# Unsupervised Detection of Distinctive Regions on 3D Shapes

XIANZHI LI, LEQUAN YU, and CHI-WING FU, The Chinese University of Hong Kong
DANIEL COHEN-OR, Tel Aviv University
PHENG-ANN HENG, The Chinese University of Hong Kong

This paper presents a novel approach to learn and detect distinctive regions on 3D shapes. Unlike previous works, which require labeled data, our method is unsupervised. We conduct the analysis on point sets sampled from 3D shapes, then formulate and train a deep neural network for an unsupervised shape clustering task to learn local and global features for distinguishing shapes with respect to a given shape set. To drive the network to learn in an unsupervised manner, we design a clustering-based nonparametric softmax classifier with an iterative re-clustering of shapes, and an adapted contrastive loss for enhancing the feature embedding quality and stabilizing the learning process. By then, we encourage the network to learn the point distinctiveness on the input shapes. We extensively evaluate various aspects of our approach and present its applications for distinctiveness-guided shape retrieval, sampling, and view selection in 3D scenes.

CCS Concepts: • **Computing methodologies** → **Neural networks**; **Shape analysis**.

Additional Key Words and Phrases: shape analysis, unsupervised, learning, neural network, distinctive regions

## 1 INTRODUCTION

Reasoning about distinctive regions on 3D shapes has a wide range of applications in computer graphics and geometric processing, *e.g.*, object retrieval [Gal and Cohen-Or 2006; Shilane and Funkhouser 2006], shape matching [Castellani et al. 2008; Shilane and Funkhouser 2007], and view selection [Lee et al. 2005; Leifman et al. 2012]. In this work, we follow the definition of *distinctive* regions proposed by Shilane and Funkhouser [2006; 2007], *i.e.*, the distinction of a surface region in an object is defined as

> *how useful the region is for distinguishing the object from others of different types*[1].

---

[1] Type refers to the specific class that an object belongs to, *e.g.*, chair, table, and car.

Authors' addresses: Xianzhi Li; Lequan Yu; Chi-Wing Fu, The Chinese University of Hong Kong, {xzli,lqyu,cwfu}@cse.cuhk.edu.hk; Daniel Cohen-Or, Tel Aviv University, dcor@mail.tau.ac.il; Pheng-Ann Heng, The Chinese University of Hong Kong, pheng@cse.cuhk.edu.hk.
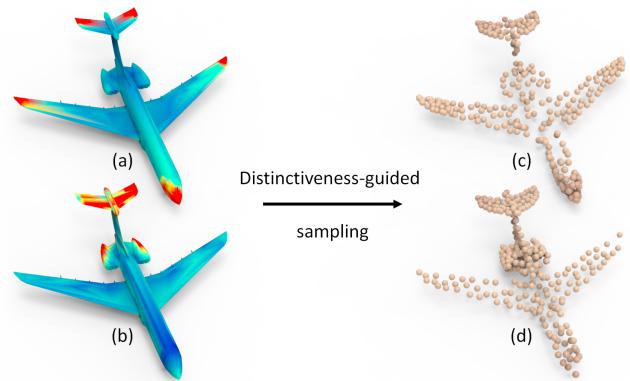
Fig. 1. Distinctive regions detected by our method on the same shape shown in (a,b), relative to (a) an inter-class dataset (i.e., the whole ModelNet40 dataset) and (b) an intra-class dataset (i.e., only the tail-engine and four-engine airplanes in ModelNet40), with their corresponding distinctiveness-guided sampling results (c,d). The colors in (a,b) indicate the region distinctiveness with red color being the most distinctive.
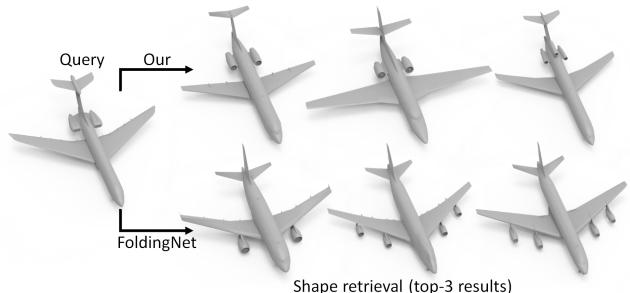


Fig. 2. Top-3 Shape retrieval results using our intra-class distinctiveness-guided retrieval method (top row) and using FoldingNet [Yang et al. 2018], a general unsupervised feature learning method (bottom row). Conventional methods can only retrieve shapes with similar overall appearance, while *guided by distinctiveness* enables the retrieval of shapes with similar distinctive features, *e.g.*, the two engines at the back.

Hence, distinctive regions of a shape should be *common and unique* in its own type, compared with shapes of other types. So, distinctive involves, and is quantified *relative to*, a given set of 3D shapes.

However, existing methods [Shilane and Funkhouser 2006, 2007; Song et al. 2018] either rely on hand-crafted local features, or detect distinctive regions in a supervised setting, meaning that they require labels on data. Hand-crafted features generally do not generalize well to other shapes, since their representation capabilities are limited by the pre-defined fixed operations. Also, in most application scenarios, it is difficult to acquire labels or pre-classify 3D shapes due to annotation efforts. In light of these limitations, the challenging problem is to explore unsupervised methods to detect distinctive

regions directly from the 3D shapes in a data-driven manner. In such a setting, the set of given shapes has *not* been pre-analyzed by any means, and no local descriptors are pre-defined on them.

In this work, we present a method to compute distinctive regions on 3D shapes *in an unsupervised setting*. Our method is based on a neural network that learns and analyzes a given set of shapes without relying on hand-crafted local features, and assigns to each point on the shapes a degree of distinctiveness. First, we sample and represent each given shape as a point cloud, a lightweight and flexible representation. We design a deep neural network and train it on the point clouds for an *unsupervised clustering task*. By training, our network can learn both per-point features and per-shape features. In particular, to drive the network to learn to cluster the shapes for detecting distinctive regions, we design a joint loss function composed of a clustering-based nonparametric softmax loss and an adapted contrastive loss. For the network to learn to cluster the shapes, it has to attend to the discriminative features among the shape clusters. Hence, by analyzing the resulting per-point features and per-shape features, we can obtain a degree of distinctiveness per point in the point sets, and further project the per-point distinctiveness in each point set back to the original shape.

Figure 1 shows two visual examples of using our method to find distinctive regions in the same shape (red being the most distinctive) relative to two different datasets: the inter-class dataset (a) is the whole ModelNet40 dataset [Wu et al. 2015], whereas the intra-class dataset (b) comprises only the tail-engine and four-engine airplanes in ModelNet40. Comparing the distinctive regions in (a) and (b) reveals an interesting phenomenon that (a) tends to focus on the contour of the airplane, while (b) tends to pay more attention to the local regions. This result corresponds to the definition of distinction, since in the inter-class dataset, the contour of the tail-engine airplane is common and unique in the airplane class compared with others, while in the intra-class dataset, the tail and engines are more helpful for distinguishing the tail-engine airplanes from the four-engine airplanes. Sampling the point sets away from the distinctive regions only has little effect on the classification (see Figure 1 (c) & (d)). Later, we shall show, in an extensive empirical experiment presented in Section 4.2, that the points located on the distinctive regions are the *key points* for classification performance. Besides, the detected distinctive regions can further facilitate the development of various applications, e.g., fine-grained shape retrieval; see an example result in Figure 2, and Sections 4 and 5 for more results.

It is worth noting that, the notion of distinction is closely related to saliency as they both measure regional importance. However, while distinction considers how common and unique a region is relative to objects of other types, saliency considers how unique and visible a region is relative to other regions within the same object.

Overall, the contributions of this work are summarized below. We develop a novel unsupervised framework to detect distinctive regions on 3D shapes that does not require hand-crafted features and labels on data. We design a new clustering-based nonparametric softmax classifier and adopt an adapted contrastive loss to encourage the network to learn in an unsupervised manner. We performed extensive experiments to evaluate the effectiveness of our method: quantitatively evaluating on how the detected distinctive regions help shape classification, a user study to compare our results with

human, etc. Further, we show how distinctiveness contributes to applications for shape retrieval, sampling, and view selection.

## 2 RELATED WORK

*Distinctive region detection.* The concept of *distinction*, or *distinctiveness*, was first proposed by Shilane and Funkhouser [2006; 2007]. The main idea of their methods is to extract local shape descriptors for local regions on each shape, then to obtain the distinctiveness of each local region by comparing the difference between all pairs of shape descriptors in the training database. To avoid the drawback of hand-crafted shape descriptors, Song *et al.* [2018] employed a classification network to consume multi-view images of given 3D shapes as input and learn view-based distinction by back-propagating the classification probability. Next, a Markov random field is employed to combine the view-based distinctions across multiple views. Despite the success in finding distinctive regions, existing approaches are *all* supervised, meaning that they all need class labels on the shapes given in the training dataset. In contrast, our method detects distinctive regions in an unsupervised manner.

Besides 3D shapes, the concept of *distinction* was also mentioned in several works on images. Given a large collection of geo-localized images, Doersch *et al.* [2012] developed a discriminative clustering approach to find visual elements that occur much more often in one geographic region than in others, *e.g.*, the kinds of windows, balconies, and street signs that are distinctive in Paris, compared with those in London. Later, several approaches were developed to extract discriminative regions from images for image classification [Juneja et al. 2013; Singh et al. 2012; Sun and Ponce 2013]. More recently, Wang *et al.* [2016] proposed a patch-based framework by introducing triplets of patches with geometric constraints to mine discriminative regions for fine-grained intra-class classification.

*Saliency region detection.* Similar to distinction, saliency also measures regional importance of a shape, but it considers how unique and visible a region is relative to other regions *within the same object*. Lee *et al.* [2005] devised a scale-dependent measure to compute the mesh saliency, while Gal and Cohen-Or [2006] developed local surface descriptors to extract salient geometric features for partial shape matching and retrieval. These techniques are typically based on curvature, or other geometric features. To alleviate the limitation of hand-crafted geometric features, several works adopt data-driven methods to effectively find the saliency for 3D surfaces [Chen et al. 2012; Shu et al. 2019]. In other aspects, Shtrom *et al.* [2013] detected saliency in large point sets, while Ponjou Tasse *et al.* [2015] detected saliency in point sets with a cluster-based approach. Very recently, Wang *et al.* [2018] developed an eye tracking system to obtain mesh saliency from human viewing behavior.

