
Cluster3D: A Dataset and Benchmark for Clustering Non-Categorical 3D CAD Models

Siyuan Xiang^{1*} Chin Tseng^{1*} Congcong Wen^{1†} Deshana Desai^{1†}
Yifeng Kou^{1†} Binil Starly² Daniele Panozzo^{1‡} Chen Feng^{1‡}
¹New York University ²North Carolina State University
<https://cluster3d.github.io/>

Abstract

1 We introduce the first large-scale dataset and benchmark for non-categorical an-
2 notation and clustering of 3D CAD models. We use the geometric data of the
3 ABC dataset, and we develop an interface to allow expert mechanical engineers
4 to efficiently annotate pairwise CAD model similarities, which we use to evaluate
5 the performance of seven baseline deep clustering methods. Our dataset contains
6 a manually annotated subset of 22,968 shapes, and 252,648 annotations. Our
7 dataset is the first to directly target deep clustering algorithms for geometric shapes,
8 and we believe it will be an important building block to analyze and utilize the mas-
9 sive 3D shape collections that are starting to appear in deep geometric computing.
10 Our results suggest that, differently from the already mature shape classification
11 algorithms, deep clustering algorithms for 3D CAD models are in their infancy and
12 there is much room for improving their performance.

13 1 Introduction

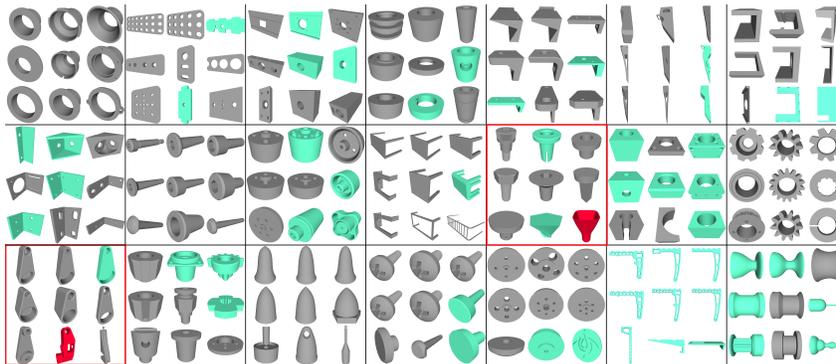


Figure 1: **The overview of Cluster3D** via a subset of clustering results from a baseline (DeepCluster), demonstrating the challenges of classification-based labeling in our task, due to many non-standard mechanical components (green). Each section shows some random CAD models in the same cluster. A red section shows a cluster with annotation violations highlighted at the red objects.

14 Shape classification is a core component in many modern 3D computer vision pipelines, and for
15 which many datasets and benchmarks have been introduced in the last decade, usually focusing on a
16 small number of object classes.

*Equal contributions.

†Equal contributions.

‡The corresponding authors are Chen Feng cfeng@nyu.edu and Daniele Panozzo panozzo@nyu.edu.

17 With the advent of large geometric collections of data, it is natural to expect that harder classification
18 tasks will be solvable using modern data-driven classification approaches. However, we discovered
19 that when the number of categorical classes becomes large and hierarchical, the task becomes very
20 challenging. Even annotation of such object classes by human experts may yield varying classification
21 labels, particularly because the necessary contextual information about the objects themselves are
22 unavailable. We attempted to annotate the ABC dataset [35], which is composed of 1 million 3D
23 CAD models, manually modeled by hobbyists and experts alike. The dataset has very complex
24 and unbalanced class distributions, i.e., essentially being non-categorical. This makes the problem
25 intractable even by our subject matter expert annotators (graduate students in mechanical engineering).
26 *One may easily see such annotation challenges when trying to name each group of objects in Figure 1.*

27 Despite the challenges, without an automatic way to classify shapes, the practical utility of large
28 geometric collections is hindered, especially considering that manual annotation is economically
29 unfeasible at large scale, if not impossible to carry out due to incomplete information on the design
30 intent of the original asset creator. Even when design information or engineering specifications are
31 available, the content is in text form that cannot be easily related to features of the 3D model. The
32 obvious alternative to supervised shape classification is the use of an unsupervised clustering method
33 to group shapes, helping in both the annotation process and in the dataset exploration.

34 Surprisingly, we could not find any large-scale dataset or benchmark for deep clustering of collections
35 of 3D shapes, especially non-categorical ones where clustering is more useful. To fill this gap, we
36 propose to construct such a dataset based on the ABC dataset. Because of the non-categorical feature,
37 we cannot annotate each object with a class label. Instead, we propose to annotate the pairwise 3D
38 shape similarity relationship. Although it might sound even more intractable, we allow experts to
39 focus on a small subset of carefully selected pairs (instead of all the pairs) to provide useful and
40 scalable annotation. We developed a web-based user interface to implement this annotation workflow
41 and to annotate 252, 648 selected pairwise similarities on 22, 968 ABC objects.

42 We benchmark seven clustering baseline methods to analyze the properties of our dataset. These
43 clustering methods can be classified into two types: 1) two-stage clustering and 2) end-to-end deep
44 clustering. For the two-stage clustering approach, we first perform deep representation learning on
45 3D mechanical components to extract high dimensional features, using pre-trained neural networks.
46 Then we apply classic clustering algorithms, like KMeans [40], to group these learned features into
47 different clusters. For end-to-end deep clustering methods, we combine representation learning,
48 dimensionality reduction, and clustering in an end-to-end framework. Based on either the ground
49 truth annotations or CAD model distance metrics, we respectively apply either external or internal
50 cluster validation indexes [53] to evaluate their performances. Particularly, since we are the first
51 to use similarity-based annotations, we need to design an external cluster validation index. We
52 propose two formulations: one measuring the pairwise similarity accuracy, and the other measuring
53 the compactness for the cluster elements.

54 We discovered that the performance of existing deep clustering methods is still insufficient for the
55 automatic clustering of large datasets, and there is a lot of room for algorithmic improvement. We
56 believe that our dataset will help by providing an objective metric on a large dataset specifically
57 designed for this task. *We plan to continue collecting data and periodically update the dataset.* The
58 dataset data, the annotation software, the implementation of all baseline methods, and scripts to run
59 the evaluations are publicly released as open source using the MIT license.

60 In summary, our contributions are the following:

- 61 • To the best of our knowledge, Cluster3D is the first dataset focusing on non-categorical
62 annotation for 3D mechanical components, which could stimulate a new direction for deep
63 clustering on large-scale mechanical component collections.
- 64 • We propose a scalable and effective pairwise similarity annotation workflow, implemented
65 in a graphical user interface, to allow experts to efficiently label a large number of object
66 pairs (for a total of 252, 648 annotation pairs per annotator).
- 67 • We design/adapt 7 clustering methods on our dataset and benchmark their performances.
- 68 • We propose 2 external cluster evaluation indexes to evaluate the clustering results, using
69 the similarity annotation. Also, we analyze our evaluation metrics, comparing them with
70 several internal cluster evaluation indexes.

