

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

ARTICLE INFO

Article history:

Received 28 February 2015

Received in revised form

10 July 2015

Accepted 23 August 2015

Available online 8 April 2016

Keywords:

3D Shape Matching

3D Shape Retrieval

Autoencoder

Shape Representation

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 combining 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.

© 2016 Published by Elsevier B.V.

1. Introduction

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

Along with the development of computer vision and machine learning, deep learning methods have been proven to be very effective for visual recognition. For example, deep convolutional neural network (CNN) [1] has achieved the state-of-the-art performance for object recognition on the ImageNet dataset [2] and for object detection on the PASCAL dataset [3]. One reason of the success of deep learning for visual recognition is that the deep learning methods can automatically learn the features with the superior discriminatory power for image representation, rather than using hand-crafted image descriptors. Currently, in the context of 3D shape recognition, shape descriptors are mainly hand-crafted and deep learning representation has not been widely applied. It seems that it is hard to directly apply deep learning methods to 3D shape representation, since deep learning methods need a large amount of data to bridge the visual gap among

training examples from the same object category; and it is unlikely to learn a good representation using a few data with large visual variation.

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

Until now, few approaches based on deep learning frameworks have been proposed to deal with 3D shape retrieval. Following [6], Fang et al. [7] trained a deep neural network using Eigen-shape descriptor and Fisher-shape descriptor as target values to guide the network. Heat shape descriptor developed from Heat Kernel Signature is fed into the network. Wu et al. [8] constructed a large-scale 3D CAD model dataset to train a convolutional deep belief network. This network learns the distribution of 3D shapes with different categories and arbitrary poses. Therefore, adopting deep

* Corresponding author. Tel.: +86 027 87543236; fax: +86 027 87543236.

E-mail addresses: zhuotun@gmail.com (Z. Zhu), xgwang@hust.edu.cn (X. Wang), songbai@hust.edu.cn (S. Bai), yaocong2010@gmail.com (C. Yao), xbai@hust.edu.cn (X. Bai).

learning approaches for 3D shape retrieval needs to be further justified.

Different from recent works in [7] and [8], we adopt view-based approaches. Motivated by other view-based 3D shape methods [9,10], in which a 3D shape can be projected into many 2D depth images, we aim to use autoencoder to learn a 3D shape representation based on the depth images obtained by projection. As shown in Fig. 1, a 3D shape is projected into many different depth images; the learnt autoencoder can reconstruct the depth images nicely. Matching 3D shape based on the autoencoder features can be converted to a set-to-set matching problem, conventional set-to-set distance, like the Hausdorff distance, can be adopted. Our autoencoder based 3D shape representation is a deep learning representation; compared to the representations based on local descriptor, e.g. SIFT, it is a global representation. This global deep learning representation and the representation based on local descriptors are complementary to each other.

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

It is worth noting that we extended the conference version [6] in this manuscript as follows: (1) in Section 1, we added the discussions on recent deep learning methods for 3D shape analysis; (2) in the Section 2, we enriched it with the detailed description of LFD descriptor; and (3) in Section 5, we added new evaluation protocols (the precision-recall curve) and experiment results on the NTU dataset.

The remainder of this paper is organised in the following part: In Section 2, we offer an overview of the previous work on the content-based 3D shape retrieval. In Section 3, we present an explicit description of our method to extract the global features of 3D shape. In Section 4, we briefly depict the local descriptor formerly implemented in [11] on 3D shape. Experimental results and extensive evaluation are then carried out in Section 5. At last, we conclude this paper in Section 6.

2. Related work

Based on the main idea that “two 3D models are similar if they look similar with each other from all viewing angles”, there are plenty of view-based approaches that have been regarded as the most discriminative methods on literature [12]. Since our shape descriptor is also view-based, we mainly discuss some effective, competing view-based approaches during the following part.