*Network explanation.* Our work is also related to the visualization of neuron activities in deep neural networks. Zeiler and Fergus [2014] devised a perturbation-based method to find the contribution of each portion of the input by removing or masking them, and then running a forward pass on the new input to contrast with the original input. Such approach tends to be slow as the number of test regions grows. Instead, backpropagation-based methods [Ancona et al. 2018; Shrikumar et al. 2017; Sundararajan et al. 2017; Zhang et al. 2018] compute the contribution of all the input regions in
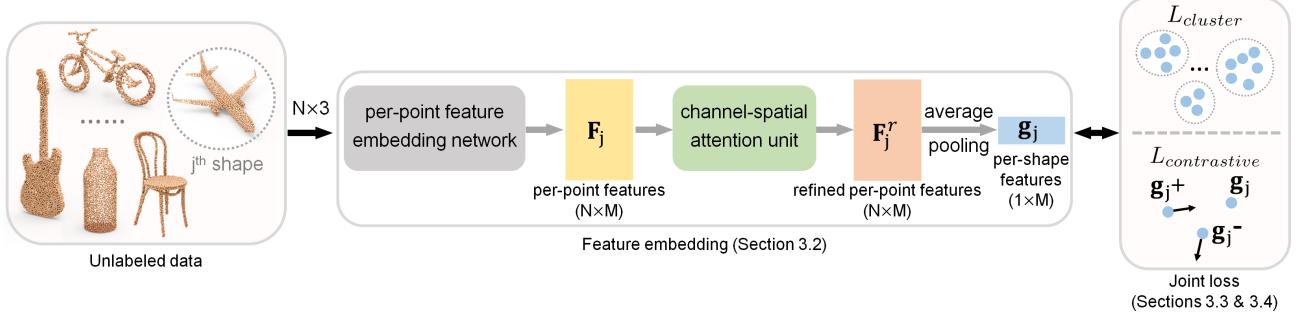
Fig. 3. The overall framework of our unsupervised learning approach to detect distinctive regions on a given set of 3D shapes.

a single forward and backward pass through the network. Unlike the backpropagation-based methods, Zhou *et al.* [2016] formulated the Class Activation Map (CAM) to localize the discriminative regions, and Selvaraju *et al.* [2017] further developed the grad-CAM for producing visual explanations for decisions from the CNN models. These methods, however, require the class labels to visualize the neuron activities, while in our work, we analyze the network activities in an unsupervised training.

*Deep neural networks for point set.* Following PointNet [Qi et al. 2017a] and PointNet++ [Qi et al. 2017b] to embed features directly from point sets, several works successively introduce methods to improve the capturing of geometric information, *e.g.*, SpiderCNN [Xu et al. 2018], KCNet [Shen et al. 2018], PointGrid [Le and Duan 2018], pointwise convolution [Hua et al. 2018], DGCNN [Wang et al. 2019], SPLATNet [Su et al. 2018], and PointCNN [Li et al. 2018]. Besides embedding point features for recognition tasks, several works propose to learn point features for registration, *e.g.*, [Aoki et al. 2019; Wang and Solomon 2019]. In our network, we adopt PointCNN as a module to extract point features, but other network models can also be used for the purpose.

## 3 METHOD

### 3.1 Overview

Given a set of shapes $\mathcal{S} = \{S_j\}_{j=1}^{N_{obj}}$, let $P_j = \{\mathbf{p}_{i,j}\}_{i=1}^{N}$ be a set of 3D points sampled on the $j$-th shape $S_j$, where $N_{obj}$ is the number of shapes in $\mathcal{S}$; $N$ is the number of points in each point set $P_j$; and $\mathbf{p}_{i,j} \in \mathbb{R}^3$ is the 3D coordinates of the $i$-th point in $P_j$. The problem of detecting distinctive regions on shape $S_j$ is

> To predict a per-point distinctiveness value $d_{i,j} \in [0, 1]$ for each point $\mathbf{p}_{i,j}$ on $S_j$ relative to the shapes in $\mathcal{S}$ that are of different types from $S_j$,

where a large $d_{i,j} \approx 1$ indicates that the associated region exists mainly in shapes of the same type as $S_j$ but not in shapes of other types, and a small $d_{i,j} \approx 0$ indicates that the associated region exists in all types of shapes. Hence, $d_{i,j}$ indicates the degree to which the associated region distinguishes shape $S_j$ from others, and there is no requirement for the size of the distinctive regions.

Such goal requires us to consider not only the object itself, but also the other objects in the given reference set. Intuitively, the designed network should contain per-point importance, while also having the ability to classify object types. Hence, we propose a novel

framework to learn both per-point features and per-shape features, as shown in Figure 3. The per-point features are for calculating the distinctiveness $d_{i,j}$, while the per-shape features are for shape classification. To perform the feature embedding in an unsupervised manner, we drive the network to learn by solving an unsupervised shape clustering task using our joint loss (see Figure 3 (right)).

In the following, we first introduce the network for feature embedding (Section 3.2). Next, we introduce a clustering-based nonparametric softmax classifier (Section 3.3) and an adapted contrastive loss (Section 3.4) to drive the unsupervised network to learn. Lastly, we give details on the end-to-end network training (Section 3.5), and describe how we obtain the per-point distinctiveness values from the embedded features (Section 3.6).

### 3.2 Feature Embedding

*Extracting the per-point features.* In this part, we aim to learn an embedding function $f_\theta$:

$$\mathbf{F}_j = f_\theta(P_j), \tag{1}$$

where $\mathbf{F}_j \in \mathbb{R}^{N \times M}$ is the set of extracted per-point features (see $\mathbf{F}_j$ in Figure 3), each row $\mathbf{f}_{i,j}$ in $\mathbf{F}_j$ is a per-point feature of $M$ channels, and $f_\theta$ is a deep neural network with parameters $\theta$. In theory, $\mathbf{f}_{i,j}$ should represent the underlying local geometric structures around each point and further reveal the point's distinctiveness. We apply a point-set-network to extract per-point features, and the choice of the point-set-network is flexible. Most recent networks on processing point sets can be employed; here, we adopt the segmentation architecture of PointCNN [Li et al. 2018] to learn $f_\theta$.

*Adaptive per-point features refinement.* So far, the per-point features are extracted locally over each given shape by the point-set network. As the distinctiveness requires shape-level context, not simply local context around each point, we further propose to refine the per-point features $\mathbf{F}_j$. To this end, we adopt [Woo et al. 2018] to formulate the channel-spatial attention unit (see Figure 4) to fuse the $M$ feature channels over the $N$ per-point features together and produce the refined per-point embedding feature $\mathbf{F}_j^r$.

*Extracting the per-shape feature.* Further, we use an average-pooling operation to obtain the global feature $\mathbf{g}_j$ from $\mathbf{F}_j^r$:

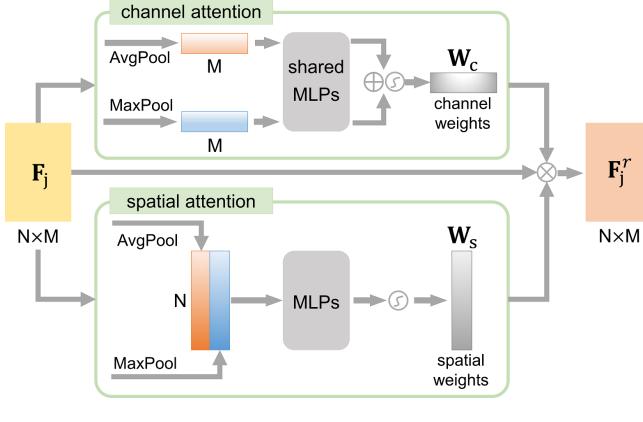$$\mathbf{g}_j = \frac{\sum_{i=1}^{N} \mathbf{f}_{i,j}^r}{N}, \tag{2}$$

Fig. 4. The channel-spatial attention unit.

where $\mathbf{f}_{i,j}^r$ denotes the $i$-th per-point feature vector in $\mathbf{F}_j^r$. As shown in Figure 3, the process of training the network to learn the local (per-point) and global (per-shape) features is driven by an unsupervised joint loss function, which we shall elaborate below.

### 3.3 Clustering-based Nonparametric Softmax

To drive the network to classify objects in an unsupervised manner, we propose a clustering-based nonparametric softmax classifier. In a typical supervised deep neural network for classification, the softmax classifier is commonly employed and the probability of the $j$-th object being recognized as the $q$-th class is

$$P(y_j = q|\mathbf{g}_j) = \frac{exp(\mathbf{w}_q^T\mathbf{g}_j)}{\sum_{k=1}^{C} exp(\mathbf{w}_k^T\mathbf{g}_j)}, \tag{3}$$

where $y_j$ is the class label of the $j$-th object; $C$ is a hyperparameter that denotes the number of classes; $q \in \{1, ..., C\}$ is the class assignment; $\mathbf{w}_k$ is the weight vector for the $k$-th class; $k \in \{1, ..., C\}$; and $\mathbf{w}_k^T\mathbf{g}_j$ measures how well $\mathbf{g}_j$ matches the $k$-th class, so $\mathbf{w}_k$ serves as the class prototype of the $k$-th class.

In supervised learning, we can leverage the class labels provided in training data to learn the class prototype $\mathbf{w}_k$ of each class. This is, however, not possible for unsupervised learning. Recently, an observation was reported by Liu *et al.* [2018] that when the network has successfully converged, the class prototype is usually consistent with the average of all the feature vectors belonging to the same class. Based on this observation, we thus approximate the class prototype by using the average of all the feature vectors belonging to the class, since no class labels are given in our setting.

