

PoinTr: Diverse Point Cloud Completion with Geometry-Aware Transformers

Xumin Yu*, Yongming Rao*, Ziyi Wang, Zuyan Liu, Jiwen Lu†, Jie Zhou

Department of Automation, Tsinghua University, China

State Key Lab of Intelligent Technologies and Systems, China

Beijing National Research Center for Information Science and Technology, China

yuxm20@mails.tsinghua.edu.cn; raoyongming95@gmail.com;

{wziyi20, liuzuyan19}@mails.tsinghua.edu.cn; {lujiwen, jzhou}@tsinghua.edu.cn

Abstract

Point clouds captured in real-world applications are often incomplete due to the limited sensor resolution, single viewpoint, and occlusion. Therefore, recovering the complete point clouds from partial ones becomes an indispensable task in many practical applications. In this paper, we present a new method that reformulates point cloud completion as a set-to-set translation problem and design a new model, called PoinTr that adopts a transformer encoder-decoder architecture for point cloud completion. By representing the point cloud as a set of unordered groups of points with position embeddings, we convert the point cloud to a sequence of point proxies and employ the transformers for point cloud generation. To facilitate transformers to better leverage the inductive bias about 3D geometric structures of point clouds, we further devise a geometry-aware block that models the local geometric relationships explicitly. The migration of transformers enables our model to better learn structural knowledge and preserve detailed information for point cloud completion. Furthermore, we propose two more challenging benchmarks with more diverse incomplete point clouds that can better reflect the real-world scenarios to promote future research. Experimental results show that our method outperforms state-of-the-art methods by a large margin on both the new benchmarks and the existing ones. Code is available at <https://github.com/yuxumin/PoinTr>.

1. Introduction

Recent developments in 3D sensors largely boost researches in 3D computer vision. One of the most commonly used 3D data format is the point cloud, which requires less memory to store but convey detailed 3D shape

*Equal contribution.

†Corresponding author.

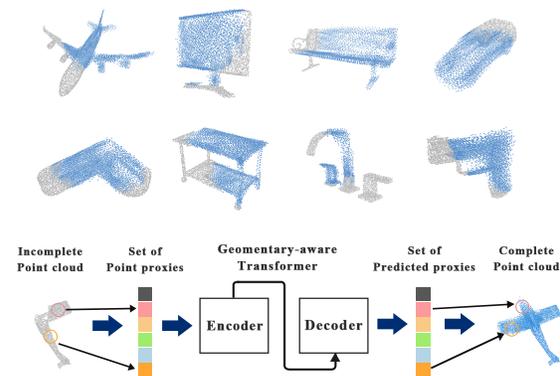


Figure 1: PoinTr is designed for point cloud completion task. It takes the downsampled partial point clouds as inputs (gray points), and predicts the missing parts and upsamples the known parts simultaneously (blue points). We propose to formulate the point cloud completion task as a set-to-set translation task and use a transformer encoder-decoder architecture to learn the complex dependencies among the point groups. Furthermore, we design two new benchmarks with more diverse tasks (*i.e.*, upsampling and completion of point cloud), more diverse categories (*i.e.*, from 8 categories to 55 categories), more diverse viewpoints (*i.e.*, from 8 viewpoints to all possible viewpoints) and more diverse levels of incompleteness (*i.e.*, missing 25% to 75% points of the ground-truth point clouds) to better reflect the real-world scenarios and promote future research.

information. However, point cloud data from existing 3D sensors are not always complete and satisfactory because of inevitable self-occlusion, light reflection, limited sensor resolution, *etc.* Therefore, recovering complete point clouds from partial and sparse raw data becomes an indispensable task with ever-growing significance.

Over the years, researchers have tried many approaches

to tackle this problem in the realm of deep learning. Early attempts on point cloud completion [6, 15, 30, 31, 23, 36, 20, 19, 51, 47, 40] try to migrate mature methods from 2D completion tasks to 3D point clouds by voxelization and 3D convolutions. However, these methods suffer from a heavy computational cost that grows cubically as the spatial resolution increases. With the success of PointNet and PointNet++ [25, 26], directly processing 3D coordinates becomes the mainstream of point cloud based 3D analysis. The technique is further applied to many pioneer works [1, 49, 35, 16, 22, 14, 29] in point cloud completion task, in which an encoder-decoder based architecture is designed to generate complete point clouds. However, the bottleneck of such methods lies in the max-pooling operation in the encoding phase, where fine-grained information is lost and can hardly be recovered in the decoding phase.

Reconstructing complete point cloud is a challenging problem since the structural information required in the completion task runs counter to the unordered and unstructured nature of point cloud data. Therefore, learning structural features and long-range correlations among local parts of the point cloud becomes the key ingredient towards better point cloud completion. In this paper, we propose to adopt Transformers [37], one of the most successful architecture in Natural Language Processing (NLP), to learn the structural information of pairwise interactions and global correlations for point cloud completion. Our model, named *PoinTr*, is characterized by five key components: 1) *Encoder-Decoder Architecture*: We adopt the encoder-decoder architecture to convert point cloud completion as a set-to-set translation problem. The self-attention mechanism of transformers models all pairwise interactions between elements in the encoder, while the decoder reasons about the missing elements based on the learnable pairwise interactions among features of the input point cloud and queries; 2) *Point Proxy*: We represent the set of point clouds in a local region as a feature vector called *Point Proxy*. The input point cloud is convert to a sequence of Point Proxies, which are used as the inputs of our transformer model; 3) *Geometry-aware Transformer Block*: To facilitate transformers to better leverage the inductive bias about 3D geometric structures of point clouds, we design a geometry-aware block that models the geometric relations explicitly; 4) *Query Generator*: We use dynamic queries instead of fixed queries in the decoder, which are generated by a query generation module that summarizes the features produced by the encoder and represents the initial sketch of the missing points; 5) *Multi-Scale Point Cloud Generation*: We devise a multi-scale point generation module to recover the missing point cloud in a coarse-to-fine manner.

