

Accurate Instance-Level CAD Model Retrieval in a Large-Scale Database

Jiaxin Wei¹, Lan Hu^{1,2}, Chenyu Wang¹, Laurent Kneip¹

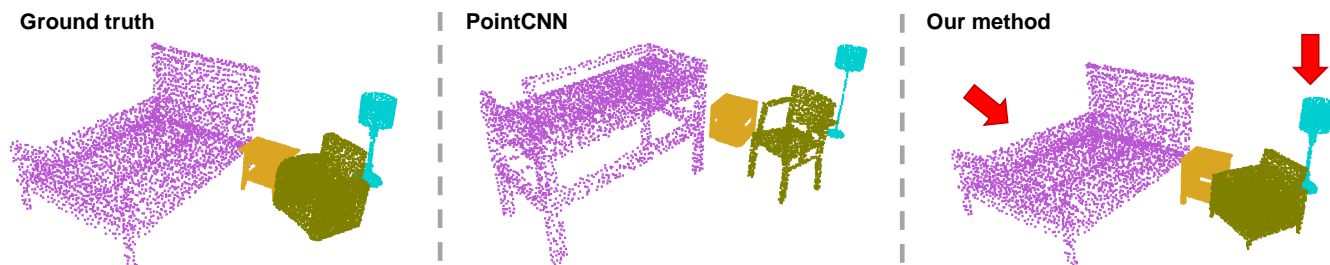


Fig. 1: Scene reconstruction of scene0143.00 from the Scan2CAD [1] dataset. This scene contains a bed, a table, a chair, and a lamp. Compared to PointCNN [2], our method not only retrieves the ground truth model for the bed and the lamp, but also gives more reasonable results for the table and the chair.

Abstract—We present a new solution to the fine-grained retrieval of clean CAD models from a large-scale database in order to recover detailed object shape geometries for RGBD scans. Unlike previous work simply indexing into a moderately small database using an object shape descriptor and accepting the top retrieval result, we argue that in the case of a large-scale database a more accurate model may be found within a neighborhood of the descriptor. More importantly, we propose that the distinctiveness deficiency of shape descriptors at the instance level can be compensated by a geometry-based re-ranking of its neighborhood. Our approach first leverages the discriminative power of learned representations to distinguish between different categories of models and then uses a novel robust point set distance metric to re-rank the CAD neighborhood, enabling fine-grained retrieval in a large shape database. Evaluation on a real-world dataset shows that our geometry-based re-ranking is a conceptually simple but highly effective method that can lead to a significant improvement in retrieval accuracy compared to the state-of-the-art.

I. INTRODUCTION

Nowadays, emerging robotics and augmented reality [3] applications have a huge demand for real-time 3D reconstructions of object-level scene models from RGBD scans. Traditional incremental fusion solutions [4][5] often fail to recover detailed scene geometries due to the noisy nature of the measurements captured by consumer-grade hardware. Instead of focusing on the exhaustive refinement of 3D scans [6], an increasing number of approaches [7][1][8] turn to an offline CAD database for better reconstructions. Given sufficient object detection, segmentation, and representation abilities, the original incomplete scans may be greatly improved by substituting corresponding measurement parts with

similar CAD models.

However, the imperfection of object detection and segmentation, together with the incapability of 3D shape descriptors to encode precise geometry information, jointly place a major restriction on CAD model retrieval at a finer level. Current methods, such as [1][8], concentrate on computing the 9DoF alignment between scanned objects and CAD models while only requiring the retrieval of a same-category model due to the usage of a small ground truth pool. Zhao *et al.* [9] and Dahnert *et al.* [10] claim retrieval of CAD instances with higher similarity, but they still constrain the model search to a relatively small database containing an augmented ground truth set for each scanned object. We instead focus on instance-level retrieval of CAD models in a large shape database, which is more practical in real-world scenarios.

CAD models from the same category usually possess a universal geometric structure. Therefore, it is quite easy for learned representations [11][12][13] to achieve outstanding performance in 3D shape classification. Yet, individual CAD models within the same class vary drastically in terms of shape details, making it hard to differentiate them from one another. To this end, we propose a two-step pipeline to tackle this challenging problem. We first retrieve multiple CAD candidates from the shape descriptor space using nearest neighbor search, and then use geometric residuals to re-rank those candidates, thereby obtaining much better results both quantitatively and qualitatively. Our method is based on two important insights. First, the discriminative power of learned features can help narrow down the search to models with the same semantic class as that of the scanned object. And a suitable hierarchical data structure on top of our large-scale database makes the search more efficiently. Second, after figuring out which part of a database is of interest, we rely on point set distance metrics to further perform a local search and refine the ranking order of initial candidates. We

¹Mobile Perception Lab, School of Information Science and Technology, ShanghaiTech University;

{weijx, wangchy4, lkneip}@shanghaitech.edu.cn

²Shanghai Institute of Microsystem and Information Technology, Chinese Academy of Sciences and University of Chinese Academy of Sciences. hulan@shanghaitech.edu.cn

also suggest a new variant of Chamfer Distance to deal with pathological cases in evaluating candidates.