We adopt the above observation to our problem by formulating the *clustering-based nonparametric softmax classifier*, where we iteratively re-cluster the per-shape feature vectors $\mathbf{g}_j$ in the network and take the average feature vector of each cluster to estimate the cluster prototype $\mathbf{w}_k$. In this way, we can approximate the probability $P(y_j = q|\mathbf{g}_j)$ for unsupervised learning as

$$P(y_j = q|\mathbf{g}_j) \approx \frac{exp(\bar{\mathbf{g}}_q^T\mathbf{g}_j/\tau)}{\sum_{k=1}^{C} exp(\bar{\mathbf{g}}_k^T\mathbf{g}_j/\tau)}, \tag{4}$$

where $\bar{\mathbf{g}}_k = \frac{1}{|\mathbb{C}_k|}\sum_{t \in \mathbb{C}_k} \mathbf{g}_t$ is the average feature vector over all per-shape global feature vectors $\mathbf{g}_t$ of cluster $\mathbb{C}_k$; we take $\bar{\mathbf{g}}_k$ (per-cluster) to approximate $\mathbf{w}_k$ for unsupervised learning; and $C$ denotes the number of clusters in our unsupervised setting. Further, we enforce $\|\mathbf{g}_j\|=1$ via an L2-normalization layer in the network and make use of $\tau$, which is a temperature parameter, to control the concentration level of the distribution [Hinton et al. 2015; Wu et al. 2018].

In our experiments, we set $\tau$ as 0.07, following the setting in [Wu et al. 2018]. Then, our learning objective is to maximize $P(y_j = q|\mathbf{g}_j)$, or equivalently, to minimize the negative log-likelihood of the probability. Therefore, our *clustering-based nonparametric softmax loss* is formulated as

$$L_{cluster} = -\sum_{j=1}^{N_{obj}} \log P(y_j = q|\mathbf{g}_j). \tag{5}$$

In our implementation, we use spectral clustering [Stella and Shi 2003; Von Luxburg 2007] to cluster the per-shape global features $\mathbf{g}_j$ in each training epoch. Experimentally, we found that our network detects similar distinctive regions when equipped with different clustering algorithms; see Supplementary Material Part A for the evaluation. Also, please see Supplementary Material Part B for the effect of having different $C$ on extracting distinctive regions.

### 3.4 Adapted Contrastive Learning

Inaccurate clustering results are inevitable, so relying only on the clustering-based loss may mislead the learning of the network. Motivated by [Bachman et al. 2019; Hénaff et al. 2019; Hjelm et al. 2019], to stabilize and enhance the feature learning in the network, we formulate an *adapted contrastive loss*, which is particularly important at the beginning of the training process when the clustering results are more random. Considering input point set $P_j$ to the network as the *anchor*, for each training epoch, we form a *positive* point set sample $P_j^+$ and a *negative* point set sample $P_j^-$ for $P_j$, such that the per-shape global feature $\mathbf{g}_j^+$ associated with $P_j^+$ is close to $\mathbf{g}_j$, while the per-shape global feature $\mathbf{g}_j^-$ associated with $P_j^-$ is far from $\mathbf{g}_j$.

- For $P_j^-$, we randomly pick a point set from the shapes in the clusters that $P_j$ does not belong to.
- For $P_j^+$, since the intra-class clustering results may not be reliable, especially at the beginning of the training process, we thus do not randomly pick from the cluster that $P_j$ belongs to. Rather, we resample another point set $P_j^+$ on the given 3D shape ($S_j$) associated with $P_j$ and pass $P_j^+$ to the network to generate $\mathbf{g}_j^+$. Note that $P_j^+$ and $P_j$ are different point sets due to randomness in the point sampling process, but essentially, they describe the same object, *i.e.*, $S_j$.

We take the above triplet $\{\mathbf{g}_j, \mathbf{g}_j^+, \mathbf{g}_j^-\}$ to form an *adapted contrastive loss* following [Hadsell et al. 2006] as

$$L_{contrastive} = D(\mathbf{g}_j, \mathbf{g}_j^+) + \max(0, \lambda - D(\mathbf{g}_j, \mathbf{g}_j^-)), \tag{6}$$

where $D$ is the Euclidean distance in feature space, and we set $\lambda = 2.0$ in our experiments. Importantly, we generate such triplet input *dynamically* for each $P_j$ in each training epoch. Using this strategy, we can increase the diversity of the training samples and produce more reliable samples as the training progresses.
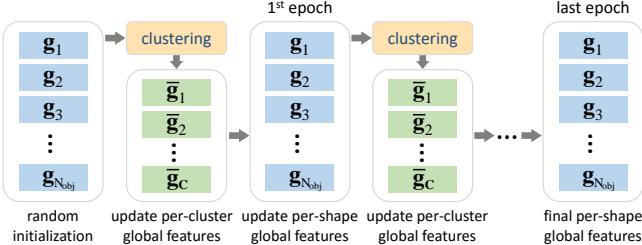
Fig. 5. Illustration of the network training process.

## 3.5 End-to-end Network Training

Overall, we end-to-end train the network to learn the features for clustering the given shapes using the joint loss function

$$L(\theta) = L_{cluster} + \alpha L_{contrastive} + \beta \|\theta\|^2, \qquad (7)$$

where $\alpha$ balances the two loss terms and $\beta$ is the multiplier of weight decay in the regularization term (see Section 4.1 for their values).

In summary, the feature embedding is conducted in a self-training way by iteratively learning the feature vectors, re-clustering them, then using the clustering results to fine-tune the model. Figure 5 illustrates the whole training process, where we first randomly initialize the per-shape global feature $\mathbf{g}_j$ of each training sample $P_j$, cluster them into $C$ classes, and generate the per-cluster global feature $\bar{\mathbf{g}}_k$ for each cluster. Early in the training, these $\mathbf{g}_j$ and $\bar{\mathbf{g}}_k$ are unlikely reliable, but as the training progresses, we iteratively update these per-shape and per-cluster features in each training epoch, these features can then gradually converge and let us further obtain the per-point distinctiveness in the given shapes. Here, $C$ remains unchanged during the training, and our method does not require even class size in the training data. The clustering method we employed, *i.e.*, spectral clustering, will divide the training samples automatically into $C$ clusters based on the feature similarity.

Figure 6 shows t-SNE visualizations that reveal the clustering of the per-shape features ($\mathbf{g}_j$) during the unsupervised training. Here, we cluster over 9000 shapes into 40 classes. The dimension of the features (*i.e.*, $M$) is 128 in our implementation.

## 3.6 Obtaining and Visualizing the Distinctiveness

As introduced earlier in Section 3.1, we design our unsupervised approach to learn both per-shape and per-point features in the given shapes. After the training to meet the shape clustering task, the response of the activation neuron associated with the per-point features should positively correlate to its confidence of the detection [Zhang et al. 2018]. Therefore, we obtain the per-point distinctiveness $d_{i,j}$ from $\mathbf{f}_{i,j}^r$ of each point $\mathbf{p}_{i,j}$ by taking the maximum value in $\mathbf{f}_{i,j}^r$ and normalizing $d_{i,j}$ between 0 and 1 for each shape. For a comparison of applying other alternatives to extract $d_{i,j}$ from $\mathbf{f}_{i,j}^r$, including the mean, $L_2$ norm, average of the three largest values, etc., please refer to Supplementary Material Part C.

Furthermore, to visualize the distinctiveness results, we project the per-point distinctiveness on point set $P_j$ back to the original shape $S_j$ and obtain a distinctiveness value for every vertex on the original shape by averaging the distinctiveness values over the
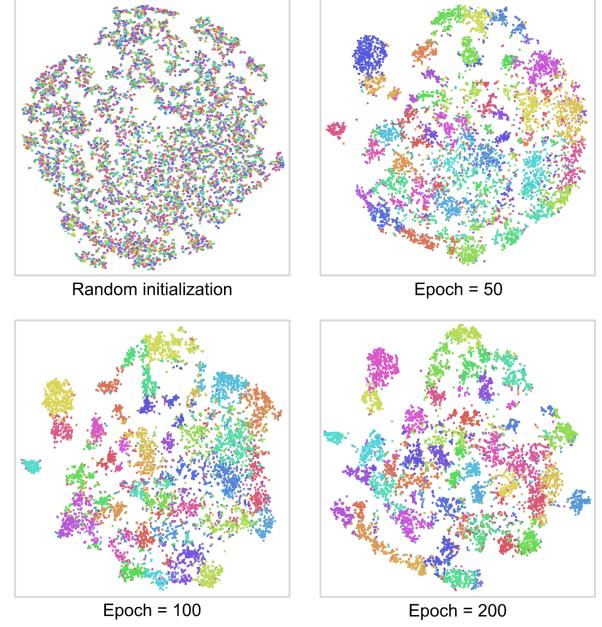


Fig. 6. t-SNE visualizations of the per-shape features clustering during the unsupervised training process.

nearby sampled points in $P_j$; see Figure 1 and Figure 7 for example results, where regions in red are the most distinctive.

## 4 EXPERIMENTS AND RESULTS

### 4.1 Implementation Details

Our method was implemented using TensorFlow [Abadi et al. 2016]. To train the network, we randomly sampled 2,048 points for each shape as input and augmented the input point sets on-the-fly, including random rotation, scaling, shifting, and jittering. Moreover, we empirically set $\alpha$ and $\beta$ in Eq. (7) as 3.0 and $10^{-5}$, respectively, and trained our network for 200 epochs using the Adam optimizer [Kingma and Ba 2014] with a learning rate of 0.01. See Supplementary Material Part J for further analysis on $\alpha$ and $\beta$. For the ModelNet40 dataset [Wu et al. 2015] with 9,839 training samples, it took about 25 hours to train our network on a 12GB TITAN Xp GPU with a batch size of 50. The network has 0.48M parameters. For inference, it took about 0.017 seconds for our method to predict the per-point distinctiveness values for a point cloud of 2,048 points.

Based on PointCNN [Li et al. 2018], we further made the following adaptations to the feature embedding component in the network architecture. First, we use a fixed-size query ball [Qi et al. 2017b] instead of KNN or geodesic-like KNN [Yu et al. 2018] to find the local neighborhood for extracting the point features; note also that according to [Qi et al. 2017b], query ball is preferred for finding the local neighborhood in tasks that require local pattern recognition, e.g., our distinctive detection task. Second, we removed the $\mathcal{X}$-conv operation in the deconvolution part of the PointCNN segmentation network and directly used feature interpolation [Qi et al. 2017b] for per-point feature restoration. In this way, we can reduce the number of network parameters and speed up the network training

Fig. 7. Distinctive regions detected by our method *unsupervisedly* on various 3D models in ModelNet40 [Wu et al. 2015]; red indicates high distinctive regions.

process with little degradation in the quality of the results. Lastly, we explored different network backbones (*i.e.*, PointNet and Point-Net++) for learning the per-point local features; see Supplementary Material Part D for the experimental results.

