

DeepShape: Deep-Learned Shape Descriptor for 3D Shape Retrieval

Jin Xie, Guoxian Dai, Fan Zhu, Edward K.Wong, and Yi Fang

Abstract—Complex geometric structural variations of 3D models usually pose great challenges in 3D shape matching and retrieval. In this paper, we propose a novel 3D shape feature learning method to extract high-level shape features that are insensitive to geometric deformations of shapes. Our method uses a discriminative deep auto-encoder to learn deformation-invariant shape features. First, a multiscale shape distribution is computed and used as input to the auto-encoder. We then impose the Fisher discrimination criterion on the neurons in the hidden layer to develop a deep discriminative auto-encoder. Finally, the outputs from the hidden layers of the discriminative auto-encoders from different scales are concatenated to form the shape descriptor. The proposed method is evaluated on four benchmark datasets that contain 3D models with large geometric variations: McGill, SHREC'10 ShapeGoogle, SHREC'14 Human and SHREC'14 Large Scale Comprehensive Retrieval Track Benchmark datasets. Experimental results on the benchmark datasets demonstrate the effectiveness of the proposed method for 3D shape retrieval.

Index Terms—3D shape retrieval, heat kernel signature, heat diffusion, auto-encoder, Fisher discrimination criterion.

I. INTRODUCTION

NOWADAYS there is an explosive growth of 3D meshed surface models in a variety of fields, such as engineering, entertainment and medical imaging [1–6]. Due to the data-richness of 3D models, shape retrieval for 3D model searching, understanding and analyzing has received more and more attention. Using a shape as a query, the shape retrieval algorithm aims to find similar shapes. The performance of a shape retrieval algorithm mainly relies on a shape descriptor that can effectively capture the distinctive properties of shape. It is preferable that a shape descriptor be deformation-insensitive and invariant to different classes of transformations. The shape descriptor should also be insensitive to both topological and numerical noise. Once the shape descriptor is formed, the similarity between two shapes can be determined by comparing their descriptors.

Shape descriptors for shape matching and retrieval have been extensively studied in the geometry community [7–11]. In the past decades, plenty of shape descriptors have been proposed; these include the $D2$ shape distribution [9], statistical moments of model [10, 12], Fourier descriptor [13] and Eigenvalue Descriptor (EVD)[14]. Although these shape descriptors can represent shapes effectively, they are either sensitive to non-rigid transformations or topological changes. To be invariant to isometric transformations, local geometric features, such as spin images [15], shape context [16] and

mesh HOG [17] are extracted to represent the shape. However, these features are sensitive to local geometric noise and they do not capture the global structure of the shape very well.

Apart from the above earlier shape descriptors, another popular approach to shape retrieval uses diffusion-based methods [1, 18, 19]. Based on the Laplace-Beltrami operator, the global point signature (GPS) [1] was proposed to represent shapes. Since the eigenfunctions of the Laplace-Beltrami operator are able to robustly characterize the points on a meshed surface, each vertex is represented by a high-dimensional vector (called GPS) of the scaled eigenfunctions of the Laplace-Beltrami operator evaluated at the vertex. Another widely used shape signature is heat kernel signature (HKS) [18], where the diagonal of the heat kernel is used as a local descriptor to represent shapes. HKS is invariant to isometric deformations and insensitive to small perturbations on the surface. Both GPS and HKS are point-based signatures that characterize vertices on the meshed surface by vectors.

In the aforementioned methods, the shape descriptors are hand-crafted rather than learned from a set of training shapes. In [20], the authors applied the bag-of-features (BoF) paradigm to learn the shape descriptor. The dictionary of words is learned by applying K -means clustering to a set of HKSs. A histogram of pairs of spatially-close words over the learned dictionary is then formed as the shape descriptor for retrieval. Using K -means clustering, Lavoué *et al.* [21] combined the standard and spatial BoF descriptors for 3D shape retrieval. EINagh *et al.* proposed the compact HKS-based BoF descriptor, i.e., CompactBoFHKS [22]. In the CompactBoFHKS method, feature point detection is employed to select critical points. For each critical point, certain scales of the HKS are selected to form a compact feature vector to describe it. The BoF method is then applied to the feature vectors to learn a shape descriptor for retrieval. Litman *et al.* [23] employed sparse coding to learn the dictionary of words instead of K -means clustering. The histogram of encoded representation coefficients over the learned dictionary is used to represent shapes for retrieval. Moreover, in order to obtain discriminative representation coefficients, a class-specific dictionary is constructed using supervised learning.