The main contributions of this paper are summarized as follows:

- We present a new pipeline for instance-level CAD model retrieval at large scale. It relies on a feature-based k -nearest neighbor search followed by a geometry-based re-ranking, which can considerably improve the quality of retrieved CAD models.
- We introduce a variant of the Single-direction Chamfer Distance [9] that is more robust to outliers when measuring point set distances between partial scans and CAD models.

As demonstrated by our experiments, the combination of geometric re-ranking outperforms the plain application of state-of-the-art feature embedding networks. From a qualitative perspective, it also produces much more visually similar results, and thus for the first time enables instance-level retrieval of CAD models within a large-scale database.

II. RELATED WORK

CAD Model Retrieval: There are mainly two directions for 3D model retrieval: view-based methods and model-based methods. View-based methods extract and merge 2D features from multiple projections of a 3D object and merge those features to generate a global shape descriptor for later recognition while model-based methods directly deal with 3D data. View-based methods are out of scope of this paper. The interested readers are referred to [14][15][16] for more details.

Most traditional 3D model retrieval algorithms [17][18][19][20] are based on hand-designed descriptors [21][22] that may easily suffer from the inconsistency between clean CAD models and noisy real-world data. Meanwhile, learning-based methods are getting more and more popular. Scan2CAD [1] carries out implicit retrieval by evaluating the semantic compatibility between CAD models and a patch from the scan. Avetisyan *et al.* [8] insert an extra object detection module such that the relevant parts of a scan can be directly used to calculate descriptors. Dahnert *et al.* [10] segment the foreground from the input scan and complete the segmented objects by stacking hourglass encoder-decoders, resulting in a joint embedding of 3D scans and CAD models. More recently, CORSAIR [9] extends FCGF [23] to simultaneously learn a global feature for retrieval as well as local point features for alignment.

The aforementioned methods have shown promising results yet are not practical. They either assume a small pool of exact CAD models as input or consider a pool with augmented ground truth models for each scanned object. In contrast, our proposed method performs retrieval in a more realistic setting where the database has thousands of models and a one-to-many relationship with scanned objects.

3D Shape Descriptors: Given that there are plenty of excellent works for instance segmentation [24][25][26], we can safely assume that objects have already been segmented from scans. 3D model retrieval can thus be achieved by

learning a global shape descriptor for object scans, which has attracted increasing attention over the past few years. For example, VoxNet [27] applies a supervised 3D CNN on volumetric occupancy grids to predict object class labels. Alternative approaches directly work on point clouds. PointNet [11] extracts point features using shared MLPs and aggregates them with a max pooling operator to obtain a permutation-invariant global feature. Later on, Qi *et al.* propose PointNet++ [12] to further capture the local geometric structures of point clouds from neighborhood points. In the spirit of graph neural networks, DGCNN [13] constructs a local graph on the point cloud which can be dynamically updated using aggregated edge features after each layer of the network. As a generalization of 2D CNN, PointCNN [2] uses \mathcal{X} -Conv to transform the input point clouds and raw features into a canonical form and then applies typical convolutions. However, those works usually do not work well with noisy or partial measurements. We resort to geometric information to reduce the uncertainty.

Set Similarity for Point Clouds: Fan *et al.* [28] propose two distance metrics to compare generated point clouds against ground truth: Chamfer Distance (CD) and Earth Mover’s Distance (EMD). EMD gives the minimum effort of transforming one point set into the other and has been proven to be more accurate than CD in point cloud generation [28][29]. The biggest issue impeaching EMD from widespread use is the high computational cost for exact solutions. CD is preferred [30][31][32] when efficiency is of major concern.

III. METHOD

This section is organized as follows. Section III-A reveals our motivation of establishing instance-level retrieval in a large CAD database and provides an overview of our method. In Section III-B and III-C, we describe our two-step pipeline in all details, comprising of feature-based ranking and geometry-based re-ranking.

A. Overview

The retrieval of a similar CAD model is tightly entangled with the cross-instance alignment between a CAD model and a scanned object. With the rise of category-level object alignment [33], it is possible to align different object instances within a category. For ease of implementation, we simply assume that all the retrieved CAD models are already aligned with the query object by leveraging such alignment methods.