71 2 Related Work

72 We cover the related works more closely related to our main contributions: (1) large datasets of 3D
73 models, (2) approaches for annotating 3D datasets, and (3) clustering algorithms and corresponding
74 evaluation metrics.

75 **3D mechanical object datasets.** Large-scale 3D object datasets are routinely used for classification,
76 instance segmentation, and shape reconstruction tasks [62, 43, 72, 10]. However, these datasets focus
77 on a small group of categories, where each object can be uniquely and reliably be classified. Recently,
78 large datasets of mechanical components have been introduced, which dramatically increase the
79 difficulty in identifying and classifying 3D parts, due to their self-similarity and, generally, larger
80 number of categories. Annotation of mechanical components requires more effort since the labeling
81 work for such a dataset needs subject matter expertise rather than just common sense [34]. MCB [34]
82 contains 58,696 components and 68 categories. The Fabwave dataset [2] contains 46 classes of
83 standardized part categories, with 4000 variations under each of the standard part classes. Smaller
84 datasets have been introduced in AAD [6] and ESB [31]. These datasets have been constructed by
85 selecting components from a certain number of classes.

86 The ABC dataset[35] is different, as it was obtained by scraping all the public data available from
87 OnShape. It contains more than one million mechanical components, and a large proportion of these
88 components are non-standard, which means that are not belonging to standardized categories. The
89 combination of massive scale and non-standard components makes it challenging to build a taxonomy
90 on the dataset, even for our expert annotators. The number of classes and objects in each class
91 does not fully represent the diversity of standard and non-standard parts seen in the product design
92 category. In certain object categories, specific annotation labels do not depend on the shape alone but
93 also its dimension and eventual intended application. For example, a fastener such as a washer and a
94 gasket may look exactly the same, but the label categories vary based on the material specification
95 and dimensions of the part. Annotation based on shape alone can lead to erroneous labeling and can
96 often confuse annotators particularly when the true design intent of the user is not known or when the
97 product assembly context is unknown.

98 Our Cluster3D dataset tackles this challenge directly, recasting the classification problem as an
99 unsupervised clustering problem. To the best of our knowledge, our dataset is the first specifically
100 designed for training and evaluating deep clustering methods on non-categorical 3D shapes.

101 **Interfaces for 3D models annotation.** Web-based platforms are commonly used for annotation
102 acquisition since they require no front-end installation by annotators and the cloud infrastructure can
103 support large 3D model datasets. MCB [34] developed a web-based platform with 3D viewers to
104 provide enough information of 3D objects for annotation. ShapeNet [10] and PartNet [43] present
105 web-based interface allows operating on 3D models and hierarchical 3D part annotation. The focus
106 of these interfaces is segmentation and classification. In this work, *we introduce a web interface that*
107 *enables efficient large-scale similarity annotation tailored for non-categorical datasets.*

108 **Unsupervised 3D representation learning.** Hand-crafted 3D descriptors has been studied as
109 geometry-based methods [56], view-based methods[45, 61, 20, 11], or hybrid methods [37]. For
110 3D deep representation learning, these methods can be classified into point cloud-based method,
111 view-based method, and volume-based method, depending on the different input data formats. Unlike
112 supervised 3D deep learning that requires class labels [50, 51, 39, 70, 66], self-supervised methods
113 are more suitable in our context. For example, Foldingnet [76], Atlasnet [24], TearingNet [47],
114 and [1, 80, 12] are a series of work investigating the autoencoder architecture to learn the latent
115 representation of point clouds. Rendered images from different views could also be used to learn
116 3D shape representations [54, 22, 26]. VConv-DAE [60] uses an autoencoder to learn the latent
117 representation of 3D objects with voxel as input. A 3D shape descriptor network was also proposed
118 to model volumetric represented objects[73]. Any of these descriptors, both hand-crafted and learned,
119 can be used with the deep clustering methods discussed next.

120 **Deep clustering methods.** Deep clustering adopts deep neural networks to learn clustering-friendly
121 representations [42] by integrating representation learning and clustering into an end-to-end model.
122 The optimizing objectives are the network parameters and the clustering results. For a deep neural
123 network, autoencoder-based models are widely used. DEC [28] is a classic deep clustering method:
124 it first pre-trains the autoencoder with network loss for a few epochs, then fine-tunes the encoder
125 network by optimizing KL divergence. DBC [38] achieves better clustering results compared to

126 DEC using a convolutional autoencoder instead of a feed-forward autoencoder. For DCN [74], the
127 autoencoder network is pre-trained first, then the autoencoder network is jointly optimized with the
128 K-means clustering results. Other methods, like DCC [58] and DEPICT [21], differ from DCN in
129 that they use different clustering loss functions. Generative model-based deep clustering, such as
130 variational autoencoder based method [33] and generative adversarial network based method [13],
131 have also been studied. SCAN [67] relies on a two-step process for representation learning and
132 clustering respectively. [32] maximizes mutual information for the clustering. DeepCluster uses the
133 initial clustering results as pseudo labels for supervised training [8].

134 All these methods have been designed for 2D images but are adaptable to work on 3D object dataset
135 by changing the features they use. We adapt 7 of these methods [24, 66, 23, 28, 67, 32, 8] to work on
136 3D models, and benchmark them on our newly introduced dataset for the first time.

137 **Clustering evaluation metrics.** To evaluate the success of a clustering algorithm and to enable
138 objective comparisons between different methods requires an evaluation metric. While the choice is
139 more obvious for supervised approaches, the situation is more challenging for clustering, especially
140 in our setting where the entities involved are 3D models, which lack a canonical representation.

141 An ideal clustering result should maximize the intercluster distance (*compactness*), and minimize
142 the intracluster distance (*separability*) [53]. The evaluation metrics can be classified into two major
143 classes: external validation indexes and internal validation indexes [25].

144 *External validation indexes* use previous knowledge about the data to evaluate the clustering re-
145 sults [25]. In the computer vision community, the most common indexes are unsupervised clustering
146 accuracy [77], normalized mutual information [69, 71], and adjusted rand index [29, 68]. Unsuper-
147 vised clustering accuracy is the equivalent of usual classification accuracy, with the difference that
148 it requires a mapping function to find the best mapping between the cluster assignment output and
149 ground truth labels. Normalized mutual information measures the mutual dependencies between the
150 cluster assignment and ground truth labels. Adjusted rand index is the corrected-for-chance version
151 of the rand index [52], which measures the similarity between two clustering by comparing the all
152 data pairs in the two clustering dataset. Besides, F-measure [65], entropy [59], purity [57] and other
153 metrics [64] can also be applied.

154 The annotation matrix in our dataset can serve as the previous knowledge of the data for external
155 evaluation. However, most of these indexes are using node labels rather than edge labels. Therefore,
156 besides using [3] as one of the evaluation metric in our dataset, we also propose another evaluation
157 metric to understand and analyze the baseline results, inspired by the purity index [57].

