

Deep Learning Representation using Autoencoder for 3D Shape Retrieval

Zhuotun Zhu, Xinggang Wang, Song Bai, Cong Yao, Xiang Bai

School of Electronic Information and Communications, Huazhong University of Science and Technology, 1037 Luoyu Road, Wuhan, Hubei Province 430074, P.R. China.

Abstract

We study the problem of how to build a deep learning representation for 3D shape. Deep learning has shown to be very effective in variety of visual applications, such as image classification and object detection. However, it has not been successfully applied to 3D shape recognition. This is because 3D shape has complex structure in 3D space and there are limited number of 3D shapes for feature learning. To address these problems, we project 3D shapes into 2D space and use autoencoder for feature learning on the 2D images. High accuracy 3D shape retrieval performance is obtained by aggregating the features learned on 2D images. In addition, we show the proposed deep learning feature is complementary to conventional local image descriptors. By combing the global deep learning representation and the local descriptor representation, our method can obtain the state-of-the-art performance on 3D shape retrieval benchmarks.

Keywords: 3D Shape Matching, 3D Shape Retrieval, Autoencoder, Shape Representation

1. Introduction

2 With the fast development of 3D printer, Microsoft Kinect sensor and laser
3 scanner, etc., there are more and more digitized 3D models that need to be
4 recognized. Thus it is critical to study how to build an efficient 3D shape search
5 engine. However, due to the intrinsic complex structure of 3D shape, it is hard
6 to handle 3D shape using a simple representation for efficient search.

7 Along with the development of computer vision and machine learning, deep
8 learning methods have been proven to be very effective for visual recognition.
9 For example, deep convolutional neural network (CNN) [1] has achieved the
10 state-of-the-art performance for object recognition on the ImageNet dataset [2]
11 and for object detection on the PASCAL dataset [3]. One reason of the success

Email addresses: zhuzhuotun@hust.edu.cn (Zhuotun Zhu), xgwang@hust.edu.cn (Xinggang Wang), songbai@hust.edu.cn (Song Bai), yaocong2010@gmail.com (Cong Yao), xbai@hust.edu.cn (Xiang Bai)

12 of deep learning for visual recognition is that the deep learning methods can au-
13 tomatically learn the features with the superior discriminatory power for image
14 representation, rather than using hand-crafted image descriptors. Currently, in
15 the context of 3D shape recognition, shape descriptors are mainly hand-crafted
16 and deep learning representation has not been widely applied. It seems that
17 it is hard to directly apply deep learning methods to 3D shape representation,
18 since deep learning methods need a large amount of data to bridge the visual
19 gap among training examples from the same object category; and it is unlikely
20 to learn a good representation using a few data with large visual variation.

21 The above developments of deep learning are in a supervised way and are
22 not suitable for retrieval task. From the aspect of unsupervised deep learning,
23 Hinton and Krizhevsky [4] proposed the autoencoder algorithm with the appli-
24 cation of image retrieval, which is then used for some other specific tasks like
25 face alignment [5]. Training autoencoder does not require any label informa-
26 tion. The autoencoder can be regarded as a multi-layer sparse coding network.
27 Each node in the autoencoder network can be regarded as a prototype of object
28 image/shape. From the bottom layer to the top layer, the prototype contain-
29 s richer semantic information and becomes a better representation. After the
30 autoencoder network is learnt, the coefficients obtained by reconstructing im-
31 age/shape based on prototypes are used as feature for 3D shape matching and
32 retrieval. Since the autoencoder can learn feature adaptively to training data,
33 it can get excellent performance for image retrieval.

34 Motivated by the view-based 3D shape methods [6, 7], in which a 3D shape
35 can be projected into many 2D depth images, we aim to use autoencoder to learn
36 a 3D shape representation based on the depth images obtained by projection. As
37 shown in Fig. 1, a 3D shape is projected into many different depth images; the
38 learnt autoencoder can reconstruct the depth images nicely. Matching 3D shape
39 based on the autoencoder features can be converted to a set-to-set matching
40 problem, conventional set-to-set distance, like the Hausdorff distance, can be
41 adopted. Our autoencoder based 3D shape representation is a deep learning
42 representation; compared to the representations based on local descriptor, e.g.
43 SIFT, it is a global representation. This global deep learning representation and
44 the representation based on local descriptors are complementary to each other.

45 In summary, the main contributions of this paper are: (1) A new method to
46 learn deep learning representation for 3D shape using autoencoder; (2) combin-
47 ing the global deep learning representation with local descriptor representation
48 and obtaining the state-of-the-art 3D shape retrieval performance.

49 The remainder of this paper is organised as follows: In Section 2, we offer
50 an overview of the previous work on the content-based 3D shape retrieval. In
51 Section 3, we present an explicit description of our method to extract the global
52 features of 3D shape. In Section 4, we briefly depict the local descriptor formerly
53 implemented in [8] on 3D shape. Experimental results and extensive evaluation
54 are then carried out in Section 5. At last, we conclude this paper in Section 6.