As another contribution, we argue that existing benchmarks are not representative enough to cover real-world scenarios of incompleting point clouds. Therefore, we intro-

duce two more challenging benchmarks that contain more diverse tasks (*i.e.*, joint upsampling and completion of point cloud), more object categories (*i.e.*, from 8 categories to 55 categories), more diverse views points (*i.e.*, from 8 viewpoints to all possible viewpoints) and more diverse level of incompleteness (*i.e.*, missing 25% to 75% points of the ground-truth point clouds). We evaluate our method on both the new benchmarks and the widely used PCN dataset [49] and KITTI benchmark [10]. Experiments demonstrate that *PoinTr* outperforms previous state-of-the-art methods on all benchmarks by a large margin. The main contributions of this paper are summarized in Figure 1.

2. Related Work

3D Shape Completion. Traditional methods for 3D shape completion tasks often adopt voxel grids or distance fields to describe 3D objects [6, 15, 31]. Based on such structured 3D representations, the powerful 3D convolutions are used and achieve a great success in the tasks of 3D reconstruction [4, 11] and shape completion [6, 15, 44]. However, this group of methods suffers from heavy memory consumption and computational burden. Although these issues are further alleviated by methods based on sparse representations [32, 39, 12], the quantization operation in these methods still cause a significant loss in detailed information. Different from the above methods, researchers gradually start to use unstructured point clouds as the representation of 3D objects, given the small memory consumption and strong ability to represent fine-grained details. Nevertheless, the migration from structured 3D data understanding to point clouds analysis is non-trivial, since the commonly used convolution operator is no longer suitable for unordered points clouds. PointNet and its variants [25, 26] are the pioneer work to directly process 3D coordinates and inspire the researches in many downstream tasks. In the realm of point cloud completion, PCN [49] is the first learning-based architecture, which proposes an Encoder-Decoder framework and adopts a FoldingNet to map the 2D points onto a 3D surface by mimicking the deformation of a 2D plane. After PCN, many other methods [35, 16, 46, 18] spring up, pursuing point clouds completion in higher resolution with better robustness.

Transformers. Transformers [37] are first introduced as an attention-based framework in Natural Language Processing (NLP). Transformer models often utilize the encoder-decoder architecture and are characterized by both self-attention and cross-attention mechanisms. Transformer models have proven to be very helpful to the tasks that involve long sequences thanks to the self-attention mechanism. The cross-attention mechanism in the decoder exploit the encoder information to learn the attention map of query features, which making transformers powerful in generation

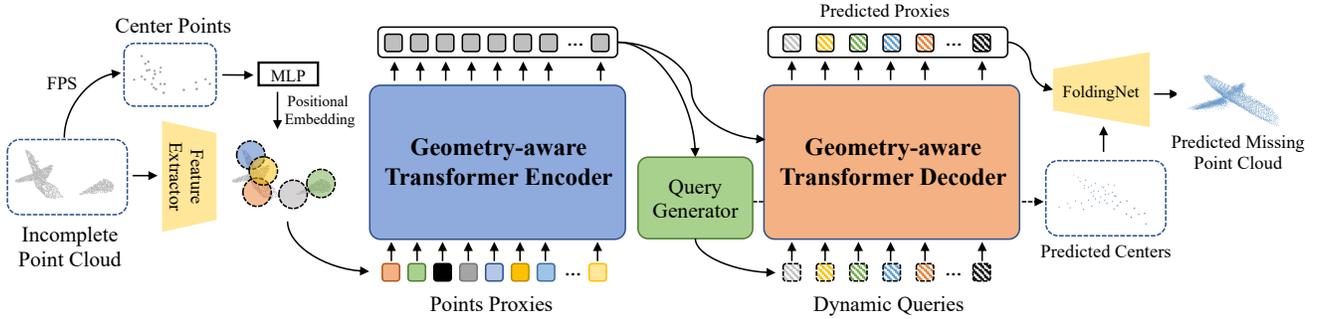


Figure 2: The Pipeline of *PointTr*. We first downsample the input partial point cloud to obtain the center points. Then, we use a lightweight DGCNN [42] to extract the local features around the center points. After adding the position embedding to the local feature, we use a transformer encoder-decoder architecture to predict the point proxies for the missing parts. A simple MLP and FoldingNet are used to complete the point cloud based on the predicted point proxies in a coarse-to-fine manner.

tasks. By taking the advantages of both self-attention and cross-attention mechanisms, transformers have a strong capability to handle long sequence input and enhance information communications between the encoder and the decoder. In the past few years, transformers have dominated the tasks that take long sequences as input and gradually replaced RNNs [38] in many domains. Now they begin their journey in computer vision [7, 21, 33, 27, 24, 28].

3. Approach

The overall framework of *PointTr* is illustrated in Figure 2. We will introduce our method in detail as follows.

3.1. Set-to-Set Translation with Transformers

The primary goal of our method is to leverage the impressive sequence-to-sequence generation ability of transformer architecture for point cloud completion tasks. We propose to first convert the point cloud to a set of feature vectors, *point proxies*, that represent the local regions in the point clouds (we will describe in Section 3.2). By analogy to the language translation pipeline, we model point cloud completion as a set-to-set translation task, where the transformers take the point proxies of the partial point clouds as the inputs and produce the point proxies of the missing parts. Specifically, given the set of point proxies $\mathcal{F} = \{F_1, F_2, \dots, F_N\}$ that represents the partial point cloud, we model the process of point cloud completion as a set-to-set translation problem:

$$\mathcal{V} = \mathcal{M}_E(\mathcal{F}), \quad \mathcal{H} = \mathcal{M}_D(\mathcal{Q}, \mathcal{V}), \quad (1)$$

where \mathcal{M}_E and \mathcal{M}_D are the encoder and decoder models, $\mathcal{V} = \{V_1, V_2, \dots, V_N\}$ are the output features of the encoder, $\mathcal{Q} = \{Q_1, Q_2, \dots, Q_M\}$ are the dynamic queries for the decoder, $\mathcal{H} = \{H_1, H_2, \dots, H_M\}$ are the predicted point proxies of the missing point cloud, and M is the number of the predicted point proxies. The recent success in

NLP tasks like text translation and question answering [8] have clearly demonstrated the effectiveness of transformers to solve this kind of problem. Therefore, we propose to adopt a transformer-based encoder-decoder architecture to solve the point cloud completion problem.