158 *Internal validation indexes* use the information intrinsic to the data [53], avoiding the need for
159 additional external information [79]. Internal criteria can be further divided into two research
160 topics [48]: 1) measurement of the fit of the cluster assignment and the inherent structure of the data
161 and 2) the stability of the clustering results [48]. To measure the fit between the cluster results and
162 structure of the data, compactness and separability is evaluated by computing the distances between
163 clustered samples using the Dunn [17], Davies-Bouldin [15], Calinski-Harabasz [7], silhouette
164 coefficient [55], and many other indexes in the literature [41]. An in-depth study on the stability of
165 the clustering results, we refer to [5, 44, 36].

166 In our dataset, by defining the distance as Chamfer distance [4] or Jaccard distance [18], we can use
167 these internal evaluation indexes. Particularly, we choose to use silhouette coefficient [55]. This
168 index does not require the cluster centroid, which is more appropriate for our clustering results.

169 3 Cluster3D Dataset

170 We view the Cluster3D dataset as an undirected complete graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$. The node set \mathcal{V} contains
171 each 3D CAD model in Cluster3D as a node v . Naturally, an edge $e_{i,j} \in \{+1, -1, 0\}$ in the edge set
172 \mathcal{E} stores the similarity annotation of two 3D CAD models v_i and v_j mentioned in the introduction.
173 The edge labels $+1$, -1 , and 0 respectively indicate similar, dissimilar, and unknown relationship
174 between two nodes. Next, we first explain how we create the Cluster3D dataset, and then discuss
175 several important design considerations.

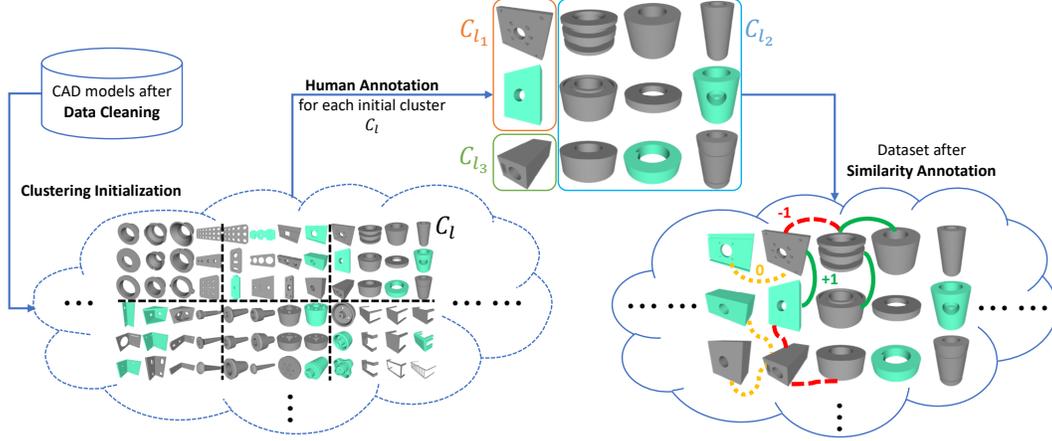


Figure 2: Cluster3D creation workflow.

176 3.1 Dataset creation workflow

177 We use the workflow illustrated in Figure 2 with the following major steps to create Cluster3D.

178 **Step 1: Data Cleaning.** We use the first four chunks of the ABC dataset [35]. We filter out all blank
 179 files and all files containing assemblies instead of single components, obtaining a dataset of 22, 968
 180 CAD models.

181 **Step 2: Similarity Annotation.** Since 3D CAD models are non-categorical, it is challenging to
 182 assign object class labels for each node in Cluster3D. Instead, we propose to manually annotate a
 183 small set of carefully sampled edges storing the pairwise CAD model similarity. This novel scalable
 184 and efficient edge-based annotation is divided into two steps.

185 *Step 2.1: Cluster Initialization.* Before manual annotation, we first grouped all the CAD models into
 186 a set of initial clusters $\{\mathcal{C}_k \mid \cup_{\forall k} \mathcal{C}_k = \mathcal{V}, \mathcal{C}_k \cap \mathcal{C}_l = \emptyset \forall k \neq l, \text{ and } |\mathcal{C}_k| \leq T, \forall k\}$, each containing
 187 no more than $T = 12$ CAD models, using a clustering method detailed in section 3.2. We then
 188 automatically assigned label 0 (meaning unknown) to all the edges across different clusters, i.e.,
 189 $e_{i,j} = 0 \iff v_i \in \mathcal{C}_k, v_j \in \mathcal{C}_l, \text{ and } k \neq l$.

190 *Step 2.2: Human Annotation.* We only manually annotate the edges residing inside the same initial
 191 cluster, i.e., $e_{i,j} \neq 0 \iff \exists k \in [1, K_I], v_i \in \mathcal{C}_k, v_j \in \mathcal{C}_k$, thus reducing considerably the
 192 human annotation cost. A number $A = 3$ of mechanical engineers served as our experts to provide
 193 their CAD model similarity annotations independently. For each one of the above initial clusters,
 194 e.g., \mathcal{C}_k , an annotator has to either *confirm* that CAD models inside \mathcal{C}_k are all similar to each
 195 other, or further *divide* the cluster into smaller clusters until such confirmations can be made for
 196 each smaller clusters. The *confirmation* of a cluster \mathcal{C}_l assigns all internal edges with the *positive*
 197 label +1, i.e., $e_{i,j} = +1 \iff v_i \in \mathcal{C}_l, v_j \in \mathcal{C}_l$. *Dividing* a cluster \mathcal{C}_l into smaller clusters
 198 $\{\mathcal{C}_{l_t} \mid t \in \mathbb{Z}^+, \cup_{\forall t} \mathcal{C}_{l_t} = \mathcal{C}_l\}$ assigns all edges across those small clusters with the *negative*
 199 label -1, i.e., $e_{i,j} = -1 \iff v_i \in \mathcal{C}_{l_t}, v_j \in \mathcal{C}_{l_s}, s \neq t$. We record each annotator independently, i.e., the a -th
 200 annotator’s annotation forms a complete edge set \mathcal{E}^a over the same node set.

201 3.2 Design Decisions on the Annotation Workflow

202 **Why manually annotate similarity?** Although it has been widely used in geometry processing and
 203 machine learning, the concept of similarity can be vague and ambiguous when applied to 3D CAD
 204 models in Cluster3D. To determine whether two 3D models are similar or not algorithmically, there
 205 are two main criteria: geometric distribution similarity [63, 27, 19, 46] and visual similarity [11].
 206 Yet for human beings, the mechanism to determine the similarity between two CAD models is often
 207 based on unconscious background knowledge [75, 16], which might be different from the similarity
 208 judgment encoded in existing algorithms [75, 14]. Therefore, even with mathematically defined 3D
 209 object similarity metrics, acquiring large-scale human annotations for pairwise CAD object similarity
 210 is still important to capture the underlying background knowledge.