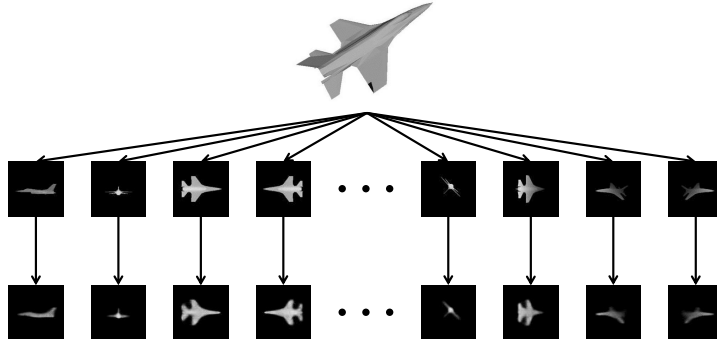


Figure 1: A specific illustration of our method to reconstruct 2D images. Note that the first row displays the original depth images in gray-scale of the 3D shape, while the second row shows the reconstructed ones corresponding to the images of the first row. And the black dots indicates those extracted from other different views.

55 2. Related Work

56 Based on the main idea that “two 3D models are similar if they look sim-
 57 ilar with each other from all viewing angles”, there are plenty of view-based
 58 approaches that have been regarded as the most discriminative methods on lit-
 59 erature [9]. Since our shape descriptor is also view-based, we mainly discuss
 60 some effective, competing view-based approaches during the following part.

61 In [10], Cyr and Kimia recognized a 3D shape by comparing a view of the
 62 shape with all views of 3D objects using shock graph matching. Osada et al. [11]
 63 proposed the shape distribution descriptor that measures properties based on
 64 area, angle, distance and volume measurements between a random set of points
 65 on the object. The similarity between two objects is defined by suitable shape
 66 functions, e.g. the D2 function. Ohbuchi et al. [8] utilized local visual features by
 67 using the Scale Invariant Feature Transform (SIFT) [12] to retrieve 3D shapes.
 68 A host of local features describing the 3D models is integrated into a histogram
 69 using Bag-of-Features [13] to reduce the computation complexity. Vranic [14] p-
 70 resented a composite 3D shape feature vector (DESIRE) which consists of depth
 71 buffer images, silhouettes and ray-extents of a polygonal mesh. The composite
 72 of various feature vectors extracted in a canonical coordinate frame generally
 73 performs better than the single method which relies on pairwise alignment of 3D
 74 objects. Later on, Papadakis et al. [15] made use of a hybrid descriptor (Hybrid)
 75 which consists of both depth buffer based 2D features and spherical harmonies
 76 based 3D features. The Hybrid adopts two alignment methods to compensate
 77 inner rotation variance and then uses Huffman coding to further compress fea-
 78 ture descriptors. Also, they presented a 3D descriptor (PANORAMA) [16] that
 79 captures the panoramic view of a 3D shape by projecting it to a lateral surface
 80 of a cylinder parallel to one of its three principal axes. By aligning its princi-

81 ple axes to capture the global information and combining 2D Discrete Fourier
 82 Transform and 2D Discrete Wavelet Transform, the PANORAMA outperforms
 83 all the other 3D shape retrieval methods on several standard 3D benchmarks.
 84 Meanwhile, Lian et al. [7] used Bag-of-Features and Clock Matching (CM-BoF)
 85 on a set of depth-buffer views obtained from the projections of the normalized
 86 object. The CM-BoF method also takes advantage of the preserved local details
 87 as well as isometry-invariant global structure to reach a competing result. Prior
 88 to that, they also proposed a shape descriptor named Geodesic Sphere based
 89 Multi-view Descriptors (GSMD) [17] measuring the extend to which a 3D poly-
 90 gon is rectilinear based on the maximum ratio of the surface area to the sum
 91 of three orthogonal projected areas. Recently, Bai et al. [18] adopted contour
 92 fragments as the input features for learning a BoW model, which is general and
 93 efficient for both 2D and 3D shape matching.

94 Among the view-based 3D shape retrieval methods, the Light Field descrip-
 95 tor (LFD) [6] may be the most famous approach. The extraction of LFD begins
 96 with the scale and translation normalization while achieving rotation invariance
 97 via an exhaustive searching from a great many of views. Then the silhouette pro-
 98 jections, rendered from uniformly sampled positions on a unit sphere, are repre-
 99 sented by 10 Fourier coefficients [19] and 35 Zernike moments coefficients [20].
 100 Finally, the dissimilarity between two objects is measured by the minimum
 101 distance of all group matching pairs. The LFD is insensitive to similarity trans-
 102 form, geometry degeneracy and noise, etc, thus shows better performance than
 103 other competing approaches.

104 3. Deep Learning Representation using Autoencoder

105 In this Section, given a 3D shape model S , we show how to perform au-
 106 toencoder initialized with deep belief network for S and then conduct 3D shape
 107 retrieval based on the calculated shape code. As shown in Fig. 2, we illustrate
 108 a specific flow chart about the whole procedure.

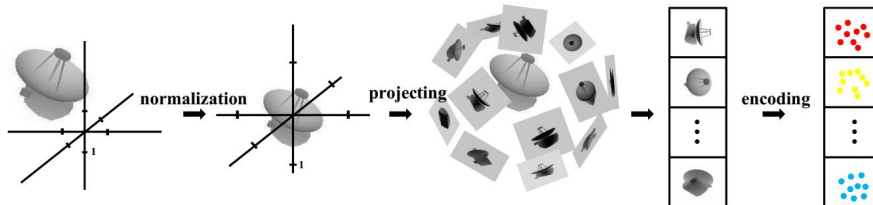


Figure 2: The flow chart of 3D shape representation using autoencoder. First, we conduct pose normalization for differences in translation and scale to each 3D model. Next, each 3D shape is represented by a set of depth-buffer images. Finally all the projections are used to train the autoencoder to acquire the code as a low-dimensional representation of the depth images, based on which to conduct 3D shape retrieval. In the last image, the colored dots indicate those features extracted from the corresponding depth images.

