative results on PCN dataset.** In Fig. 7, we provide more visual results compared with PCN [43], GRNet [40] and SnowflakeNet [38]. All the results are obtained from their released pretrained models. We can see that SeedFormer performs clearly better than previous methods.

**More results on ShapeNet-55/34.** We report complete results of our method on ShapeNet-55 in Tab. 11 and results of novel categories on ShapeNet-34 in Tab. 10. The models are tested under three difficulty levels: simple, moderate and hard. For novel objects of ShapeNet-34, we can see that SeedFormer achieves best scores on all categories.

Table 11: Detailed results on ShapeNet-55 dataset. *S.*, *M.* and *H.* stand for the simple, moderate and hard difficulty levels.

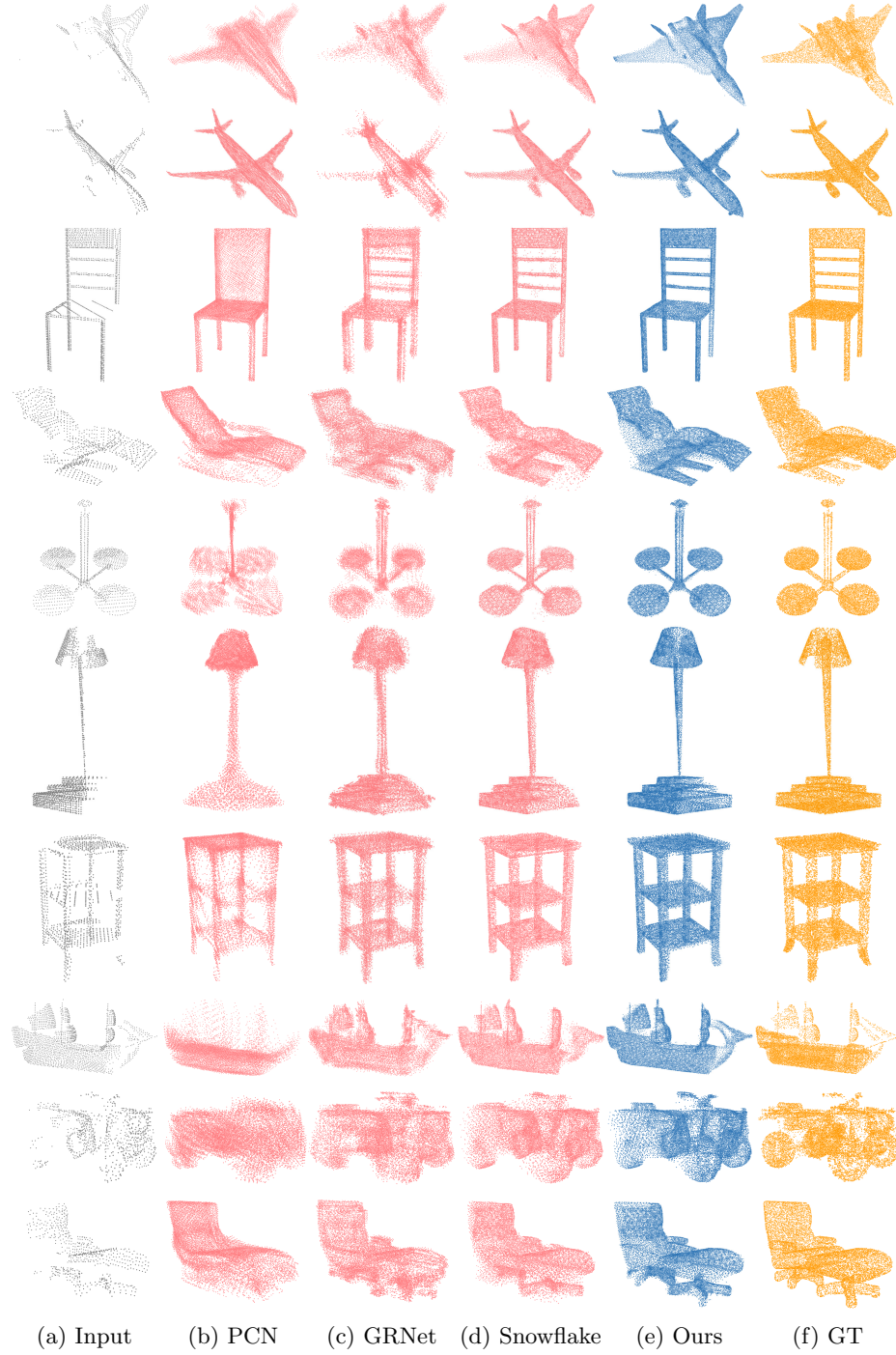| | PCN [43] | | | GRNet [40] | | | PoinTr [42] | | | SeedFormer | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | S. | M. | H. | S. | M. | H. | S. | M. | H. | S. | M. | H. |
| airplane | 0.90 | 0.89 | 1.32 | 0.87 | 0.87 | 1.27 | 0.27 | 0.38 | 0.69 | **0.23** | **0.35** | **0.61** |
| trash-bin | 2.16 | 2.18 | 5.15 | 1.69 | 2.01 | 3.48 | 0.80 | 1.15 | 2.15 | **0.73** | **1.08** | **1.94** |
| bag | 2.11 | 2.04 | 4.44 | 1.41 | 1.70 | 2.97 | 0.53 | 0.74 | 1.51 | **0.43** | **0.67** | **1.28** |
| basket | 2.21 | 2.10 | 4.55 | 1.65 | 1.84 | 3.15 | 0.73 | 0.88 | 1.82 | **0.65** | **0.83** | **1.54** |
| bathtub | 2.11 | 2.09 | 3.94 | 1.46 | 1.73 | 2.73 | 0.64 | 0.94 | 1.68 | **0.52** | **0.82** | **1.45** |
| bed | 2.86 | 3.07 | 5.54 | 1.64 | 2.03 | 3.70 | 0.76 | 1.10 | 2.26 | **0.63** | **0.91** | **1.89** |
| bench | 1.31 | 1.24 | 2.14 | 1.03 | 1.09 | 1.71 | 0.38 | 0.52 | 0.94 | **0.32** | **0.42** | **0.84** |
| birdhouse | 3.29 | 3.53 | 6.69 | 1.87 | 2.40 | 4.71 | 0.98 | 1.49 | 3.13 | **0.76** | **1.30** | **2.46** |