## 4.2 Detecting Distinctive Regions

*Distinctiveness visualization.* We employed the ModelNet40 training split dataset [Wu et al. 2015] and trained our network in an un-supervised manner to sort the models in the dataset into 40 clusters. Figure 7 shows the distinctive regions detected by our method on a variety of models in the dataset, where red color indicates high distinctive regions. When we determine if a region is distinctive, instead of looking just at the shape itself, we consider all shapes of different classes (but not shapes of its own class) in the dataset. Taking the Person shape in Figure 7 (right) as an example, the head, feet, and hand are found to be more distinctive (red), while the body part is less (blue). Since we compare the Person shapes with shapes of other types, these human parts are distinctive for the network to recognize the Person shapes relative to others. For the two Lamp shapes in Figure 7 (middle & top-right), our method detects the bulb as distinctive, since it is common and unique in this class compared with shapes in other classes. For the Chair shapes, not only the legs are detected, the back is also detected as distinctive. Similarly, our method detects as distinctive the handle of the Cup, the leaves of the Plant, the struts of the Guitar, the handle of the Door, etc. For more results, please refer to Supplementary Material Part K.

In particular, our network does not simply locate high-curvature regions and extremities as distinctive. We show several another examples in Figure 8, where the detected distinctive regions are not extreme regions. For example, for the four shapes shown in the top row, the detected distinctive regions are not extreme regions at all. Concerning the other four shapes on the bottom, for the Chair shape, its back (see the region marked by the black arrow), which is not
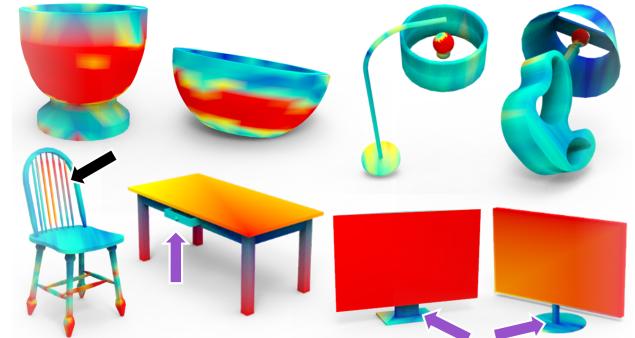


Fig. 8. Distinctive regions detected by our method may not simply lie on high-curvature and extreme regions, as shown in these examples. Also, not all extreme regions are detected as distinctive (see purple arrows).

extreme region, is also detected. For the Table and Monitor shapes, some obvious extreme regions are not detected (see the regions marked by the purple arrows). For more non-extreme examples, please see Supplementary Material Part M.

*Quantitative evaluation.* Next, we quantitatively evaluated *how helpful the detected distinctive regions are to shape classification.* Here, we employed totally 2,468 models in the ModelNet40 testing dataset as test shapes, and used three different preferences to downsample points from pre-sampled point sets ($N = 2,048$) on the test shapes: (i) probability to preserve a point based on the distinctiveness at the point; (ii) probability to preserve a point based on the curvature at the point; and (iii) ignore the point importance and downsample points at random. Therefore, preference (i) leads to the production of more points on our detected distinctive regions, while preference (ii) leads to the production of more points on high-curvature regions. Then, we fed the downsampled points into a classification network

Two different training datasets
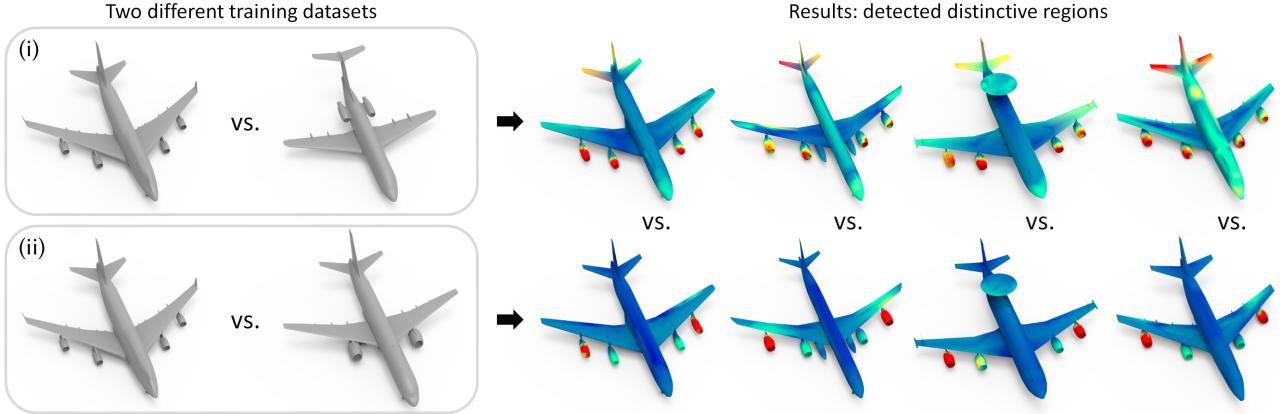
Results: detected distinctive regions



Fig. 9. Our network, when trained with different datasets (left), detects different distinctive regions on the same given objects (right): (i) a training set of four-engine and tail-engine airplanes in the top row, and (ii) a training set of four-engine and two-engine airplanes in the bottom row.
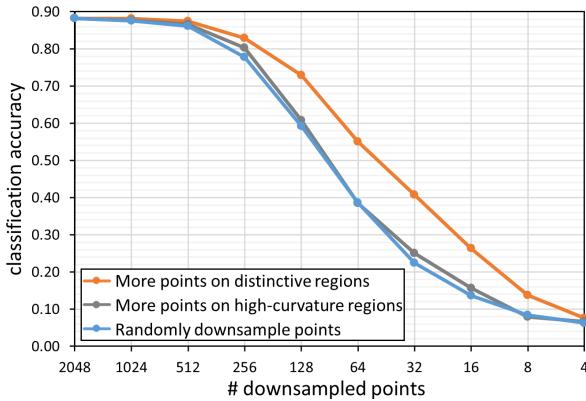


Fig. 10. Overall shape classification accuracy when we downsample points on test shapes using three different preferences. Having more distinctive points detected by our method (orange plot) leads to higher accuracy.

(*i.e.*, PointNet [Qi et al. 2017a]), which has been pre-trained on the ModelNet40 training split dataset with 2,048 points on each shape for a shape classification task, and computed the overall classification accuracy averaged over all the test shapes.

Figure 10 plots the overall classification accuracy for point sets downsampled with the three different preferences using decreasing number of downsampled points. Comparing the orange plot with the gray and blue plots in Figure 10, we can see that having more points on the distinctive regions can better preserve the discrimination of the shapes, leading to higher shape classification accuracy, particularly for results with fewer points. Hence, this quantitative comparison shows that the distinctive regions detected by our method are helpful to the classification of 3D shapes. Additionally, since the orange plot is above the gray plot, this indicates that not all high-curvature regions are important for recognition and our method does not simply take all regions of specific curvature profiles, *e.g.*, sharp corners and extremities, as distinctive.

### 4.3 Effect of using Different Training sets

A notable characteristic of detecting distinctive regions is that the results highly relate to the given set of shapes, or the *training set*. To verify this, we conducted an experiment to explore the effect of training set as follows. First, we collected two training sets: (i) 500 four-engine airplanes and 250 tail-engine airplanes (see Figure 9 (top-left) for examples), and (ii) we kept the four-engine airplanes but replaced the tail-engine airplanes with around 1,000 two-engine airplanes (see Figure 9 (bottom-left) for examples). Please see Supplementary Material Part E for more examples in the datasets. Further, we trained our network on each training set separately with $C$ set to two, and employed the two trained network models to detect distinctive regions on four-engine airplanes.

Figure 9 (right) shows the distinctive regions detected on four different four-engine airplanes. The interesting observation is that when trained with the four-engine vs. tail-engine airplane dataset (top row), our network tends to highlight all the four engines on the test airplanes. On the other hand, when trained with the four-engine vs. two-engine airplane dataset (bottom row), our network tends to only highlight the outer two engines on the airplanes as distinctive regions. In particular, the four test airplanes have different size, engine shape, and wing shape, where the middle two have special structures, *i.e.*, extra fuel tanks and a radar on top. Yet, given these shapes as inputs, our method still detects consistent distinctive regions. Further, if we look closer to the training set on top left, the two kinds of airplanes not only have different engine locations but also different tail shapes. Even though these tail regions are not as dominant as the engines, our trained network can still weakly highlight the tails on the test airplanes; compare the tails of the test airplanes on top-right vs. bottom-right in Figure 9. Please see Supplementary Material Part L for results on other datasets.

### 4.4 Comparing with Other Methods

Figure 11 shows the distinctive regions detected by a traditional method [Shilane and Funkhouser 2007] (without deep neural networks), by a weakly-supervised deep learning method [Song et al. 2018], and by our method. We train our network using the same

Fig. 11. Distinctive regions detected by Shilane et al. [2007] (top row), by Song et al. [2018] (middle row), and by our method (bottom row). For all the three methods, distinctive regions on the person shapes are detected by comparing the person shapes with shapes of other different class types.

training set, *i.e.*, the Princeton Shape Benchmark [Shilane et al. 2004] as in [Shilane and Funkhouser 2007], and most class types of training samples in [Song et al. 2018] are the same as ours. For all the three methods, distinctive regions on the person shapes are detected by comparing the person shapes with shapes of other different class types in the data. The results of [Shilane and Funkhouser 2007] and [Song et al. 2018] were directly acquired from their papers.

From the results, we can see that [Shilane and Funkhouser 2007] tends to highlight elbows as distinctive regions and ignore semantic parts such as heads and feet, due to the limited representation capability of the hand-crafted features. For [Song et al. 2018], it is able to highlight heads and hands as distinctive. Compared with [Song et al. 2018], even though our method is unsupervised, it can detect not only the heads and hands but also the feet as distinctive. In particular, our method detects consistent distinctive regions on these Person shapes, even they have different poses and shapes.