109 *3.1. Depth Projection Image*

Different from shapes of 2D images, 3D models represent the 3D objects using a collection of points in 3D space, connected by various geometric entities such as lines, curved surfaces, etc. In our method, the autoencoder initialized by a DBN described in Section 3.2 is used to reconstruct the gray-scale depth 2D images as input and acts as a low-dimensional coding method. Thus, projecting a 3D model to a collection of 2D images is required to make it possible. For a 3D shape model S preprocessed by scale and translation normalization, from a host of angles of view, we collect 2D projections set of S defined as

$$\mathbf{P}(S) = \{V_1, V_2, \dots, V_{Np}\}, \quad (1)$$

110 where Np denotes the number of projections for each model.

111 More specifically, Fig. 3 illustrates how we obtain a series of projections for
 112 the shape S viewed from different angles both in azimuth and elevation.

113 *3.2. Deep Belief Network*

114 The deep belief network (DBN) [21, 22, 23] is a generative graphical model,
 115 or alternatively a type of deep neural network, composed of multiple layers
 116 of latent variables (“hidden units”), with connections between the layers but
 117 not between units within each layer. When trained on plenty of examples in
 118 an unsupervised way, a DBN can probabilistically reconstruct the inputs by
 119 learning a stack of Restricted Boltzmann Machines (RBMs), where each of the
 120 previous RBM’s hidden layer serves as the visible layer for the next. That is
 121 to say, each time a new RBM is added to the stacked structure of DBN, then
 122 the new DBN has a better variational lower bound in the log probability of the
 123 data than the previous DBN [4].

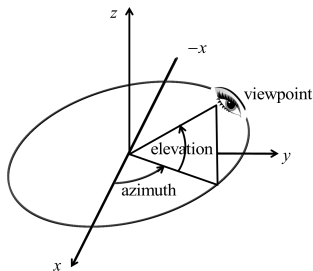


Figure 3: The illustration of how we get the projections of a 3D shape model S . Azimuth is the polar angle in the x - y plane, with positive number indicating anticlockwise rotation of the viewpoint. As for elevation, positive and negative numbers are the angle above and below the x - y plane respectively.

We introduce the “pretraining” procedure as shown in Fig. 4 for binary units, then generalize to real-valued units and show that it works well. The pixels correspond to the “visible units” since their states can be observed; as

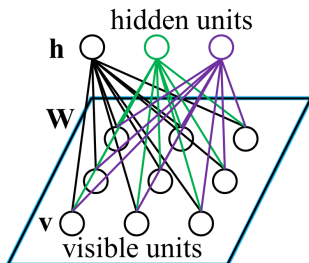


Figure 4: A graphical description of RBM. Note that a standard type of RBM has binary-valued visible and hidden units with weights of the connection between them. What needs to be specially emphasized is that there are none connections within visible units or hidden ones, which leads to a property that the hidden unit activations are mutually independent given the activations of visible units and conversely.

for the feature detectors, they correspond to the “hidden units”. The energy of a joint configuration (\mathbf{v}, \mathbf{h}) for the visible and hidden units is defined in [24] as

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} w_{ij} v_i h_j, \quad (2)$$

124 where v_i, h_j denote the binary states of visible unit i and hidden unit j respec-
 125 tively; a_i, b_j are their biases and w_{ij} is the connection weight between them.

The network assigns a probability to every possible couple of a visible vector and a hidden one by the following function

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}, \quad (3)$$

where the “partition function” Z is given by the sum of all possible pairs between visible and hidden vectors

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}. \quad (4)$$

The probability that the network assigns to a visible vector, is defined as the sum of all possible hidden vectors

$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}. \quad (5)$$

The probability of a training image can be increased by adjusting the biases and weights to lower the energy of that image but to increase the energy of the rest, especially for these that own low energy and thus are assigned high probability by the network and make great contribution to the partition function. The mathematically derived derivative of the log probability of a visible vector to a weight is simple:

$$\frac{\partial \log p(\mathbf{v})}{\partial w_{ij}} = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}}, \quad (6)$$

where the angle brackets denote expectations under the exact distribution specified by the subscript that follows. Thus, utilizing stochastic gradient descent as the learning approach is a very simple way in the log probability of training data

$$\Delta w_{ij} = \epsilon(\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}), \quad (7)$$

126 where the ϵ is the learning rate.

Because of the RBM's restricted structure that there are no direct connections within hidden units, it is pretty easy to obtain an unbiased sample of $\langle v_i h_j \rangle_{data}$. Given a training image as the visible vector \mathbf{v} , the binary state h_j of every hidden unit j is set to 1 with the probability

$$p(h_j = 1 | \mathbf{v}) = S(b_j + \sum_{i \in visible} w_{ij} v_i), \quad (8)$$

127 where $S(x)$ denotes the sigmoid function defined by the formula $1/[1 + \exp(-x)]$.

Given a hidden vector \mathbf{h} , it is also quite easy to obtain an unbiased sample of a visible unit's state as a consequence of no connections within visible units. The first equation corresponds with the construction of binary visible units and the second one with linear visible units, where $N(\mu, \sigma)$ is a Gaussian with mean value μ and standard deviation σ .

$$p(v_i = 1 | \mathbf{h}) = S(a_i + \sum_{j \in hidden} w_{ij} h_j), \text{ or} \quad (9)$$

$$v_i = N(a_i + \sum_{j \in hidden} w_{ij} h_j, 1).$$

128 Obtaining an unbiased sample of $\langle v_i h_j \rangle_{model}$, however, is much more tough.
 129 It can be done by beginning with any random state of a visible vector and
 130 performing alternating Gibbs sampling for quite a long time. One iteration of
 131 Gibbs sampling is used to update all the hidden units in parallel applying (8)
 132 followed by updating all the visible units in parallel applying (9).