| bookshelf | 2.70 | 2.70 | 4.61 | 1.42 | 1.71 | 2.78 | 0.71 | 1.06 | 1.93 | **0.57** | **0.84** | **1.57** |
| bottle | 1.25 | 1.43 | 4.61 | 1.05 | 1.44 | 2.67 | 0.37 | 0.74 | 1.50 | **0.31** | **0.63** | **1.21** |
| bowl | 2.05 | 1.83 | 3.66 | 1.60 | 1.77 | 2.99 | 0.68 | 0.78 | 1.44 | **0.56** | **0.65** | **1.18** |
| bus | 1.20 | 1.14 | 2.08 | 1.06 | 1.16 | 1.48 | 0.42 | 0.55 | 0.79 | **0.42** | **0.55** | **0.73** |
| cabinet | 1.60 | 1.49 | 3.47 | 1.27 | 1.41 | 2.09 | **0.55** | **0.66** | 1.16 | 0.57 | 0.69 | **1.05** |
| camera | 4.05 | 4.54 | 8.27 | 2.14 | 3.15 | 6.09 | 1.10 | 2.03 | 4.34 | **0.83** | **1.68** | **3.45** |
| can | 2.02 | 2.28 | 6.48 | 1.58 | 2.11 | 3.81 | 0.68 | 1.19 | 2.14 | **0.58** | **1.03** | **1.79** |
| cap | 1.82 | 1.76 | 4.20 | 1.17 | 1.37 | 3.05 | 0.46 | 0.62 | 1.64 | **0.33** | **0.45** | **1.18** |
| car | 1.48 | 1.47 | 2.60 | 1.29 | 1.48 | 2.14 | **0.64** | 0.86 | 1.25 | 0.65 | **0.86** | **1.17** |
| cellphone | 0.80 | 0.79 | 1.71 | 0.82 | 0.91 | 1.18 | 0.32 | **0.39** | 0.60 | **0.31** | 0.40 | **0.54** |
| chair | 1.70 | 1.81 | 3.34 | 1.24 | 1.56 | 2.73 | 0.49 | 0.74 | 1.63 | **0.41** | **0.65** | **1.38** |
| clock | 2.10 | 2.01 | 3.98 | 1.46 | 1.66 | 2.67 | 0.62 | 0.84 | 1.65 | **0.53** | **0.74** | **1.35** |
| keyboard | 0.82 | 0.82 | 1.04 | 0.74 | 0.81 | 1.09 | 0.30 | 0.39 | **0.45** | **0.28** | **0.36** | **0.45** |
| dishwasher | 1.93 | 1.66 | 4.39 | 1.43 | 1.59 | 2.53 | **0.55** | 0.69 | 1.42 | 0.56 | **0.69** | **1.30** |
| display | 1.56 | 1.66 | 3.26 | 1.13 | 1.38 | 2.29 | 0.48 | 0.67 | 1.33 | **0.39** | **0.59** | **1.10** |
| earphone | 3.13 | 2.94 | 7.56 | 1.78 | 2.18 | 5.33 | 0.81 | 1.38 | 3.78 | **0.64** | **1.04** | **2.75** |