The encoder-decoder architecture consists of L_E and L_D multi-head self-attention layers [37] in the encoder and decoder, respectively. The self-attention layer in the encoder first updates proxy features with both long-range and short-range information. Then a feed forward network (FFN) further updates the proxy features with an MLP architecture. The decoder utilizes self-attention and cross-attention mechanisms to learn structural knowledge. The self-attention layer enhances the local features with global information, while the cross-attention layer explores the relationship between queries and outputs of the encoder. To predict the point proxies of the missing parts, we propose to use dynamic query embeddings, which makes our decoder more flexible and adjustable for different types of objects and their missing information. More details about the transformer architecture can be found in the supplementary material and [8, 37].

Note that benefiting from the self-attention mechanism in transformers, the features learned by the transformer network are invariant to the order of point proxies, which is also the basis of using transformers to process point clouds. Considering the strong ability to capture data relationships, we expect the transformer architecture to be a promising alternative for deep learning on point clouds.

3.2. Point Proxy

The Transformers in NLP take as input a 1D sequence of word embeddings [37]. To make 3D point clouds suitable for transformers, the first step is to convert the point cloud to a sequence of vectors. A trivial solution is directly feeding the sequence of xyz coordinates to the transformers. However, since the computational complexity of the trans-

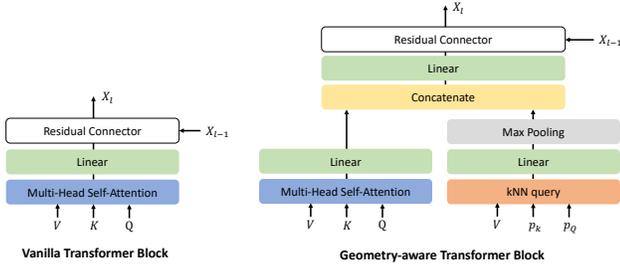


Figure 3: Comparisons of the vanilla transformer block and the proposed geometry-aware transformer block.

formers is quadratic to the sequence length, this solution will lead to an unacceptable cost. Therefore, we propose to represent the original point cloud as a set of *point proxies*. A point proxy represents a local region of the point clouds. Inspired by the set abstraction operation in [26], we first conduct *furthest point sample (FPS)* to locate a fixed number N of point centers $\{q_1, q_2, \dots, q_N\}$ in the partial point cloud. Then, we use a light-weight DGCNN [42] with hierarchical downsampling to extract the feature of the point centers from the input point cloud. The point proxy F_i is a feature vector that captures the local structure around q_i , which can be computed as:

$$F_i = F'_i + \varphi(q_i), \quad (2)$$

where F'_i is the feature of point q_i that is extracted using the DGCNN model, and φ is another MLP to capture the location information of the point proxy. The first term represents the semantic patterns of the local region, and the second term is inspired by the position embedding [3] operation in transformers, which explicitly encodes the global location of the point proxy. The detailed architecture of the feature extraction model can be found in Supplementary Material.

3.3. Geometry-aware Transformer Block

One of the key challenges of applying transformers for vision tasks is the self-attention mechanism in transformers lacks some inductive biases inherent to conventional vision models like CNNs and point cloud networks which explicitly model the structures of vision data. To facilitate transformers to better leverage the inductive bias about 3D geometric structures of point clouds, we design a geometry-aware block that models the geometric relations, which can be a plug-and-play module to incorporate with the attention blocks in any transformer architectures. The details of the proposed block are shown in Figure 3. Different from the self-attention module that uses the feature similarity to capture the semantic relation, we propose to use kNN model to capture the geometric relation in the point cloud. Given the query coordinates p_Q , we query the features of the nearest keys according to the key coordinates p_k . Then we fol-

low the practice of DGCNN [42] to learn the local geometric structures by feature aggregation with a linear layer followed by the max pooling operation. The geometric feature and semantic feature are then concatenated and mapped to the original dimensions to form the output.

3.4. Query Generator

The queries \mathcal{Q} serve as the initial state of the predicted proxies. To make sure the queries correctly reflect the sketch of the completed point cloud, we propose a query generator module to generate the query embeddings dynamically conditioned on the encoder outputs. Specifically, we first summarize \mathcal{V} with a linear projection to higher dimensions followed by the max pooling operation. Similar to [49], we use a linear projection to directly generate $M \times 3$ dimension features that can be reshaped as the M coordinates $\{c_1, c_2, \dots, c_M\}$. Lastly, we concatenate the global feature of the encoder and the coordinates, and use an MLP to produce the query embeddings.

3.5. Multi-Scale Point Cloud Generation

The goal of our encoder-decoder network is to predict the missing parts of incomplete point clouds. However, we can only get predictions for missing proxies from the transformer decoder. Therefore, we propose a multi-scale point cloud generation framework to recover missing point clouds at full resolution. To reduce redundant computations, we reuse the M coordinates produced by the query generator as the local centers of the missing point cloud. Then, we utilize a FoldingNet [48] f to recover detailed local shapes centered at the predicted proxies:

$$\mathcal{P}_i = f(H_i) + c_i, \quad i = 1, 2, \dots, M. \quad (3)$$

where \mathcal{P}_i is the set of neighboring points centered at c_i . Following previous work [16], we only predict the missing parts of the point cloud and concatenate them with the input point cloud to obtain the complete objects. Both predicted proxies and recovered point clouds are supervised during the training process, and the detailed loss function will be introduced in the following section.