## 4.5 Unsupervised vs. Weakly-supervised Learning

To explore if our unsupervised network can meaningfully cluster the shapes for detecting distinctive regions, we further compared it with a weakly-supervised version of our method. Specifically, we used the class labels provided in ModelNet40, added a fully-connected layer with 40 output neurons after the global feature (see Figure 3) to regress the class scores, then used the cross entropy loss to replace the unsupervised loss to train this weakly-supervised network. From the results presented in Figure 12, we can see that most distinctive regions detected by the weakly-supervised network (top) can also be found by the unsupervised network (bottom); the unsupervised network only misses a few of them, *e.g.*, some leaves in Plant. This comparison result gives evidence that even without the class labels, the performance of our unsupervised method is still comparable to that of the weakly-supervised version of our method.
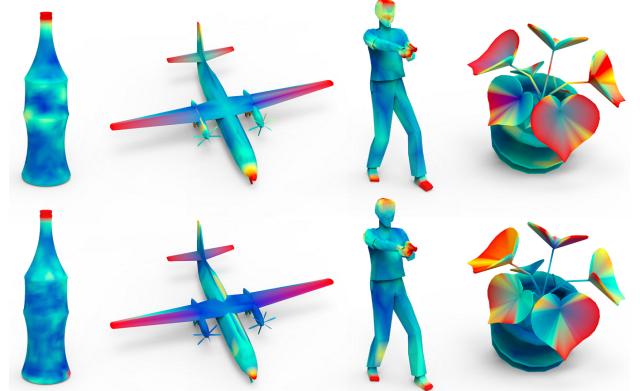


Fig. 12. Distinctive regions detected by our method (bottom row) and a weakly-supervised version of our method (top row). Even without using any class label, our method, which is unsupervised, can still detect regions that are similar to those by a weakly-supervised version of our method.

Further, we quantitatively compare the two results by following [Dutagaci et al. 2012] to compute the False Negative Error (FNE) and False Positive Error (FPE). Specifically, given distinctiveness values $d_{i,j}$ and $\widehat{d}_{i,j}$ at point $\mathbf{p}_{i,j}$ detected by the unsupervised and weakly-supervised networks, respectively, we first located two sets of more distinctive points per shape $S_j$ by a threshold $d_t$: $Q_j = \{\mathbf{p}_{i,j} | d_{i,j} > d_t\}$ and $\widehat{Q_j} = \{\mathbf{p}_{i,j} | \widehat{d}_{i,j} > d_t\}$. By regarding $\widehat{Q_j}$ as the *ground truth*, a point $\widehat{\mathbf{q}} \in \widehat{Q_j}$ is said to be *covered* by $Q_j$, if there exists point $\mathbf{q} \in Q_j$, such that $||\widehat{\mathbf{q}} - \mathbf{q}||_2 \leq r D_j$ and $\mathbf{q}$ is not closer to any other point in $\widehat{Q_j}$, where $D_j$ is the bounding sphere diameter of shape $S_j$ and $r$ is a parameter (ratio) to control the localization tolerance. Then, we compute $\text{FNE}_j = (|\widehat{Q_j}| - N_c)/|\widehat{Q_j}| = 1 - N_c/|\widehat{Q_j}|$, where $N_c$ is the number of points in $\widehat{Q_j}$ that are covered by $Q_j$. On the other hand, each covered point in $\widehat{Q_j}$ corresponds to a unique point
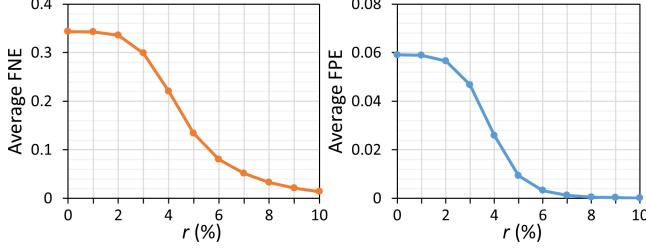
Fig. 13. Average false negative error (FNE) and false positive error (FPE) over 400 models plotted against tolerance parameter $r$.

Table 1. Comparing the overall shape classification accuracy on ModelNet40 when downsampling the test shapes with more points on the distinctive regions detected under different settings.

| Different settings | Number of downsampled points | | | | | |
|---|---|---|---|---|---|---|
| | 1024 | 512 | 256 | 128 | 64 | 32 |
| w/o Atten | 0.877 | 0.871 | 0.806 | 0.648 | 0.402 | 0.244 |
| w/o Cont | 0.878 | 0.871 | 0.808 | 0.653 | 0.413 | 0.245 |
| A-Center-Cont | 0.876 | 0.865 | 0.818 | 0.658 | 0.405 | 0.235 |
| Our full pipeline | **0.881** | **0.874** | **0.829** | **0.729** | **0.550** | **0.408** |

Figure 14 shows the distinctive regions detected by our method under the four different settings. Comparing the results produced with the *w/o Atten* setting (left-most) and with our full pipeline (right-most), we can see that the attention unit helps locate distinctive regions that span over a larger spatial areas, *e.g.*, the Chair's back. Looking at the results produced with *w/o Cont* and *A-Center-Cont*, we can see that they tend to miss some important regions, *e.g.*, the Table's desktop, or contain noise, *e.g.*, the Chair's back.

Besides visual comparison, we performed the same quantitative evaluation on the results here, as in Section 4.2. Table 1 shows the shape classification accuracy on the ModelNet40 test dataset for the four different settings in terms of decreasing number of downsampled points. From the table, we can observe that our full pipeline leads to the highest classification accuracy. Since we preserve more points on distinctive regions, this means that the three network elements being explored in this experiment all contribute to improve the detection of the distinctive regions.

### 4.7 User Studies

To obtain a sense of how consistent our results are with humans, we conducted two user studies, which we shall elaborate below. The key idea behind the studies is that we try to simulate the network clustering process with humans to obtain the distinctive regions of some test shapes. Specifically, we started by introducing the definition of distinction to each participant to confirm that all participants understood the meaning of distinction. Then, we showed the training dataset to the participants. However, to avoid fatigue, we randomly selected a subset of 3D shapes from different classes in the training dataset, and showed these shapes to each participant on a computer display, on which the participant can rotate each shape and explore its details. Please refer to Supplementary Material Part F for some screenshots. After that, each participant was given a set of test shapes, and asked to cluster the shapes and label the distinctive regions on the shapes that affect how they cluster. All the shapes are presented simply in a colorless manner to the participants in both of the studies.

*Intra-class prediction.* The first user study explores how humans find distinctive regions on shapes of the same class. Here, we employed (i) the dataset of four-engine vs. two-engine airplanes; and (ii) the dataset of four-engine vs. tail-engine airplanes, as presented in Section 4.3; see again Figure 9 (left). To avoid bias due to the dataset similarity, we randomly divided the participants into two groups, one for each set. For the first group, we randomly selected 10 four-engine airplanes and 22 two-engine airplanes, and presented them in random order on a computer display. Then, the participants



w/o Atten    w/o Cont    A-Center-Cont    Our full pipeline

Fig. 14. Visual results in ablation study.

in $Q_j$, so the points in $Q_j$ that are without any correspondence in $\widehat{Q}_j$ are regarded as false positives. Hence, $\text{FPE}_j = (|Q_j| - N_c)/|Q_j| = 1 - N_c/|Q_j|$. Here, $\text{FNE} \in [0, 1]$, $\text{FPE} \in [0, 1]$, and a small value indicates high consistency between $Q_j$ and $\widehat{Q}_j$.

Figure 13 plots the FNE and FPE values averaged over 400 randomly-selected objects (from the 40 different classes in the ModelNet40 testing split) against $r$. We can see that when $r$ is just around 8% to 10%, both FNE and FPE are very close to zero, thus demonstrating the *high consistency between the distinctive points* detected by the two networks. Particularly, the average FPE is very low even when $r = 0$, meaning that most of the distinctive regions detected by our unsupervised method are also the distinctive regions detected by the weakly-supervised version of our method.

### 4.6 Ablation Study

Next, we analyzed the major elements in our network by removing or replacing each of them when we train our network on the ModelNet40 models: (i) *w/o Atten* – removing the channel-spatial attention unit (see Figure 4 for details); (ii) *w/o Cont* – removing the $L_{contrastive}$ term from the joint loss in Eq. (7); and (iii) *A-Center-Cont* – replacing $L_{cluster}$ in Eq. (4) with an adapted center loss [Wen et al. 2016] for unsupervised clustering-based learning: $\frac{1}{2} \sum_{j=1}^{N_{\text{obj}}} \|\mathbf{g}_j - \bar{\mathbf{g}}_q\|^2$ with a goal of minimizing the intra-cluster variations, while keeping the features of different clusters separable.
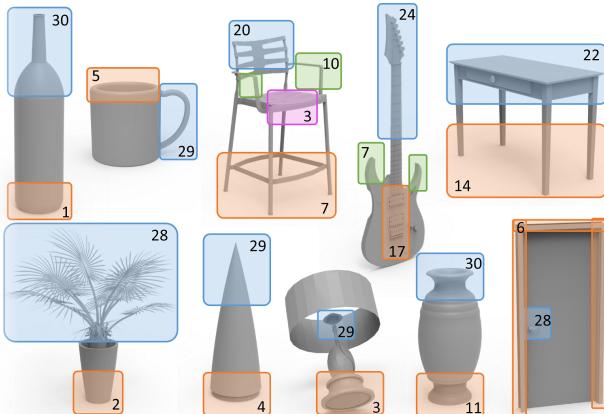
Fig. 15. The ten objects on which the participants mark distinctive regions. The above boxes reveal the participant-marked distinctive regions and the corresponding number of participants marked on each region. Note that each participant may mark more than one region on the same object.