211 **Reasons for cluster initialization.** We acknowledge that the choice of cluster initialization could
212 introduce bias in our data collection, as a different clustering method would lead to a different subset
213 of annotated edges. However, we argue that it is unavoidable, due to the sheer size of the dataset: it is
214 impossible for human experts to annotate all the similarity relationships between one CAD model
215 and all the remaining CAD models, as the quadratic number of edges is intractable. It would take 16
216 years of annotation time, assuming a 1 second time to annotate each edge in Cluster3D.

217 The most obvious, and unbiased approach, to restrict the manual annotation to a subset of the edges
218 would be to do random sampling. This is however not an option for Cluster3D, as the sampling would
219 be too unbalanced, as most of the edges indicate dissimilarity between objects. We needed a strategy
220 that would give us a more or less even split between similar and dissimilar edges so that we could use
221 the annotations to evaluate clustering methods in a balanced way.

222 After experimenting with different approaches, we found a method that, in our dataset, leads to a
223 reasonable 1 to 1 ratio of similar and dissimilar edges: we use a clustering algorithm to overcluster
224 the dataset in small clusters of 12 objects, and then ask users to annotate all edges within each cluster.
225 Overall, this approach allowed us to get a good distribution of edge labels while annotating only
226 0.5% of the entries in our similarity matrix, making the annotation problem tractable with our budget
227 and resources.

228 **Details for cluster initialization.** For cluster initialization, we use the MVCNN [66] based method.
229 We opted for this method as it is the only image-based clustering method in our baselines, and in
230 this way, we can do a fair comparison of the remaining six methods that are all using a point cloud
231 representation.

232 Specifically, we first generate 12 images for all 22,968 CAD models in our dataset, following the
233 original settings in [66]. We use all these $12 \times 22,968$ images to train a convolutional auto-encoder
234 network. Then the trained encoder is used to extract features for all these images. For each CAD
235 model, we concatenate the twelve latent vectors from its corresponding 12 images to represent
236 its features. Finally, 22,968 features representing all the CAD models are clustered by KMeans
237 algorithm. With the K number in KMeans setting to be 2,000, we have two thousand initial clusters
238 for the human annotators. We continue to split the clusters with more than 12 models in the class
239 using KMeans, until the contained number of models is not greater than 12. These clustered CAD
240 models can be loaded into our database for annotation.

241 **Annotation interface.** We developed a web-based annotation application. The interface shows CAD
242 models of one cluster at a time: it shows the 12 CAD models with checkbox, and initially all 12
243 checkbox are set to be checked. The annotators manually unmark the models which are considered
244 dissimilar from the others, effectively annotating all edges linking the 12 models in the cluster. After
245 confirmation, a new set of 12 models is shown.

246 **Conflicts in annotations.** We use one single annotated similarity matrix as the final outcome of our
247 annotation procedure. In case of conflicts between different annotators, the majority wins. Note that
248 for the final evaluation we also consider the individual similarity matrices of the different annotators.

249 **Data statistics.** Cluster3D has 22,968 number of CAD models; therefore, totally there will be
250 $22,968 \times 22,968$ number of similar or dissimilar edges between every two CAD models. Among
251 them, 275,616 edges are labeled by three human experts respectively. For the first annotator, 155,960
252 edges are labeled as 1, representing these CAD model pairs are similar, and 119,656 are labeled
253 as -1 , meaning these pairs are dissimilar. For the second and third annotator, they have labeled
254 130,442, 145,174 similar edges, and 205,582, 70,034 dissimilar edges respectively. We also check
255 the consistency of the three annotators' labeling. Here consistency means the three annotator's label
256 for a specific edge is the same. The total number of consistent label is 172,554, occupying 62.6%
257 of the labeled edges. In the next release, we will double the number A and we expect to see higher
258 consistency among the annotations.

259 4 Cluster3D Benchmark

260 4.1 Baseline methods

261 We adapt seven baseline methods to establish a benchmark for clustering algorithms. We divide
262 these baseline methods into two types: 1) two-stage clustering, and 2) end-to-end deep clustering.

263 Two-stage clustering methods use a deep neural network to extract features for all CAD models,
 264 then apply a traditional clustering algorithm, such as KMeans. End-to-end deep clustering baseline
 265 methods integrate feature extraction and clustering in one framework: during the training process
 266 both network loss and clustering loss are minimized. Note that all these methods are considered as
 267 partitional clustering [9], i.e. one CAD model will only fall into one cluster.

268 Since some of the baseline methods are designed for 2D images, we adapt them for 3D CAD models.
 269 We use either point cloud or multi-view images as the representation format, and select suitable deep
 270 neural networks. A detailed description of all networks can be found in the supplementary.

271 **Two-stage clustering.**

272 *MVCNN*: We describe this algorithm in Section 3.

273 *AtlasNet*: To compute cluster of 3D point clouds using Atlasnet, we follow the original auto-encoder
 274 architecture to reconstruct 3D point cloud for each input 3D CAD object, and then predict latent
 275 vectors based on the encoder of the trained model. The CAD objects are clustered by using KMeans
 276 on the obtained latent features.

277 *BYOL*: BYOL is proposed to compute self-supervised image representation learning. We replace the
 278 image encoder (ResNet) with a point cloud encoder (PointNet) to learn a representation of a 3D CAD
 279 shape. We then apply the KMeans algorithm on the learned latent representation to cluster CAD
 280 objects.

281 **End-to-end deep clustering.**

282 *DEC*: To adapt the DEC algorithm, we initialize DEC with the AtlasNet architecture to auto-encode
 283 3D point clouds as the input data. The deep auto-encoder is trained to minimize Chamfer loss and
 284 learns representations of the 3D shapes. We then follow the DEC algorithm by discarding the decoder
 285 layers and use the encoder layers as the initial mapping between the data and feature space. This
 286 is followed by joint optimization of the cluster centers and encoder parameters using SGD with
 287 momentum.

288 *DeepCluster*: We replace the convolution networks trained by the DeepClustering algorithm to use
 289 PointNet instead for encoding the point cloud data to predict cluster assignments. The algorithm is
 290 followed by alternating between clustering of the point cloud feature descriptors using K-Means and
 291 training the PointNet network using the multinomial logistic loss function.

292 *IIC*: Instead of the original IIC method for unsupervised image semantic task, we first randomly
 293 transform a CAD model to a pair of point clouds, and use PointNet as encoder to maximize mutual
 294 information between the class assignments of each pair. The trained model directly outputs class
 295 labels for each 3D CAD model.

296 *SCAN*: We adjust the pretext stage: Instead of using noise contrastive estimation (NCE) to determine
 297 the nearest neighbors, we use the auto-encoder of AtlasNet we have trained to output the feature
 298 vectors to generate the nearest neighbors set.

299 **4.2 Evaluation Metrics**

300 As discussed in Section 2, external and internal indexes are used for evaluating clustering results.

301 **External validation indexes.** We evaluate the clustering results of the baseline methods using three
 302 external validation indexes: *pair-wise accuracy*, *intra-cluster purity*, and *inter-cluster purity*.

