

---

# Density-aware Chamfer Distance as a Comprehensive Metric for Point Cloud Completion

---

Tong Wu<sup>1</sup>, Liang Pan<sup>2</sup>, Junzhe Zhang<sup>2,4</sup>, Tai Wang<sup>1,3</sup>, Ziwei Liu<sup>2</sup>, Dahua Lin<sup>1,3,5</sup>

<sup>1</sup>SenseTime-CUHK Joint Lab, The Chinese University of Hong Kong,

<sup>2</sup>S-Lab, Nanyang Technological University, <sup>3</sup>Shanghai AI Laboratory, <sup>4</sup>SenseTime Research,

<sup>5</sup>Centre of Perceptual and Interactive Intelligence

{wt020, wt019, dhlin}@ie.cuhk.edu.hk, junzhe001@e.ntu.edu.sg,  
{liang.pan, ziwei.liu}@ntu.edu.sg

## Abstract

Chamfer Distance (CD) and Earth Mover’s Distance (EMD) are two broadly adopted metrics for measuring the similarity between two point sets. However, CD is usually insensitive to mismatched local density, and EMD is usually dominated by global distribution while overlooks the fidelity of detailed structures. Besides, their unbounded value range induces a heavy influence from the outliers. These defects prevent them from providing a consistent evaluation. To tackle these problems, we propose a new similarity measure named **Density-aware Chamfer Distance** (DCD). It is derived from CD and benefits from several desirable properties: **1)** it can detect disparity of density distributions and is thus a more intensive measure of similarity compared to CD; **2)** it is stricter with detailed structures and significantly more computationally efficient than EMD; **3)** the bounded value range encourages a more stable and reasonable evaluation over the whole test set. We adopt DCD to evaluate the point cloud completion task, where experimental results show that DCD pays attention to both the overall structure and local geometric details and provides a more reliable evaluation even when CD and EMD contradict each other. We can also use DCD as the training loss, which outperforms the same model trained with CD loss on all three metrics. In addition, we propose a novel point discriminator module that estimates the priority for another guided down-sampling step, and it achieves noticeable improvements under DCD together with competitive results for both CD and EMD. We hope our work could pave the way for a more comprehensive and practical point cloud similarity evaluation. Our code will be available at [https://github.com/wutong16/Density\\_aware\\_Chamfer\\_Distance](https://github.com/wutong16/Density_aware_Chamfer_Distance).

## 1 Introduction

Point cloud as one of the fundamental 3D representations is attracting increasing attention due to its efficiency, flexibility, and direct connection to real-world objects through 3D scanning devices. It has been employed in a wide range of application scenarios and studied for various tasks [1, 4, 7, 11, 28, 36, 39, 30, 24, 40, 42]. A proper similarity measure between two point clouds is always a crucial aspect for both guiding the training process and providing a fair and reasonable evaluation. However, this is a challenging design considering the unordered and irregular data form and varying point numbers.

Chamfer Distance (CD) and Earth Mover’s Distance (EMD) are two of the most universally acknowledged metrics in various point cloud tasks. CD is a nearest-neighbour-based method and benefits from its efficient computation and flexible applicability for point sets with different point numbers. EMD

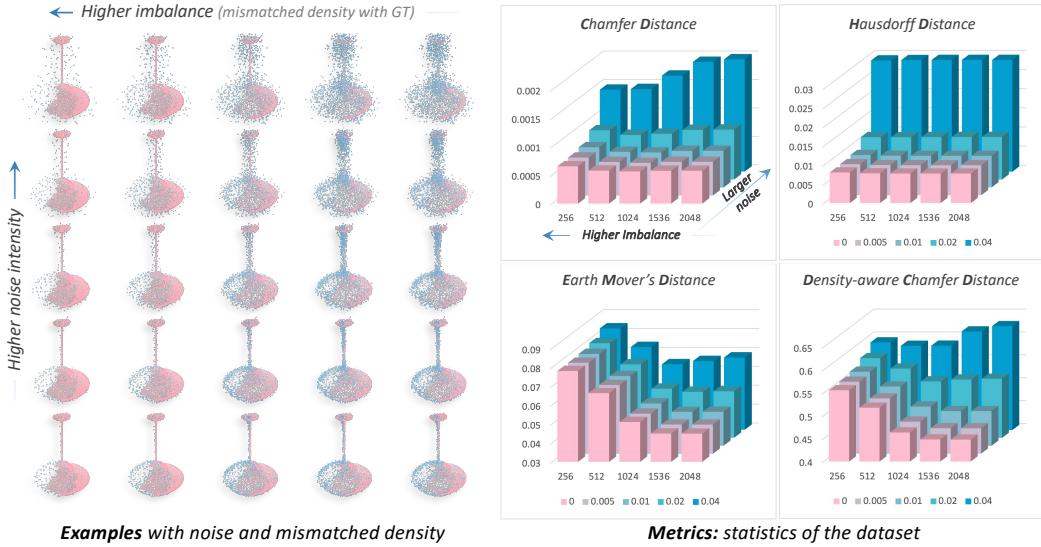


Figure 1: To generate the examples on the left, a complete yet noisy shape (blue) with  $n$  points is combined with a partial yet clean shape (pink) with 2048 points and then down-sampled back to 2048 points via FPS. Four point set distances are calculated between the generated shape and the ground truth. The results are averaged over the whole dataset to get the 3D histograms on the right: CD and HD are not sensitive to the mismatched density while highly influenced by noise when the intensity achieves a certain level; EMD and DCD share a similar pattern as  $n$  changes, but DCD is more sensitive to the noise which also represents detailed structures.

relies on solving an optimization problem to find the least expensive one-to-one transportation flow between two point sets. Although sometimes considered to be more faithful to visual quality than CD [1, 14], it is significantly more computationally expensive. The desirable properties of a good similarity metric can be different across application scenarios. 1) For perception and registration, it should be robust to sampling strategies and noisy points on detecting the similarity between the continuous surfaces represented by the discrete points. 2) For modelling and generation, it should be stricter with the quality of local point distributions, which is also crucial for visual quality. We mainly focus on the second aspect in this paper.

Take a closer look at the two metrics above. The formulation of CD sometimes suffers from its intrinsic deficiencies: being insensitive to different density distributions while significantly influenced by outliers [25]. In particular, we specialize one concept closely related to visual quality, namely *a matched density distribution between two point sets*. This can also be denoted as 'balance', assuming the ground truth is uniformly distributed. We visualize the issue in Fig. 1 through the task of point cloud completion: 1) under a low noise level, CD hardly changes with the imbalance ratio; 2) it increases dramatically once the noise intensity rises to some extent; 3) under a high noise level, higher imbalance even results in a lower CD, which leads to a tricky operation to reduce CD for the task: increasing the point density in the seen area with assured accuracy while reducing it in the unseen area to lower down the risk from abnormal points. But severe imbalance also significantly affects the global appearance. On the contrary, EMD can steadily detect the change of distribution in Fig. 1. However, the requirement for one-to-one mapping is usually over harsh. Consequently, the optima of the transportation problem is dominated by the global distribution while ignoring the fine-structured local details [5], as to be described in Sec. 3. Therefore, both CD and EMD are not ideally suitable for evaluating the quality of the generated shapes.

In this work, we mainly focus on CD's injustice as an evaluation metric and propose a new similarity measure named **D**ensity-**A**ware **C**hamfer **D**istance (DCD) to tackle the challenges above. Specifically, DCD is derived from the original CD, while it benefits from a higher sensitivity to distribution quality through a fraction term of query frequency and a higher tolerance to outliers through an approximation of Taylor Expansion. It shares a similar trend with EMD under a varying point distribution while being more computationally efficient and better at capturing the details. Moreover, due to their different focus, CD and EMD often encounter divergence when evaluating different methods, which makes them less reliable as consistent metrics, as to be demonstrated in Sec. 5. We empirically observe that DCD usually provides a more consistent and reliable evaluation, especially when the

Table 1: Comparison of properties among different metrics.

Metrics	assignment	efficient	bounded	density-aware	detail-aware
CD	nearest neighbour	✓	✗	✗	✓
EMD	optimization	✗	✗	✓	✗
DCD	nearest neighbour	✓	✓	✓	✓

results of CD and EMD contradict each other. Note that the proposed metric is also beneficial at dealing with ground truth with a non-uniform distribution such as curvature-based sampling, as to be discussed in the supplementary material.

Furthermore, we analyze the property of DCD as a loss function and compare it with L1 and L2 versions of CD. We introduce minor adjustments to it that are critical for a better training process. We then propose to make better use of the information from the query frequency mentioned above and design an MLP-based point discriminator inspired by the recent success of implicit functions [19, 15, 21]. It is further integrated to our balanced design for a two-stage completion framework, where the output of the module can be viewed as the “importance” of each point and serve as priority for a following step of down-sampling. Finally, the guided down-sampling operation benefits from removing outliers and preserving the critical points.

Extensive investigations are provided for the comparison among different metrics and methods for the task of point cloud completion. Experimental results validate that the proposed metric, Density-aware Chamfer Distance, successfully overcomes the aforementioned issues of CD. DCD can provide a more reliable evaluation when CD and EMD contradict each other, and it is proved to be more faithful to visual quality in Sec. 5 and a user study in the supplementary material. We validate its capacity as a loss function on PCN [39] and VRCNet [18], showing that it not only helps reduce DCD itself but also significantly lowers down the EMD metric and surprisingly reduces CD as well compared with the network trained with CD. Our proposed balanced design also gains noticeable improvement under the new metric, competitive results for both CD and EMD, and superior visual quality in the experiments.