Avetisyan *et al.* [1][8] retrieve CAD models from a small pool tailored for each scene. The model pool contains only the ground truth CAD models corresponding to the objects that are effectively present in the scene. On the contrary, we perform retrieval in a large-scale CAD database, which is more closer to the real situation. We furthermore focus on instance-level rather than category-level retrieval. For category-level retrieval, the matching of categories between the retrieved CAD and the query object is enough to be regarded as a success. However, real-world applications such

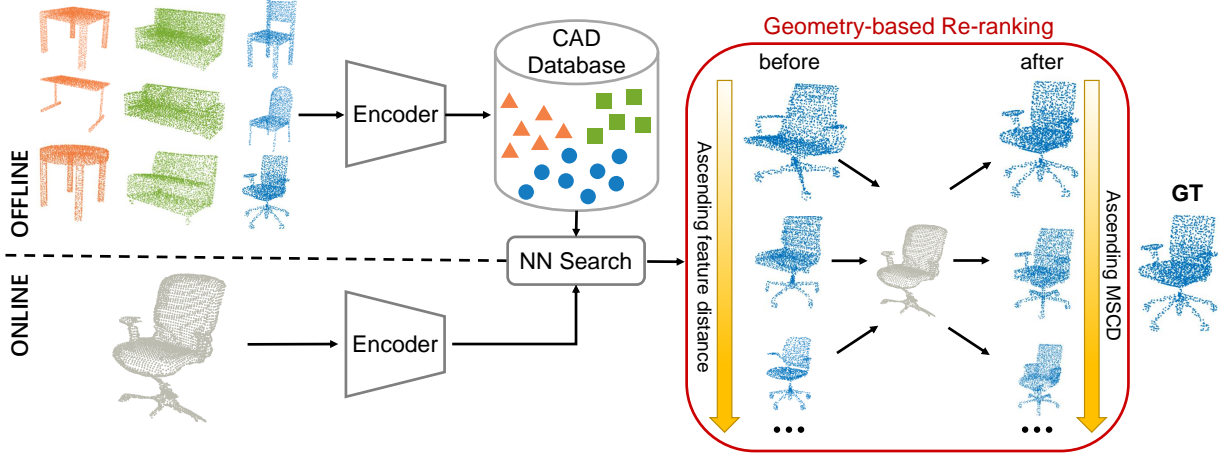


Fig. 2: An illustration of our pipeline. CAD models are processed offline to form a large database in which we can do k -nearest neighbor search for each query object. Based on the coarsely retrieved CAD candidates, we further perform geometry-based re-ranking using Modified Single-Direction Chamfer distance (MSCD) to make sure that higher-ranked candidates are generally more similar to the query object.

as robot interactions with surrounding environment require a detailed understanding of object shapes. In other words, the retrieved CAD model should replicate the object shape as closely as possible. In contrast to the works of Ishimtsev *et al.* [34] and Uy *et al.* [35], our method avoids the involvement of expensive CAD model deformation, and makes full use of a large-scale database to identify the model that is most similar to the scanned object.

Our method breaks down into two procedures. First, we calculate a shape descriptor using a neural network and then perform coarse retrieval by means of k -nearest neighbour search to obtain an initial set of CAD candidates. It points out regions of interest in the database that are highly likely to contain the best shape. Second, we apply geometric verification to assess and re-rank the initial candidates. As it is hard to evaluate the similarity between retrieved CAD models and incomplete scanned objects, we introduce a modified version of Chamfer Distance to better handle the outliers. This step can significantly eliminate the uncertainty incurred by neural networks. In addition, the refinement is quite efficient given that the number of CAD candidates is far fewer than that of models in the whole database. Our complete pipeline is shown in Figure 2.

B. Feature-Based Ranking

Although neural networks may not yet be able to capture the exact geometry of an object, it is still easy for them to categorize shapes into the right semantic classes. Consequently, we utilize the discriminative power of learned global features to characterize inter-class dissimilarity. All the query objects and CAD models are processed through the same backbone encoder. Note that there are little constraints on the network architecture as long as it outputs a global feature for each shape. Here we adopt neural networks originally designed for 3D shape classification.

Given a query object $\mathbf{X} \in \mathbb{R}^{N \times 3}$ and a CAD database $\mathcal{Y} = \{\mathbf{Y}_i | \mathbf{Y}_i \in \mathbb{R}^{M_i \times 3}, 1 \leq i \leq S\}$ consisting of S models,

all represented in the form of point clouds, our goal is to retrieve a CAD model $\mathbf{Y} \in \mathcal{Y}$ that is not only semantically compatible to \mathbf{X} but also geometrically close to it. Let $f(\cdot)$ denote the backbone encoder. It takes as input a point cloud \mathbf{P} and outputs a d -dimensional shape descriptor