303 *Pair-wise accuracy*. It is natural to compare the pair-wise clustering results with the annotated simi-
 304 larity matrix, which is the evaluation metric in correlation clustering [3]. We note that the similarity
 305 matrix is an undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$ on N nodes. Let e_{ij} denote the label of the edge relationship
 306 between object i, j , and $e_{ij} = e_{ji}$. $E = \{e_{ij}\}$ denote all the edges. $G' = (V', E')$ is the subgraph of
 307 of G , which is only composed of the known labels. $E' = \{e_{ij} | e_{ij} = 1 \vee e_{ij} = -1, e_{ij} \in E\}$. For
 308 the clustering results obtained from the baseline methods, \hat{e}_{ij} denote the clustered edge relationship
 309 between object i, j . If objects i, j are grouped into the same cluster, we assume the two objects
 310 are similar, therefore $\hat{e}_{ij} = 1$. Otherwise $\hat{e}_{ij} = -1$. The pair-wise accuracy is defined as: $\text{acc} =$
 311 $\sum_{e_{ij} \in E'} \frac{|\hat{e}_{ij} - e_{ij}|}{2n(E')}$, where $n(E')$ is the number of elements in E' . The range of the pair-wise accuracy
 312 is $[0, 1]$.

313 *Intra-cluster purity.* The pair-wise accuracy might miss information, since the clustering performance
 314 can be also evaluated cluster-wise [78]. The concept of purity [57] in node-labeled clustering
 315 evaluation inspired us, as it is used to evaluate the extent at which one cluster contains one single
 316 class. Similarly, we propose intra-cluster purity to detect the false positive edges in each cluster.
 317 Intuitively, the intra-cluster purity metric measures the extent at which each cluster contains similar
 318 edges. Let K denote the number of cluster, $S = \{s_1, s_2, \dots, s_K\}$ denote the set of the number of
 319 labeled similar edges in each cluster, $T = \{t_1, t_2, \dots, t_K\}$ denote the set of the number of all known
 320 edges in each cluster. $T' = \{t'_i | t_i \neq 0, t_i \in T\}$ denote the subset of T , where only those clusters
 321 with at least one labeled edge are considered. Intra-cluster purity is defined as: $\frac{1}{n(T')} \sum_{i=1}^{n(T')} \left| \frac{s_i}{t'_i} \right|$.
 322 The range of inter-cluster purity is $[0, 1]$; 0 means the worst clustering result, and 1 means the best
 323 clustering result.

324 **Internal validation indexes.** It is not meaningful to use cluster centroids in our case, since the
 325 features of the CAD models are in high dimension. Therefore, we opt for the silhouette coefficient [55]
 326 method, which is widely used and does not require cluster centroids. Based on its definition, we need
 327 to determine the distance between every two objects. In our dataset, the objects are 3D CAD models
 328 which can be represented as point cloud or voxel. Therefore, we choose to use Chamfer distance [4]
 329 and Jaccard distance [18] as two distances between every two CAD models.

330 5 Benchmark Results and Discussions

331 **Experiment settings.** All the baseline methods are implemented using PyTorch [49] and run
 332 on an NVIDIA GeForce GTX 1080 Ti GPU. For hyperparameter settings, we tune learning
 333 rate and batch size for each baseline method. The learning rates for MVCNN-based
 334 method, Atlasnet-based method, BYOL-based method, DEC, DeepCluster, IIC, and SCAN are
 335 0.0001, 0.001, 0.0003, 0.00001, 0.05, 0.0003, 0.0001 respectively. The batch sizes for these methods
 336 are 60, 11, 10, 128, 50, 10, 96 respectively.

337 **Clustering results using external evaluation metrics.** Figure 3-(Pair-wise accuracy) shows the
 338 benchmark results on the Cluster3D dataset, using the *pair-wise accuracy* and *intra-cluster purity*
 339 as external evaluation metrics. Since we do not know a priori the number of clusters in ABC, we
 340 test our baseline methods on seven different number of K , from 32 to 2,000, following exponential
 341 growth.

342 *All baseline methods perform well with respect to intra-cluster accuracy, but their pair-wise accuracy*
 343 *is much lower.* Figure 3-(Intra-cluster purity) shows that all baseline methods can achieve *intra-cluster*
 344 *accuracy* higher than 0.8 for all K . Some of the baseline methods can even achieve accuracy higher
 345 than 0.9 when K is no less than 128. However, for *pair-wise accuracy*(Figure 3-(pair-wise accuracy))
 346 results, it reveals that all these deep neural networks do not obtain enough ability to group similar
 347 CAD models, with their accuracy lower than 0.7.

348 *Sensitivity to the cluster number K .* Figure 3-(Pair-wise accuracy) shows that most of the baseline
 349 methods' performances decrease as K increases. We hypothesize that it is due to the imbalance of
 350 the annotated similarity matrix. There are a total of 431,576 edges labeled as similar, and 395,272
 351 edges labeled as dissimilar for the three annotators. Therefore, we think the annotated similarity
 352 matrix might be biased towards the clustering results which have more similar pair predictions. For
 353 each baseline method, if the number of clusters K increases, the CAD models will be more separate,
 354 causing more dissimilar pair predictions.

355 Figure 3-(Intra-cluster purity) shows that larger number of K increase most baseline methods'
 356 performances. We believe it is because of the definition of *intra-cluster accuracy*. When K becomes
 357 larger, the CAD models will be grouped into more clusters, which causes each cluster to be more
 358 pure. By definition, the more pure the clusters are, the higher the *intra-cluster accuracy* will be.

359 *Surprisingly, end-to-end deep clustering methods do not outperform two-stage clustering methods.*
 360 As shown in Figure 3-(Pairwise accuracy), there is no obvious evidence showing higher performances
 361 of end-to-end methods (DeepCluster, DEC, IIC, SCAN), compared to two-stage clustering methods
 362 (AtlasNet-based method, BYOL-based method). Therefore, we believe it is necessary to study how
 363 to take the advantage of the clustering loss when we are training a deep neural network.

364 *MVCNN-based method is the only image-based method, and it was used during annotation: we*
 365 *believe these are the reason why it behaves noticeably differently than the other methods based on*
 366 *point clouds.* Figure 3-(Intra-cluster purity) shows the performance of MVCNN-based method is
 367 significantly lower than all other methods. Also, the MVCNN-based method is not as sensitive to the
 368 K value as other methods.

369 *It still requires effort to study why some baseline methods perform differently than the overall trend*
 370 *with other methods.* Although most of the baseline methods show the same trend when the number of
 371 K increases, SCAN and IIC are different. For SCAN, the input K is used as the maximum number
 372 of clusters. Indeed, the actual number of cluster is often smaller than K , and it might be the reason
 373 that SCAN performances differently.