## 2 Related Works

**Point Cloud Completion.** Point cloud completion aims to recover a complete shape based on a partial observation. Earlier works represent shapes with voxels [3, 6, 23], while PCN [39] first proposes to use raw point data and leverages an encoder-decoder structure to generate a global-feature-based coarse shape followed by the folding-based up-sampling [37]. Following works enhance the feature representation by techniques like attention mechanisms [31, 18, 17], hierarchical aggregation [8, 41], and grid structure for cubic feature sampling [35], etc.; the decoding process can also be examined to leverage, for example, a tree-structure [26], iterative refinement [29, 34], multiple patches generation [14], or separated prediction for the seen and unseen [41], etc. These works use either CD or EMD for evaluation, yet the two may not be consistently satisfied due to their different concerns, as to be shown in Sec. 5, and thus a more comprehensive metric is essential for a fair and reliable comparison.

**Point Cloud Distance.** The term “distance” refers to a non-negative function that measures the dissimilarity between two point sets. Considering the unordered structure of point clouds, the shape-level distance usually comes from statistics of pair-wise point-level distances based on certain assignment strategy. Chamfer Distance (CD) is one of the most widely used metrics based on the nearest neighbour (Eqn. 1). There are variants of it used for training [37, 4] and some other similarly formed distances like Hausdorff [9, 2, 32]. Another generally adopted metric is Earth Mover’s Distance (EMD), which relies on solving an optimization problem to find the best mapping function from one set to the other. It is sometimes considered to be more rational than CD [1, 14], but is much more computationally expensive. Recently, Urbach *et al* [27] propose DPDist, which compares point clouds by measuring the distance between the surfaces that they were sampled on. However, it is estimated by a network rather than a mathematical formulation, making it inconvenient and potentially unstable to be adopted into various tasks. Another closely related work is by Nguyen *et al* [16], who propose the sliced Wasserstein distance that has equivalent properties as EMD and similar computational complexity to CD, with the Monte Carlo method involved for approximation. In comparison, we start from the opposite point of view by deriving a new formulation based on CD

and result in a clean and explicit expression. Similarly, our DCD also shares the properties of EMD in many cases (Fig. 1), while it detects the detail preserving issue better than EMD.

### 3 Density-aware Chamfer Distance for Point Sets

#### 3.1 Preliminaries

Chamfer Distance between two point sets  $S_1$  and  $S_2$  is defined as:

$$d_{CD}(S_1, S_2) = \frac{1}{|S_1|} \sum_{x \in S_1} \min_{y \in S_2} \|x - y\|_2 + \frac{1}{|S_2|} \sum_{y \in S_2} \min_{x \in S_1} \|y - x\|_2. \quad (1)$$

Each point  $x \in S_1$  finds its nearest neighbour in  $S_2$  and vice versa; all the point-level pair-wise distances are averaged to produce the shape-level distance. The simple and flexible formulation generalizes well across many tasks. Earth Mover’s Distance is defined as:

$$d_{EMD}(S_1, S_2) = \min_{\phi: S_1 \rightarrow S_2} \sum_{x \in S_1} \|x - \phi(x)\|_2. \quad (2)$$

It relies on solving an optimization problem that finds a one-to-one bijection mapping  $\phi : S_1 \rightarrow S_2$ , thus only applicable when  $|S_1| = |S_2|$ . The pair-wise distances are then calculated between  $x$  and  $\phi(x)$ . As computing the optimal mapping is computationally expensive and even hardly affordable, several approximation schemes [14, 13] have been developed to relieve the computation burden.

#### 3.2 Density-aware Chamfer Distance

**Formulation and Interpretation.** As discussed in Sec. 1, CD is not a comprehensive metric for evaluating visual quality for the generation tasks, e.g., point cloud completion. We explain this from its formulation: **1)** the square operation makes it intensively influenced by outliers, and the evaluation results have a huge varying range across the dataset; **2)** the nearest point query operation makes it less sensitive to the issue of mismatched density distribution, and hence resulting in a less discriminative evaluation of visual quality. Therefore, we aim to propose a new metric based on the original formulation that not only preserves the capacity of similarity measure but also highly alleviates the problems above, namely Density-aware Chamfer Distance (DCD).

Firstly, CD grows quadratically with point pair distances, which can be dominated by the worst cases and hence overlooking the others. To address this problem, we introduce the first order approximation of Taylor Expansion  $e^z = \sum_{n=0}^{\infty} z^n / n!$ , i.e.  $e^z \approx 1 + z$  where  $z = -\|x - y\|_2$ . Thanks to the nearest neighbour assignment, the condition that  $z \approx 0$  is usually satisfied and the approximation is reasonable. Thus we have:

$$d_{CD}(S_1, S_2) \approx \frac{1}{|S_1|} \sum_{x \in S_1} \min_{y \in S_2} (1 - e^{-\|x - y\|_2}) + \frac{1}{|S_2|} \sum_{y \in S_2} \min_{x \in S_1} (1 - e^{-\|y - x\|_2}). \quad (3)$$

Considering the property of the exponential function  $e^z (z < 0)$ , each point-level distance is mapped to a value between  $[0, 1]$ . As a result, the formulation also sets a natural boundary of  $[0, 1]$  for the overall shape distance. The approximation would be less accurate as the point-level distance gets away from zero, yet it exactly helps to mitigate the over-sensitivity of the outliers by suppressing the square growth. We add another scale factor  $\alpha$  as in Eqn. 4 to adjust the sensitivity. The absolute distance value depends on this factor, and we reveal its relative consistency across different choices of  $\alpha$  in the supplementary material, and we fix  $\alpha = 1000$  for evaluation in this paper.

Secondly, the ambiguity of CD is partially due to its “blindness” that each point only considers its nearest neighbour in the other set while ignoring the surroundings. We denote the nearest neighbour assignment process by “query” for simplicity and present a simple example here: assume  $S_1$  to be a uniform point cloud and  $S_2$  an in-homogeneous one; consider two points  $y_1, y_2 \in S_2$  with  $y_1$  located at a sparse area and  $y_2$  in a relatively dense area. The calculation of  $d_{CD}(S_1, S_2)$  would likely get  $y_1$  frequently queried by points in  $S_1$  due to the local sparsity, and the case for  $y_2$  can be exactly the opposite. Thus we denote that  $y_1$  and  $y_2$  are not equally critical in representing the shape. Furthermore, assume a subset  $S_1^y \subseteq S_1$  so that each point in  $S_1^y$  queries  $y$  and that  $n_y = |S_1^y|$ , points in this set are unaware of each other under the formulation of CD, and the contribution of  $y$  to each point in  $S_1^y$  is not affected as  $n_y$  gets larger, which is unreasonable.

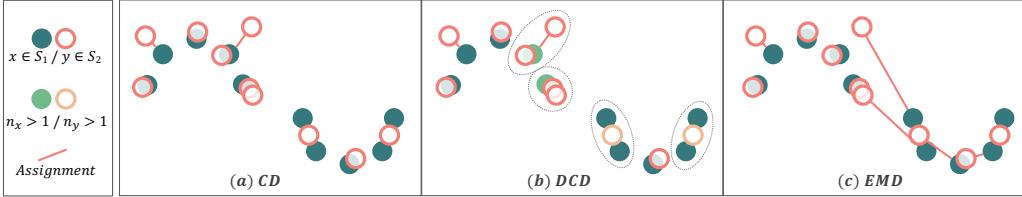


Figure 2: Comparison of assignment strategies and distance calculation. CD and DCD take the nearest neighbour locally, and DCD further considers the point-specific query frequency; EMD forces a one-to-one mapping, and the assigned pair of points may locate far from each other with weaker physical meaning.

Alternatively, we denote that the overall contribution of each point to the evaluation system shall be normalized, and we introduce  $1/n_y$  to deal with the case where  $y$  is being shared by multiple  $x$ s:

$$d_{DCD}(S_1, S_2) = \frac{1}{2} \left( \frac{1}{|S_1|} \sum_{x \in S_1} \left( 1 - \frac{1}{n_{\hat{y}}} e^{-\alpha ||x - \hat{y}||_2} \right) + \frac{1}{|S_2|} \sum_{y \in S_2} \left( 1 - \frac{1}{n_{\hat{x}}} e^{-\alpha ||y - \hat{x}||_2} \right) \right), \quad (4)$$

where  $\hat{y} = \min_{y \in S_2} ||x - y||_2$ ,  $\hat{x} = \min_{x \in S_1} ||y - x||_2$ , and  $\alpha$  denotes a temperature scalar. Finally, consider the first term without loss of generality, each  $y$  contributes  $| - \frac{1}{n_y} \sum_{x \in S_1^y} e^{-||x - y||_2} | \in [0, 1]$  to the overall distance metric (before averaging). A variants for DCD to deal with point sets with different number of points will be discussed in the supplementary material.

**Comparison among the Distance Metrics.** We conduct a brief comparison of the properties among the three metrics in Table 1, sketching their assignment strategies and computation schemes in Fig. 2. As discussed above, both CD and DCD assign point pairs by querying the nearest neighbours, while DCD further considers the point-specific query frequency  $n_y$ , and incorporates the property of density distribution into the measurement. In comparison, EMD naturally forces an equal query frequency of 1 for each point via the mapping function, and the metric is highly sensitive to the global point distribution. However, the harsh constraint not only imposes a significant increase in computational cost but also tends to sacrifice its attention to visual quality for optimal mapping. As shown in Fig. 2, the assigned pair of points could be located far from each other, and the distance is less physically meaningful. Experimental results in Sec. 5 would also show that an over compacted and smoothed shape can be favored by EMD despite its loss of detailed structure.