In [13], Cyr and Kimia recognized a 3D shape by comparing a view of the shape with all views of 3D objects using shock graph matching. Osada et al. [14] proposed the shape distribution descriptor that measures properties based on area, angle, distance and volume measurements between a random set of points on the object. The similarity between two objects is defined by suitable shape functions, e.g. the D2 function. Ohbuchi et al. [11] utilized local visual features by using the Scale Invariant Feature Transform (SIFT) [15] to retrieve 3D shapes. A host of local features describing the 3D models is integrated into a histogram using Bag-of-Features [16] to reduce the computation complexity. Vranic [17] presented a composite 3D shape feature vector (DESIRE) which consists of depth buffer images, silhouettes and ray-extents of a polygonal mesh. The composite of various feature vectors extracted in a canonical coordinate frame generally performs better than the single method which relies on pairwise alignment of 3D objects. Later on, Papadakis et al. [18] made use of a hybrid descriptor (Hybrid) which consists of both depth buffer based 2D features and spherical harmonies based 3D features. The Hybrid adopts two alignment methods to compensate inner rotation variance and then uses Huffman coding to further compress feature descriptors. Also, they presented a 3D descriptor (PANORAMA) [19] that captures the panoramic view of a 3D shape by projecting it to a lateral surface of a cylinder parallel to one of its three principal axes. By aligning its principle axes to capture the global information and combining 2D Discrete Fourier Transform and 2D Discrete Wavelet Transform, the PANORAMA outperforms all the other 3D shape retrieval methods on several standard 3D benchmarks. Meanwhile, Lian et al. [10] used Bag-of-Features and Clock Matching (CM-BoF) on a set of depth-buffer views obtained from the projections of the normalized object. The CM-BoF method also takes advantage of the preserved local details as well as isometry-invariant global structure to reach a competing result. Prior to that, they also proposed a shape descriptor named Geodesic Sphere based Multi-view Descriptors (GSMD) [20] measuring the extend to which a 3D polygon is rectilinear based on the maximum ratio of the surface area to the sum of three orthogonal projected areas. Recently, Bai et al. [21] adopted contour fragments as the input features for learning a BoW model, which is general and efficient for both 2D and 3D shape matching.

Among the view-based 3D shape retrieval methods, the Light Field descriptor (LFD) [9] may be the most famous approach. The extraction of LFD begins with the scale and translation normalization while achieving rotation invariance via an exhaustive searching from a great many of views. Then the silhouette projections, rendered from uniformly sampled positions on a unit sphere, are represented by 10 Fourier coefficients [22] and 35

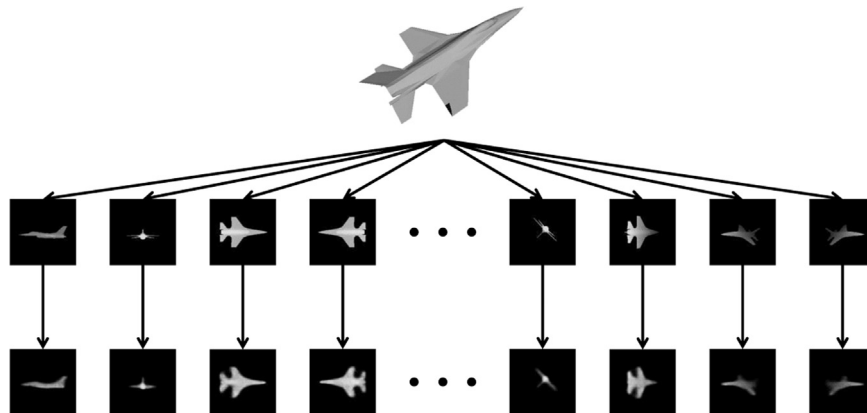


Fig. 1. A specific illustration of our method to reconstruct 2D images. Note that the first row displays the original depth images in gray-scale of the 3D shape, while the second row shows the reconstructed ones corresponding to the images of the first row. And the black dots indicates those extracted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

Zernike moments coefficients [23]. Finally, the dissimilarity between two objects is measured by the minimum distance of all group matching pairs. The LFD is insensitive to similarity transform, geometry degeneracy and noise, etc, thus shows better performance than other competing approaches.

3. Deep learning representation using autoencoder

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

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

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

where Np denotes the number of projections for each model.

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

3.2. Deep belief network

The deep belief network (DBN) [24–26] is a generative graphical model, or alternatively a type of deep neural network, composed of multiple layers of latent variables (“hidden units”), with connections between the layers but not between units within each layer. When trained on plenty of examples in an unsupervised way, a DBN can probabilistically reconstruct the inputs by learning a stack of Restricted Boltzmann Machines (RBMs), where each of the previous RBM’s hidden layer serves as the visible layer for the next. That is to say, each time a new RBM is added to the stacked structure of DBN, then the new DBN has a better variational lower bound in the log probability of the data than the previous DBN [4].