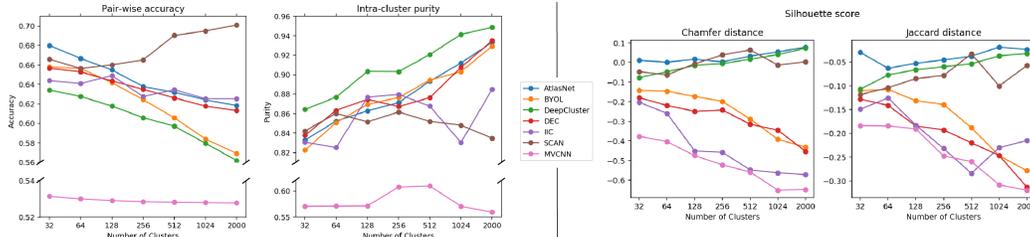


Figure 3: Cluster3D benchmark results.

374 **Clustering results using internal evaluation metric.** Figure 3-(Chamfer distance) shows the bench-
 375 mark results using the *silhouette score* as internal evaluation metric. Using Chamfer distance and
 376 Jaccard distance lead to similar performances.

377 For all methods, we find that the clusters are not obvious different or even wrongly assigned,
 378 since most of the silhouette score is below 0. Second, the AtlasNet-based method, DEC, SCAN
 379 methods perform better when the number of cluster K increases, while other baseline methods show
 380 the opposite trend. Future investigation should be conducted to further understand this peculiar
 381 phenomenon.

382 5.1 Limitations and discussion

383 The major challenge in our study is the very high cost of annotating a similarity matrix which has a
 384 quadratic number of entries with respect to the number of objects in the dataset. We introduced a
 385 technique to reduce the annotation cost, but it is possible that the filtering introduced a bias in the
 386 annotations. This bias could be reduced by picking multiple initial clustering methods, which we
 387 plan to explore in the future.

388 The different evaluation metrics lead to different ranking for these baseline methods, suggesting that
 389 they evaluate different criteria. Identifying which metric is best for specific applications would be
 390 crucial to guide the development of clustering algorithms, and we believe it is an interesting venue
 391 for future work in deep clustering of 3D CAD models.

392 6 Conclusion

393 Cluster3D is a manually annotated dataset for the development and evaluation of clustering
 394 algorithms on 3D CAD models. We introduce the dataset, two external evaluation metrics
 395 based on the matrix, and benchmarked seven state-of-the-art clustering methods. Our conclu-
 396 sion is that the gap between human annotators and state-of-the-art methods is large: we be-
 397 lieve our dataset will be an important resource to improve clustering methods for 3D geome-
 398 try.

399 References

400 [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and
 401 generative models for 3d point clouds. In *International conference on machine learning*, pages 40–49.
 402 PMLR, 2018. 3

- 403 [2] Binil Starly Atin Angrish, Akshay Bharadwaj. Mvncnn++: Computer-aided design model shape classifica-
404 tion and retrieval using multi-view convolutional neural networks. *J. Comput. Inf. Sci. Eng.*, 21(1): 011001,
405 2021. 3
- 406 [3] Nikhil Bansal, Avrim Blum, and Shuchi Chawla. Correlation clustering. *Machine learning*, 56(1):89–113,
407 2004. 4, 7
- 408 [4] Harry G Barrow, Jay M Tenenbaum, Robert C Bolles, and Helen C Wolf. Parametric correspondence and
409 chamfer matching: Two new techniques for image matching. Technical report, SRI INTERNATIONAL
410 MENLO PARK CA ARTIFICIAL INTELLIGENCE CENTER, 1977. 4, 8
- 411 [5] Asa Ben-Hur and Isabelle Guyon. Detecting stable clusters using principal component analysis. In
412 *Functional genomics*, pages 159–182. Springer, 2003. 4
- 413 [6] Dmitriy Bespalov, Cheuk Yiu Ip, William C Regli, and Joshua Shaffer. Benchmarking cad search
414 techniques. In *Proceedings of the 2005 ACM symposium on Solid and physical modeling*, pages 275–286,
415 2005. 3
- 416 [7] Tadeusz Caliński and Jerzy Harabasz. A dendrite method for cluster analysis. *Communications in*
417 *Statistics-theory and Methods*, 3(1):1–27, 1974. 4
- 418 [8] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised
419 learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages
420 132–149, 2018. 4
- 421 [9] M Emre Celebi. *Partitional clustering algorithms*. Springer, 2014. 7
- 422 [10] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio
423 Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository.
424 *arXiv preprint arXiv:1512.03012*, 2015. 3
- 425 [11] Ding-Yun Chen, Xiao-Pei Tian, Yu-Te Shen, and Ming Ouhyoung. On visual similarity based 3d model
426 retrieval. In *Computer graphics forum*, volume 22, pages 223–232. Wiley Online Library, 2003. 3, 5
- 427 [12] Siheng Chen, Chaojing Duan, Yaoqing Yang, Duanshun Li, Chen Feng, and Dong Tian. Deep unsupervised
428 learning of 3d point clouds via graph topology inference and filtering. *IEEE Transactions on Image*
429 *Processing*, 29:3183–3198, 2019. 3
- 430 [13] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel. Infogan:
431 Interpretable representation learning by information maximizing generative adversarial nets. *arXiv preprint*
432 *arXiv:1606.03657*, 2016. 4
- 433 [14] Florin Cutzu and Shimon Edelman. Representation of object similarity in human vision: psychophysics
434 and a computational model. *Vision research*, 38(15-16):2229–2257, 1998. 5
- 435 [15] David L Davies and Donald W Bouldin. A cluster separation measure. *IEEE transactions on pattern*
436 *analysis and machine intelligence*, (2):224–227, 1979. 4
- 437 [16] Hans P Op de Beeck, Katrien Torfs, and Johan Wagemans. Perceived shape similarity among unfamiliar
438 objects and the organization of the human object vision pathway. *Journal of Neuroscience*, 28(40):10111–
439 10123, 2008. 5
- 440 [17] Joseph C Dunn. Well-separated clusters and optimal fuzzy partitions. *Journal of cybernetics*, 4(1):95–104,
441 1974. 4
- 442 [18] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew
443 Zisserman. The pascal visual object classes challenge: A retrospective. *International journal of computer*
444 *vision*, 111(1):98–136, 2015. 4, 8
- 445 [19] Thomas Funkhouser, Patrick Min, Michael Kazhdan, Joyce Chen, Alex Halderman, David Dobkin, and
446 David Jacobs. A search engine for 3d models. *ACM Transactions on Graphics (TOG)*, 22(1):83–105, 2003.
447 5
- 448 [20] Yue Gao, Qionghai Dai, Meng Wang, and Naiyao Zhang. 3d model retrieval using weighted bipartite graph
449 matching. *Signal Processing: Image Communication*, 26(1):39–47, 2011. 3
- 450 [21] Kamran Ghasedi Dizaji, Amirhossein Herandi, Cheng Deng, Weidong Cai, and Heng Huang. Deep
451 clustering via joint convolutional autoencoder embedding and relative entropy minimization. In *Proceedings*
452 *of the IEEE international conference on computer vision*, pages 5736–5745, 2017. 4