Fortunately, a much faster learning algorithm was proposed in [25]. This algorithm begins by setting the visible units' states to a training vector. Then the whole hidden units' binary states are calculated in parallel applying (8). After those binary states have been probabilistically chosen for the hidden units, a "confabulation" is produced via setting each visible unit v_i to 1 with probability as in (9). Update the states of the hidden units once more in order that they can represent features of the confabulation. Then the adjustment of the weight is formulated by

$$\Delta w_{ij} = \epsilon(\langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{recon}), \quad (10)$$

133 where the $\langle v_i h_j \rangle_{data}$ is the fraction of times that the visible unit i and the hidden
 134 unit j are on together when the hidden units are driven by data, and $\langle v_i h_j \rangle_{recon}$
 135 is the corresponding part given by the confabulation. A same learning rules is
 136 used to adjust the biases.

137 In our experiments, this fast learning procedure works out well even though
 138 it is just approximating the derivative of the log probability with respect to the
 139 training data.

140 *3.3. Fine-tuning the Autoencoder*

141 After pretraining a DBN which acts as initialization of an autoencoder, a
 142 global fine-tuning procedure replaces the former stochastic, binary activities
 143 with crucial, real-valued probabilities and uses backpropagation through the
 144 whole structure of autoencoder to adjust the weights as well as biases for a
 145 reconstruction model. By minimizing the root mean squared reconstruction
 146 error $\sqrt{\sum_i (\langle v_i \rangle_{data} - \langle v_i \rangle_{recon})^2}$, we finally obtain a deep-structured, optimal
 147 reconstruction model of the 2D depth images as input.

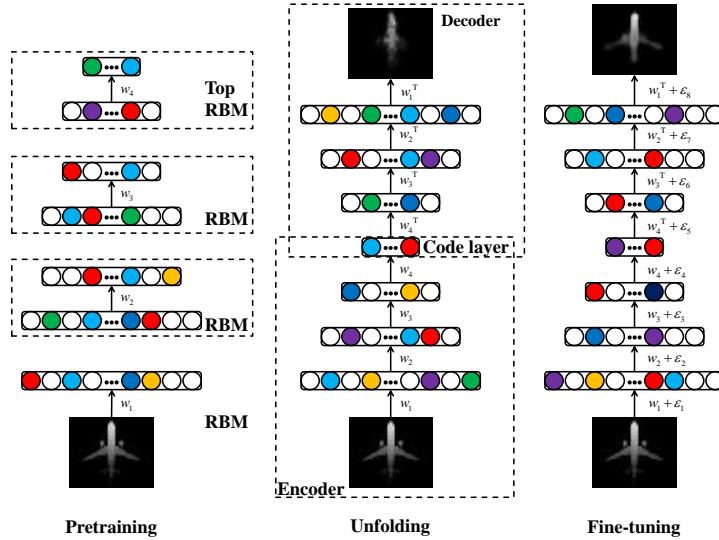


Figure 5: Details of autoencoder implemented on depth images. The circles enclosed by rectangle in each layer denote the units with various filling colour indicating different probability that the network assigns to them, and the rectangle’s length corresponds to the relative size of dimension on that layer. As we can see, the reconstruction performance becomes much better after doing the fine-tuning procedure compared to the only pretraining procedure done, which ensures the low-dimensional code layer being a good representation of the 2D image and has a great influence on the retrieval results.

148 To sum up, the whole autoencoder system is depicted in Fig. 5. Pretraining
 149 consists of a stacked RBMs where the hidden units in the previous layer acts
 150 as the visible units of the next layer. Then the “unfolded” autoencoder initialized
 151 by DBN is fine-tuned to obtain a better reconstruction performance. Finally,
 152 the code layer that is an efficient representation of the input image is utilized
 153 to conduct 3D retrieval.

154 *3.4. Set-to-Set Distance*

After projecting 3D model and then reconstructing 2D depth images, we get a low-dimensional representation of S with a code set C

$$C(S) = \{\vec{C}_1, \vec{C}_2, \dots, \vec{C}_{N_p}\}, \quad (11)$$

where Np denotes the number of projection images of each model; and \vec{C}_i ($i = 1, 2, \dots, Np$) denotes the coding vector corresponds to the projection V_i with respect to that shape model S , defined by

$$\vec{C}_i = (c_{i1}, c_{i2}, \dots, c_{iNc}), \quad (12)$$

155 where Nc denotes the dimensionality of every code vector; c_{ij} is the value of
 156 j -th dimensionality corresponding to code vector \vec{C}_i .

Based on the effective and efficient autoencoder, we can obtain the quantified distance within each 3D model by defining specific distance method given any two shape model S_A and S_B , whose code sets are as follows

$$\begin{aligned} \mathbf{C}(S_A) &= \{\vec{C}_{A_1}, \vec{C}_{A_2}, \dots, \vec{C}_{A_{Np}}\} \\ \mathbf{C}(S_B) &= \{\vec{C}_{B_1}, \vec{C}_{B_2}, \dots, \vec{C}_{B_{Np}}\}, \end{aligned} \quad (13)$$

157 where A_i and B_i denote the i -th projection index of model S_A , S_B respectively.