In brief, DCD takes a step from CD and attempts to provide a rationale bridge towards EMD for a better sense of point distribution rather than being blinded by its nearest neighbour. Compared with EMD, it is not only more efficient but also stricter with local structures. A balanced distribution and good preservation of detailed structures are both important factors for the visual quality of the completion result. More examples are to be shown in Sec. 5 for better illustration.

### 3.3 Application as an Objective Function

**Gradient Analysis and Comparisons.** Besides its usage as an evaluation metric, DCD is also expected to serve as an objective function to guide the training process. Considering the same nearest neighbour assignment, our analysis below will mainly focus on comparing CD (denoted by CD-T) and its L1 version (denoted by CD-P), which is also widely adopted for training. The exponential formulation in DCD takes in the L2 distance and modifies the gradient curve while bounding the loss value between  $[0, 1]$  at the same time.

We visualize the loss value and gradient curves of the three in Fig. 3 (a and b). The gradient by CD-T grows linearly with the distance  $l$  of a single point pair: it becomes rather small when  $l$  is close to zero and performs a heavy punishment on points with a large  $l$ . The CD-P loss, on the contrary, produces a constant gradient. For the DCD loss, it would first rise and then approach 1 as  $l$  increases, and its gradient can be calculated as:  $\delta \hat{d}_{DCD}(l) / \delta l = 2\alpha l e^{-\alpha l^2} / n$ , where  $\hat{d}_{DCD}(l)$  denotes the contribution of one point pair with a distance of  $l$ , and  $n = n_{\hat{y}}$ . The gradient first rises and then falls to zero, indicating that the loss is only effective to point pairs whose distance lies within a certain range. Similar to CD-T, it becomes small when  $l$  reduces to zero, which is more reasonable than the

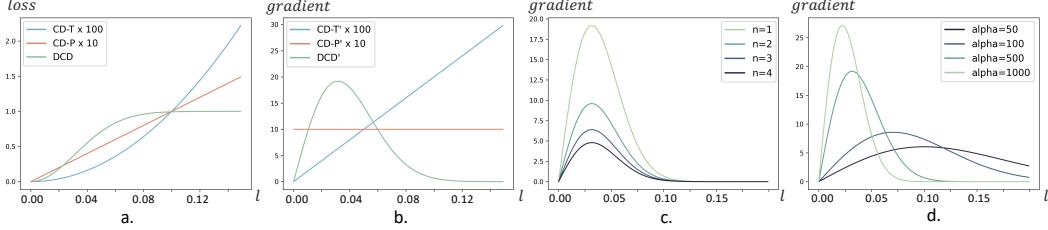


Figure 3: **a** presents the loss curves of CD-T, CD-P, and DCD and **b** presents the gradient curves of them; **c** and **d** visualizes the gradient for DCD with different  $n$  and  $\alpha$ , respectively.

constant gradient by CD-P, while it is also bounded by a maximum and will not be so harsh to the points with an extremely large  $l$  and thus stabilizes the training process.

**Adjustment on DCD for Training.** There are two characteristic components that decide the property of the curve, namely the hyper-parameter  $\alpha$  and the query frequency  $n$ . As shown in Fig. 3(c), a larger  $\alpha$  promotes a higher peak and a smaller range for non-zero gradient, and we need to set a proper  $\alpha$  during training based on the practical distance distribution. We empirically find that setting  $\alpha \in [40, 100]$  best promotes the training. A larger  $n$  indicates a more serious mismatching of density in a local region and leads to a higher loss value. When viewing both  $n$  and  $l$  as variables, they are encouraged to be smaller to lower down the loss, but  $n$  is not differentiable and cannot be directly changed by backpropagation. Meanwhile, a higher  $n$  reduces the gradient to  $l$  linearly (Fig. 3(d)), which prevents one point from moving towards an over-dense region, while it may also block the training process when the gradient to  $l$  gets too small. As a result, we introduce another hyper-parameter  $\lambda \in [0, 1]$  and replace  $n_{\hat{x}}$  and  $n_{\hat{y}}$  in Eqn. 4 with  $n_{\hat{x}}^\lambda$  and  $n_{\hat{y}}^\lambda$  to get the practical loss function for training.

## 4 Incorporating Balanced Design in Point Cloud Completion

### 4.1 Point Cloud Completion Framework

We adopt a typical two-stage coarse-to-fine completion pipeline [39, 14, 29, 18] (Fig. 4(d)). In the first stage, we extract a global feature from the partial observation and generate a complete yet coarse point cloud; in the second stage, we then introduce local features with abundant geometry information and obtain the final output with more details, higher visual quality, and high fidelity to the input. More details of the network architecture and training loss are included in the supplementary material.

We would like to highlight two observations regarding the point distributions here. First, a mean shape usually exists for each category, and there is an obvious imbalance of density for different regions according to how commonly they are shared across the dataset (Fig. 4(a)); second, the evaluation results by CD does not perfectly align with human assessment: the trick that place more points in the seen region with high confidence usually boosts the CD performance while it introduces imbalanced distribution and hurts the benign global distribution at the same time.

To address the problems above, we propose a simple yet effective method based on a current SoTA approach [18] that significantly boosts both the qualitative and quantitative results. Specifically, we tried our new metric as the objective function to replace  $L_{CD}$ ; we then propose a novel point discriminator which is supervised by a carefully designed density-aware signal and estimates the importance of each point; the output from the point discriminator is used as a priority for a final guided down-sampling process that helps remove outliers and maintain a balanced distribution.

### 4.2 Point Discriminator

The generated points in a shape are usually not equally important. We can roughly group them into three categories with the help of query frequency  $n$  introduced in Sec. 3: 1) when  $n > 0$ , the points with a higher  $n$  often lie in a sparse area and play a critical role in representing the shape; 2) when  $n = 0$ , the point can lie on the ground truth surface while in an over-populated region and become relatively unimportant; 3) another case with  $n = 0$  is that the point locates far from the underlying surface and can be regarded as an outlier that hurts the overall appearance. There are

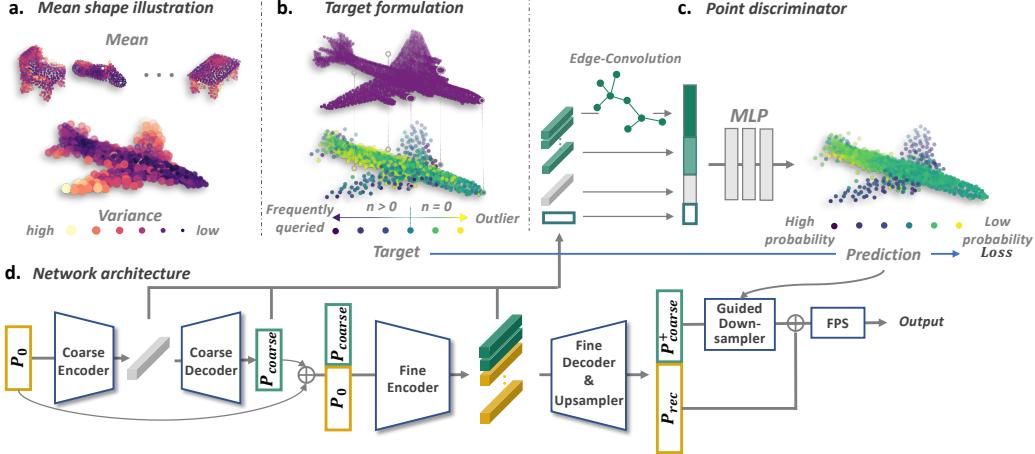


Figure 4: **a.** visualization of mean shape over the dataset; **b.** visualization of  $g(x)$  in Eqn. 5 at instance level; **c.** module structure of the point discriminator; **d.** two-stage framework with guided down-sampling.

no hard boundaries among these cases above, but we carefully design a formulation  $g(x)$  in Eqn. 5 to map them into a proper value range for learning. Specifically, suppose we train the module with the ground truth  $P_{gt}$  and the coarse output  $P_{coarse}$ , when  $n > 0$ , we adopt a logarithmic function for  $n$  and introduce the expected frequency  $|P_{gt}|/|P_{coarse}|$  as the denominator; we add a bias of 1 to ensure  $g(0) = 0$  at the boundary and finally invert the sign to distinguish from the following cases. When  $n = 0$ , we denote that the distance from a point  $x$  to  $P_{gt}$  together with a scaling factor  $t$  can serve as a suitable indicator to distinguish between case 2) and 3). To this end, we formulate a conjoint objective  $g(x)$  as below:

$$g(x) = \begin{cases} \min_{y \in P_{gt}} \|x - y\|_2 \cdot t & n_x = 0, \\ -\log_2\left(\frac{|P_{coarse}|}{|P_{gt}|} n_x + 1\right) & n_x > 0. \end{cases} \quad (5)$$

Encouraged by the recent success of implicit functions [19, 15, 21], we train an MLP module  $h$  along with the main network that learns to predict the target function  $g(x)$ , namely point discriminator. Specifically, it takes the local features of a point and its neighbours for edge convolution; the results are concatenated with the local feature  $f_x^l$ , global feature  $f_x^g$ , and point coordinates  $c_x$ , denoted by  $z = h(f_x^l, f_x^g, f_x^l, c_x)$ , where  $x' \in \mathcal{N}_x$  and  $\mathcal{N}_x$  denotes the neighbours for point  $x$ . It outputs a single scalar  $z$  which is supervised by the target function  $g(x)$  with a regression loss:

$$L_h = \frac{1}{|P_{coarse}|} \sum_{x \in P_{coarse}} \|z - g(x)\|_2. \quad (6)$$

The discriminator would be used for deciding point sampling privilege at inference time.