- 453 [22] Rohit Girdhar, David F Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and
454 generative vector representation for objects. In *European Conference on Computer Vision*, pages 484–499.
455 Springer, 2016. 3
- 456 [23] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H Richemond, Elena Buchatskaya,
457 Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap
458 your own latent: A new approach to self-supervised learning. *arXiv preprint arXiv:2006.07733*, 2020. 4
- 459 [24] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché
460 approach to learning 3d surface generation. In *Proceedings of the IEEE conference on computer vision
461 and pattern recognition*, pages 216–224, 2018. 3, 4
- 462 [25] Maria Halkidi, Yannis Batistakis, and Michalis Vazirgiannis. Cluster validity methods: part i. *ACM Sigmod
463 Record*, 31(2):40–45, 2002. 4
- 464 [26] Zhizhong Han, Mingyang Shang, Yu-Shen Liu, and Matthias Zwicker. View inter-prediction gan: Unsu-
465 pervised representation learning for 3d shapes by learning global shape memories to support local view
466 predictions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8376–8384,
467 2019. 3
- 468 [27] Masaki Hilaga, Yoshihisa Shinagawa, Taku Kohmura, and Toshiyasu L Kunii. Topology matching for fully
469 automatic similarity estimation of 3d shapes. In *Proceedings of the 28th annual conference on Computer
470 graphics and interactive techniques*, pages 203–212, 2001. 5
- 471 [28] Peihao Huang, Yan Huang, Wei Wang, and Liang Wang. Deep embedding network for clustering. In *2014
472 22nd International conference on pattern recognition*, pages 1532–1537. IEEE, 2014. 3, 4
- 473 [29] Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of classification*, 2(1):193–218, 1985.
474 4
- 475 [30] Krishna Murthy Jatavallabhula, Edward Smith, Jean-Francois Lafleche, Clement Fuji Tsang, Artem
476 Rozantsev, Wenzheng Chen, Tommy Xiang, Rev Lebededian, and Sanja Fidler. Kaolin: A pytorch library
477 for accelerating 3d deep learning research. *arXiv:1911.05063*, 2019. 15
- 478 [31] Subramaniam Jayanti, Yagnanarayanan Kalyanaraman, Natraj Iyer, and Karthik Ramani. Developing an
479 engineering shape benchmark for cad models. *Computer-Aided Design*, 38(9):939–953, 2006. 3
- 480 [32] Xu Ji, João F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image
481 classification and segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer
482 Vision*, pages 9865–9874, 2019. 4
- 483 [33] Zhuxi Jiang, Yin Zheng, Huachun Tan, Bangsheng Tang, and Hanning Zhou. Variational deep embedding:
484 An unsupervised and generative approach to clustering. *arXiv preprint arXiv:1611.05148*, 2016. 4
- 485 [34] Sangpil Kim, Hyung-gun Chi, Xiao Hu, Qixing Huang, and Karthik Ramani. A large-scale annotated
486 mechanical components benchmark for classification and retrieval tasks with deep neural networks. In
487 *Proceedings of 16th European Conference on Computer Vision (ECCV)*, 2020. 3
- 488 [35] Sebastian Koch, Albert Matveev, Zhongshi Jiang, Francis Williams, Alexey Artemov, Evgeny Burnaev,
489 Marc Alexa, Denis Zorin, and Daniele Panozzo. Abc: A big cad model dataset for geometric deep
490 learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
491 9601–9611, 2019. 2, 3, 5
- 492 [36] Tilman Lange, Volker Roth, Mikio L Braun, and Joachim M Buhmann. Stability-based validation of
493 clustering solutions. *Neural computation*, 16(6):1299–1323, 2004. 4
- 494 [37] Bo Li and Henry Johan. 3d model retrieval using hybrid features and class information. *Multimedia tools
495 and applications*, 62(3):821–846, 2013. 3
- 496 [38] Fengfu Li, Hong Qiao, and Bo Zhang. Discriminatively boosted image clustering with fully convolutional
497 auto-encoders. *Pattern Recognition*, 83:161–173, 2018. 3
- 498 [39] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on
499 x-transformed points. *Advances in neural information processing systems*, 31:820–830, 2018. 3
- 500 [40] James MacQueen et al. Some methods for classification and analysis of multivariate observations. In
501 *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages
502 281–297. Oakland, CA, USA, 1967. 2

- 503 [41] Glenn W Milligan and Martha C Cooper. An examination of procedures for determining the number of
504 clusters in a data set. *Psychometrika*, 50(2):159–179, 1985. 4
- 505 [42] Erxue Min, Xifeng Guo, Qiang Liu, Gen Zhang, Jianjing Cui, and Jun Long. A survey of clustering with
506 deep learning: From the perspective of network architecture. *IEEE Access*, 6:39501–39514, 2018. 3
- 507 [43] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su.
508 Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In
509 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 909–918,
510 2019. 3
- 511 [44] G Bel Mufti, P Bertrand, and EL Moubarki. Determining the number of groups from measures of cluster
512 stability. In *Proceedings of international symposium on applied stochastic models and data analysis*, pages
513 17–20, 2005. 4
- 514 [45] Ryutarou Ohbuchi, Kunio Osada, Takahiko Furuya, and Tomohisa Banno. Salient local visual features
515 for shape-based 3d model retrieval. In *2008 IEEE International Conference on Shape Modeling and*
516 *Applications*, pages 93–102. IEEE, 2008. 3
- 517 [46] Robert Osada, Thomas Funkhouser, Bernard Chazelle, and David Dobkin. Shape distributions. *ACM*
518 *Transactions on Graphics (TOG)*, 21(4):807–832, 2002. 5
- 519 [47] Jiahao Pang, Duanshun Li, and Dong Tian. Tearingnet: Point cloud autoencoder to learn topology-friendly
520 representations. *arXiv preprint arXiv:2006.10187*, 2020. 3
- 521 [48] Damaris Pascual, Filiberto Pla, and J Salvador Sánchez. Cluster validation using information stability
522 measures. *Pattern Recognition Letters*, 31(6):454–461, 2010. 4
- 523 [49] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
524 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang,
525 Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie
526 Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In
527 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in*
528 *Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. 8
- 529 [50] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d
530 classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern*
531 *recognition*, pages 652–660, 2017. 3
- 532 [51] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on
533 point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017. 3
- 534 [52] William M Rand. Objective criteria for the evaluation of clustering methods. *Journal of the American*
535 *Statistical association*, 66(336):846–850, 1971. 4
- 536 [53] Eréndira Rendón, Itzel Abundez, Alejandra Arizmendi, and Elvia M Quiroz. Internal versus external
537 cluster validation indexes. *International Journal of computers and communications*, 5(1):27–34, 2011. 2, 4
- 538 [54] Danilo Jimenez Rezende, SM Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas
539 Heess. Unsupervised learning of 3d structure from images. *arXiv preprint arXiv:1607.00662*, 2016. 3
- 540 [55] Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis.
541 *Journal of computational and applied mathematics*, 20:53–65, 1987. 4, 8
- 542 [56] Dietmar Saupe and Dejan V Vranić. 3d model retrieval with spherical harmonics and moments. In *Joint*
543 *Pattern Recognition Symposium*, pages 392–397. Springer, 2001. 3
- 544 [57] Hinrich Schütze, Christopher D Manning, and Prabhakar Raghavan. *Introduction to information retrieval*,
545 volume 39. Cambridge University Press Cambridge, 2008. 4, 8
- 546 [58] Sohil Atul Shah and Vladlen Koltun. Deep continuous clustering. *arXiv preprint arXiv:1803.01449*, 2018.
547 4
- 548 [59] Claude E Shannon. A mathematical theory of communication. *The Bell system technical journal*,
549 27(3):379–423, 1948. 4
- 550 [60] Abhishek Sharma, Oliver Grau, and Mario Fritz. Vconv-dae: Deep volumetric shape learning without
551 object labels. In *European Conference on Computer Vision*, pages 236–250. Springer, 2016. 3