$$f : \mathbf{P} \mapsto f(\mathbf{P}) \in \mathbb{R}^d.$$

Out of consideration for efficiency, features for all CAD models in the database are pre-computed offline. Let $\mathbf{F}^{\mathbf{X}} \in \mathbb{R}^d$ and $\mathbf{F}^{\mathcal{Y}} = \{\mathbf{F}_i^{\mathcal{Y}} | \mathbf{F}_i^{\mathcal{Y}} \in \mathbb{R}^d, 1 \leq i \leq S\}$ be the extracted global features for the query object and the models, respectively. Each time a new query arrives, we perform a k -nearest neighbor search using the object feature to obtain an ordered candidate set:

$$\mathcal{C} = \{\mathbf{C}_j | \mathbf{C}_j \in \mathcal{Y}, 1 \leq j \leq k\}$$

such that $d(\mathbf{F}^{\mathbf{X}}, \mathbf{F}_j^{\mathcal{C}})$ is the j -th smallest element in $d(\mathbf{F}^{\mathbf{X}}, \mathbf{F}^{\mathcal{Y}})$. Here $d(\cdot)$ represents the Euclidean distance between feature vectors. To reduce the search complexity, we explore the database via hierarchical data structures, such as a kd-tree.

C. Geometry-Based Re-ranking

In most cases, the initial set may contain a very similar model but its ranking order is unsatisfactory owing to the limited ability of learned features to encode precise shape information. The distinctiveness can be further brought down when some parts of the query object are occluded or unobserved. Meanwhile, the bounded number of CAD models in the database means that there usually do not exist identical matches. Therefore, we exploit geometric information to address the intra-class similarity. Specifically, we re-rank the initial CAD candidates by measuring the geometric difference between point clouds.

Chamfer Distance (CD) is a popular point set distance metric used in many learning based point cloud generation

frameworks [28][29]. For two point clouds P and Q , a common definition of CD is to find the nearest neighbor of P in Q and vice versa, and then add up all squared distances. That is,

$$d_{CD}(P, Q) = \frac{1}{|P|} \sum_{p \in P} \min_{q \in Q} \|p - q\|_2^2 + \frac{1}{|Q|} \sum_{q \in Q} \min_{p \in P} \|p - q\|_2^2, \quad (1)$$

where $|P|$ and $|Q|$ are the number of points in P and Q , respectively. The division by $|P|$ and $|Q|$ acts as a normalization which can balance the influence caused by different point densities in P and Q . Note that the term *Chamfer Distance* is a slight abuse of notation as it does not actually satisfy the triangle inequality required by a proper distance metric [28].

An interesting variant of CD has been introduced in [9], called the Single-direction Chamfer Distance (SCD). It simply drops one side of squared distances in (1), resulting in

$$d_{SCD}(P, Q) = \frac{1}{|P|} \sum_{p \in P} \min_{q \in Q} \|p - q\|_2^2. \quad (2)$$

Hence, SCD only maintains the distances from scanned object to CAD model. It is intuitively clear that this simplification helps SCD to concentrate on the visible parts of a scanned object while neglecting unobserved parts.

For a better understanding of (2), let us define a distance vector

$$\mathbf{d} = \begin{bmatrix} \min_{q \in Q} \|p_1 - q\|_2 \\ \min_{q \in Q} \|p_2 - q\|_2 \\ \dots \\ \min_{q \in Q} \|p_{|P|} - q\|_2 \end{bmatrix} \in \mathbb{R}^{|P|}.$$

Thus, SCD can be reformulated as

$$d_{SCD}(P, Q) = \frac{1}{|P|} \|\mathbf{d}\|_2^2, \quad (3)$$

the L_2 -norm of the distance vector \mathbf{d} . However, it is well known from the literature of robust optimization that the L_2 -norm is sensitive to outliers as it emphasises on residuals by squaring them. In practical scenarios, object scans often contain outlier measurements owing to the existence of segmentation errors and other artifacts such as floating point measurements. In order to improve the robustness of the above distance against such artifacts, we propose to use the L_1 -norm of \mathbf{d} instead. We denote this distance the Modified Single-direction Chamfer Distance (MSCD), and it is given by

$$d_{MSCD}(P, Q) = \frac{1}{|P|} \sum_{p \in P} \min_{q \in Q} \|p - q\|_2. \quad (4)$$

Though the change seems subtle, our experimental results demonstrate that MSCD achieves better results in practice.

To conclude, we take each candidate CAD model from \mathcal{C} and calculate the MSCD with respect to our object scan, yielding a new ordered set

$$\hat{\mathcal{C}} = \{\hat{\mathbf{C}}_j | \mathbf{C}_j \in \mathcal{C}, 1 \leq j \leq k\}$$

such that $d_{MSCD}(\mathbf{X}, \hat{\mathbf{C}}_j)$ is the j -th smallest element in $d_{MSCD}(\mathbf{X}, \mathcal{C})$.