In recent years, due to the success of deep neural networks in different application domains, deep learning based 3D shape features have been proposed for 3D shape analysis. Wu *et al.* [24] proposed to represent 3D shapes as a probability distribution of binary variables on a 3D voxel grid. Then a convolutional deep belief network is developed to learn the joint probabilistic distribution of the voxel data and the category label. Boscaini *et al.* [25] employed windowed Fourier transform to points on the meshed surface to form a local frequency representation. These local frequency rep-

Jin Xie, Guoxian Dai, Fan Zhu and Yi Fang are with the Department of Electrical and Computer Engineering, New York University Abu Dhabi, UAE. Edward K.Wong is with Tandon School of Engineering, New York University, New York, USA (e-mail: {jin.xie, guoxian.dai, fan.zhu, ewong, yfang}@nyu.edu).

representations are passed through a bank of filters to form a deep representation for 3D shapes. The filter coefficients can be learned by using a task-specific cost. Masci *et al.* [26] generalized the convolutional neural network to non-Euclidean manifolds for 3D shape retrieval and correspondence. Also, Dosovitskiy *et al.* [27] trained a generative convolutional neural network to simulate the viewpoints of the given 3D shape. The trained convolutional neural network can also be used to find correspondences between the 3D shape models.

In this paper, we develop a novel deep neural network based method for learning shape descriptors for retrieval applications. Our method uses the Fisher discrimination criterion on the hidden layer to make the shape features discriminative and insensitive to geometric structural variations. The neurons in the hidden layer have small within-class scatter but large between-class scatter. To more effectively represent shape, we use multiscale shape distribution as inputs to the auto-encoders. We train a discriminative auto-encoder at each scale and concatenate the outputs from the hidden layers of different scales as the shape descriptor. We tested the proposed shape descriptor on several benchmark shape datasets and promising results have been obtained.

The rest of the paper is organized as follows. In Section II, we briefly introduce the HKS and auto-encoder. In Section III, we present the proposed shape descriptor with the discriminative auto-encoder. We describe our experimental results in Section IV and we conclude the paper in Section V.

II. BACKGROUND

A. Heat Kernel Signature

A 3D model is represented as a graph $G = (V, E, W)$, where V is the set of vertices, E is the set of edges and W is the set of weights for the edges. Given a graph constructed by connecting pairs of data points with weighted edges, the heat kernel $H_t(g_0, g_1)$ measures the heat flow across the graph, defined as the amount of heat passing from vertex g_0 to vertex g_1 within a certain amount of time. The heat flow across the surface is governed by a heat equation. Provided that there is an initial Dirac delta heat distribution on the meshed surface at $t = 0$, the heat kernel provides the fundamental solution of the heat equation, which is associated with the Laplace-Beltrami operator Ψ by:

$$\frac{\partial H_t}{\partial t} = -\Psi H_t \quad (1)$$

where H_t denotes the heat kernel and t is the diffusion time. The solution of Eq. (1) can be obtained by the eigenfunction expansion of the Laplace-Beltrami operator described as:

$$H_t = \exp(-t\Psi). \quad (2)$$

By the spectral theorem, the heat kernel can be expressed as:

$$H_t(g_0, g_1) = \sum_i e^{-\lambda_i t} \phi_i(g_0) \phi_i(g_1) \quad (3)$$

where λ_i is the i^{th} eigenvalue of the Laplacian-Beltrami operator and ϕ_i is the i^{th} eigenfunction. The heat kernel signature (HKS) [18] of vertex g_0 at time t can be defined

as the diagonal of the heat kernel of vertex g_0 taken at time t :

$$H_t(g_0, g_0) = \sum_i e^{-\lambda_i t} \phi_i(g_0)^2. \quad (4)$$

The defined HKS $H_t(g_0, g_0)$, is a point signature that captures the neighborhood information at point g_0 and scale t on the surface.

B. Auto-encoder

An auto-encoder neural network [28, 29] consists of two parts, i.e., encoder and decoder. The encoder, denoted by F , maps the input $\mathbf{h} \in \mathcal{R}^{d \times 1}$ to the hidden layer, denoted by $\mathbf{z} \in \mathcal{R}^{r \times 1}$, where d is the dimension of the input and r is the number of neurons in the hidden layer. In the auto-encoder neural network, a neuron in layer l is connected to all neurons in the next layer $l + 1$. We denote the weight and bias connecting layers l and $l + 1$ by \mathbf{W}^l and \mathbf{b}^l , respectively. A non-linear activation function, such as the sigmoid function $\sigma(\mathbf{h}) = \frac{1}{1+e^{-\mathbf{h}}}$ or tanh function $\sigma(\mathbf{h}) = \frac{e^{\mathbf{h}} - e^{-\mathbf{h}}}{e^{\mathbf{h}} + e^{-\mathbf{h}}}$ is usually used to produce the output at each neuron. The output at layer $l + 1$ can be represented as

$$f_{l+1}(\mathbf{a}^l) = \sigma(\mathbf{W}^l \mathbf{a}^l + \mathbf{b}^l) \quad (5)$$

where $f_{l+1}(\mathbf{a}^l)$ is the activation function for layer $l + 1$ and \mathbf{a}^l is the neurons in layer l . Thus, the encoder $F(\mathbf{h})$ can be represented as

$$F(\mathbf{h}) = f_k(f_{k-1}(\cdots, f_2(\mathbf{h}))). \quad (6)$$

The decoder, denoted by G , maps the hidden layer representation \mathbf{z} back to the input \mathbf{h} . It is defined as