We use one variant of ‘‘Hausdorff Distance’’ to define the distance of S_A to S_B , given by

$$D(S_A, S_B) = \frac{1}{Np} \sum_{i=1}^{Np} \min_j d\{\vec{C}_{A_i}, \vec{C}_{B_j}\}, \quad (14)$$

158 where $d\{\vec{C}_{A_i}, \vec{C}_{B_j}\}$ denotes one specific distance function between two vector,
 159 such as p-norm distance in ‘‘Euclidean Space’’, algebraic distance, etc. Depend-
 160 ing on the distance of any two models, shape retrieval could be directly done
 161 according to the ranked list.

162 4. Bag of Features Representation

163 In this Section, we describe the local descriptor formerly implemented by
 164 Ohbuchi et al. [8] on 3D shape. Considering that our method autoencoder
 165 mentioned above is a global descriptor, it is much reasonable to boost a better
 166 performance if combining with a local descriptor. Bag-of-Features using Scale-
 167 invariant feature transform (BoF-SIFT) model is selected as the local description
 168 for a 3D model. Different from previous work in [8] that considers the SIFTs
 169 of each depth image separately, we put all SIFTs in a single bag, *i.e.*, rotation
 170 normalization is not conducted.

171 We first learn the visual word vocabulary with size of 1500 in a randomly
 172 selected subset of all features via K-means off-line. In order to encode the set of
 173 SIFTs in each 3D model, we conduct Vector Quantization proposed in [26] to get
 174 a histogram representation that counts the number of SIFTs belonging to each
 175 visual word. Before computing the pairwise distance among the models, all the
 176 histogram is L_1 normalized. We will display the good property of extraordinary
 177 complementarity between autoencoder and BoF-SIFT in Section 5.

178 **5. Experiments**

179 In this Section, we test our method on two widely used, standard datasets
180 of 3D shapes and compare our results with the-state-of-the-art approaches for
181 3D shape retrieval. The algorithm is implemented in MATLAB and experi-
182 ments are carried out on a laptop machine with Intel(R) Core(TM) i5-3210M
183 CPU(2.5GHz) and 4GB memory.

184 *5.1. Princeton Shape Benchmark (PSB)*

185 The *Princeton Shape Benchmark* [9] dataset provides a repository of 3D
186 models and software tools for comparing different shape-based models. It's
187 created to promote the use of standardized datasets and evaluate methods for
188 research in matching, classification, clustering, and recognition of 3D models.
189 Each model of the 3D shape consists of the polygonal geometry surface of the
190 corresponding shape. There are totally 1814 models and the base classification
191 is partitioned equally into training and testing sets. The training set with
192 90 classes, 907 models is used to attain parameters of shape models through
193 training procedure, while the other with 92 classes, equal number of models for
194 comparison with other algorithm. In addition, the number of models belonging
195 to the same class in the base classification varies from each class and ranges from
196 4 to 50. Some 3D models from the PSB are randomly selected to be exhibited
197 in Fig. 6.

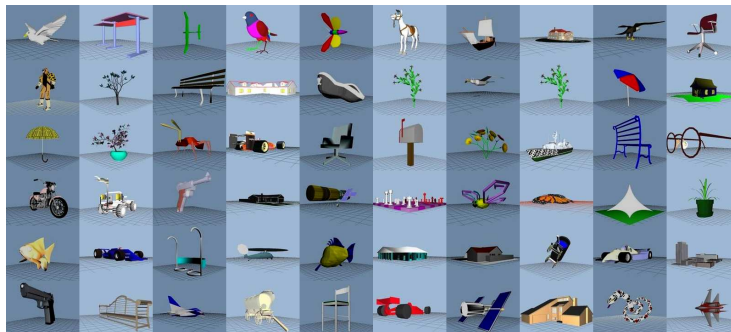


Figure 6: Exemplar images randomly chosen from the PSB dataset. The base classification spans a large various of classes including animals, buildings, etc.

198 *5.2. Engineering Shape Benchmark (ESB)*

199 The *Engineering Shape Benchmark* [27] is particularly proposed to evalu-
200 ate shape-based searching methods relevant to the mechanical engineering do-
201 main. More specifically, the ESB dataset has totally 867 3D CAD models clas-
202 sified into 45 classes with the number of models ranging from 4 to 58 in a class.
203 The 3D models contained in the ESB cover a wide variety of real-world engi-
204 neering models so that different methods can compete with each other more
205 fairly. As shown in Fig. 7, we randomly select some models in the ESB to show
206 engineering properties of the models.

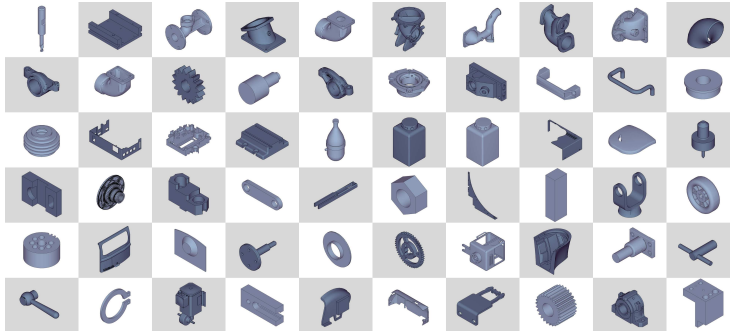


Figure 7: Exemplar images randomly chosen from the ESB dataset. Compared to the PSB dataset, all of the models contained in the ESB are mechanical engineering objects (parts) such as bearing assemblies, spacer, spinner, etc.