IV. EXPERIMENTS

This section is organized as follows. Section IV-A introduces the dataset on which we conduct all our experiments. Section IV-B describes the alternative state-of-the-art methods and Section IV-C defines three evaluation metrics. Discussions of the performance and ablation studies are presented in Section IV-D and Section IV-E, respectively. Finally, in Section IV-F, we show some qualitative results to further demonstrate the superiority of our method.

A. Dataset

We test our method on the real-world dataset Scan2CAD [1], which contains 14225 (3049 unique) CAD models from ShapeNet [37] and corresponding scanned objects from ScanNet [38]. For CAD model retrieval, we assume that instance segmentations are known, and apply the segmentation masks over the 3D scans to obtain object surfaces. We furthermore sample points on the surface of the CAD models to generate point clouds for the database. Note that we filter out some poor quality point clouds with too few points. Finally, we work with 12598 object scans spreading over 1506 scenes, and a database of 2827 CAD models.

B. Shape Descriptor Alternatives

For comparison purposes, we choose several state-of-the-art 3D shape descriptors. This includes the handcrafted feature VFH [36] as well as the learning-based features DGCNN [13], VoxNet [27], PointNet2 (or PointNet++) [12], and PointCNN [2]. Viewpoint Feature Histogram (VFH) is a global descriptor based on FPFH [21]. We compute the 308-dimensional VFH feature using the PCL implementation [39]. As for DGCNN, VoxNet, PointNet2 and PointCNN, they are representatives of graph-based, volumetric-based, point-based, and convolution-based methods for 3D shape classification, respectively. All networks are trained on ModelNet40 [40] and we extract the global feature from the input of the final classification layer. Both DGCNN and PointNet2 output 256-dimensional feature vectors while VoxNet and PointCNN produce 128-dimensional features.

C. Evaluation Metrics

We use three different metrics to evaluate the retrieval performance on Scan2CAD:

Topk Retrieval Accuracy. Suppose we have T query objects in \mathcal{X} and there is a ground truth CAD model $\mathbf{Y}^* \in \mathcal{Y}$ for each query object $\mathbf{X} \in \mathcal{X}$. The Topk retrieval accuracy (RA) is defined as

$$\text{RA} = \frac{1}{T} \sum_{\mathbf{X} \in \mathcal{X}} \mathbf{1}_{\mathcal{A}_k}(\mathbf{Y}^*), \quad (5)$$

where \mathcal{A}_k can be replaced either by the top-k subset of the initial ordered set \mathcal{C} , or the top-k subset of the re-ordered set $\hat{\mathcal{C}}$. Note that $\mathbf{1}_{\mathcal{A}_k} : \mathcal{Y} \rightarrow \{0, 1\}$ is an indicator function:

$$\mathbf{1}_{\mathcal{A}_k}(\mathbf{Y}^*) = \begin{cases} 1 & \text{if } \mathbf{Y}^* \in \mathcal{A}_k, \\ 0 & \text{otherwise.} \end{cases}$$

TABLE I: Top1/Top5 retrieval accuracy

Method	bed	bookshelf	cabinet	chair	display	sofa	table	other	class avg.	instance avg.
VFH [36]	0.00/0.00	0.00/0.02	0.00/0.00	0.01/0.02	0.01/0.03	0.00/0.01	0.01/0.04	0.00/0.00	0.00/0.01	0.01/0.02
DGCNN [13]	0.04/0.17	0.02/0.10	0.00/0.06	0.01/0.04	0.01/0.04	0.01/0.07	0.01/0.04	0.06/0.24	0.25/0.47	0.02/0.09
VoxNet [27]	0.12/0.37	0.06/0.17	0.02/0.06	0.02/0.05	0.03/0.08	0.04/0.12	0.03/0.10	0.07/0.22	0.22/0.50	0.04/0.11
PointNet2 [12]	0.10/0.28	0.02/0.08	0.01/0.05	0.01/0.05	0.01/0.07	0.02/0.08	0.02/0.08	0.09/0.28	0.29/0.51	0.03/0.11
PointCNN [2]	0.05/0.29	0.16/0.31	0.04/0.11	0.02/0.08	0.04/0.11	0.02/0.10	0.05/0.15	0.11/0.32	0.29/0.57	0.06/0.16
VFH+MSCD	0.01/0.01	0.05/0.14	0.02/0.03	0.02/0.08	0.03/0.07	0.03/0.05	0.08/0.11	0.01/0.01	0.02/0.04	0.03/0.07
DGCNN+MSCD	0.42/0.74	0.28/0.53	0.18/0.31	0.08/0.17	0.10/0.26	0.18/0.33	0.13/0.21	0.30/0.63	0.39/0.70	0.17/0.33
VoxNet+MSCD	0.41/ 0.75	0.26/0.51	0.21/0.43	0.08/0.18	0.11/0.29	0.25/0.40	0.17/0.31	0.30/0.62	0.43/0.69	0.18/0.36
PointNet2+MSCD	0.41/0.72	0.28/0.53	0.23/0.44	0.10/0.21	0.12/0.29	0.22/0.37	0.19/0.32	0.30/0.63	0.43/0.70	0.19/0.37
PointCNN+MSCD	0.41/0.72	0.29/0.56	0.27/0.54	0.11/0.26	0.16/0.35	0.25/0.42	0.27/0.47	0.30/0.63	0.43/0.71	0.22/0.43