3.6. Optimization

The loss function for point cloud completion should provide a quantitative measurement for the quality of output. However, since the point clouds are unordered, many loss functions that directly measure the distance between two points (*i.e.* ℓ_2 distance) are unsuitable. Fan *et al.* [9] introduce two metrics that are invariant to the permutation of points, which are Chamfer Distance (CD) and Earth Mover’s Distance (EMD). We adopt Chamfer Distance as our loss function for its $\mathcal{O}(N \log N)$ complexity. We use \mathcal{C} to represent the n_C local centers and \mathcal{P} to represent n_P points of the completed point cloud. Give the ground-truth

completed point cloud \mathcal{G} , the loss functions for these two predictions can be written as:

$$J_0 = \frac{1}{n_C} \sum_{c \in \mathcal{C}} \min_{g \in \mathcal{G}} \|c - g\| + \frac{1}{n_G} \sum_{g \in \mathcal{G}} \min_{c \in \mathcal{C}} \|g - c\|,$$

$$J_1 = \frac{1}{n_P} \sum_{p \in \mathcal{P}} \min_{g \in \mathcal{G}} \|p - g\| + \frac{1}{n_G} \sum_{g \in \mathcal{G}} \min_{p \in \mathcal{P}} \|g - p\|.$$

Note that we also concatenate the predicted local centers and the centers of the input point cloud to form the local centers of the whole object \mathcal{C} . We directly use the high-resolution point cloud \mathcal{G} to supervise the sparse point cloud \mathcal{C} to encourage them to have similar distributions. Our final objective function is the sum of these two objectives $J = J_0 + J_1$.

4. Experiments

In this section, we first introduce the new benchmarks for diverse point cloud completion and the evaluation metric. Then, we show the results of both our method and several baseline methods on our new benchmarks. Lastly, we demonstrate the effectiveness of our model on the widely used PCN dataset and KITTI benchmark. We also provide ablation study and visual analysis of our method.

4.1. Benchmarks for Diverse Point Completion

We choose to generate the samples in our benchmarks based on the synthetic dataset, ShapeNet [45], because it contains the complete object models that cannot be obtained from real-world datasets like ScanNet [5] and S3DIS [2]. What makes our benchmarks distinct is that our benchmarks contain more object categories, more incomplete patterns and more viewpoints. Besides, we pay more attention to the ability of networks to deal with the objects from novel categories that do not appear in the training set.

ShapeNet-55 Benchmark: In this benchmark, we use all the objects in ShapeNet from 55 categories. Most existing datasets for point cloud completion like PCN [49] only consider a relatively small number of categories (*e.g.*, 8 categories in PCN). However, the incomplete point clouds from the real-world scenarios are much more diverse. Therefore, we propose to evaluate the point cloud completion models on all 55 categories in ShapeNet to more comprehensively test the ability of models with a more diverse dataset. We split the original ShapeNet using the 80-20 strategy: we randomly sample 80% objects from each category to form the training set and use the rest for evaluation. As a result, we get 41,952 models for training and 10,518 models for testing. For each object, we randomly sample 8,192 points from the surface to obtain the point cloud.

ShapeNet-34 Benchmark: In this benchmark, we want to explore another important issue in point cloud completion:

the performance on novel categories. We believe it is necessary to build a benchmark for this task to better evaluate the generalization performance of models. We first split the origin ShapeNet into two parts: 21 unseen categories and 34 seen categories. In the seen categories, we randomly sample 100 objects from each category to construct a test set of the seen categories (3,400 objects in total) and leave the rest as the training set, resulting in 46,765 object models for training. We also construct another test set consisting of 2,305 objects from 21 novel categories. We evaluate the performance on both the seen and unseen categories to show the generalization ability of models.

Training and Evaluation: In both benchmarks, the partial point clouds for training are generated online. We sample 2048 points from the object as the input and 8192 points as the ground truth. In order to mimic the real-world situation, we first randomly select a viewpoint and then remove the n furthest points from the viewpoint to obtain a training partial point cloud. Although the projection method proposed in [49] is a better approximation to real scans, our strategy is more flexible and efficient. Our experiments on KITTI also show the model learned on our dataset works well when finetuning to real-world scans. Besides, our strategy also ensures the diversity of our training samples in the aspect of viewpoints. During training, n is randomly chosen from 2048 to 6144 (25% to 75% of the complete point cloud), resulting in different level of incompleteness. We then down-sample the remaining point clouds to 2048 points as the input data for models.

During evaluation, we fix 8 view points and n is set to 2048, 4096 or 6144 (25%, 50% or 75% of the whole point cloud) for convenience. According to the value of n , we divide the test samples into three difficulty degrees, *simple*, *moderate* and *hard* in our experiments. In the following experiments, we will report the performance for each method in *simple*, *moderate* and *hard* to show the ability of each network to deal with tasks at difficulty levels. In addition, we use the average performance under three difficulty degrees to report the overall performance (*Avg*).

4.2. Evaluation Metric

We follow the existing works [49, 35, 16, 46] to use the mean Chamfer Distance as the evaluation metric, which can measure distance between the prediction point cloud and ground-truth in set-level. For each prediction, the Chamfer Distance between the prediction point set \mathcal{P} and the ground-truth point set \mathcal{G} is calculated by:

$$d_{CD}(\mathcal{P}, \mathcal{G}) = \frac{1}{|\mathcal{P}|} \sum_{p \in \mathcal{P}} \min_{g \in \mathcal{G}} \|p - g\| + \frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \min_{p \in \mathcal{P}} \|g - p\|$$

Following the previous methods, we use two versions of Chamfer distance as the evaluation metric to compare the performance with existing works. $CD-\ell_1$ uses L1-norm to

Table 1: Results of our methods and state-of-the-art methods on ShapeNet-55. We report the detailed results for each method on 10 categories and the overall results on 55 categories for three difficulty degrees. We use CD-S, CD-M and CD-H to represent the CD results under the *Simple*, *Moderate* and *Hard* settings. We also provide results under the F-Score@1% metric.