| faucet | 3.21 | 3.48 | 7.52 | 1.81 | 2.32 | 4.91 | 0.71 | 1.42 | 3.49 | **0.55** | **1.15** | **2.63** |
| filecabinet | 2.02 | 1.97 | 4.14 | 1.46 | 1.71 | 2.89 | **0.63** | **0.84** | 1.69 | **0.63** | **0.84** | **1.49** |
| guitar | 0.42 | 0.38 | 1.23 | 0.44 | 0.48 | 0.76 | 0.14 | 0.21 | 0.42 | **0.13** | **0.19** | **0.32** |
| helmet | 3.76 | 4.18 | 7.53 | 2.33 | 3.18 | 6.03 | 0.99 | 1.93 | 4.22 | **0.79** | **1.52** | **3.61** |
| jar | 2.57 | 2.82 | 6.00 | 1.72 | 2.37 | 4.37 | 0.77 | 1.33 | 2.87 | **0.63** | **1.13** | **2.36** |
| knife | 0.94 | 0.62 | 1.37 | 0.72 | 0.66 | 0.96 | 0.20 | 0.33 | 0.56 | **0.15** | **0.28** | **0.45** |
| lamp | 3.10 | 3.45 | 7.02 | 1.68 | 2.43 | 5.17 | 0.64 | 1.40 | 3.58 | **0.45** | **1.06** | **2.67** |
| laptop | 0.75 | 0.79 | 1.59 | 0.83 | 0.87 | 1.28 | **0.32** | **0.34** | 0.60 | **0.32** | 0.37 | **0.55** |
| loudspeaker | 2.50 | 2.45 | 5.08 | 1.75 | 2.08 | 3.45 | 0.78 | 1.16 | 2.17 | **0.67** | **1.01** | **1.80** |
| mailbox | 1.66 | 1.74 | 5.18 | 1.15 | 1.59 | 3.42 | 0.39 | 0.78 | 2.56 | **0.30** | **0.67** | **2.04** |
| microphone | 3.44 | 3.90 | 8.52 | 2.09 | 2.76 | 5.70 | 0.70 | 1.66 | 4.48 | **0.62** | **1.61** | **3.66** |
| microwaves | 2.20 | 2.01 | 4.65 | 1.51 | 1.72 | 2.76 | 0.67 | 0.83 | 1.82 | **0.63** | **0.79** | **1.47** |
| motorbike | 2.03 | 2.01 | 3.13 | 1.38 | 1.52 | 2.26 | 0.75 | 1.10 | 1.92 | **0.68** | **0.96** | **1.44** |
| mug | 2.45 | 2.48 | 5.17 | 1.75 | 2.16 | 3.79 | 0.91 | 1.17 | 2.35 | **0.79** | **1.03** | **2.06** |
| piano | 2.64 | 2.74 | 4.83 | 1.53 | 1.82 | 3.21 | 0.76 | 1.06 | 2.23 | **0.62** | **0.87** | **1.79** |
| pillow | 1.85 | 1.81 | 3.68 | 1.42 | 1.67 | 3.04 | 0.61 | 0.82 | 1.56 | **0.48** | **0.75** | **1.41** |
| pistol | 1.25 | 1.17 | 2.65 | 1.11 | 1.06 | 1.76 | 0.43 | 0.66 | 1.30 | **0.37** | **0.56** | **0.96** |