Top1 Chamfer Distance. With a slight modification of the metric defined in [9], the Top1 Chamfer Distance is computed by measuring MSCD between the ground truth CAD model and the retrieved top 1 CAD model:

$$d = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{X}} d_{MSCD}(\mathcal{A}_1, \mathbf{Y}^*). \quad (6)$$

Ground truth Ranking. The ranking order of the ground truth CAD model in the returned CAD candidate set:

$$R = \frac{1}{T} \sum_{\mathbf{x} \in \mathcal{X}} \text{ranking}_{\mathcal{A}}(\mathbf{Y}^*) \quad (7)$$

where $\text{ranking}_{\mathcal{A}}(\cdot)$ denotes the ranking position in \mathcal{A} .

D. Results

We use the methods mentioned in Section IV-B as our front-end, followed by MSCD-based re-ranking to improve the CAD candidate set returned by the initial feature-based search. Note that the nearest neighbour search is configured to find 90 nearest neighbors. Table I indicates Top1 and Top5 retrieval accuracy of different methods and demonstrates the impact of adding MSCD-based re-ranking. It is obvious that the proposed geometric re-ranking strategy leads to a significant performance boost over the original feature-based ranking results. In addition, our method achieves much better results when combined with PointCNN.

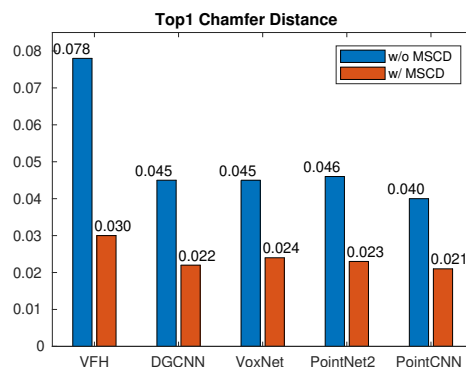
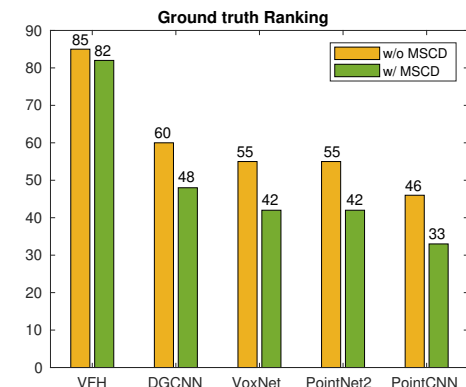
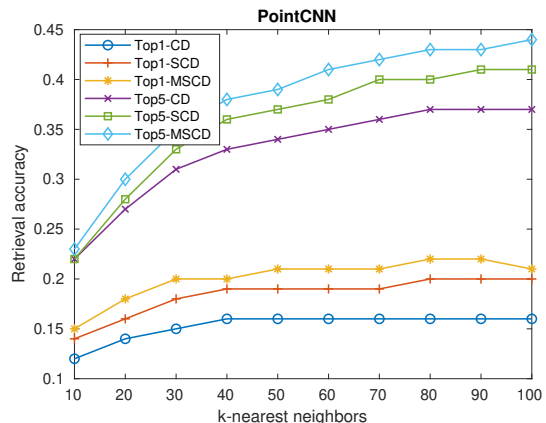
E. Ablation Studies

In Table II, we show that the top 5 CAD models produced by 3D shape descriptors almost always have the same category as that of the query object, only with the exception of VFH. This proves the usefulness of neural networks in handling inter-class dissimilarity and also explains the poor performance of VFH.

We further test our MSCD-based re-ranking method using Top1 Chamfer Distance. As shown in Figure 3, the MSCD between the ground truth and the retrieved top 1 CAD model declines for all 3D shape descriptors when combining with a

TABLE II: Top5 Category Ratio

Method	Top5 Category Ratio
VFH [36]	0.32
DGCNN [13]	0.99
VoxNet [27]	1.00
PointNet2 [12]	0.99
PointCNN [2]	0.99


Fig. 3: Comparative results of Top1 Chamfer Distance on Scan2CAD.