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 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 [27] as

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

where v_i, h_j denote the binary states of visible unit i and hidden unit j respectively; 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

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

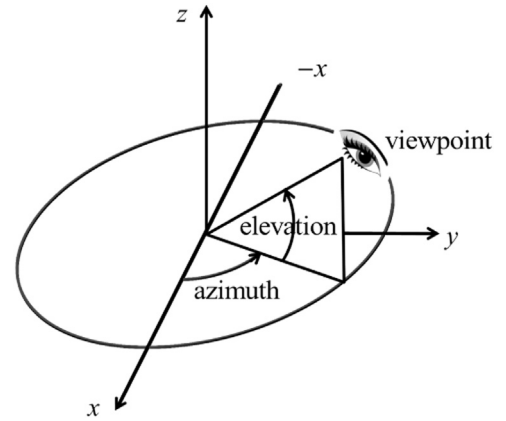


Fig. 3. The illustration of how we get the projections of a 3D shape model S . Azimuth is the polar angle in the x - y plane, with positive number indicating anticlockwise rotation of the viewpoint. As for elevation, positive and negative numbers are the angle above and below the x - y plane respectively.

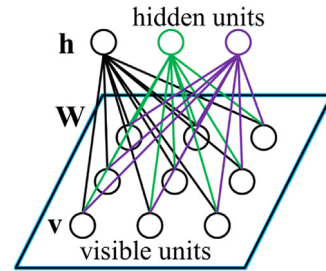


Fig. 4. A graphical description of RBM. Note that a standard type of RBM has binary-valued visible and hidden units with weights of the connection between them. What needs to be specially emphasized is that there are none connections within visible units or hidden ones, which leads to a property that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

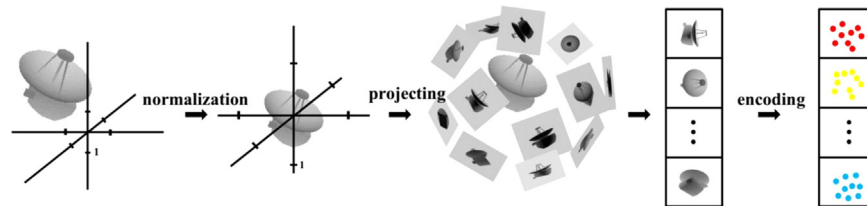


Fig. 2. The flow chart of 3D shape representation using autoencoder. First, we conduct pose normalization for differences in translation and scale to each 3D model. Next, each 3D shape is represented by a set of depth-buffer images. Finally all the projections are used to train the autoencoder to acquire 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.

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:

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

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

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

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

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

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 .

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

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

Fortunately, a much faster learning algorithm was proposed in [28]. 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

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

where the $\langle v_i h_j \rangle_{data}$ is the fraction of times that the visible unit i and the hidden unit j are on together when the hidden units are driven by data, and $\langle v_i h_j \rangle_{recon}$ is the corresponding part given by

the confabulation. A same learning rules is used to adjust the biases.

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

3.3. Fine-tuning the autoencoder

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

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

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 \mathbf{C}

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

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

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

where Nc denotes the dimensionality of every code vector; c_{ij} is the value of 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

$$\mathbf{C}(S_A) = \{ \vec{C}_{A1}, \vec{C}_{A2}, \dots, \vec{C}_{ANp} \} \\ \mathbf{C}(S_B) = \{ \vec{C}_{B1}, \vec{C}_{B2}, \dots, \vec{C}_{BNp} \}, \quad (13)$$

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

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

where $d(\vec{C}_{Ai}, \vec{C}_{Bj})$ denotes one specific distance function between two vector, such as p-norm distance in "Euclidean Space", algebraic distance, etc. Depending on the distance of any two models, shape retrieval could be directly done according to the ranked list.

4. Bag of features representation

In this Section, we describe the local descriptor formerly implemented by Ohbuchi et al. [11] on 3D shape. Considering that our method autoencoder mentioned above is a global descriptor, it