| flowerpot | 3.32 | 3.39 | 6.04 | 2.02 | 2.48 | 4.19 | 1.01 | 1.51 | 2.77 | **0.93** | **1.30** | **2.32** |
| printer | 2.90 | 3.19 | 5.84 | 1.56 | 2.38 | 4.24 | 0.73 | 1.21 | 2.47 | **0.58** | **1.11** | **2.13** |
| remote | 0.99 | 0.97 | 2.04 | 0.89 | 1.05 | 1.29 | 0.36 | 0.53 | 0.71 | **0.29** | **0.46** | **0.62** |
| rifle | 0.98 | 0.80 | 1.31 | 0.83 | 0.77 | 1.16 | 0.30 | 0.45 | 0.79 | **0.27** | **0.41** | **0.66** |
| rocket | 1.05 | 1.04 | 1.87 | 0.78 | 0.92 | 1.44 | 0.23 | 0.48 | 0.99 | **0.21** | **0.46** | **0.83** |
| skateboard | 1.04 | 0.94 | 1.68 | 0.82 | 0.87 | 1.24 | 0.28 | 0.38 | 0.62 | **0.23** | **0.32** | **0.62** |
| sofa | 1.65 | 1.61 | 2.92 | 1.35 | 1.45 | 2.32 | 0.56 | 0.67 | 1.14 | **0.50** | **0.62** | **1.02** |
| stove | 2.07 | 2.02 | 4.72 | 1.46 | 1.72 | 3.22 | 0.63 | 0.92 | 1.73 | **0.59** | **0.87** | **1.49** |
| table | 1.56 | 1.50 | 3.36 | 1.15 | 1.33 | 2.33 | 0.46 | 0.64 | 1.31 | **0.41** | **0.58** | **1.18** |
| telephone | 0.80 | 0.80 | 1.67 | 0.81 | 0.89 | 1.18 | **0.31** | **0.38** | 0.59 | **0.31** | 0.39 | **0.55** |
| tower | 1.91 | 1.97 | 4.47 | 1.26 | 1.69 | 3.06 | 0.55 | 0.90 | 1.95 | **0.47** | **0.84** | **1.65** |
| train | 1.50 | 1.41 | 2.37 | 1.09 | 1.14 | 1.61 | **0.50** | 0.70 | 1.12 | 0.51 | **0.66** | **1.01** |
| watercraft | 1.46 | 1.39 | 2.40 | 1.09 | 1.12 | 1.65 | 0.41 | 0.62 | 1.07 | **0.35** | **0.56** | **0.92** |
| washer | 2.42 | 2.31 | 6.08 | 1.72 | 2.05 | 4.19 | 0.75 | 1.06 | 2.44 | **0.64** | **0.91** | **2.04** |
| mean | 1.96 | 1.98 | 4.09 | 1.35 | 1.63 | 2.86 | 0.58 | 0.88 | 1.80 | **0.50** | **0.77** | **1.49** |

(a) Input     (b) PCN     (c) GRNet     (d) Snowflake     (e) Ours     (f) GT

Fig. 7: Visual comparisons on PCN dataset.