Fig. 4: Comparative results of Ground truth Ranking on Scan2CAD.

Fig. 5: Top1 and Top5 retrieval accuracy using different metrics and different search ranges based on PointCNN.

Query	Gt	Top5 Retrieved CAD (w/o MSCD)					Top5 Retrieved CAD (w/ MSCD)				
		Ascending feature distance →					Ascending MSCD →				

Fig. 6: Qualitative results of CAD model retrieval on Scan2CAD [1] based on PointCNN [2]. Rows 1 to 5 show different kinds of examples for classes chair, table, bed, lamp and trash bin, respectively. The retrieved CAD model identical to ground truth is highlighted with corresponding color while the models of clearly inferior quality are highlighted with red. The last two rows are failure cases.

re-ranking method. We also compare the rank of the ground truth CAD model with and without re-ranking to validate its effectiveness (see Figure 4). After re-ranking, the ground truth CAD model is promoted to a higher rank.

Our method has very few hyper-parameters. To figure out the influence of hyper-parameters on retrieval performance, we run a series of experiments based on PointCNN using different nearest neighbor search ranges and different point cloud distance metrics (introduced in Section III-C). Figure 5 shows that MSCD outperforms both CD and SCD in terms of Top1 and Top5 retrieval accuracy. The performance gradually improves as the search range becomes larger.

F. Qualitative Results

Qualitative results for CAD model retrieval on Scan2CAD are shown in Figure 6. Our method takes good advantage of visible parts in the query object and retrieves more reasonable CAD models with similar structure. For example, let us consider the first row in Figure 6. The wheel and armrest fragments in the query indicate that this is a swivel chair with armrests. It is obvious that our proposed re-ranking method returns CAD models with both of these properties while the feature-based search struggles to get reasonable results. Furthermore, our method successfully rejects models

of clearly inferior quality (cf. row 3 and row 4 of Fig. 6).

The last two rows in Figure 6 show failure cases of our approach. Most of them are caused by the inherent ambiguity in query objects and inaccurate ground truth annotations. Taking the bottom row as an example, the query barely captures the back and the seat of a chair, which creates possibilities for arbitrary types of legs. Though our method retrieves some similar CAD models, it fails to find the exact model of the ground truth.

V. CONCLUSION

We have presented a two-step pipeline for instance-level CAD model retrieval in a large-scale database. Since state-of-the-art learned representations still have limitations on understanding all aspects of 3D shapes, we have proposed a geometry-based re-ranking method to facilitate fine-grained retrieval. We have furthermore introduced a Modified Single-direction Chamfer Distance to measure the point set distance between partial objects and CAD models. Thorough experiments on real-world test cases validate the superiority of our method. Future work will investigate the combination of instance-level retrieval and category-level alignment in order to simultaneously perform CAD model retrieval and alignment with high accuracy.