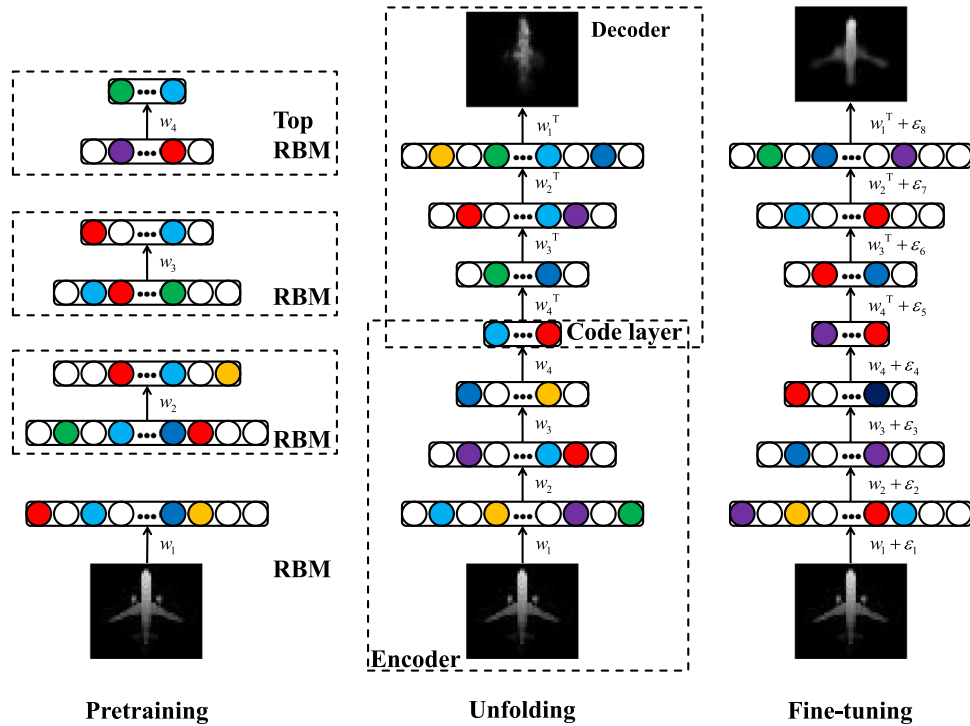


Fig. 5. Details of autoencoder implemented on depth images. The circles enclosed by rectangle in each layer denote the units with various filling colour indicating different probability that the network assigns to them, and the rectangle's length corresponds to the relative size of dimension on that 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.

is much reasonable to boost a better performance if combining with a local descriptor. Bag-of-Features using Scale-invariant feature transform (BoF-SIFT) model is selected as the local description for a 3D model. Different from previous work in [11] that considers the SIFTs of each depth image separately, we put all SIFTs in a single bag, i.e., rotation normalization is not conducted.

We first learn the visual word vocabulary with size of 1500 in a randomly selected subset of all features via K-means off-line. In order to encode the set of SIFTs in each 3D model, we conduct Vector Quantization proposed in [29] to get a histogram representation that counts the number of SIFTs belonging to each visual word. Before computing the pairwise distance among the models, all the histogram is L_1 normalized. We will display the good property of extraordinary complementarity between autoencoder and BoF-SIFT in Section 5.

5. Experiments

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

5.1. Princeton shape benchmark (PSB)

The *Princeton Shape Benchmark* [12] dataset provides a repository of 3D models and software tools for comparing different shape-based models. It's created to promote the use of standardized datasets and evaluate methods for research in matching, classification, clustering, and recognition of 3D models. Each

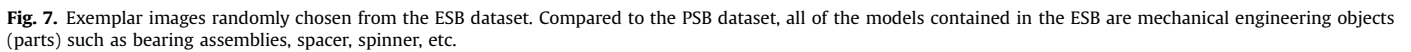
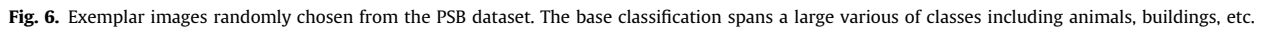
model of the 3D shape consists of the polygonal geometry surface of the corresponding shape. There are totally 1814 models and the base classification is partitioned equally into training and testing sets. The training set with 90 classes, 907 models is used to attain parameters of shape models through training procedure, while the other with 92 classes, equal number of models for comparison with other algorithm. In addition, the number of models belonging to the same class in the base classification varies from each class and ranges from 4 to 50. Some 3D models from the PSB are randomly selected to be exhibited in Fig. 6.

5.2. Engineering shape benchmark (ESB)

The *Engineering Shape Benchmark* [30] is particularly proposed to evaluate shape-based searching methods relevant to the mechanical engineering domain. More specifically, the ESB dataset has totally 867 3D CAD models classified into 45 classes with the number of models ranging from 4 to 58 in a class. The 3D models contained in the ESB cover a wide variety of real-world engineering models so that different methods can compete with each other more fairly. As shown in Fig. 7, we randomly select some models in the ESB to show engineering properties of the models.

5.3. National taiwan university benchmark (NTU)

The *National Taiwan University Benchmark* [9] provides 3D models for the purpose of 3D shape retrieval, matching, recognition and classification. Based on functional similarities, 549 3D models mainly for vehicles and household items are classified into 47 categories. As shown in Fig. 8, we give illustrations of the randomly selected models from NTU.