**Guided Down-sampling.** In view of the imbalanced population issue discussed above, an intuitive idea to alleviate it would be to allow a larger number of points generated via up-sampling before reducing to the desired number via Furthest Point Sampling (FPS). However, FPS tends to select more points from outer regions, which increases the risk of including more outliers, and thus this operation usually results in a noisier output. As a result, we aim to take advantage of the point discriminator above for the down-sampling stage at inference time together with FPS. Specifically, for each point  $x \in P_{coarse}$ , we define a point-wise existing probability by:

$$p(x) = \sigma(-\beta \cdot z - \gamma) = \frac{1}{1 + e^{(\beta \cdot z + \gamma)}}. \quad (7)$$

Let  $s$  denote the scale for up-sampling, then  $x$  is corresponding to  $s$  points in the  $P_{coarse}^+$  (up-sampled from  $P_{coarse}$ ) and they all share the same  $p(x) \in (0, 1)$ , an independent probability of each point being sampled.  $P_{coarse}^+$  is first down-sampled in this manner and then combined with  $P_{rec}$  (up-sampled from  $P_0$ ) before the final FPS, which ensures that we remain a pre-defined number of points. This process leverages the prior learned during training to effectively reduce the unreasonably located points at inference time (Fig. 7), relieving the side effect brought by FPS operation.

Table 2: Point cloud completion results in terms of CD  $\times 10^4$ , EMD  $\times 10^2$ , and DCD, lower is better.

Methods	Metrics	airplane	cabinet	car	chair	lamp	sofa	table	watrcraft	bed	bench	bookshelf	bus	guitar	motorbike	pistol	skateboard	Avg.
		CD	EMD	CD	EMD	CD	EMD	CD	EMD	CD	EMD	CD	EMD	CD	EMD	CD	EMD	
PCN	CD	4.50	8.83	6.41	13.01	21.33	9.90	12.86	9.46	20.00	10.26	14.63	4.94	1.73	6.17	5.84	5.76	9.78
	EMD	4.70	7.99	5.75	6.90	11.99	5.32	6.60	5.40	9.84	4.85	7.87	5.24	10.56	4.93	4.86	5.59	6.80
	DCD	0.478	0.519	0.490	0.617	0.710	0.552	0.559	0.580	0.662	0.562	0.608	0.429	0.446	0.548	0.491	0.445	0.553
PCN++	CD	4.06	9.08	6.64	13.11	19.25	9.78	14.36	9.66	22.33	9.73	15.51	5.13	1.86	6.25	5.81	4.99	10.29
	EMD	3.44	3.75	<b>3.15</b>	4.65	8.00	3.56	4.69	4.22	6.13	3.85	4.39	<b>2.62</b>	2.78	3.60	3.71	3.07	4.27
	DCD	0.428	0.464	0.451	0.574	0.661	0.504	0.517	0.540	0.617	0.524	0.563	0.389	0.369	0.527	0.447	0.393	0.508
TopNet	CD	4.12	9.84	7.44	13.26	18.64	10.77	12.95	8.98	19.99	9.21	16.06	5.47	2.36	7.06	7.04	4.68	10.30
	EMD	4.89	6.30	4.07	7.01	10.75	6.47	7.50	4.68	8.09	6.27	6.80	3.50	4.21	4.26	6.02	3.49	6.18
	DCD	0.536	0.558	0.548	0.650	0.711	0.598	0.599	0.600	0.678	0.588	0.622	0.492	0.487	0.572	0.542	0.496	0.598
MSN	CD	2.73	8.92	6.50	10.75	13.37	9.26	10.17	7.70	17.27	6.64	12.10	5.21	1.37	4.59	4.62	3.38	7.99
	EMD	2.75	4.02	3.47	4.44	6.28	3.74	4.46	3.82	5.27	3.34	4.28	2.92	<b>2.07</b>	3.30	3.62	2.21	3.94
	DCD	0.404	0.509	0.516	0.537	0.539	0.532	0.498	0.515	0.574	0.471	0.541	0.458	0.388	0.491	0.463	0.422	0.499
VRC	CD	2.20	<b>7.92</b>	5.60	<b>7.49</b>	<b>8.15</b>	<b>7.45</b>	<b>7.52</b>	<b>5.20</b>	<b>11.90</b>	<b>4.88</b>	<b>7.39</b>	4.53	<b>1.15</b>	3.90	3.44	3.22	<b>6.09</b>
	EMD	3.03	7.57	6.14	5.49	6.15	5.80	4.65	4.97	6.58	3.45	5.28	6.59	3.08	4.45	4.56	3.20	5.27
	DCD	0.374	0.509	0.499	0.488	0.475	0.515	0.438	0.478	0.527	0.401	0.470	0.462	0.349	0.452	0.443	0.363	0.462
VRC-EMD	CD	2.72	9.03	6.58	9.93	11.53	9.38	9.80	6.71	17.22	6.88	10.34	5.32	1.39	4.47	4.62	4.69	7.87
	EMD	2.50	<b>3.65</b>	3.23	<b>4.15</b>	5.31	<b>3.61</b>	3.93	3.58	<b>5.17</b>	3.19	3.97	2.69	2.08	3.06	3.48	<b>2.29</b>	<b>3.62</b>
	DCD	0.369	0.483	0.473	0.502	0.499	0.509	0.450	0.478	0.547	0.423	0.487	0.424	0.349	0.445	0.430	0.370	0.461
Our	CD	<b>2.22</b>	8.00	<b>5.41</b>	7.88	8.28	7.94	8.89	5.46	14.76	5.78	9.37	<b>4.44</b>	1.30	<b>3.59</b>	<b>3.43</b>	<b>2.39</b>	6.51
	EMD	<b>2.29</b>	4.43	3.46	<b>3.92</b>	<b>4.98</b>	3.98	<b>3.89</b>	<b>3.51</b>	5.34	<b>3.13</b>	<b>3.91</b>	3.29	2.21	<b>3.02</b>	<b>3.38</b>	2.39	3.67
	DCD	0.335	<b>0.447</b>	0.427	<b>0.451</b>	<b>0.445</b>	0.469	0.423	0.426	<b>0.504</b>	0.399	0.453	<b>0.382</b>	<b>0.336</b>	<b>0.401</b>	0.365	<b>0.345</b>	<b>0.420</b>

## 5 Experiments

**Dataset.** We use the recently proposed MVP Dataset [18] for our study and experiments. It is a multi-view partial point cloud dataset covering 16 categories with 62,400 and 41,600 pairs for training and testing, respectively. It renders the partial 3D shapes from 26 uniformly distributed camera poses for each 3D CAD model selected from ShapeNet [33], and the ground truth point cloud is sampled via Poisson Disk Sampling (PDS).

**Comparison Methods and Metrics.** We include the following methods for comparison in the main experiments: PCN [39], TopNet [26], MSN [14], and VRCNet [18] (with the PSK module discarded to improve efficiency while slightly sacrificing performance). PCN++ is a simple extension of PCN [39] that generates the double number of points for training and down-sampled to the required point number at inference time. We report per-class results on CD, EMD, and DCD for a clear comparison across methods and a view of the different properties of these metrics.

**Implementation Details.** All the models are trained using the Adam optimizer [10] with the learning rate initialized at  $1e^{-4}$  and decayed by 0.7 every 40 epochs. We use a batch size of 32 and a total epoch of 80. We set  $\alpha = 1000$  for the evaluation of DCD, and  $\alpha \in [40, 100]$  for training. We set  $\lambda \in [0, 0.5]$  and  $\beta = 9, \gamma = 1$  for our approach in the main experiments. Our work is implemented with PyTorch and is run on a Tesla V100 GPU.

### 5.1 Comparison of the Methods

The main experimental results for point cloud completion on MVP dataset are reported in Table 2. Note that most networks are trained with the CD loss except when specified. Early works, PCN [39] and TopNet [26] have relatively high loss for all the evaluated metrics, CD, EMD, and DCD. PCN with an additional up-sampling step, namely PCN++, notably lowers down EMD and DCD, while marginally raising CD for sacrifice; MSN [14] benefits from the multi-surface design and obtains low EMD and DCD results, while its CD performance is not that satisfying. The previous SoTA method, VRCNet [18] outperforms the other methods under the metric of CD and DCD, but we observe that its EMD loss is surprisingly high. We further replace the CD loss with EMD when training VRCNet (denoted as VRC-EMD), and it achieves the lowest EMD yet obviously higher CD than the original version. In comparison, our method reports the lowest DCD, second-lowest CD, and comparable EMD with VRC-EMD. Qualitative results (Fig. 5) show that our results are apparently superior in both the global balanced point distributions and local structures. A user study in the supplementary material will further validate that our method benefits from a higher visual quality.