Table 2. Quantitative evaluation on the inter-class prediction consistency between our method and the participants.

|     | Bottle | Chair | Cone | Cup | Door | Guitar | Lamp | Plant | Table | Vase | Avg |
|-----|--------|-------|------|-----|------|--------|------|-------|-------|------|-----|
| FNE | 0.02 | 0.40 | 0.08 | 0.10 | 0 | 0.51 | 0.07 | 0.07 | 0 | 0 | 0.13 |
| FPE | 0 | 0.53 | 0.03 | 0.03 | 0.40 | 0.20 | 0.03 | 0.07 | 0.35 | 0.30 | 0.19 |
| WME | 0.03 | 0.33 | 0.12 | 0.15 | 0.18 | 0.50 | 0.09 | 0.06 | 0 | 0 | 0.15 |

were asked to divide the 32 airplanes into two clusters and label the distinctive regions on the four-engine airplanes. For the other group of participants, we randomly selected 10 four-engine and 22 tail-engine airplanes from the other dataset, and performed the same procedure with the participants.

All ten participants (both groups) recruited in this study clustered the airplanes in the same way as our network. For the four-engine vs. two-engine dataset, all participants marked the outer two engines as distinctive in the four-engine airplanes. Their results are the same as our network predictions; see Figure 9 (bottom-right). For the four-engine vs. tail-engine dataset, all participants marked the four engines as distinctive in the four-engine airplanes, and their results are almost the same as our network (see Figure 9 (top-right)), except for the tails of the airplanes; since the tail-engine airplanes mostly have T-shaped tails, which are generally absent in four-engine airplanes. Without our reminder, only one participant noticed the T-shaped tails, but when we asked the other participants whether such tails are also distinctive, they all strongly agreed. This study shows that our network is able to attend to large and small distinctive regions, which may even be overlooked by humans.

*Inter-class prediction.* The second study explores how humans find distinctive regions between shapes of various kinds. Here, we employed the models from ModelNet40 and recruited 30 participants. To avoid fatigue, we randomly selected 75 shapes evenly from 15 different classes, and further selected ten objects of different classes from the set; see Figure 15. Then, we showed the 75 shapes to each participant and asked him/her to mark distinctive regions on the ten selected objects. Particularly, we explained the definition of

distinctiveness, *i.e.*, the distinctive regions should be common in each specific class, while being unique *relative* to other classes.

During the study, we found that the size of the marked region varies among the participants, even at the same object location. Also, they focus more on the shape features than on the size of the marked regions. Figure 15 shows the summary of human marked regions on the ten objects. In each marked region, the corresponding number indicates how many participants marked on the same area.

To quantitatively compare the distinctive regions marked by the participants and detected by our method, we adapted the FNE and FPE metrics (see Section 4.5) as follows. First, for regions detected by our network, we set a threshold to keep only the high-distinctive regions as the final detected distinctive regions; see Figure 7 for the distinctiveness distribution on the ten objects. Intuitively, we keep only the yellow and red regions on the objects. Next, we compared these regions with the regions marked by each participant, and regard his/her marked regions as the ground truth, where a marked region is said to be *covered*, if both the participant-marked region and the network-detected region cover almost the same structures. On the other hand, the network-detected regions that are covered by the participant-marked regions are said to be the true positives. Hence, for each ($k$-th) participant, we define $N_{j,k}^c$ as the number of participant-marked regions that are covered and $T_{j,k}^h$ as the total number of participant-marked regions, on the $j$-th object. Also, we define $T_j^d$ as the total number of network-detected regions on the $j$-th object. Then, similar to Section 4.5, we compute the corresponding $\text{FNE}_{j,k} = 1 - N_{j,k}^c/T_{j,k}^h$ and $\text{FPE}_{j,k} = 1 - N_{j,k}^c/T_j^d$, and further compute the FNE and FPE values averaged over all the participants per object. Additionally, to account for the frequently-marked regions, we adopted the Weighted Miss Error (WME) metric [Dutagaci et al. 2012], *i.e.*, $\text{WME} = 1 - \sum_k N_{j,k}^c/\sum_k T_{j,k}^h$.

Table 2 shows the per-object FNE, FPE, and WME values, as well as their overall averages over the ten objects. All three metrics range [0, 1], where a low value indicates high consistency between the participant-marked and network-detected regions. From the table, we can see that most values are very low and several are even zeros, indicating that most participant-marked distinctive regions can be detected by the network, and vice versa. However, values for some complex objects are a bit higher, *e.g.*, the Chair and the Guitar, since different participants may mark on different structures. Yet, the overall values are very close to zeros, meaning that the distinctive regions detected by our network are highly consistent with the distinctive regions marked by the participants. Furthermore, we also conducted a statistical test to show that 30 participants are sufficient to produce stable human labels in the analysis. Please refer to Supplementary Material Part F for the details.

## 5 APPLICATIONS

Being able to discover distinctive regions on 3D shapes enables us to support and enhance various applications, such as shape retrieval [Gal and Cohen-Or 2006; Shilane and Funkhouser 2006], shape simplification [Garland and Heckbert 1997], remeshing [Alliez et al. 2002], best view selection [Shilane and Funkhouser 2007], perception-aware 3D printing [Zhang et al. 2015]. In this section, we
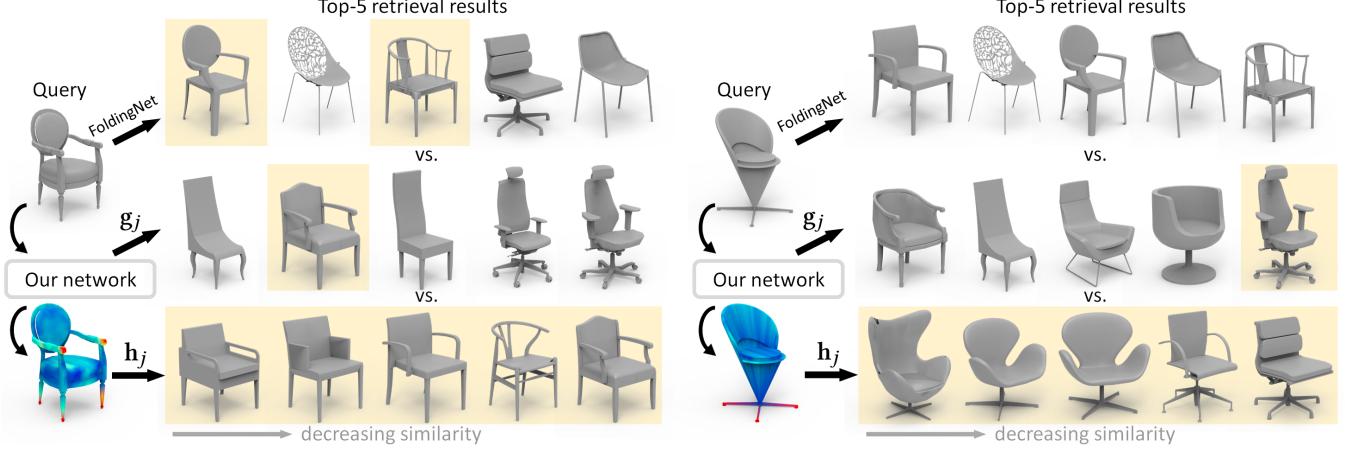
Fig. 16. Two sets of shape retrieval results (left & right). In each set, we show the top-five similar shapes retrieved by using the per-shape global feature from FoldingNet [Yang et al. 2018] (top) and from our network, *i.e.*, $\mathbf{g}_j$ (middle), and by using our distinctiveness-guided global feature $\mathbf{h}_j$ (bottom). Retrieved shapes of substructures (armrest (left) & swivel base (right)) similar to the query shape are marked over a yellow background.

present three typical applications to demonstrate the applicability of our technique in distinctiveness-guided shape retrieval, sampling, and view selection in 3D scenes.

## 5.1 Distinctiveness-guided Shape Retrieval

Conventional approaches for shape retrieval first extract representative shape descriptors then retrieve similar shapes based on the distance between the extracted descriptors; see [Tangelder and Veltkamp 2004] for a survey. Though these approaches have achieved promising performance for inter-class retrieval, they tend to have limited ability for fine-grained intra-class retrieval, since the extracted descriptors are global. With our detected distinctive regions, we can perform *distinctiveness-guided shape retrieval*, which enables fine-grained intra-class shape retrieval, *e.g.*, retrieving swivel chairs from a large collection of chairs; see Figure 16 (right).

The main idea is to only consider these distinctive point features rather than regarding all the point features equally. In detail, we randomly sample $N$ (*e.g.*, 2,048) points on each shape $S_j$ and feed these points into our trained network to obtain local per-point features $\mathbf{F}_j^r$ and global per-shape features $\mathbf{g}_j$; see Section 3.2. As described in Section 3.6, we further obtain per-point distinctiveness $d_{i,j}$ from $\mathbf{f}_{i,j}^r$. To facilitate fine-grained intra-class retrieval, instead of directly using $\mathbf{g}_j$ as the representative descriptor, we select only the more distinctive per-point features and average over them to obtain the distinctiveness-guided global feature $\mathbf{h}_j$:

$$\mathbf{h}_j = \frac{\sum_{i=1}^{N} \mathbb{I}\{d_{i,j} > \Delta_d\} \mathbf{f}_{i,j}^r}{\sum_{i=1}^{N} \mathbb{I}\{d_{i,j} > \Delta_d\}}, \tag{8}$$

where $\mathbb{I}$ is the indicator function and $\Delta_d$ is the threshold. Lastly, we measure the similarity between shapes by computing the Euclidean distance between $\mathbf{h}_j$ of the shapes.