	Table	Chair	Airplane	Car	Sofa	Bird house	Bag	Remote	Key board	Rocket	CD-S	CD-M	CD-H	CD-Avg	F1
FoldingNet [48]	2.53	2.81	1.43	1.98	2.48	4.71	2.79	1.44	1.24	1.48	2.67	2.66	4.05	3.12	0.082
PCN [49]	2.13	2.29	1.02	1.85	2.06	4.50	2.86	1.33	0.89	1.32	1.94	1.96	4.08	2.66	0.133
TopNet [35]	2.21	2.53	1.14	2.18	2.36	4.83	2.93	1.49	0.95	1.32	2.26	2.16	4.3	2.91	0.126
PFNet [16]	3.95	4.24	1.81	2.53	3.34	6.21	4.96	2.91	1.29	2.36	3.83	3.87	7.97	5.22	0.339
GRNet [46]	1.63	1.88	1.02	1.64	1.72	2.97	2.06	1.09	0.89	1.03	1.35	1.71	2.85	1.97	0.238
PoinTr	0.81	0.95	0.44	0.91	0.79	1.86	0.93	0.53	0.38	0.57	0.58	0.88	1.79	1.09	0.464

Table 2: Results of our methods and state-of-the-art methods on ShapeNet-34. We report the results of 34 seen categories and 21 unseen categories in three difficulty degrees. We use CD-S, CD-M and CD-H to represent the CD results under the *Simple*, *Moderate* and *Hard* settings. We also provide results under the F-Score@1% metric.

	34 seen categories					21 unseen categories				
	CD-S	CD-M	CD-H	CD-Avg	F1	CD-S	CD-M	CD-H	CD-Avg	F1
FoldingNet [48]	1.86	1.81	3.38	2.35	0.139	2.76	2.74	5.36	3.62	0.095
PCN [49]	1.87	1.81	2.97	2.22	0.154	3.17	3.08	5.29	3.85	0.101
TopNet [35]	1.77	1.61	3.54	2.31	0.171	2.62	2.43	5.44	3.50	0.121
PFNet [16]	3.16	3.19	7.71	4.68	0.347	5.29	5.87	13.33	8.16	0.322
GRNet [46]	1.26	1.39	2.57	1.74	0.251	1.85	2.25	4.87	2.99	0.216
PoinTr	0.76	1.05	1.88	1.23	0.421	1.04	1.67	3.44	2.05	0.384

calculate the distance between two points, while CD- ℓ_2 uses L2-norm instead. We also follow [34] to adopt F-Score as another evaluation metric.

4.3. Results on ShapeNet-55

We first conduct experiments on ShapeNet-55, which consists of objects from 55 categories. To compare with existing methods, We implement FoldingNet [48], PCN [49], TopNet [35], PFNet [16] and GRNet [46] on our benchmark according to their open-source code and use the best hyper-parameters in their papers for fair comparisons. We first investigate how the existing methods and our method perform when there are objects from more categories. The last four columns in Table 1 show that our PoinTr can better cope with different situations with diverse viewpoints, diverse object categories, diverse incomplete patterns and diverse incompleteness levels. We achieve 0.58, 0.6 and 0.69 improvement in CD- ℓ_2 (multiplied by 1000) under three settings (simple, moderate and hard) comparing with the SOTA method GRNet [46]. PFNet [16], which proposes to directly predict the missing parts of objects, fail in our benchmarks due to the high diversity. We further report the performance on categories with sufficient and insufficient samples. We only sample 10 categories out from 55 categories to show the results due to the limited space, in which *Table*, *chair*, *Airplane*, *Car* and *Sofa* contain more than 2500 samples in the training set while *Birdhouse*, *Bag*, *Remote*,

Keyboard and *Rocket* contain less than 80 samples. We also provide the detailed results for all 55 categories in our supplementary material. We place the categories with sufficient samples at the first five columns and the categories with insufficient samples in the following five columns in Table 1. The average CD results for three difficulty degrees are also reported. Surprisingly, there is no obvious difference between the results on these two kinds of categories. However, except for our PoinTr and SOTA method GRNet, the imbalance in the number of training samples lead to a relatively high CD in the categories with insufficient samples. Besides, our model achieves 0.46 F-Score on ShapeNet-55, while the state-of-the-art GRNet only obtain 0.24 F-Score. These results clearly demonstrate the effectiveness of PoinTr under the more diverse setting.

4.4. Results on ShapeNet-34

On ShapeNet-34, we also conduct experiments for our method and other five state-of-the-art methods. The results are shown in Table 2. For the 34 seen categories, we can see our method outperforms all the other methods. For the 21 unseen categories, we using the networks that are trained on the 34 seen categories to evaluate the performance on the novel objects from the other 21 categories that do not appear in the training phase. We see our method also achieves the best performance in this more challenging setting. Comparing with the results of seen categories, we see in the simple

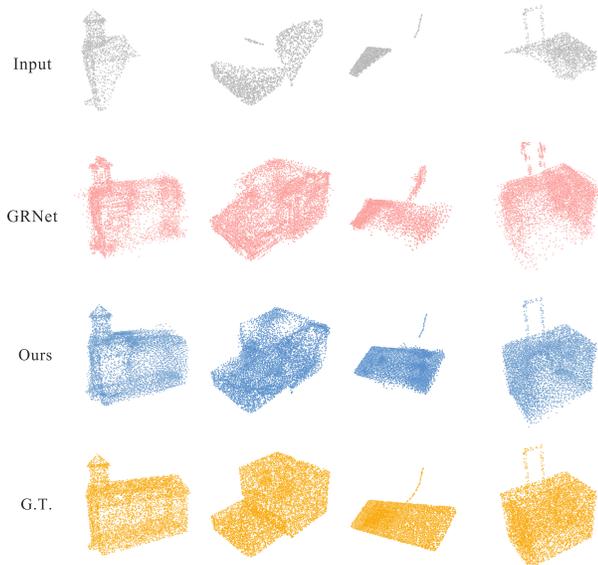


Figure 4: Point cloud completion results on some objects from novel categories. We show the input point cloud and the ground truth as well as the predictions of GRNet and our model.