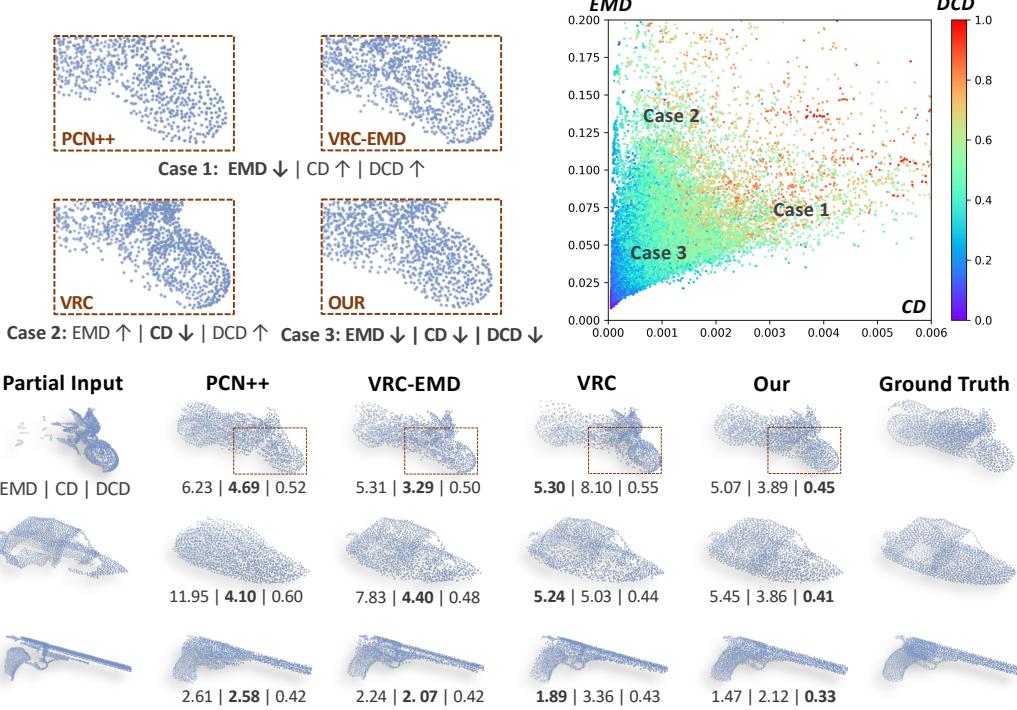


Figure 5: Comparison of CD, EMD, and DCD through examples and statistics. **Below:** examples from different methods that represent three typical cases: 1) EMD is low while CD is high; 2) CD is low while EMD is high; 3) both CD and EMD are low. DCD can only be lowered down in the third case. **Upper left:** A close shot at the cases above. **Upper right:** visualization of the three metrics on the test set, we can observe that 1) the positive correlation between CD and EMD is weak and the scatter points form a fan-shaped area, and 2) DCD (denoted by the point color) becomes lower towards the original point where both CD and EMD are low.

## 5.2 Comparison of the Metrics

**Consistency and Reliability.** When taking a closer look at all the CD-EMD-DCD tuples (either for each category or for the averaged results) in Table 2, the results by different metrics do not exhibit a clear positive correlation. But we also observe that for those that have a similar value of CD, DCD is usually dominated by their EMD performance, and for those with similar EMD results, DCD is highly correlated with CD. Although the law is not strictly held or theoretically proved, the scatter plot by each instance from different methods in Fig. 5 (upper right) also supports our empirical observation. It reveals that since CD and EMD focus on different aspects of the point cloud, inconsistency and confusion may exist for the similarity measurement, which prevents either of them from being a comprehensive metric. In comparison, DCD reflects the behaviors of them both and could only be reduced when both CD and EMD are relatively low so that it serves as a more consistent, stable, and comprehensive metric, especially when CD and EMD encounter disparity. Our user study in the supplementary material further indicates that DCD is a more faithful metric to visual quality.

**Bounded Distance.** One advantage of DCD is its bounded value range, which promotes an equal consideration for all samples in the dataset at the evaluation stage. In contrast, CD and EMD are dominated by the worst cases in the dataset due to their unbounded nature, as shown in Fig. 6. Specifically, for CD, 80% of the normalized loss accumulation is contributed by only top 50% of the samples ranked by the shape-level distance, and 50% contributed by the top 25%; a similar case is also observed for EMD. This property makes them fail to comprehensively evaluate the quality of all samples in the dataset. Our DCD overcomes this problem by assuring a clear [0,1] boundary for the distance value, enabling more comprehensive and stable statistical results over the whole test set.

## 5.3 Ablation Study

We study the effectiveness of each component in our method separately, including DCD loss, the additional point up-sampling (PU), and the guided point down-sampling (GPD). As shown in Table 3,

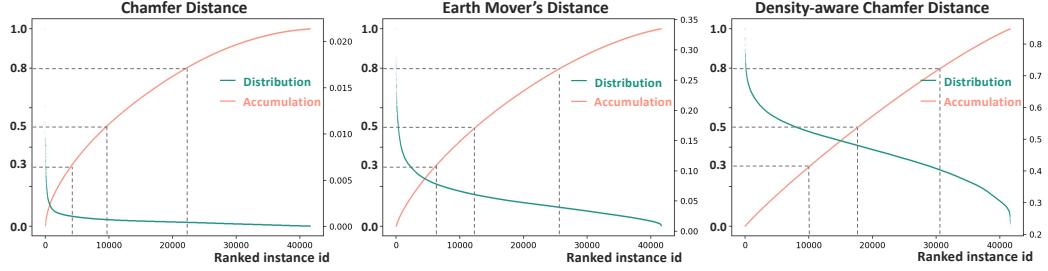


Figure 6: Distribution and accumulation (normalized to 1) of evaluation results per-shape in the test set.

Table 3: Ablation Study.

DCD	PU	GPD	Metric		
			CD	EMD	DCD
✓			6.09	5.27	0.462
	✓		<b>5.85</b>	5.14	0.457
✓	✓		6.91	3.78	0.425
✓	✓	✓	6.65	3.68	0.422
✓	✓	✓	6.51	<b>3.67</b>	<b>0.420</b>

Table 4: Evaluation results on three metrics when trained with each of them. \* denotes applying  $L_{EMD}$  on the final outputs while  $L_{CD}$  on intermediate ones.

Model	PCN [39]			VRC [18]		
	Metric	$L_{CD}$	$L_{EMD}^*$	$L_{DCD}$	$L_{CD}$	$L_{EMD}^*$
$CD$	9.78	10.70	<b>9.36</b>	6.09	7.87	<b>5.85</b>
$EMD$	6.80	<b>3.97</b>	4.71	5.27	<b>3.62</b>	5.14
$DCD$	0.553	0.537	<b>0.526</b>	0.462	0.461	<b>0.457</b>

leveraging  $L_{DCD}$  loss usually outperforms the same model trained with  $L_{CD}$  in all the three metrics; the additional point up-sampling (PU) significantly reduces EMD while it increases CD at the same time; after guided down-sampling (GPD) is applied to the model with both  $L_{DCD}$  and PU, we achieve the lowest EMD and DCD and a relatively low CD.

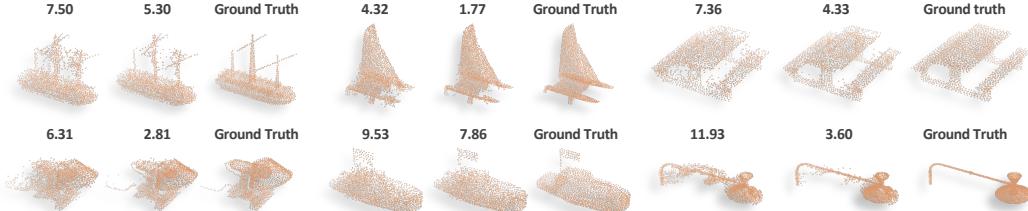


Figure 7: Examples of guided down sampling (from left to right) and evaluation in terms of  $CD \times 10^4$ .

#### 5.4 Performance as the Loss Function

We evaluate the effectiveness of DCD as a loss function by training networks with each of the three metrics (denoted by  $L_{CD}$ ,  $L_{EMD}$ , and  $L_{DCD}$ ) and evaluate on all of them (denoted by CD, EMD, and DCD). We conduct experiments on the baseline model PCN [39] and the SOTA VRCNet [18], as shown in Table 4. Compared with  $L_{CD}$ -trained networks, training with  $L_{EMD}$  produces the lowest EMD and a lower  $d_{DCD}$ , yet it suffers from non-negligible scarification of  $CD$  and much heavier computational burden, which is to be discussed in the supplementary material; training with  $L_{DCD}$  will produce the lowest DCD, significantly reduce EMD, and can even slightly reduce CD than the  $L_{CD}$ -trained models (especially on PCN). The time consumption is also comparable with  $L_{CD}$ , which further validates DCD’s convenience and superiority as an objective function. We set  $\alpha = 50$  or 100 here with  $\beta = 0$  and replace all the  $L_{CD}$  with  $L_{DCD}$ . Experiments on how  $\alpha$  and  $\beta$  affect the performance are included in the supplementary material.

## 6 Conclusion

In this work, we propose a new similarity measure for point clouds named Density-aware Chamfer Distance (DCD). It is bounded in value, effective in computation, and faithful to visual quality by considering both the density distribution and detailed structures. Our method achieves noticeable improvements under DCD and superior visual quality compared with previous works.