Figure 16 shows two sets of results, each using a different chair as the query shape. As a comparison, we applied a recent unsupervised network FoldingNet [Yang et al. 2018] to extract per-shape global feature for similarity computation, and the results are shown in the
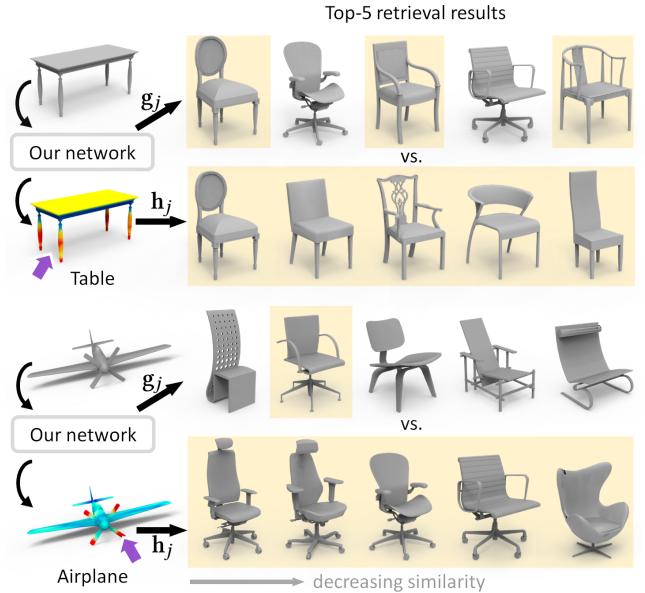


Fig. 17. Top-five shape retrieval results by using a Table shape with four legs (top) or using an Airplane shape with a propeller (bottom) as the query to search over the Chair dataset. Retrieved shapes with similar substructures are marked over a yellow background.

top row. Besides, we also directly employed our per-shape feature $\mathbf{g}_j$ for retrieval; see the middle row. In each set, we retrieved the top-five similar shapes from a data pool of 100 different chairs, and manually marked in yellow (see the figure) the results with substructures similar to the query shape: for the left example in Figure 16, we consider a retrieval result as similar to the query shape, if it also has an armrest and four legs as its substructures, while for the right example in Figure 16, we consider a retrieval result as similar to the query shape, if it also has a swivel base as its substructures. From Figure 16, we can see that using the global descriptors, no matter extracted by FoldingNet or our network, the

Best views (a & b) vs. worst views (c & d) for the whole 3D scene
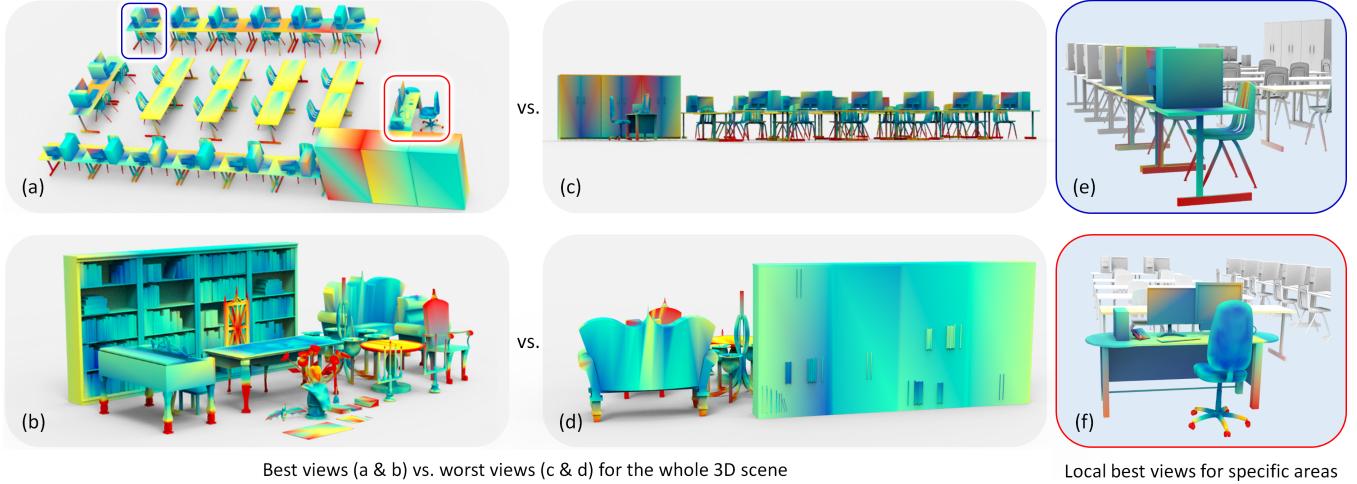
Local best views for specific areas

Fig. 18. Using the network-predicted distinctiveness values over a scene, we can find best views with maximized distinctiveness.
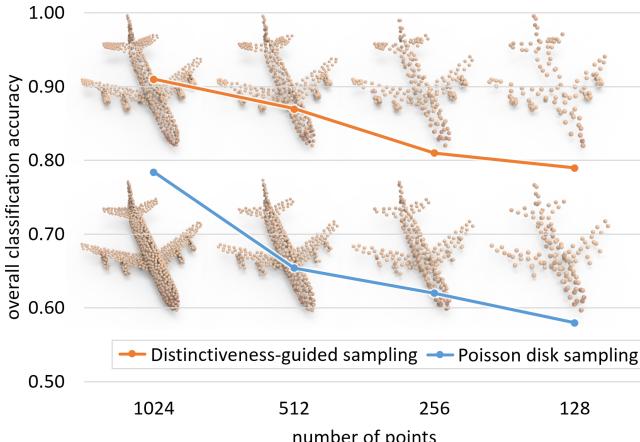


Fig. 19. Shape classification accuracy using distinctiveness-guided sampling (top row) vs. conventional Poisson disk sampling (bottom row).

results have very different substructures from the query shapes. Guided by the network-predicted distinctiveness values, $\mathbf{h}_j$ can help retrieve shapes with more similar substructures (marked in yellow). Please see Supplementary Material Part G for more results.

To further explore the ability of our distinctiveness-guided shape retrieval, we perform shape retrieval on the Chair dataset using a Table and an Airplane as the query. Figure 17 shows the results. For the top case, the four legs in the Table shape are found to be distinctive, so when using it as the query, four-leg chairs are obtained. For the bottom case, the Airplane shape has a propeller, so chairs with substructures that look like the propeller (distinctive regions) can be retrieved as the results. These results demonstrate that our distinctiveness-guided shape retrieval can pay more attention to the local distinctive regions, instead of simply the overall structure.

## 5.2 Distinctiveness-guided Sampling

Sampling is a common task in computer graphics, as well as in many domains, for generating point samples to represent a continuous shape. Using the distinctiveness detected on 3D shapes, we can guide the point sampling process by emphasizing the distinctive regions on the shapes. In this way, the sampled points can more effectively describe the shapes in terms of discrimination ability.

In our implementation, we take the point distinctiveness values detected by our network to control the local sampling density in an adaptive Poisson disk sampling process. That is, we set higher sampling density (or equivalent, smaller Poisson disks) for more distinctive regions, and vice versa. Figure 19 presents the sampling results on a four-engine airplane object with decreasing number of points, where (i) the top row shows the results produced using an adaptive Poisson disk sampling guided by the network-predicted distinctiveness values, and (ii) the bottom row shows the results produced by the conventional Poisson disk sampling, which randomly but uniformly samples the given object. From the results, we can see that distinctiveness-guided sampling (top) arranges more points in high distinctive regions, such as the engines, thereby enhancing the preservation of the shape's characteristics. Please refer to Supplementary Material Part H for more sampling results.

Furthermore, we compare the two-class classification accuracy for the sampled point sets using our unsupervised network, and present the overall shape classification accuracy as plots in Figure 19. In this quantitative comparison, we can further show that the points produced from distinctiveness-guided sampling lead to higher classification accuracy, even with fewer points.

### 5.3 View Selection in 3D Scenes

Given a 3D scene, where are the distinctive regions to attend to? Using our unsupervised framework, we can find the distinctive regions in an input 3D scene and locate the best views to look at these regions. Here, we define the best views as those with maximized distinctiveness displayed in the views.

To find such views, we first sample local patches of 2,048 points in the input scene, and feed these patches as inputs to our network to predict per-point distinctiveness. Since our network is trained on individual objects, we crop regions of around one-meter diameter
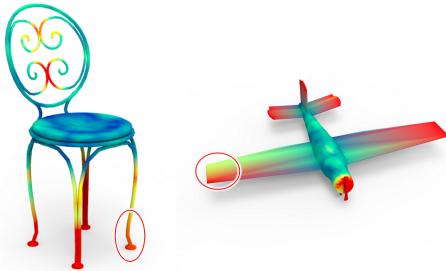
Fig. 20. Asymmetry in the detected distinctiveness.

in the scene for point sampling. Then, we can combine the results from the local patches to obtain distinctiveness over the scene. Next, we generate a set of candidate views by uniformly sampling 50 different views on the upper hemisphere that bounds the scene, and evaluate the quality of each view by averaging the distinctiveness over the visible points in the view. Lastly, we choose the view with the highest averaged distinctiveness as the best view.

Figures 18 (a) & (b) show the best views that were automatically selected for two different scenes (courtesy of 3D Warehouse [2019]), where the camera was set to look at the center of the scene when sampling the candidate views. As a comparison, we also show the corresponding worst views for the two scenes; see (c) & (d). From these results we can see that, the selected best views can reasonably present most distinctive regions in the views and contain rich information of the scene, while the worst views exhibit very limited information. Besides setting the camera to look at the whole scene, we can set it to look at specific areas or distinctive regions in the scene, and find *local best views*, meaning that we consider only the distinctive regions in the user-specified area when searching for the best view with maximized distinctiveness. The red and blue boxes in Figure 18 (a) mark two example areas, while Figures 18 (e) & (f) present the corresponding local best views found in the areas. Please refer to Supplementary Material Part I for more view selection results.

## 6 CONCLUSION AND FUTURE WORK

We presented a technique to learn and detect distinctive regions on 3D shapes. The technique analyzes a given set of objects *without any supervision*. The shapes are represented by point clouds, and the analysis is performed by a deep neural network that learns per-point and per-shape features from the point clouds. Further, we formulate the unsupervised joint loss for a shape clustering task of the per-shape features, thereby implicitly encouraging the network to learn the distinctiveness of the per-point features relative to the shapes in different clusters. We demonstrated the effectiveness of our method via extensive experiments, and presented several applications based on the network-predicted distinctiveness.