207 *5.3. Implementation Details*

208 As described in Section 3.1, we set the number of each model’s projection to
 209 64 (8×8) on the dataset. Then the total raw gray-scale images with real value
 210 in the range of $[0, 1]$, preprocessed by transform invariant low-rank textures
 211 (TILT) [28] to eliminate the large orientation variance, served as the visible
 212 units of the DBN’s first layer.

213 More specifically, the visible units of the first RBM layer were the normalized
 214 value of the depth images’ pixels. When training higher level layer, the visible
 215 units of a RBM were set to the activation probabilities of the previous RBM’s
 216 hidden units. As for the hidden units, they had stochastic binary value except
 217 the top layer’s hidden units, which had stochastic real-valued states calculated
 218 from the unit standard deviation Gaussian whose mean value was defined by
 219 the input from that RBM’s logistic visible units. The real-valued states are in
 220 the range $[0, 1]$, compared to the binary states either 0 or 1, allowed the low-
 221 dimensional codes to take good advantage of continuous data and could avoid
 222 unnecessary sampling noise. Note that we trained each RBM for 40 epochs using
 223 mini-batches of size 100 and adopted a learning rate of 0.1 for the linear-binary
 224 RBMs, 0.001 for the top layer RBM.

225 With the DBN structure constructed, we initialized an autoencoder with
 226 the weights trained from the DBN and fine-tuned them using backpropagation
 227 as described in Section 3.3. The autoencoder consisted of an encoder with the
 228 designed layers and a symmetric structure for the decoder. The hidden units
 229 in the last layer were linear while all the other units were logistic. The deep,
 230 well-trained autoencoder was able to find how to convert each depth image
 231 into low-dimensional code that leads to a discriminative description and well
 232 reconstruction.

233 Then all the parameters including weights and biases are well-trained in
 234 an unsupervised way, we used them to obtain the low-dimensional code for
 235 projections of 3D models on the dataset. For the PSB, we constructed an
 236 encoder with the layers structure of 5184 (72×72)-1000-500-250-30 while a

237 structure of 5184 (72×72)-2000-500-100-20 for the ESB. In addition, we only
 238 used the testing set to both train the parameters and evaluate our results for
 239 the PSB while experiments were done on the whole dataset of the ESB since it
 240 provides no training set or testing set.

Finally, we define the distance function as mentioned in Section 3.3 as

$$d\{\overrightarrow{C_{A_i}}, \overrightarrow{C_{B_j}}\} = \|\overrightarrow{C_{A_i}} - \overrightarrow{C_{B_j}}\|_p, \quad p = 2, \quad (15)$$

241 where $\|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}}$, please note that x is a vector in the
 242 n -dimensional real vector space \mathfrak{R}^n .

243 5.4. Evaluation Methods

244 In this Section, we introduce statistical description for the retrieval perfor-
 245 mance of a specific algorithm. The PSB provides open source code for evaluating
 246 different algorithms and judging how well one algorithm is compared to others.
 247 Thus, the performance can be fairly judged by the same evaluation tools in
 248 varieties of perspectives. When any doubt comes to you, please refer to [9] for
 249 more details about definition of every evaluation method.

250 **Nearest Neighbor (NN):** the percentage of the closest matches that be-
 251 long to the same class as the query. This statistic offers us an indication of how
 252 well a nearest neighbor classifier could perform. As we can see, higher score
 253 represents better performance.

254 **First-Tier (FT) and Second-Tier (ST):** the percentage of models in the
 255 query’s class that appear within the top M matches. where M is determined
 256 by the size of the query’s class. Given that the query’s class owns C models,
 257 $M = C - 1$ for the first-tier and $M = 2(C - 1)$ for the second tier.

258 The three statistics mentioned above put emphases upon different aspects.
 259 The Nearest Neighbor (NN) evaluation merely lays emphasis on the discrimina-
 260 tive ability since it only accounts for the most similar one in the retrieved, sorted
 261 list. However, the First-Tier (FT) and Second-Tier (ST) indicate how well the
 262 average performance of an algorithm taking into consideration the tradeoff be-
 263 tween intra-class variation and inter-class discrepancy.

264 5.5. Retrieval Results

265 As shown in Table 1 and 2, we compare the global-feature-based autoencoder
 266 with the other global descriptors on the two standard datasets to explore the
 267 efficacy of using autoencoder to tackle 3D shape retrieval.

268 We come to a solid conclusion that the autoencoder is more efficient than
 269 the other global-features-based methods for 3D shape retrieval.