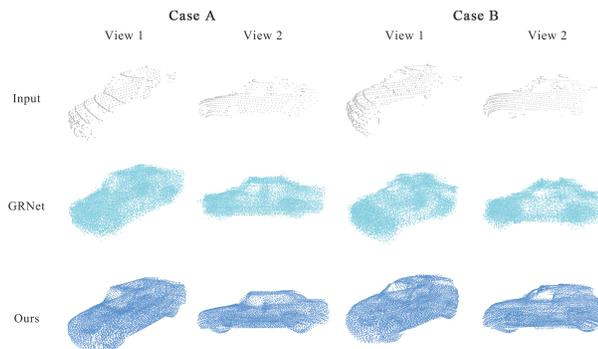


Figure 5: Qualitative results on the KITTI dataset. In order to better show the shape of the car, we provide two views of the same point cloud in each case. Our method can recover the complete point cloud of a car with more accurate boundaries and details (e.g. tires of cars).

setting (25% of point cloud will be removed), the performance drop of our method is less than 0.3. But when the difficulty level increases, the performance gap between seen categories and unseen categories significantly increases. We also visualize the results in Figure 4 to show the effectiveness of our method on the unseen categories.

4.5. Results on the Existing Benchmarks

Apart from the experiments on the two newly proposed challenging benchmarks, we also conduct the experiments on the existing benchmarks including the PCN dataset [49]

Table 3: Results on the PCN dataset. We use the L1 Chamfer Distance to compare with other methods.

CD- ℓ_1 ($\times 1000$)	Avg	Air	Cab	Car	Cha	Lam	Sof	Tab	Wat
FoldingNet [48]	14.31	9.49	15.80	12.61	15.55	16.41	15.97	13.65	14.99
AtlasNet [13]	10.85	6.37	11.94	10.10	12.06	12.37	12.99	10.33	10.61
PCN [49]	9.64	5.50	22.70	10.63	8.70	11.00	11.34	11.68	8.59
TopNet [35]	12.15	7.61	13.31	10.90	13.82	14.44	14.78	11.22	11.12
MSN [17]	10.0	5.6	11.9	10.3	10.2	10.7	11.6	9.6	9.9
GRNet [46]	8.83	6.45	10.37	9.45	9.41	7.96	10.51	8.44	8.04
PMP-Net [43]	8.73	5.65	11.24	9.64	9.51	6.95	10.83	8.72	7.25
CRN [41]	8.51	4.79	9.97	8.31	9.49	8.94	10.69	7.81	8.05
PoinTr	8.38	4.75	10.47	8.68	9.39	7.75	10.93	7.78	7.29

Table 4: Ablation study on the PCN dataset. We investigate different designs including query generator (Query), DGCNN feature extractor (DGCNN) and Geometry-aware Blocks (Geometry).

Model	Query	DGCNN	Geometry	CD- ℓ_1	F-Score@1%
A				9.43	67.82
B	✓			9.09	0.713
C	✓	✓		8.69	0.736
D	✓	✓	all	8.44	0.741
E	✓	✓	1 st	8.38	0.745

and KITTI benchmark [10].

Results on the PCN Dataset. The PCN dataset [49] is one of the most widely used benchmark datasets for the point cloud completion task. To verify the effectiveness of our method on existing benchmarks and compare it with more state-of-the-art methods, we conducted experiments on this dataset following the standard protocol and evaluation metric used in previous work [49, 17, 46, 43, 41]. The results are shown in Table 3. We see our method largely improves the previous methods and establishes the new state-of-the-art on this dataset.

Results on KITTI Benchmark. To show the performance of our method in real-world scenarios, we follow [46] to finetune our trained model on ShapeNetCars [49] and evaluate the performance of our model on KITTI dataset, which contains the incomplete point clouds of cars in the real-world scenes from LiDAR scans. We report the Fidelity and MMD metrics in Table 5 and show some reconstruction results in Figure 5. Our method achieves better qualitative and quantitative performance.

4.6. Model Design Analysis

To examine the effectiveness of our designs, we conduct a detailed ablation study on the key components of PoinTr. The results are summarized in Table 4. The baseline model A is the vanilla transformer model for point cloud completion, which uses the encoder-decoder architecture with the standard transformer blocks. In this model, we form the point proxies directly from the point cloud using a single-layer DGCNN model. We then add the query generator

Table 5: Results on LiDAR scans from KITTI dataset under the Fidelity and MMD metrics.