REFERENCES

- [1] A. Avetisyan, M. Dahnert, A. Dai, M. Savva, A. X. Chang, and M. Nießner, "Scan2cad: Learning cad model alignment in rgb-d scans," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2609–2618, 2019.
- [2] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, "Pointcnn: Convolution on x-transformed points," in *NeurIPS*, 2018.
- [3] A. J. Davison, "Futuremapping: The computational structure of spatial ai systems," *arXiv preprint arXiv:1803.11288*, 2018.
- [4] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon, "Kinectfusion: Real-time dense surface mapping and tracking," *2011 10th IEEE International Symposium on Mixed and Augmented Reality*, pp. 127–136, 2011.
- [5] J. McCormac, R. Clark, M. Bloesch, A. Davison, and S. Leutenegger, "Fusion++: Volumetric object-level slam," *2018 International Conference on 3D Vision (3DV)*, pp. 32–41, 2018.
- [6] J. Hou, A. Dai, and M. Nießner, "Revealnet: Seeing behind objects in rgb-d scans," *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2095–2104, 2020.
- [7] S. Gupta, P. Arbeláez, R. Girshick, and J. Malik, "Aligning 3d models to rgb-d images of cluttered scenes," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 4731–4740.
- [8] A. Avetisyan, A. Dai, and M. Nießner, "End-to-end cad model retrieval and 9dof alignment in 3d scans," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2551–2560, 2019.
- [9] T. Zhao, Q. Feng, S. Jadhav, and N. A. Atanasov, "Corsair: Convolutional object retrieval and symmetry-aided registration," *ArXiv*, vol. abs/2103.06911, 2021.
- [10] M. Dahnert, A. Dai, L. Guibas, and M. Nießner, "Joint embedding of 3d scan and cad objects," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8748–8757, 2019.
- [11] C. Qi, H. Su, K. Mo, and L. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 77–85, 2017.
- [12] C. Qi, L. Yi, H. Su, and L. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in *NIPS*, 2017.
- [13] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. Bronstein, and J. Solomon, "Dynamic graph cnn for learning on point clouds," *ACM Transactions on Graphics (TOG)*, vol. 38, pp. 1–12, 2019.
- [14] A. Grabner, P. M. Roth, and V. Lepetit, "Location field descriptors: Single image 3d model retrieval in the wild," in *2019 International Conference on 3D Vision (3DV)*. IEEE, 2019, pp. 583–593.
- [15] A. Bansal, B. Russell, and A. Gupta, "Marr revisited: 2d-3d alignment via surface normal prediction," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5965–5974.
- [16] X. He, Y. Zhou, Z. Zhou, S. Bai, and X. Bai, "Triplet-center loss for multi-view 3d object retrieval," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1945–1954.
- [17] Y. M. Kim, N. J. Mitra, D.-M. Yan, and L. Guibas, "Acquiring 3d indoor environments with variability and repetition," *ACM Transactions on Graphics (TOG)*, vol. 31, no. 6, pp. 1–11, 2012.
- [18] Y. Li, A. Dai, L. Guibas, and M. Nießner, "Database-assisted object retrieval for real-time 3d reconstruction," in *Computer Graphics Forum*, vol. 34, no. 2. Wiley Online Library, 2015, pp. 435–446.
- [19] Y. M. Kim, N. J. Mitra, Q. Huang, and L. Guibas, "Guided real-time scanning of indoor objects," in *Computer Graphics Forum*, vol. 32, no. 7. Wiley Online Library, 2013, pp. 177–186.
- [20] B. Drost and S. Ilic, "3d object detection and localization using multimodal point pair features," in *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission*. IEEE, 2012, pp. 9–16.
- [21] R. Rusu, N. Blodow, and M. Beetz, "Fast point feature histograms (fpfh) for 3d registration," *2009 IEEE International Conference on Robotics and Automation*, pp. 3212–3217, 2009.
- [22] A. E. Johnson, "Spin-images: a representation for 3-d surface matching," 1997.
- [23] C. Choy, J. Park, and V. Koltun, "Fully convolutional geometric features," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 8957–8965, 2019.
- [24] H. Chen, K. Sun, Z. Tian, C. Shen, Y. Huang, and Y. Yan, "Blendmask: Top-down meets bottom-up for instance segmentation," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8573–8581.
- [25] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "Yolact: Real-time instance segmentation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 9157–9166.
- [26] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [27] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 922–928, 2015.
- [28] H. Fan, H. Su, and L. Guibas, "A point set generation network for 3d object reconstruction from a single image," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2463–2471, 2017.
- [29] P. Achlioptas, O. Diamanti, I. Mitliagkas, and L. Guibas, "Learning representations and generative models for 3d point clouds," in *ICML*, 2018.
- [30] H. Deng, T. Birdal, and S. Ilic, "Ppf-foldnet: Unsupervised learning of rotation invariant 3d local descriptors," *ArXiv*, vol. abs/1808.10322, 2018.
- [31] —, "3d local features for direct pairwise registration," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3239–3248, 2019.
- [32] C. Duan, S. Chen, and J. Kovacevic, "3d point cloud denoising via deep neural network based local surface estimation," *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 8553–8557, 2019.
- [33] Q. Feng and N. A. Atanasov, "Fully convolutional geometric features for category-level object alignment," *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 8492–8498, 2020.
- [34] V. Ishimtsev, A. Bokhovkin, A. Artemov, S. Ignatyev, M. Nießner, D. Zorin, and E. Burnaev, "Cad-deform: Deformable fitting of cad models to 3d scans," in *ECCV*, 2020.
- [35] M. A. Uy, V. G. Kim, M. Sung, N. Aigerman, S. Chaudhuri, and L. Guibas, "Joint learning of 3d shape retrieval and deformation," *ArXiv*, vol. abs/2101.07889, 2021.
- [36] R. Rusu, G. Bradski, R. Thibaux, and J. M. Hsu, "Fast 3d recognition and pose using the viewpoint feature histogram," *2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 2155–2162, 2010.
- [37] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu, "Shapenet: An information-rich 3d model repository," *ArXiv*, vol. abs/1512.03012, 2015.
- [38] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "ScanNet: Richly-annotated 3d reconstructions of indoor scenes," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 2432–2443, 2017.
- [39] R. B. Rusu and S. Cousins, "3D is here: Point Cloud Library (PCL)," in *IEEE International Conference on Robotics and Automation (ICRA)*. Shanghai, China: IEEE, May 9-13 2011.
- [40] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1912–1920, 2015.