**Acknowledgements.** This study is supported in part by the NTU NAP, and under the RIE2020 Industry Alignment Fund – Industry Collaboration Projects (IAF-ICP) Funding Initiative, as well as cash and in-kind contribution from the industry partner(s). It is also supported in part by Centre for Perceptual and Interactive Intelligence Limited, in part by the GRF through the Research Grants Council of Hong Kong under Grants (Nos. 14208417, 14207319 and 14203518) and ITS/431/18FX, in part by CUHK Strategic Fund and CUHK Agreement TS1712093, in part by the Shanghai Committee of Science and Technology, China (Grant No. 20DZ1100800).

## A Further Analysis of Balanced Chamfer Distance

### A.1 The consistency and choice of hyper-parameter $\alpha$

We introduce a temperature scalar  $\alpha$  in Eqn. (4) so that  $e^{-\alpha \cdot \|x-y\|^2}$  can have a relatively wider varying range. We set  $\alpha = 1000$  in the paper and would briefly show why it is a proper value here. Notice that  $\|x-y\|^2$  for each nearest point pair is usually about  $10^{-4}$  or  $10^{-3}$ , thus setting  $\alpha = 1000$  maps it to around  $10^{-1}$  or  $10^0$  where the exponential term has a large gradient. As visualized in Fig. S8(a), either setting  $\alpha$  too large or too small would not result in an ideally shaped function. We also regenerate the DCD value matrix under the same settings as in Fig. 1 with different  $\alpha$  values, and we calculate the mean and variance accordingly, as shown in Fig. S8(b). A larger  $\alpha$  results in a higher mean which is reasonable according to Eqn. (4), and we observe the variance is at its largest with  $\alpha = 1000$ , which also aligns well with the theoretical analysis.

Furthermore, we track the PCN training loss calculated by DCD with different  $\alpha$ , as shown in Fig. S8(c), and we visualize the per-instance evaluation results with different  $\alpha$  for a well-trained model in Fig. S8(d). The statistical results show that the relative value and trend of DCD is relatively consistent with different data distributions when  $\alpha$  changes, while their absolute values are different. At evaluation time, we can use the same  $\alpha$  (e.g., 1000) for all the methods for a fair comparison.

### A.2 Dealing with mismatched point numbers.

We consider the case where the two point sets  $S_1$  and  $S_2$  do not have the same number of points, suppose  $|S_1| = \eta \cdot |S_2|$ ,  $\eta > 1$ . A naive extension of DCD (Eqn. (4) in Sec. 3) is presented as Eqn. 8, where we add  $\eta$  or  $1/\eta$  to indicate the one-to-many mapping in this case:

$$d_{DCD}(S_1, S_2) = \frac{1}{2|S_1|} \sum_{x \in S_1} \left( 1 - \frac{\eta}{n_{\hat{y}}} e^{-\alpha \|x-\hat{y}\|_2} \right) + \frac{1}{2|S_2|} \sum_{y \in S_2} \left( 1 - \frac{1}{\eta \cdot n_{\hat{x}}} e^{-\alpha \|y-\hat{x}\|_2} \right), \quad (8)$$

The formulation above usually works well in practice, but it may also lead to negative results in the first term when  $n_{\hat{y}} < \eta$  and  $\frac{\eta}{n_{\hat{y}}} e^{-\alpha \|x-\hat{y}\|_2} > 1$ . We thus propose another variant of DCD in Eqn. 9. Considering **the first term**, since  $|S_1| > |S_2|$  and each  $y \in S_2$  should naturally be assigned to more than one  $x \in S_1$ , the decaying term should not follow the tendency of  $1/n_{\hat{y}}$ , but rather updated to  $\max(\eta/n_{\hat{y}}, 1)$ . On the one hand, the contribution of  $\hat{y}$  would not be reduced before the querying frequency of it reaches  $\eta$ ; on the other hand, it should not exceed 1 either, which is important for keeping a non-negative result. As for **the second term**, each  $x \in S_1$  is corresponding for more than one point in  $S_2$ . We take  $\bar{\eta} = \text{ceiling}(\eta)$  and find  $\bar{\eta}$ -nearest neighbours for  $x$ , denoted by  $N(x)\bar{\eta}$ . And the overall formulation of the variant is:

$$d_{DCD-E}(S_1, S_2) = \frac{1}{2|S_1|} \sum_{x \in S_1} \left( 1 - \frac{1}{\max(\eta/n_{\hat{y}}, 1)} e^{-\alpha \|x-\hat{y}\|_2} \right) + \frac{1}{2|S_2|} \sum_{y \in S_2} \left( 1 - \frac{1}{\bar{\eta} \cdot n_{\hat{x}}} \sum_{\hat{x} \in N(y)\bar{\eta}} e^{-\alpha \|y-\hat{x}\|_2} \right). \quad (9)$$

We use Eqn. 8 in training for the loss between the coarse shape with 1024 points and the ground truth with 2048 points for simplicity.

### A.3 Time Complexity.

EMD relies on solving the linear assignment problem in an iterative approximation manner with a practical time complexity between  $O(n^2)$  [13] and  $O(n^3)$  [20] and usually  $O(n^2)$  memory footprints. We adopt an  $O(n)$  memory-efficient implementation by [14] and with the error rate  $\epsilon = 0.004$  and an iteration of 3000. The most computationally expensive part for CD and DCD is the nearest neighbour

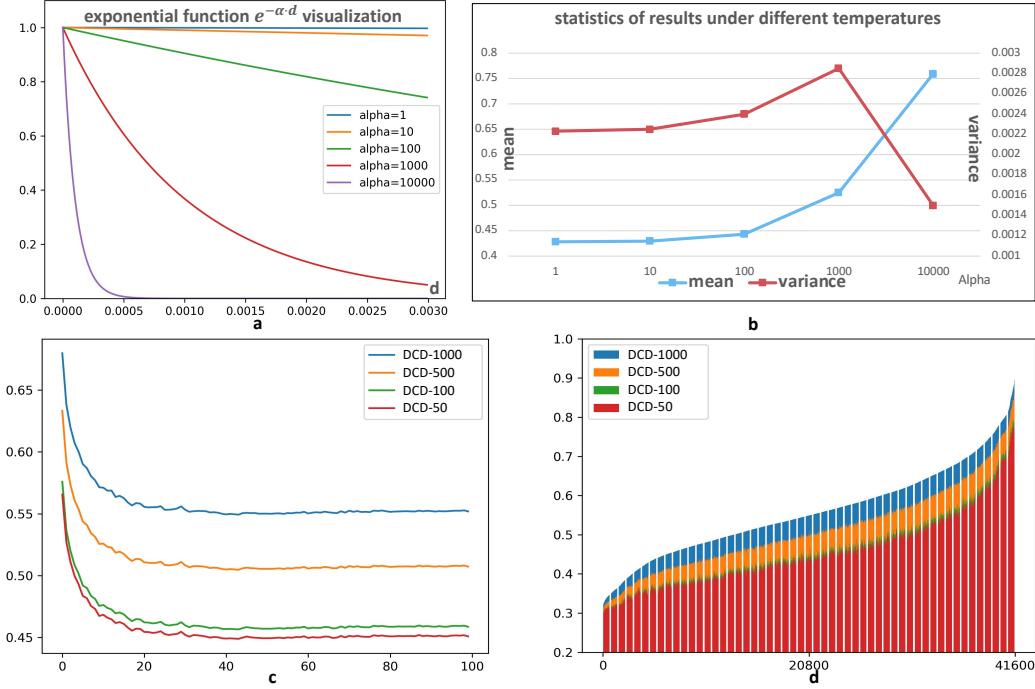


Figure S8: **a** visualization of the exponential function  $e^{-\alpha \cdot d}$  with a varying  $\alpha$ ; by empirically setting the value range of  $d$  to be similar with  $\|x - y\|^2$  for paired point-wise distance, we found that  $\alpha = 1000$  gains a desirable value range by the function. **b** visualization of the statistical results of DCD under different temperatures; a larger  $\alpha$  results in a larger mean value and  $\alpha = 1000$  gets the highest variance. **c** we track the loss with different  $\alpha$  when training a PCN network. **d** the per-instance evaluation with different  $\alpha$  for a well-trained network.

selection, which is usually  $O(n)$  for time complexity and can be accelerated by special data structures like KD-tree. We observe that both CD and DCD are significantly efficient to be computed compared with EMD; the running time of EMD is also affected by the distributions of the two point sets, which decides the difficulty of finding the optimal assignment. As shown in Table R5, with the same setting as in Fig. 1, a higher mismatched distribution results in obviously heavier time consumption, and the noise intensity also influences the results.

Table R5: Time consumption evaluation.

	256	512	1024	1536	2048
CD	0.006	0.007	0.008	0.010	0.012
EMD	0.362	0.327	0.267	0.241	0.239
DCD	0.013	0.013	0.013	0.013	0.013
	0	0.005	0.01	0.02	0.04
CD	0.008	0.008	0.008	0.008	0.008
EMD	0.272	0.271	0.269	0.264	0.256
DCD	0.013	0.013	0.013	0.013	0.013

#### A.4 Evaluation on other tasks and ground truth distributions.

Apart from the task of completion, DCD is also a suitable evaluation metric for tasks like point cloud upsampling or denoising, where the ground truth with desirable point distribution is provided, and the model is expected to generate high-quality point cloud outputs. Taking upsampling as an example, a desirable dense output is expected to be uniform, clean, and faithfully located on the underlying surface, and thus metrics like NUC [38] and uniform loss [12] were proposed to evaluate the distribution uniformity besides Chamfer Distance. However, these metrics usually make strong

assumptions that points in a small patch lie on a surface or that there should be an expected number of points in a ball anywhere with a certain radius. And the metrics are sometimes sensitive with the choice of hyper-parameters and the geometry itself. Moreover, they always encourage uniformity rather than the specific distribution of the ground truth, which also limits the application scope.