$CD-\ell_2 (\times 1000)$	AtlasNet [13]	PCN [49]	FoldingNet [48]	TopNet [35]	MSN [17]	NSFA [50]	PFNet [16]	CRN [41]	GRNet [46]	PoinTr
Fidelity ↓	1.759	2.235	7.467	5.354	0.434	1.281	1.137	1.023	0.816	0.000
MMD ↓	2.108	1.366	0.537	0.636	2.259	0.891	0.792	0.872	0.568	0.526

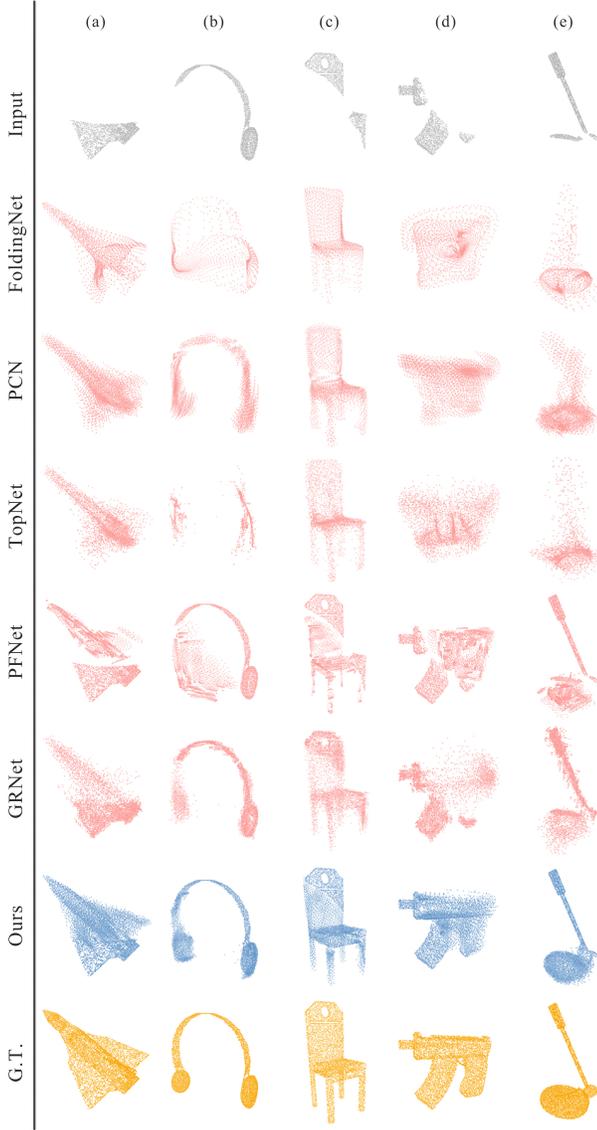


Figure 6: Qualitative results on ShapeNet-55. All methods above takes the point clouds in the first line as inputs and generate complete point clouds. Our methods can complete the point clouds with higher fidelity, which clearly shows the effectiveness of our method.

between the encoder and decoder (model B). We see the query generator improve the baseline by 0.34 in Chamfer distance. When using DGCNN to extract features from the input point cloud (model C), we observe a significant im-

provement to 8.69. By adding the geometric block to all the transformer blocks (model D), we see the performance can be further improved, which clearly demonstrates the effectiveness of the geometric structures learned by the block. We find that only adding the geometric block to the first transformer block in both encoder and decoder can lead to a slightly better performance (model E), which indicates the role of geometric block is to introduce the inductive bias and a single layer is sufficient while adding more blocks may result in over-fitting. Besides, we see our method can achieve over 0.74 F-Score on the PCN dataset while obtaining only 0.46 F-Score on our ShapeNet-55, which also suggests our new dataset is much more challenging.

4.7. Qualitative Results

In Figure 6, we show some completion results for all methods and find our method perform better. For example, the input data in (a) nearly lose all the geometric information and can be hardly recognized as an airplane. In this case, other methods can only roughly complete the shape with unsatisfactory geometry details, while our method can still complete the point cloud with higher fidelity. These results show our method has a stronger ability to recover details and is more robust to various incomplete patterns. More results can be found in the supplementary material.

5. Conclusion

In this paper, we have proposed a new architecture, PoinTr, to convert the point cloud completion task into a set to set translation tasks. With several technical innovations, we successfully applied the transformer model to this task and achieved state-of-the-art performance. Moreover, we proposed two more challenging benchmarks for more diverse point cloud completion. Extending our transformer architecture to other 3D tasks can an interesting future direction.