270 5.6. Complementary Property

271 We adopt autoencoder described in Section 3 to obtain the distance between
 272 any two shape models on the dataset. Then we get the retrieval results eval-
 273 uated by the source code provided in the PSB [9]. Furthermore, based on the
 274 knowledge that autoencoder reconstructs global information while BoF-SIFT

Table 1: Statistic evaluation of global descriptors on Princeton Shape Benchmark

| Algorithm | NN(%) | FT(%) | ST(%) |
|-------------|-------------|-------------|-------------|
| Autoencoder | 72.4 | 43.3 | 54.6 |
| GSMD [17] | 67.1 | 41.8 | 52.0 |
| DESIRE [14] | 65.8 | 40.4 | 51.3 |
| LFD [6] | 65.7 | 38.0 | 48.7 |

Table 2: Statistic evaluation of global descriptors on Engineering Shape Benchmark

| Algorithm | NN(%) | FT(%) | ST(%) |
|-------------|-------------|-------------|-------------|
| Autoencoder | 85.7 | 47.9 | 63.1 |
| Hybrid [15] | 82.9 | 46.5 | 60.5 |
| DESIRE[14] | 82.3 | 41.7 | 55.0 |
| LFD [6] | 82.0 | 40.4 | 53.9 |

275 described in Section 4 captures the local details, a linear combination of them is
 276 proposed to boost the retrieval performance. More specifically, we empirically
 277 choose the weights as $W_{global} = W_{local}$ for global and local descriptors.

Table 3: Statistic evaluation on Princeton Shape Benchmark

| Algorithm | NN(%) | FT(%) | ST(%) |
|----------------------|-------------|-------------|-------------|
| Autoencoder+BoF-SIFT | 77.5 | 52.4 | 65.4 |
| BoF-SIFT [8] | 71.4 | 45.1 | 57.6 |
| CM-BoF+GSMD [7] | 75.4 | 50.9 | 64.0 |
| PANORAMA [16] | 75.3 | 47.9 | 60.3 |
| CM-BoF [7] | 73.1 | 47.0 | 59.8 |

278 We compare our hybrid method (Autoencoder+BoF-SIFT) with the previ-
 279 ous state-of-the-art methods including PANORAMA, CM-BoF and CM-BoF+GSMD,
 280 which are able to capture both the global and local information of a 3D shape.
 281 For the retrieval results displayed in Table 3 on the PSB dataset, we can find
 282 that: our autoencoder shows pretty well complementary property with the ex-
 283 isting local-features-based method BoF-SIFT, whose retrieval results of FT and
 284 ST are both improved by more than 7 percent.

285 Furthermore, Fig. 8 shows a precision-recall plot of six methods on the PSB
 286 and ESB dataset respectively. Among all the methods, the composite of autoen-
 287 coder and BoF-SIFT achieves the best performance. The proposed autoencoder
 288 consistently outperforms other global-descriptors-based methods.

289 6. Conclusions

290 In this paper, we present a novel view-based 3D shape retrieval method using
 291 autoencoder, which is firstly utilized to 3D shape retrieval. A set of experiments

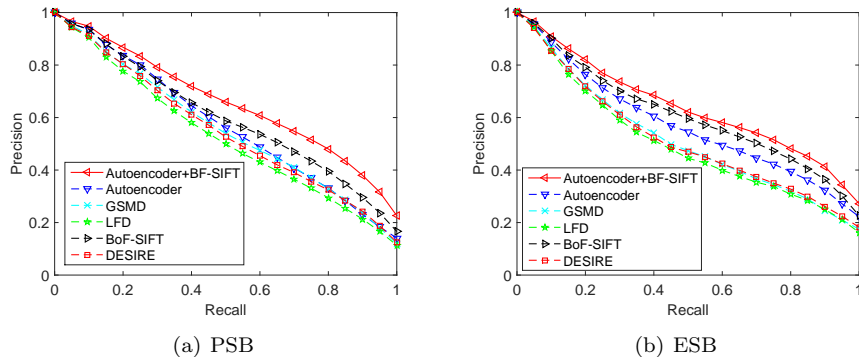


Figure 8: Precision-recall curves of nine methods implemented on two standard benchmarks. (a), (b) illustrates the results evaluated on the PSB and ESB respectively.

292 were carried out to investigate the effectiveness and efficiency of our method
 293 on two standard datasets, which shows that the autoencoder outperforms other
 294 global descriptor on retrieval results. Furthermore, the experiments demonstrate
 295 that the autoencoder displays good complementarity with the local descriptor,
 296 for linearly combining them achieves the state-of-the-art performance. Our future
 297 work might focus on studying the effect of the proposed representation with
 298 context-based shape similarity method [29].

299 7. Acknowledgement

300 This work was primarily supported by National Natural Science Foundation
 301 of China (NSFC) (No.61222308), and Program for New Century Excellent Tal-
 302 ents in University (No.NCET-12-0217), Fundamental Research Funds for the
 303 Central Universities (No.HUST 2013TS115). Xinggong Wang was supported
 304 by Microsoft Research Fellow Award 2012 and Excellent Ph.D Thesis Founding
 305 of HUST 2014.

306 References

- 307 [1] Yann LeCun, Bernhard Boser, John S Denker, Donnie Henderson,
 308 Richard E Howard, Wayne Hubbard, and Lawrence D Jackel. Backpropa-
 309 gation applied to handwritten zip code recognition. *Neural Computation*,
 310 1(4):541–551, 1989.
- 311 [2] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classifi-
 312 cation with deep convolutional neural networks. In *NIPS*, volume 1, page 4,
 313 2012.
- 314 [3] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich fea-
 315 ture hierarchies for accurate object detection and semantic segmentation.
 316 *arXiv preprint arXiv:1311.2524*, 2013.