On the contrary, the density-aware DCD would focus on the faithfulness of the output to the ground truth distribution without any strong assumptions; it is not sensitive to the choice of hyper-parameters; it is beneficial at reflecting the mismatched density in local areas in scenarios where the ground truth is not uniformly distributed for specific purposes (*e.g.*, curvature-based sampling), as shown in Fig. S9. We also perform quantitative evaluation by applying a mixture of curvature-based sampling (for a ratio of  $R_c$ ) and the standard Poisson-Disk Sampling (PDS), while the output is noisy and basically uniform. When  $R_c$  changes, DCD and EMD can reflect the increasingly mismatched density, while CD and the uniform loss are not sensitive to it. (Table R6).

Table R6: Evaluation on ground truth with non-uniform sampling. We apply a mixture of curvature-based sampling (for a ratio of  $R_c$ ) and the standard PDS, while the output is noisy and basically uniform.  $CD(gt)$  denotes the averaged L2 distance from the ground truth points to their nearest neighbor, and the definition can be extended to  $CD(x)$ ,  $DCD(gt)$ , and  $DCD(x)$ . When  $R_c$  changes, DCD and EMD can reflect the increasingly mismatched density, while CD and uniform loss are not sensitive to it.

$R_c$	$CD(gt)$	$CD(x)$	$DCD(gt)$	$DCD(x)$	$EMD$	$L_{uni}(x)$
0%	1.62	2.06	3.97	3.54	2.13	19.96
25%	1.68	2.21	4.24	4.00	5.64	19.96
50%	1.68	2.41	4.62	4.34	7.25	19.96

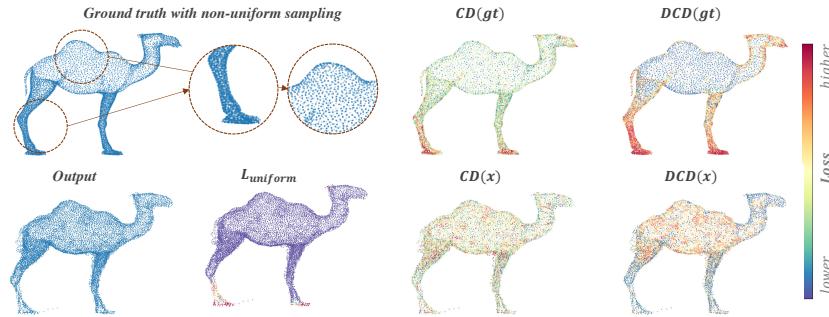


Figure S9: Visualization of the distance contributed by each point. DCD is better at reflecting the mismatched density in local areas between the two sets.

### A.5 Effect of hyper-parameters in $L_{DCD}$ .

We introduce two hyper-parameters in Sec. 3.3 that would affect the performance of DCD as a loss function, and we conduct experiments accordingly to explore the pattern. As shown in Table R7, a larger  $\alpha$  usually promotes a lower DCD, while the performance of CD worsens when  $\alpha$  is as large as 1000.  $\lambda = 0$  usually results in the best CD results yet with sub-optimal EMD and DCD, while  $\lambda = 1$  would obviously hurt CD;  $0 < \lambda < 1$  exhibits a trade-off among the metrics.

## B User Study on Visual Quality

In Sec. 5.1 and 5.2, we conduct comparisons on 1) the completion performance among different methods and 2) the characteristics of three metrics. We show that our method gains the best performance under the metric of DCD and that DCD is a more comprehensive measure through several examples. In this section, we further validate the conclusions above via user study. Specifically, we randomly select 25 partial inputs and select five methods (*i.e.*, PCN++ [39], MSN [14], VRC [18], VRC-EMD [18], and our method) to generate 125 completion outputs in total. We exclude the

Table R7: The effect of  $\alpha$  and  $\lambda$  when applying  $L_{DCD}$  on PCN [39].

$\alpha$	$\lambda$	CD	EMD	DCD
50	0.0	9.56	4.92	0.533
	0.5	9.86	4.68	0.529
	1.0	10.21	4.63	0.527
100	0.0	9.33	4.92	0.535
	0.5	9.85	4.68	0.526
	1.0	9.93	4.69	0.525
200	0.0	9.36	4.71	0.526
	0.5	9.82	4.59	0.520
	1.0	10.16	4.64	0.521
1000	0.0	10.14	4.72	0.516
	0.5	10.56	4.72	0.516
	1.0	11.12	4.96	0.519

Table R8: The average ratio of each method being considered to have produced the best results.

Methods	PCN++ [39]	MSN [14]	VRC [18]	VRC-EMD [18]	Ours
Average Ratio	0.280	0.123	0.286	0.151	<b>0.460</b>

original PCN [39] and TopNet [26] since their performance is not desirable on any of the three metrics; each of the five methods above is able to achieve favorable results under at least one metric, which provides a good playground to evaluate the quality of different metrics.

We invite 15 volunteers to complete the study. For each shape with five generated point clouds in a random order, the volunteers are asked to select a single or multiple (for at most two) items with the highest comprehensive score according to the following evaluation indicators: **1)** The similarity between each option and the ground truth, including the similarity of global shape and the fidelity of local details. **2)** The quality of the point cloud distribution, considering whether there exists a significant shift in the center of gravity and whether there exists obvious clustering or sparseness. Once we get the statistical results, we analyze the data from two points of view as follows.

**Comparison of the methods.** For the output of each method on each shape, we calculate the average ratio of it being selected as the favorite option by the volunteers (e.g., if three people out of fifteen select one method as the best for one shape, the ratio is 0.20). And then, the results are averaged for all the shapes, as shown in Table. R8. Our method outperforms the others by a large margin, indicating that it benefits from high visual quality.

**Comparison of the metrics.** We also leverage the data to evaluate how faithful each metric is to human vision. Specifically, according to each of the metrics, we can get the top-1 completion result for each shape; similarly, we collect the ratio of it being selected by volunteers as the best option and average the ratio over all the shapes. As shown in Table. R9, DCD gains the highest alignment with human vision.

## C Network architecture and training details.

### C.1 Review of a Typical Two-Stage Pipeline.

The two-stage coarse-to-fine completion pipeline was first proposed by PCN [39] and improved by a series of its following works [14, 29, 18]. The first stage takes a partial point cloud as input and extracts a global feature  $f^g$  by an MLP, e.g., PointNet [22], followed by a decoder to generate a coarse point cloud  $P_{coarse}$ . Though it maintains a relatively reasonable global shape,  $P_{coarse}$  usually fails to capture and depict the details. Therefore, the second stage follows up, which aims **1)** to precisely reconstruct the input point cloud without loss of details, denoted by  $P_{rec}$ , **2)** to improve the quality of  $P_{coarse}$  especially in the unseen part via coordinates adjusting, up-sampling, and probably

Table R9: The faithfulness of best-result-selection between each metric and human users.

Metrics	CD	EMD	DCD
Average Ratio	0.574	0.437	<b>0.623</b>

detail transferring, denoted by  $P_{coarse}^+$ , and finally 3) to obtain the final output with desirable visual quality and high fidelity to the original input. Point-wise local features  $f^l$  with abundant geometry information encoded are usually involved in this stage. We leverage the same network architecture as VRCNet [18] and adopt the official implementation from their public code <sup>1</sup>.

## C.2 Interpretation of the mean shape.

Although CD does not involve any hard assignment between prediction and ground truth point sets, it is observed that the output coordinates from each node of the last Fully Connected (FC) layer have a relatively convergent local distribution. By assuming a Gaussian-like distribution, we visualize their **mean** and **variance** for each category separately across the test set in Fig. 4, where the mean coordinates form the category-specific **mean shapes**, and the color of each node denotes the node variance <sup>2</sup>. This observation indicates that for the coarse shape, there exists an obvious imbalance of density across different regions according to how commonly they are shared across the dataset. This pattern occurs not only for the *statistical results* but also for *each shape instance* [1, 5].

Another thing to emphasize is the trade-off between accuracy and distribution balance. The first few points located in an unseen area significantly reduce the CD loss, while the marginal gain soon decays with more points predicted there. Considering the side effect that inaccurate points would bring in extra corruption to the overall CD loss, lying more points in the seen region with high confidence usually boosts the CD performance. This trick benefits the CD metric, while it violates the overall distribution and hurts the visual quality at the same time.

## C.3 Training Details

We introduce the regression loss  $L_h$  for training the point discriminator in Eqn. (6) (Sec. 4.2), and we would clarify the full training loss to train the two-stage framework here. Specifically, the loss function includes another two parts despite  $L_h$ :  $L_d$  involves multiple paired point cloud distances, and  $L_{KL}$  indicates the KL divergence for the dual-path VAE architecture following [18]. These two terms are formulated as follows:

$$L_d = \lambda_1 \cdot d_{CD}(P_{coarse}, P_{gt}) + \lambda_2 \cdot d_{CD}(P_{coarse}^+, P_{gt}) + \lambda_3 \cdot d_{CD}(P_{fine}, P_{gt}), \quad (10)$$

where  $P_{fine}$  denotes the final output after sampling, and  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda_3$  denote the loss weights. Note that we still use the Chamfer Distance as training loss in this paper, and the reason why would be partially explained in Sec. D. We thus have:

$$L_{KL} = -\lambda_{KL} \cdot (\text{KL}(q_\phi(f^g | P_{gt}), \mathcal{N}(0, I)) + \text{KL}(q_\psi(f^g | P_0), q_\phi(f^g | P_{gt}))), \quad (11)$$

where  $q$  denotes the encoder for latent distributions with network weights denoted by  $\phi$  and  $\psi$ , **KL** denotes the calculation of KL divergence with a loss weight of  $\lambda_{KL}$ . Finally, the overall loss function is formulated as:

$$L = L_h + L_d + L_{KL}. \quad (12)$$

We set  $\lambda_1 = 10$ ,  $\lambda_2 = 0.5$ ,  $\lambda_3 = 1$  and  $\lambda_{KL} = 20$  in the experiments.