Acknowledgements

This work was supported in part by the National Key Research and Development Program of China under Grant 2017YFA0700802, in part by the National Natural Science Foundation of China under Grant 61822603, Grant U1813218, and Grant U1713214, in part by a grant from the Beijing Academy of Artificial Intelligence (BAAI), and in part by a grant from the Institute for Guo Qiang, Tsinghua University.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Learning representations and generative models for 3d point clouds. In *ICLR*, 2018. 2
- [2] Iro Armeni, Sasha Sax, Amir R Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017. 5
- [3] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. In *ICCV*, pages 3286–3295, 2019. 4
- [4] Christopher B. Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *ECCV*, pages 628–644, 2016. 2
- [5] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, pages 5828–5839, 2017. 5
- [6] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *CVPR*, pages 6545–6554, 2017. 2
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *NAACL*, pages 4171–4186, 2019. 3
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. 3
- [9] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *CVPR*, pages 605–613, 2017. 4
- [10] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *international journal of robotics research (ijrr)*. 2013. 2, 7
- [11] Rohit Girdhar, David F. Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling, editors, *ECCV*, pages 484–499, 2016. 2
- [12] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, pages 9224–9232, 2018. 2
- [13] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. Atlasnet: A papier-mâché approach to learning 3d surface generation. *CoRR*, 2018. 7, 8
- [14] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *CVPR*, pages 216–224, 2018. 2
- [15] 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 *ICCV*, pages 85–93, 2017. 2
- [16] Zitian Huang, Yikuan Yu, Jiawen Xu, Feng Ni, and Xinyi Le. PF-net: Point fractal network for 3d point cloud completion. In *CVPR*, pages 7659–7667, 2020. 2, 4, 5, 6, 8
- [17] Minghua Liu, Lu Sheng, Sheng Yang, Jing Shao, and Shi-Min Hu. Morphing and sampling network for dense point cloud completion. *arXiv preprint arXiv:1912.00280*, 2019. 7, 8
- [18] Minghua Liu, Lu Sheng, Sheng Yang, Jing Shao, and Shi-Min Hu. Morphing and sampling network for dense point cloud completion. In *AAAI*, pages 11596–11603, 2020. 2
- [19] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *CVPR*, pages 8895–8904, 2019. 2
- [20] Zhijian Liu, Haotian Tang, Yujun Lin, and Song Han. Point-voxel CNN for efficient 3d deep learning. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *NeurIPS*, pages 963–973, 2019. 2
- [21] Christoph Lüscher, Eugen Beck, Kazuki Irie, Markus Kitzka, Wilfried Michel, Albert Zeyer, Ralf Schlüter, and Hermann Ney. RWTH ASR systems for librispeech: Hybrid vs attention - w/o data augmentation. *CoRR*, abs/1905.03072, 2019. 3
- [22] Priyanka Mandikal and R. Venkatesh Babu. Dense 3d point cloud reconstruction using a deep pyramid network. *CoRR*, abs/1901.08906, 2019. 2
- [23] Duc Thanh Nguyen, Binh-Son Hua, Minh-Khoi Tran, Quang-Hieu Pham, and Sai-Kit Yeung. A field model for repairing 3d shapes. In *CVPR*, pages 5676–5684, 2016. 2
- [24] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Lukasz Kaiser, Noam Shazeer, and Alexander Ku. Image transformer. *CoRR*, abs/1802.05751, 2018. 3
- [25] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 77–85, 2017. 2
- [26] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *NeurIPS*, pages 5099–5108, 2017. 2, 4
- [27] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, page 9, 2019. 3
- [28] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. Dynamicvit: Efficient vision transformers with dynamic token sparsification. *arXiv preprint arXiv:2106.02034*, 2021. 3
- [29] Muhammad Sarmad, Hyunjoo Jenny Lee, and Young Min Kim. RL-gan-net: A reinforcement learning agent controlled gan network for real-time point cloud shape completion. In *CVPR*, pages 5898–5907, 2019. 2
- [30] Abhishek Sharma, Oliver Grau, and Mario Fritz. Vconv-dae: Deep volumetric shape learning without object labels. *CoRR*, abs/1604.03755, 2016. 2

- [31] David Stutz and Andreas Geiger. Learning 3d shape completion from laser scan data with weak supervision. In *CVPR*, pages 1955–1964, 2018. [2](#)
- [32] Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. Splatnet: Sparse lattice networks for point cloud processing. In *CVPR*, pages 2530–2539, 2018. [2](#)
- [33] Gabriel Synnaeve, Qiantong Xu, Jacob Kahn, Tatiana Likhomanenko, Edouard Grave, Vineel Pratap, Anuroop Sriram, Vitaliy Liptchinsky, and Ronan Collobert. End-to-end ASR: from Supervised to Semi-Supervised Learning with Modern Architectures. *arXiv e-prints*, 2019. [3](#)
- [34] Maxim Tatarchenko, Stephan R. Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? *CoRR*, abs/1905.03678, 2019. [6](#)
- [35] Lyne P. Tchapmi, Vineet Kosaraju, Hamid Rezaatofghi, Ian D. Reid, and Silvio Savarese. Topnet: Structural point cloud decoder. In *CVPR*, pages 383–392, 2019. [2](#), [5](#), [6](#), [7](#), [8](#)
- [36] Jacob Varley, Chad DeChant, Adam Richardson, Joaquín Ruales, and Peter K. Allen. Shape completion enabled robotic grasping. In *IROS*, pages 2442–2447, 2017. [2](#)
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *NeurIPS*, pages 5998–6008, 2017. [2](#), [3](#)
- [38] Oriol Vinyals, Samy Bengio, and Manjunath Kudlur. Order matters: Sequence to sequence for sets. *arXiv preprint arXiv:1511.06391*, 2015. [3](#)
- [39] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-CNN: octree-based convolutional neural networks for 3d shape analysis. *ACM Trans. Graph.*, 2017. [2](#)
- [40] Weiyue Wang, Qiangui Huang, Suyu You, Chao Yang, and Ulrich Neumann. Shape inpainting using 3d generative adversarial network and recurrent convolutional networks. In *ICCV*, pages 2298–2306, 2017. [2](#)
- [41] Xiaogang Wang, Marcelo H Ang Jr, and Gim Hee Lee. Cascaded refinement network for point cloud completion. In *CVPR*, pages 790–799, 2020. [7](#), [8](#)
- [42] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *TOG*, 38(5):1–12, 2019. [3](#), [4](#)
- [43] Xin Wen, Peng Xiang, Zhizhong Han, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Yu-Shen Liu. Pmp-net: Point cloud completion by learning multi-step point moving paths. In *CVPR*, 2021. [7](#)
- [44] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, pages 1912–1920, 2015. [2](#)
- [45] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, pages 1912–1920, 2015. [5](#)
- [46] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Jiageng Mao, Shengping Zhang, and Wenxiu Sun. Grnet: Gridding residual network for dense point cloud completion. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *ECCV*, pages 365–381, 2020. [2](#), [5](#), [6](#), [7](#), [8](#)
- [47] Bo Yang, Hongkai Wen, Sen Wang, Ronald Clark, Andrew Markham, and Niki Trigoni. 3d object reconstruction from a single depth view with adversarial learning. In *ICCVW*, pages 679–688, 2017. [2](#)
- [48] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Interpretable unsupervised learning on 3d point clouds. *CoRR*, abs/1712.07262, 2017. [4](#), [6](#), [7](#), [8](#)
- [49] Wentao Yuan, Tejas Khot, David Held, Christoph Mertz, and Martial Hebert. PCN: point completion network. In *3DV*, pages 728–737, 2018. [2](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [50] Wenxiao Zhang, Qingan Yan, and Chunxia Xiao. Detail preserved point cloud completion via separated feature aggregation. In *ECCV*, 2020. [8](#)
- [51] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *CVPR*, pages 4490–4499, 2018. [2](#)