Despite the promising performance that our method has achieved, one limitation is that insufficient training samples or severely uneven class size would certainly affect the network's classification capability. However, such requirement on training data also appears in typical deep-learning-based methods. On the other hand, our method may detect asymmetrical distinctive regions on shapes; see

Figure 20. This may be related to our current distinctiveness extraction method, which is rather local when extracting the per-point distinctiveness. In the future, we plan to integrate the prior knowledge of shape symmetry into our current analysis framework when extracting the distinctiveness. While distinction and saliency are different (as explained in the introduction), we may further explore their relationship by taking a machine learning approach. Besides, armed with the distinctiveness analysis, in our future work, we are considering generating novel shapes with control on their distinctive features, making them more inter-class distinct and possibly more intra-class distinct, thereby enriching the variability within the class, while remaining distinct to the other classes. Generally speaking, we believe that a stronger and better set analysis will lead, in the future, to better synthesis of 3D shapes.

## REFERENCES

Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. TensorFlow: A system for large-scale machine learning. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. 265–283. https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf

Pierre Alliez, Mark Meyer, and Mathieu Desbrun. 2002. Interactive geometry remeshing. *ACM Transactions on Graphics (SIGGRAPH)* 21, 3 (2002), 347–354.

Marco Ancona, Enea Ceolini, Cengiz Oztireli, and Markus Gross. 2018. Towards better understanding of gradient-based attribution methods for deep neural networks. In *International Conference on Learning Representations (ICLR)*.

Yasuhiro Aoki, Hunter Goforth, Rangaprasad Arun Srivatsan, and Simon Lucey. 2019. PointNetLK: Robust & efficient point cloud registration using PointNet. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 7163–7172.

Philip Bachman, R Devon Hjelm, and William Buchwalter. 2019. Learning representations by maximizing mutual information across views. In *International Conference on Neural Information Processing Systems (NIPS)*. 15509–15519.

Umberto Castellani, Marco Cristani, Simone Fantoni, and Vittorio Murino. 2008. Sparse points matching by combining 3D mesh saliency with statistical descriptors. *Computer Graphics Forum (Eurographics)* 27, 2 (2008), 643–652.

Xiaobai Chen, Abulhair Saparov, Bill Pang, and Thomas Funkhouser. 2012. Schelling points on 3D surface meshes. *ACM Transactions on Graphics (SIGGRAPH)* 31, 4 (2012), 29:1–29:12.

Carl Doersch, Saurabh Singh, Abhinav Gupta, Josef Sivic, and Alexei Efros. 2012. What makes Paris look like Paris? *ACM Transactions on Graphics (SIGGRAPH)* 31, 4 (2012), 101:1–101:9.

Helin Dutagaci, Chun Pan Cheung, and Afzal Godil. 2012. Evaluation of 3D interest point detection techniques via human-generated ground truth. *The Visual Computer* 28, 9 (2012), 901–917.

Ran Gal and Daniel Cohen-Or. 2006. Salient geometric features for partial shape matching and similarity. *ACM Transactions on Graphics* 25, 1 (2006), 130–150.

Michael Garland and Paul S. Heckbert. 1997. Surface simplification using quadric error metrics. In *Proceedings of SIGGRAPH*. 209–216.

Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1735–1742.

Olivier J. Hénaff, Ali Razavi, Carl Doersch, S. M. Eslami, and Aaron van den Oord. 2019. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint*

*arXiv:1905.09272* (2019).

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531* (2015).

R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. 2019. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations (ICLR)*.

Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. 2018. Pointwise convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 984–993.

Mayank Juneja, Andrea Vedaldi, C.V. Jawahar, and Andrew Zisserman. 2013. Blocks that shout: Distinctive parts for scene classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 923–930.

Diederik P. Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

Truc Le and Ye Duan. 2018. PointGrid: A deep network for 3D shape understanding. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 9204–9214.

Chang Ha Lee, Amitabh Varshney, and David W. Jacobs. 2005. Mesh saliency. *ACM Transactions on Graphics (SIGGRAPH)* 24, 3 (2005), 659–666.

George Leifman, Elizabeth Shtrom, and Ayellet Tal. 2012. Surface regions of interest for viewpoint selection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 414–421.

Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. 2018. PointCNN: Convolution on $\mathcal{X}$-transformed points. In *International Conference on Neural Information Processing Systems (NIPS)*. 828–838.

Yu Liu, Guanglu Song, Jing Shao, Xiao Jin, and Xiaogang Wang. 2018. Transductive centroid projection for semi-supervised large-scale recognition. In *European Conference on Computer Vision (ECCV)*. 70–86.

Flora Ponjou Tasse, Jiri Kosinka, and Neil Dodgson. 2015. Cluster-based point set saliency. In *IEEE International Conference on Computer Vision (ICCV)*. 163–171.

Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. 2017a. PointNet: Deep learning on point sets for 3D classification and segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 652–660.

Charles R. Qi, Li Yi, Hao Su, and Leonidas J. Guibas. 2017b. PointNet++: Deep hierarchical feature learning on point sets in a metric space. In *International Conference on Neural Information Processing Systems (NIPS)*. 5099–5108.

Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-CAM: Visual explanations from deep networks via gradient-based localization. In *IEEE International Conference on Computer Vision (ICCV)*. 618–626.

Yiru Shen, Chen Feng, Yaoqing Yang, and Dong Tian. 2018. Mining point cloud local structures by kernel correlation and graph pooling. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 4548–4557.

Philip Shilane and Thomas Funkhouser. 2006. Selecting distinctive 3D shape descriptors for similarity retrieval. In *IEEE Intl. Conf. on Shape Modeling and Applications (SMI)*. 18:1–18:10.

Philip Shilane and Thomas Funkhouser. 2007. Distinctive regions of 3D surfaces. *ACM Transactions on Graphics* 26, 2 (2007), 7:1–7:15.

Philip Shilane, Patrick Min, Michael Kazhdan, and Thomas Funkhouser. 2004. The Princeton shape benchmark. In *IEEE Intl. Conf. on Shape Modeling and Applications (SMI)*. 167–178.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning Important Features Through Propagating Activation Differences. In *Proceedings of International Conference on Machine Learning (ICML)*. 3145–3153.

Elizabeth Shtrom, George Leifman, and Ayellet Tal. 2013. Saliency detection in large point sets. In *IEEE International Conference on Computer Vision (ICCV)*. 3591–3598.

Zhenyu Shu, Shiqing Xin, Xin Xu, Ligang Liu, and Ladislav Kavan. 2019. Detecting 3D points of interest using multiple features and stacked auto-encoder. *IEEE Transactions Visualization & Computer Graphics* 25, 8 (2019), 2583–2596.

Saurabh Singh, Abhinav Gupta, and Alexei A. Efros. 2012. Unsupervised discovery of mid-level discriminative patches. In *European Conference on Computer Vision (ECCV)*. 73–86.

Ran Song, Yonghuai Liu, and Paul Rosin. 2018. Distinction of 3D objects and scenes via classification network and Markov random field. *IEEE Transactions Visualization & Computer Graphics* (2018), To appear.

X. Yu Stella and Jianbo Shi. 2003. Multiclass spectral clustering. In *IEEE International Conference on Computer Vision (ICCV)*. 313–320.

Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. 2018. SPLATNet: Sparse lattice networks for point cloud processing. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2530–2539.

Jian Sun and Jean Ponce. 2013. Learning discriminative part detectors for image classification and cosegmentation. In *IEEE International Conference on Computer Vision (ICCV)*. 3400–3407.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic Attribution for Deep Networks. In *Proceedings of International Conference on Machine Learning (ICML)*.

3319–3328.

Johan W. H. Tangelder and Remco C. Veltkamp. 2004. A survey of content based 3D shape retrieval methods. In *IEEE Intl. Conf. on Shape Modeling and Applications (SMI)*. 145–156.

Ulrike Von Luxburg. 2007. A tutorial on spectral clustering. *Statistics and computing* 17, 4 (2007), 395–416.

Xi Wang, Sebastian Koch, Kenneth Holmqvist, and Marc Alexa. 2018. Tracking the gaze on objects in 3D: how do people really look at the Bunny? *ACM Transactions on Graphics (SIGGRAPH Asia)* 37, 6 (2018), 188:1–188:18.

Yaming Wang, Jonghyun Choi, Vlad Morariu, and Larry S. Davis. 2016. Mining discriminative triplets of patches for fine-grained classification. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1163–1172.

Yue Wang and Justin M. Solomon. 2019. Deep Closest Point: Learning Representations for Point Cloud Registration. In *IEEE International Conference on Computer Vision (ICCV)*. 3523–3532.

Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. 2019. Dynamic graph CNN for learning on point clouds. *ACM Transactions on Graphics* 38, 5 (2019), 146:1–146:12.

3D Warehouse. 2019. https://3dwarehouse.sketchup.com/ [Online; accessed 02-Jan-2019].

Yandong Wen, Kaipeng Zhang, Zhifeng Li, and Yu Qiao. 2016. A discriminative feature learning approach for deep face recognition. In *European Conference on Computer Vision (ECCV)*. 499–515.

Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. 2018. CBAM: Convolutional block attention module. In *European Conference on Computer Vision (ECCV)*. 3–19.

Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 2015. 3D ShapeNets: A deep representation for volumetric shapes. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1912–1920.

Zhirong Wu, Yuanjun Xiong, X. Yu Stella, and Dahua Lin. 2018. Unsupervised feature learning via non-parametric instance discrimination. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 3733–3742.

Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. 2018. SpiderCNN: Deep learning on point sets with parameterized convolutional filters. In *European Conference on Computer Vision (ECCV)*. 90–105.

Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. 2018. FoldingNet: Point cloud auto-encoder via deep grid deformation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 206–215.

Lequan Yu, Xianzhi Li, Chi-Wing Fu, Daniel Cohen-Or, and Pheng-Ann Heng. 2018. EC-Net: an Edge-aware Point set Consolidation Network. In *European Conference on Computer Vision (ECCV)*. 398–414.

Matthew D. Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision (ECCV)*. 818–833.

Jianming Zhang, Sarah Adel Bargal, Zhe Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. 2018. Top-down neural attention by excitation backprop. *International Journal Computer Vision* 126, 10 (2018), 1084–1102.

Xiaoting Zhang, Xinyi Le, Athina Panotopoulou, Emily Whiting, and Charlie C. L. Wang. 2015. Perceptual models of preference in 3D printing direction. *ACM Transactions on Graphics (SIGGRAPH Asia)* 34, 6 (2015), 215:1–215:12.

Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. 2016. Learning deep features for discriminative localization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2921–2929.