When we adopt  $L_{DCD}$  during network training, we can simply replace all the occurrence of  $L_{CD}$  with  $L_{DCD}$ , while there are also tricks for better performance: 1) we can use different  $\alpha$  for different terms in Eqn. 10, e.g.,  $\alpha = 50$  for the first term and  $\alpha = 100$  for the others, which empirically works slightly better than using  $\alpha = 50$  or  $\alpha = 100$  for all the terms; 2) we can add another L1-version CD (CD-P) along with  $L_{DCD}$  for training, which is our implementation to gain the results reported in Table 2.

<sup>1</sup><https://github.com/paul007pl/VRCNet>

<sup>2</sup>The category information is not provided during training, while only used for statistics.

## D Limitations and Future Work

This approach still has some limitations that can be further explored in the future: we investigate the properties of the proposed metric on the task of point cloud completion, while it is actually applicable in many other tasks and scenarios as both evaluation metric and training loss. In the future, we will conduct more experiments to validate the generalization ability of DCD across different tasks.

## E Broader Impact

A comprehensive, reliable, and effective similarity measure is critical to point cloud analysis. It not only provides a fair comparison among different methods but also encourages the design of algorithms to take more critical factors into consideration, such as preserving accurate local details, keeping a uniform global distribution, and avoiding outliers. As shown in the paper, the broadly used Chamfer Distance and Earth Mover’s Distance usually encounter obvious disparity due to their different focus, making it hard to provide a consistent evaluation. It reveals the necessity and importance of formulating a more comprehensive metric to close the gap. We hope that the Balanced Chamfer Distance we propose in this paper can better serve the demands above than the existing metrics, so that hopefully it can encourage a more reasonable evaluation and influence method designs for tasks in point cloud analysis.

## References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pages 40–49. PMLR, 10–15 Jul 2018. [1](#), [2](#), [3](#), [16](#)
- [2] Xuelin Chen, Baoquan Chen, and Niloy J Mitra. Unpaired point cloud completion on real scans using adversarial training. *The International Conference on Learning Representations*, 2019. [3](#)
- [3] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5868–5877, 2017. [3](#)
- [4] Haowen Deng, Tolga Birdal, and Slobodan Ilic. 3d local features for direct pairwise registration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3244–3253, 2019. [1](#), [3](#)
- [5] Haoqiang Fan, Hao Su, and Leonidas J Guibas. 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*, pages 605–613, 2017. [2](#), [16](#)
- [6] Xiaoguang Han, Zhen Li, Haibin Huang, Evangelos Kalogerakis, and Yizhou Yu. High-resolution shape completion using deep neural networks for global structure and local geometry inference. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 85–93, 2017. [3](#)
- [7] Pedro Hermosilla, Tobias Ritschel, and Timo Ropinski. Total denoising: Unsupervised learning of 3d point cloud cleaning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 52–60, 2019. [1](#)
- [8] Zitian Huang, Yikuan Yu, Jiawen Xu, Feng Ni, and Xinyi Le. Pf-net: Point fractal network for 3d point cloud completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7662–7670, 2020. [3](#)
- [9] Daniel P Huttenlocher, Gregory A. Klanderman, and William J Rucklidge. Comparing images using the hausdorff distance. *IEEE Transactions on pattern analysis and machine intelligence*, 15(9):850–863, 1993. [3](#)
- [10] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. [8](#)
- [11] Ruihui Li, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. Pu-gan: a point cloud upsampling adversarial network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7203–7212, 2019. [1](#)
- [12] Ruihui Li, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. Pu-gan: a point cloud upsampling adversarial network. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. [13](#)

- [13] Tianyi Lin, Nhat Ho, and Michael Jordan. On efficient optimal transport: An analysis of greedy and accelerated mirror descent algorithms. In *International Conference on Machine Learning*, pages 3982–3991. PMLR, 2019. 4, 12
- [14] Minghua Liu, Lu Sheng, Sheng Yang, Jing Shao, and Shi-Min Hu. Morphing and sampling network for dense point cloud completion. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11596–11603, 2020. 2, 3, 4, 6, 8, 12, 14, 15
- [15] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. 3, 7
- [16] Trung Nguyen, Quang-Hieu Pham, Tam Le, Tung Pham, Nhat Ho, and Binh-Son Hua. Point-set distances for learning representations of 3d point clouds. *arXiv preprint arXiv:2102.04014*, 2021. 3
- [17] Yinyu Nie, Yiqun Lin, Xiaoguang Han, Shihui Guo, Jian Chang, Shuguang Cui, and Jian Jun Zhang. Skeleton-bridged point completion: From global inference to local adjustment. *Advances in Neural Information Processing Systems*, 2020. 3
- [18] Liang Pan, Xinyi Chen, Zhongang Cai, Junzhe Zhang, Haiyu Zhao, Shuai Yi, and Ziwei Liu. Variational relational point completion network. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 3, 6, 8, 10, 14, 15, 16
- [19] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 3, 7
- [20] Ofir Pele and Michael Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th international conference on computer vision*, pages 460–467. IEEE, 2009. 12
- [21] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 523–540. Springer, 2020. 3, 7
- [22] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 1(2):4, 2017. 15
- [23] Edward J Smith and David Meger. Improved adversarial systems for 3d object generation and reconstruction. In *Conference on Robot Learning*, pages 87–96, 2017. 3
- [24] Yongbin Sun, Yue Wang, Ziwei Liu, Joshua Siegel, and Sanjay Sarma. Pointgrow: Autoregressively learned point cloud generation with self-attention. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 61–70, 2020. 1
- [25] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3405–3414, 2019. 2
- [26] Lyne P Tchapmi, Vineet Kosaraju, Hamid Rezatofighi, Ian Reid, and Silvio Savarese. Topnet: Structural point cloud decoder. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 383–392, 2019. 3, 8, 15
- [27] Dahlia Urbach, Yizhak Ben-Shabat, and Michael Lindenbaum. Dpdist: Comparing point clouds using deep point cloud distance. *The European Conference on Computer Vision (ECCV)*, 2020. 3
- [28] Tai Wang, Xinge Zhu, and Dahua Lin. Reconfigurable voxels: A new representation for lidar-based point clouds. In *Conference on Robot Learning*, 2020. 1
- [29] Xiaogang Wang, Marcelo H Ang Jr, and Gim Hee Lee. Cascaded refinement network for point cloud completion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 790–799, 2020. 3, 6, 15
- [30] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019. 1
- [31] Xin Wen, Tianyang Li, Zhizhong Han, and Yu-Shen Liu. Point cloud completion by skip-attention network with hierarchical folding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1939–1948, 2020. 3
- [32] Rundi Wu, Xuelin Chen, Yixin Zhuang, and Baoquan Chen. Multimodal shape completion via conditional generative adversarial networks. In *The European Conference on Computer Vision (ECCV)*, August 2020. 3

- [33] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 8
- [34] Chulin Xie, Chuxin Wang, Bo Zhang, Hao Yang, Dong Chen, and Fang Wen. Style-based point generator with adversarial rendering for point cloud completion. *arXiv preprint arXiv:2103.02535*, 2021. 3
- [35] Haozhe Xie, Hongxun Yao, Shangchen Zhou, Jiageng Mao, Shengping Zhang, and Wenxiu Sun. Grnet: Gridding residual network for dense point cloud completion. *arXiv preprint arXiv:2006.03761*, 2020. 3
- [36] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4541–4550, 2019. 1
- [37] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 206–215, 2018. 3
- [38] Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. Pu-net: Point cloud upsampling network. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 13
- [39] Wentao Yuan, Tejas Khot, David Held, Christoph Mertz, and Martial Hebert. Pcn: Point completion network. In *International Conference on 3D Vision (3DV)*, pages 728–737. IEEE, 2018. 1, 3, 6, 8, 10, 14, 15
- [40] Junzhe Zhang, Xinyi Chen, Zhongang Cai, Liang Pan, Haiyu Zhao, Shuai Yi, Chai Kiat Yeo, Bo Dai, and Chen Change Loy. Unsupervised 3d shape completion through gan inversion. In *CVPR*, 2021. 1
- [41] Wenxiao Zhang, Qingan Yan, and Chunxia Xiao. Detail preserved point cloud completion via separated feature aggregation. *arXiv preprint arXiv:2007.02374*, 2020. 3
- [42] Xinge Zhu, Hui Zhou, Tai Wang, Fangzhou Hong, Yuexin Ma, Wei Li, Hongsheng Li, and Dahua Lin. Cylindrical and asymmetrical 3d convolution networks for lidar segmentation. In *Proceedings of the European Conference on Computer Vision*, 2021. 1