As described in [Section 3.1](#), we set the number of each model's projection to 64 (8×8) on the dataset. Then the total raw gray-scale images with real value in the range of $[0, 1]$, preprocessed by transform invariant low-rank textures (TILT) [\[31\]](#) to eliminate the large orientation variance, served as the visible units of the DBN's first layer.

units. The real-valued states are in the range $[0, 1]$, compared to the binary states either 0 or 1, allowed the low-dimensional codes to take good advantage of continuous data and could avoid unnecessary sampling noise. Note that we trained each RBM for 40 epochs using mini-batches of size 100 and adopted a learning rate of 0.1 for the linear-binary RBMs, 0.001 for the top layer RBM.

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

Then all the parameters including weights and biases are well-trained in an unsupervised way, we used them to obtain the low-dimensional code for projections of 3D models on the dataset. For the PSB and NTU datasets, we constructed an encoder with the layer structure of 5184 (72×72)-1000-500-250-30 while a structure of 5184 (72×72)-2000-500-100-20 for the ESB. In addition, we only used the testing set to both train the parameters and evaluate our results for the PSB while experiments were done on the whole dataset of the ESB and NTU since they are not divided into training and testing sets.

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

$$d(\vec{C}_{A_i}, \vec{C}_{B_j}) = \|\vec{C}_{A_i} - \vec{C}_{B_j}\|_p, \quad p=2, \quad (15)$$

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 n -dimensional real vector space \mathbb{R}^n .

5.5. Evaluation methods

In this Section, we introduce statistical description for the retrieval performance of a specific algorithm. The PSB provides open source code for evaluating different algorithms and judging how well one algorithm is compared to others. Thus, the performance can be fairly judged by the same evaluation tools in varieties of perspectives. When any doubt comes to you, please refer to [12] for more details about definition of every evaluation method.

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

First-Tier (FT) and Second-Tier (ST): the percentage of models in the query's class that appear within the top M matches, where M is determined by the size of the query's class. Given that

the query's class owns C models, $M = C - 1$ for the first-tier and $M = 2(C - 1)$ for the second tier.

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

5.6. Retrieval results

We adopt autoencoder described in Section 3 to obtain the distance between any two shape models from datasets. Then we get the retrieval results evaluated by the source code provided in [12]. As shown in Table 1, 2 and 3, we compare the

Table 1
Statistic evaluation of global descriptors on PSB.

Algorithm	NN(%)	FT(%)	ST(%)
Autoencoder	72.4	43.3	54.6
GSMD [20]	67.1	41.8	52.0
DESIRE [17]	65.8	40.4	51.3
LFD [9]	65.7	38.0	48.7
SH-GEDT [32]	55.3	31.0	41.4

Table 2
Statistic evaluation of global descriptors on ESB.

Algorithm	NN(%)	FT(%)	ST(%)
Autoencoder	85.7	47.9	63.1
DESIRE [17]	82.3	41.7	55.0
LFD [9]	82.0	40.4	53.9
SH-GEDT [32]	80.3	40.1	53.6



Fig. 8. Exemplar images randomly chosen from the NTU dataset. There are 3D models mainly for vehicles and household items such as trucks, motorcars, bottles, beds, etc.

global-feature-based autoencoder with the other global descriptors on the three standard datasets to explore the efficacy of using autoencoder to tackle 3D shape retrieval. Compared with other competing approaches including GSMD [20],

Table 3
Statistic evaluation of global descriptors on NTU.

Algorithm	NN(%)	FT(%)	ST(%)
Autoencoder	74.8	43.6	55.4
DESIRE [17]	71.9	42.7	55.4
LFD [9]	70.0	39.0	50.1
SH-GEDT [32]	58.8	33.9	46.3

Table 4
Statistic evaluation of complementarity on PSB.

Algorithm	NN(%)	FT(%)	ST(%)
Autoencoder + BoF-SIFT	77.5	52.4	65.4
BoF-SIFT [11]	71.4	45.1	57.6
CM-BoF + GSMD [10]	75.4	50.9	64.0
PANORAMA [19]	75.3	47.9	60.3
CM-BoF [10]	73.1	47.0	59.8
Hybrid [18]	74.2	47.3	60.6

Table 5
Statistic evaluation of complementarity on ESB.

Algorithm	NN(%)	FT(%)	ST(%)
Autoencoder + BoF-SIFT	88.1	55.2	70.2
BoF-SIFT	88.0	52.4	65.4
PANORAMA [19]	86.5	49.4	64.1
Hybrid [18]	82.9	46.5	60.5

Table 6
Statistic evaluation of complementarity on NTU.