- 552 [61] Jau-Ling Shih, Chang-Hsing Lee, and Jian Tang Wang. A new 3d model retrieval approach based on the
553 elevation descriptor. *Pattern Recognition*, 40(1):283–295, 2007. 3
- 554 [62] Philip Shilane, Patrick Min, Michael Kazhdan, and Thomas Funkhouser. The princeton shape benchmark.
555 In *Proceedings Shape Modeling Applications, 2004.*, pages 167–178. IEEE, 2004. 3
- 556 [63] Heung-Yeung Shum, Martial Hebert, and Katsushi Ikeuchi. On 3d shape similarity. In *Proceedings CVPR
557 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 526–531. IEEE,
558 1996. 5
- 559 [64] Satya Chaitanya Sripada and M Sreenivasa Rao. Comparison of purity and entropy of k-means clustering
560 and fuzzy c means clustering. *Indian journal of computer science and engineering*, 2(3):343–346, 2011. 4
- 561 [65] Michael Steinbach, George Karypis, and Vipin Kumar. A comparison of document clustering techniques.
562 2000. 4
- 563 [66] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional
564 neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer
565 vision*, pages 945–953, 2015. 3, 4, 6
- 566 [67] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool.
567 Scan: Learning to classify images without labels. In *European Conference on Computer Vision*, pages
568 268–285. Springer, 2020. 4
- 569 [68] Nguyen Xuan Vinh, Julien Epps, and James Bailey. Information theoretic measures for clusterings
570 comparison: Variants, properties, normalization and correction for chance. *The Journal of Machine
571 Learning Research*, 11:2837–2854, 2010. 4
- 572 [69] Paul Viola and William M Wells III. Alignment by maximization of mutual information. *International
573 journal of computer vision*, 24(2):137–154, 1997. 4
- 574 [70] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon.
575 Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019.
576 3
- 577 [71] William M Wells III, Paul Viola, Hideki Atsumi, Shin Nakajima, and Ron Kikinis. Multi-modal volume
578 registration by maximization of mutual information. *Medical image analysis*, 1(1):35–51, 1996. 4
- 579 [72] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao.
580 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on
581 computer vision and pattern recognition*, pages 1912–1920, 2015. 3
- 582 [73] Jianwen Xie, Zilong Zheng, Ruiqi Gao, Wenguan Wang, Song-Chun Zhu, and Ying Nian Wu. Learning
583 descriptor networks for 3d shape synthesis and analysis. In *Proceedings of the IEEE conference on
584 computer vision and pattern recognition*, pages 8629–8638, 2018. 3
- 585 [74] Bo Yang, Xiao Fu, Nicholas D Sidiropoulos, and Mingyi Hong. Towards k-means-friendly spaces:
586 Simultaneous deep learning and clustering. In *international conference on machine learning*, pages
587 3861–3870. PMLR, 2017. 4
- 588 [75] Cong Yang and Marcin Grzegorzec. Object similarity by humans and machines. In *2014 AAAI Fall
589 Symposium Series*, 2014. 5
- 590 [76] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid
591 deformation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 206–215, 2018. 3
- 592 [77] Yi Yang, Dong Xu, Feiping Nie, Shuicheng Yan, and Yueting Zhuang. Image clustering using local
593 discriminant models and global integration. *IEEE Transactions on Image Processing*, 19(10):2761–2773,
594 2010. 4
- 595 [78] Yujing Zeng, Jianshan Tang, Javier Garcia-Frias, and Guang R Gao. An adaptive meta-clustering approach:
596 combining the information from different clustering results. In *Proceedings. IEEE Computer Society
597 Bioinformatics Conference*, pages 276–287. IEEE, 2002. 8
- 598 [79] Qinpei Zhao. *Cluster validity in clustering methods*. PhD thesis, Itä-Suomen yliopisto, 2012. 4
- 599 [80] Yongheng Zhao, Tolga Birdal, Haowen Deng, and Federico Tombari. 3d point capsule networks. In
600 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1009–1018,
601 2019. 3

602 **Checklist**

- 603 1. For all authors...
- 604 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
605 contributions and scope? [Yes]
- 606 (b) Did you describe the limitations of your work? [Yes] See discussions in Section 5.
- 607 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 608 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
609 them? [Yes]
- 610 2. If you are including theoretical results...
- 611 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 612 (b) Did you include complete proofs of all theoretical results? [N/A]
- 613 3. If you ran experiments...
- 614 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
615 perimental results (either in the supplemental material or as a URL)? [Yes] See the
616 provided URL under the title.
- 617 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
618 were chosen)? [Yes] See Section 5.
- 619 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
620 ments multiple times)? [No] We did not repeat experiments for multiple times since
621 experiments for all benchmark methods is time-consuming.
- 622 (d) Did you include the total amount of compute and the type of resources used (e.g., type
623 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5.
- 624 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 625 (a) If your work uses existing assets, did you cite the creators? [Yes] See Section 3.
- 626 (b) Did you mention the license of the assets? [Yes] The assets will be released under the
627 MIT license. We adopted the same license as the ABC dataset.
- 628 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 629 (d) Did you discuss whether and how consent was obtained from people whose data you're
630 using/curating? [No] The data is public.
- 631 (e) Did you discuss whether the data you are using/curating contains personally identifiable
632 information or offensive content? [No] The data we are using do not contain personally
633 identifiable information or offensive content.
- 634 5. If you used crowdsourcing or conducted research with human subjects...
- 635 (a) Did you include the full text of instructions given to participants and screenshots, if
636 applicable? [Yes] See Section A.
- 637 (b) Did you describe any potential participant risks, with links to Institutional Review
638 Board (IRB) approvals, if applicable? [No]
- 639 (c) Did you include the estimated hourly wage paid to participants and the total amount
640 spent on participant compensation? [Yes] See Section A.