- 317 [4] Alex Krizhevsky and Geoffrey E Hinton. Using very deep autoencoders for
318 content-based image retrieval. In *ESANN*. Citeseer, 2011.
- 319 [5] Jie Zhang, Shiguang Shan, Meina Kan, and Xilin Chen. Coarse-to-fine
320 auto-encoder networks (cfan) for real-time face alignment. In *European*
321 *Conference on Computer Vision*, 2014.
- 322 [6] Dingyun Chen, Xiaopei Tian, Yute Shen, and Ming Ouhyoung. On visual
323 similarity based 3D model retrieval. *Computer Graphics Forum*, 22:223–
324 232, 2003.
- 325 [7] Zhouhui Lian, A. Godil, and Xianfang Sun. Visual similarity based 3d shape
326 retrieval using bag-of-features. In *Shape Modeling International Conference*
327 *(SMI), 2010*, pages 25–36, June 2010.
- 328 [8] Ryutarou Ohbuchi, Kunio Osada, Takahiko Furuya, and Tomohisa Banno.
329 Salient local visual features for shape-based 3D model retrieval. In *Shape*
330 *Modeling International*, pages 93–102, 2008.
- 331 [9] Philip Shilane, Patrick Min, Michael Kazhdan, and Thomas Funkhouser.
332 The princeton shape benchmark. In *Shape Modeling Applications, 2004.*
333 *Proceedings*, pages 167–178. IEEE, 2004.
- 334 [10] Christopher M. Cyr and Benjamin B. Kimia. 3D object recognition us-
335 ing shape similarity-based aspect graph. In *International Conference on*
336 *Computer Vision*, pages 254–261, 2001.
- 337 [11] Robert Osada, Thomas A. Funkhouser, Bernard Chazelle, and David P.
338 Dobkin. Shape distributions. *ACM Transactions on Graphics*, 21:807–832,
339 2002.
- 340 [12] David G Lowe. Distinctive image features from scale-invariant keypoints.
341 *International Journal of Computer Vision*, 60(2):91–110, 2004.
- 342 [13] Feifei Li and Perona Pietro. A bayesian hierarchical model for learning
343 natural scene categories. In *Computer Vision and Pattern Recognition,*
344 *2005. IEEE Computer Society Conference on*, volume 2, pages 524–531.
345 IEEE, 2005.
- 346 [14] Dejan V. Vranic. DESIRE: a composite 3D-shape descriptor. In *Inter-*
347 *national Conference on Multimedia Computing and Systems/International*
348 *Conference on Multimedia and Expo*, pages 962–965, 2005.
- 349 [15] Panagiotis Papadakis, Ioannis Pratikakis, Theoharis Theoharis, Georgios
350 Passalis, Stavros J. Perantonis, and Agia Paraskevi. 3D object retrieval
351 using an efficient and compact hybrid shape descriptor. pages 9–16, 2008.
- 352 [16] Panagiotis Papadakis, Ioannis Pratikakis, Theoharis Theoharis, and S-
353 tavros J. Perantonis. PANORAMA: A 3D shape descriptor based on
354 panoramic views for unsupervised 3D object retrieval. *International Jour-*
355 *nal of Computer Vision*, 89:177–192, 2010.

- 356 [17] Zhouhui Lian, Paul L. Rosin, and Xianfang Sun. Rectilinearity of 3D
357 meshes. *International Journal of Computer Vision*, 89:130–151, 2010.
- 358 [18] Xiang Bai, Cong Rao, and Xinggang Wang. Shape vocabulary: A robust
359 and efficient shape representation for shape matching. *IEEE Transactions*
360 *on Image Processing*, 29:3935–3949, 2014.
- 361 [19] Dengsheng Zhang and Guojun Lu. Shape-based image retrieval using gener-
362 ic Fourier descriptor. *Signal Processing-image Communication*, 17:825–848,
363 2002.
- 364 [20] Alireza Khotanzad and Yaw Hua Hong. Invariant image recognition by
365 Zernike moments. *IEEE Transactions on Pattern Analysis and Machine*
366 *Intelligence*, 12:489–497, 1990.
- 367 [21] Yoshua Bengio, Pascal Lamblin, Dan Popovici, and Hugo Larochelle.
368 Greedy layer-wise training of deep networks. In *NIPS*, 2006.
- 369 [22] Geoffrey E. Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning
370 algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.
- 371 [23] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hub-
372 bard, and L. D. Jackel. Backpropagation applied to handwritten zip code
373 recognition. *Neural Computation*, 1:541–551, 1989.
- 374 [24] J J Hopfield. Neural networks and physical systems with emergent collective
375 computational abilities. *Proceedings of the National Academy of Sciences*,
376 79(8):2554–2558, 1982.
- 377 [25] Geoffrey E Hinton. Training products of experts by minimizing contrastive
378 divergence. *Neural Computation*, 14(8):1771–1800, 2002.
- 379 [26] Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. Beyond bags of
380 features: spatial pyramid matching for recognizing natural scene categories.
381 In *Computer Vision and Pattern Recognition*, volume 2, pages 2169–2178,
382 2006.
- 383 [27] Subramaniam Jayanti, Yagnanarayanan Kalyanaraman, Natraj Iyer, and
384 Karthik Ramani. Developing an engineering shape benchmark for CAD
385 models. *Computer-aided Design*, 38:939–953, 2006.
- 386 [28] Zhengdong Zhang, Xiao Liang, Arvind Ganesh, and Yi Ma. TILT: trans-
387 form invariant low-rank textures. In *Asian Conference on Computer Vision*
388 *2010*.
- 389 [29] Xiang Bai, Xingwei Yang, L.J. Latecki, Wenyu Liu, and Zhuowen Tu.
390 Learning context-sensitive shape similarity by graph transduction. *IEEE*
391 *Transactions on Pattern Anal. Mach. Intell.*, 32:861–874, 2010.