Algorithm	NN(%)	FT(%)	ST(%)
Autoencoder + BoF-SIFT	78.9	49.0	60.9
BoF-SIFT [11]	73.1	43.6	56.3
PANORAMA [19]	78.6	49.0	61.5
Hybrid [18]	76.2	46.6	59.1

DESIRE [17], LFD [9] and SH-GEDT [32] on PSB, ESB and NTU datasets, Autoencoder consistently displays better performance in the evaluations of NN, FT and ST. Especially on the ESB, the proposed Autoencoder brings about at least 6% performance gain in FT and ST, which reaches a great margin. Therefore, the autoencoder is more efficient than other global-features-based methods for 3D shape retrieval.

5.7. Complementary property

Furthermore, based on the knowledge that autoencoder reconstructs global information while BoF-SIFT described in Section 4 mainly captures the local details, a linear combination of them is proposed to boost the retrieval performance. More specifically, we empirically choose the weights as $W_{global} = W_{local}$ for global and local descriptors.

We compare our hybrid method (Autoencoder + BoF-SIFT) with the previous state-of-the-art methods including PANORAMA [19], CM-BoF and CM-BoF + GSMD [10], and Hybrid [18], which are able to capture both the global and local information of a 3D shape. For the retrieval results displayed in Table 4, 5 and 6, we can find that: our autoencoder shows pretty well complementary property with the existing local-features-based method BoF-SIFT since all retrieval results are more or less improved. It's worth noting that our hybrid method reaches at least 6% evaluation increment of NN, FT and ST on the PSB, and at least 4% increment on the NTU compared with BoF-SIFT [11] alone. Moreover, our hybrid method gets the state-of-the-art results on the three datasets except the evaluation of ST on NTU, which is slightly worse than PANORAMA [19].

Furthermore, Fig. 9 shows a precision-recall plot of six methods on the PSB and ESB dataset. The graphic illustrations confirm the conclusions drawn based on former statistic results. The proposed autoencoder consistently outperforms other global-descriptors-based methods. Among all methods, the composite of autoencoder and BoF-SIFT achieves relatively the best performance.

6. Conclusions

In this paper, we present a novel view-based 3D shape retrieval method using autoencoder, which is firstly utilized to 3D shape retrieval. A set of experiments were carried out to investigate the effectiveness and efficiency of our method on two standard datasets, which shows that the autoencoder outperforms other global descriptors on retrieval results. Furthermore, the experiments demonstrate that the autoencoder displays good complementarity with the local descriptor, since linear combination achieves the

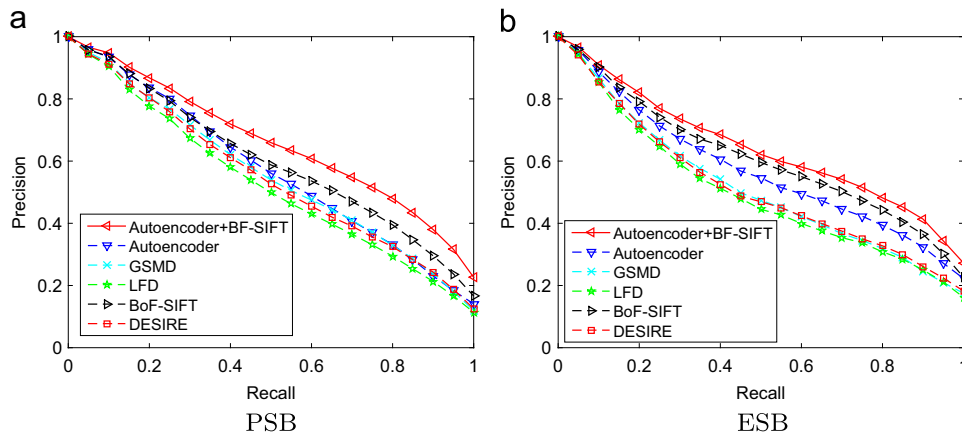


Fig. 9. Precision-recall curves of nine methods implemented on two standard benchmarks. (a), (b) illustrate the results evaluated on the PSB and ESB respectively.

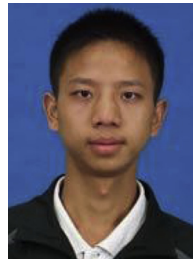
state-of-the-art performance. Our future work might focus on studying the effect of the proposed representation with context-based shape similarity method [33].

Acknowledgement

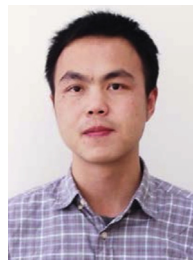
This work was primarily supported by National Natural Science Foundation of China (NSFC) (No.61222308), and Program for New Century Excellent Talents in University (No.NCET-12-0217), Fundamental Research Funds for the Central Universities (No.HUST 2013TS115). Xinggang Wang was supported by Microsoft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

References

- [1] Yann LeCun, Bernhard Boser, John S. Denker, Donnie Henderson, Richard E. Howard, Wayne Hubbard, Lawrence D. Jackel, Backpropagation applied to handwritten zip code recognition, *Neural Computation* 1 (4) (1989) 541–551.
- [2] Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton, Imagenet classification with deep convolutional neural networks, In NIPS, volume 1, 2012, pp. 4.
- [3] Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, arXiv preprint [arXiv:1311.2524](https://arxiv.org/abs/1311.2524), 2013.
- [4] Alex Krizhevsky, Geoffrey E. Hinton, Using very deep autoencoders for content-based image retrieval, In ESANN. Citeseer, 2011.
- [5] Jie Zhang, Shiguang Shan, Meina Kan, Xilin Chen, Coarse-to-fine auto-encoder networks (cfan) for real-time face alignment, In European Conference on Computer Vision, 2014.
- [6] Zhuotun Zhu, Xinggang Wang, Song Bai, Cong Yao, Xiang Bai, Deep learning representation using autoencoder for 3d shape retrieval, In Security, Pattern Analysis, and Cybernetics (SPAC), 2014 International Conference on, IEEE, 2014, pp. 279–284.
- [7] Yi Fang, Jin Xie, Guoxian Dai, Meng Wang, Fan Zhu, Tiantian Xu, Edward Wong, 3d deep shape descriptor, In IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 2319–2328.
- [8] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, Jianxiong Xiao, 3d shapenets: A deep representation for volumetric shapes, In IEEE Conference on Computer Vision and Pattern Recognition, 2015.
- [9] Dingyun Chen, Xiaopei Tian, Yute Shen, Ming Ouhyoung, On visual similarity based 3D model retrieval, *Computer Graphics Forum* 22 (2003) 223–232.
- [10] Zhouhui Lian, A. Godil, Xianfang Sun, Visual similarity based 3d shape retrieval using bag-of-features, In Shape Modeling International Conference (SMI), 2010, June 2010, pp. 25–36.
- [11] Ryutarou Ohbuchi, Kunio Osada, Takahiko Furuya, Tomohisa Banno, Salient local visual features for shape-based 3D model retrieval, In Shape Modeling International, 2008, pp. 93–102.
- [12] Philip Shilane, Patrick Min, Michael Kazhdan, Thomas Funkhouser, The princeton shape benchmark, In Shape Modeling Applications, 2004. Proceedings, IEEE, 2004, pp. 167–178.
- [13] Christopher M. Cyr, Benjamin B. Kimia, 3D object recognition using shape similarity-based aspect graph, In International Conference on Computer Vision, 2001, pp. 254–261.
- [14] Robert Osada, Thomas A. Funkhouser, Bernard Chazelle, David P. Dobkin, Shape distributions, *ACM Transactions on Graphics* 21 (2002) 807–832.
- [15] David G. Lowe, Distinctive image features from scale-invariant keypoints, *International Journal of Computer Vision* 60 (2) (2004) 91–110.
- [16] Feifei Li, Perona Pietro, A bayesian hierarchical model for learning natural scene categories, In Computer Vision and Pattern Recognition, volume 2, IEEE, 2005, pp. 524–531.
- [17] Dejan V. Vranic, DESIRE: a composite 3D-shape descriptor, In International Conference on Multimedia Computing and Systems/International Conference on Multimedia and Expo, 2005, pp. 962–965.
- [18] Panagiotis Papadakis, Ioannis Pratikakis, Theoharis Theoharis, Georgios Pas-salis, Stavros J. Perantonis, Agia Paraskevi, 3D object retrieval using an efficient and compact hybrid shape descriptor, 2008, pp. 9–16.
- [19] Panagiotis Papadakis, Ioannis Pratikakis, Theoharis Theoharis, Stavros J. Perantonis, PANORAMA: a 3D shape descriptor based on panoramic views for unsupervised 3D object retrieval, *International Journal of Computer Vision* 89 (2010) 177–192.
- [20] Zhouhui Lian, Paul L. Rosin, Xianfang Sun, Rectilinearity of 3D meshes, *International Journal of Computer Vision* 89 (2010) 130–151.
- [21] Xiang Bai, Cong Rao, Xinggang Wang, Shape vocabulary: a robust and efficient shape representation for shape matching, *IEEE Transactions on Image Processing* 29 (2014) 3935–3949.
- [22] Dengsheng Zhang, Guojun Lu, Shape-based image retrieval using generic Fourier descriptor, *Signal Processing-image Communication* 17 (2002) 825–848.
- [23] Alireza Khotanzad, Yaw Hua Hong, Invariant image recognition by Zernike moments, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 12 (1990) 489–497.
- [24] Yoshua Bengio, Pascal Lamblin, Dan Popovici, Hugo Larochelle, Greedy layer-wise training of deep networks, In NIPS, 2006.
- [25] Geoffrey E. Hinton, Simon Osindero, Yee-Whye Teh, A fast learning algorithm for deep belief nets, *Neural Computation* 18 (2006) 1527–1554.
- [26] Y. LeCun, B. Boser, J.S. Denker, D. Henderson, R.E. Howard, W. Hubbard, L. D. Jackel, Backpropagation applied to handwritten zip code recognition, *Neural Computation* 1 (1989) 541–551.
- [27] J.J. Hopfield, Neural networks and physical systems with emergent collective computational abilities, *Proceedings of the National Academy of Sciences* 79 (8) (1982) 2554–2558.
- [28] Geoffrey E. Hinton, Training products of experts by minimizing contrastive divergence, *Neural Computation* 14 (8) (2002) 1771–1800.
- [29] Svetlana Lazebnik, Cordelia Schmid, Jean Ponce, Beyond bags of features: spatial pyramid matching for recognizing natural scene categories, In IEEE Conference on Computer Vision and Pattern Recognition, volume 2, 2006.
- [30] Subramaniam Jayanti, Yagnanarayanan Kalyanaraman, Natraj Iyer, Karthik Ramani, Developing an engineering shape benchmark for CAD models, *Computer-aided Design* 38 (2006) 939–953.
- [31] Zhengdong Zhang, Xiao Liang, Arvind Ganesh, Yi Ma, TILT: transform invariant low-rank textures, In Asian Conference on Computer Vision 2010.
- [32] Michael Kazhdan, Thomas Funkhouser, Szymon Rusinkiewicz, Rotation invariant spherical harmonic representation of 3 d shape descriptors, In Symposium on geometry processing, volume 6, 2003, pp. 156–164.
- [33] Xiang Bai, Xingwei Yang, L.J. Latecki, Wenyu Liu, Zhuowen Tu, Learning context-sensitive shape similarity by graph transduction, *IEEE Transactions on Pattern Anal. Mach. Intell.* 32 (2010) 861–874.



Zhuotun Zhu is currently working toward his bachelor degree in Department of Electronics and Information Engineering in Huazhong University of Science and Technology (HUST), Wuhan, China. He was awarded Young Microsoft Fellow in 2014, National Scholarship twice in 2013 and 2012. His research interests include deep learning and shape retrieval.



Xinggang Wang received the B.E. degree in electronic information engineering from Huazhong University of Technology and Science, Wuhan, China, in 2009. He is currently working toward the Ph.D. degree in Department of Electronic Information Engineering from Huazhong University of Technology and Science. From May 2010 to July 2011, he was with the Department of Computer and Information Science, Temple University, Philadelphia, PA, as a visiting scholar. From February 2013 to September 2013, he was with the University of California, Los Angeles, as a visiting graduate researcher. He is a reviewer of IEEE Transaction on Cybernetics, Pattern Recognition, Computer Vision and Image Understanding, CVPR and ECCV etc. His research interests include computer vision and machine learning.



Song Bai received the B.S. degree in Electronics and Information Engineering from Huazhong University of Science and Technology (HUST), Wuhan, China 2009. He is currently working toward the Ph.D. degree in Department of Electronic Information Engineering from Huazhong University of Technology and Science. His research interests include shape analysis, image classification and retrieval, object detection.



Cong Yao received the B.S. and Ph.D. degrees in electronics and information engineering from the Huazhong University of Science and Technology (HUST), Wuhan, China, in 2008 and 2014, respectively. He was a research intern at Microsoft Research Asia (MSRA), Beijing, China, from 2011 to 2012. He was a Visiting Research Scholar with Temple University, Philadelphia, PA, USA, in 2013. His research has focused on computer vision and machine learning, in particular, the area of text detection and recognition in natural images.



Xiang Bai received the B.S., M.S., and Ph.D. degrees from Huazhong University of Science and Technology (HUST), Wuhan, China, in 2003, 2005, and 2009, respectively, all in Electronics and Information Engineering. He is currently a Professor with the Department of Electronics and Information Engineering, HUST. He is also the Vice-Director of National Center of Anti-Counterfeiting Technology, HUST. His research interests include object recognition, shape analysis, scene text recognition and intelligent systems.