$$\mathbf{h} = f_L(f_{L-1}(\cdots, f_{k+1}(\mathbf{z}))) \quad (7)$$

where L is the layer number of the auto-encoder neural network. The matrices \mathbf{W} and \mathbf{b} represent the weights and biases of the neural network with $\mathbf{W} = [\mathbf{W}^1, \mathbf{W}^2, \cdots, \mathbf{W}^{L-1}]$ and $\mathbf{b} = [\mathbf{b}^1, \mathbf{b}^2, \cdots, \mathbf{b}^{L-1}]$. To optimize parameters \mathbf{W} and \mathbf{b} , the standard auto-encoder minimizes the following cost function:

$$\begin{aligned} \langle \hat{\mathbf{W}}, \hat{\mathbf{b}} \rangle = & \argmin_{\mathbf{W}, \mathbf{b}} \frac{1}{2} \sum_{i=1}^M \|\mathbf{h}_i - G(F(\mathbf{h}_i))\|_2^2 \\ & + \frac{1}{2} \lambda \|\mathbf{W}\|_F^2 \end{aligned} \quad (8)$$

where \mathbf{h}_i represents the i^{th} training sample, M represents the total number of training samples, and parameter λ is a positive scalar. In Eq. (8), the first term is the reconstruction error and the second term is the regularization term that prevents overfitting. An efficient optimization method can be implemented by the restricted Boltzman machine and back-propagation framework. The reader can refer to [28] for more details.

III. DISCRIMINATIVE AUTO-ENCODER BASED SHAPE DESCRIPTOR

In this section, we describe our proposed discriminative auto-encoder based shape descriptor. As depicted in Fig. 1, our proposed framework comprises three components, namely, multiscale shape distribution, discriminative auto-encoder and 3D shape descriptor. In the multiscale shape distribution component, the distributions of heat kernel signatures of shape at different scales are extracted as low-level features and used as input to the discriminative auto-encoder. In the second component, we train a discriminative auto-encoder to learn high-level shape features (embedded in the hidden layer of the discriminative auto-encoder.) Finally, the 3D shape descriptor is formed by concatenating the outputs from the hidden layers of the discriminative auto-encoders from different scales.

A. Multiscale Shape Distribution

The HKS describes the amount of heat that remains at the vertex within a time interval. And it is highly related to the curvature of the meshed surface. Thus, as a point signature, the HKS can characterize the intrinsic geometry structure of the neighborhood of the shape well. It has attractive geometric properties that include invariance to isometric transformation, robustness against other geometric changes and local numerical noise, and multiscale representation with diffusion time [18]. Therefore, as a shape function, the HKS can capture information about the intrinsic geometry of the shape and is robust to geometric changes. Compared to the voxelization method [24], the HKS does not need shape alignment. And unlike the parameterized local surface patch method in [25], it does not need to calculate the complex patch operator by constructing the local geodesic polar coordinates on the surface.

Shape distribution [2] refers to a probability distribution sampled from a shape function describing the 3D model. In this work, we use histogram to estimate the probability distribution of the HKS to form the shape distribution. Suppose there are C shape classes, each of which has O samples. We use $y_{i,j}$ to index the j^{th} sample of the i^{th} shape class. For each shape $y_{i,j}$, we extract HKS feature $S_{i,j} \in \mathcal{R}^{N \times T}$, where $S_{i,j} = [S_{i,j}^1, S_{i,j}^2, \dots, S_{i,j}^T]$, $S_{i,j}^t$ denotes the HKS vector of N vertices of shape $y_{i,j}$ at the t^{th} scale, $t = 1, 2, \dots, T$, and T is the number of scales. For scale t , we calculate the histogram of $S_{i,j}^t$ to form the shape distribution $h_{i,j}^t$. By considering probability distributions of shape functions derived from HKS at different scales, a multiscale shape distribution can be obtained. In addition, we normalize the shape distribution, which is centralized by the mean and variance of the shape distributions over all training samples from C classes, namely,

$$h_{i,j}^t = \frac{h_{i,j}^t - \bar{h}^t}{v^t} \quad (9)$$

where \bar{h}^t and v^t are the mean and variance of all training shape distributions $h_{i,j}^t$.

Fig. 2 shows the multiscale shape distributions of the Centaur and Human models with different poses. From this

figure, we can see that the multiscale shape distributions are different for the Centaur and Human shapes. Also, the three centaur models with isometric geometric transformations have consistent multiscale shape distributions. This demonstrates the invariance of the multiscale shape distribution to isometric transformations. For the three human models with structural variations, their multiscale shape distributions capture their common geometric characteristics despite the inconsistency in their detailed descriptions.

B. Discriminative Auto-encoder

In this subsection, we propose a discriminative auto-encoder to extract high-level features for 3D shape retrieval. To boost the discriminative power of the hidden layer features, we impose a Fisher discrimination criterion [30] on them. Given the shape distribution input x_i^t of shape class i at scale t , $x_i^t = [h_{i,1}^t, h_{i,2}^t, \dots, h_{i,O}^t]$, we denote by z^t the features in the hidden layer of the auto-encoder. We can write z^t as $z^t = [z_1^t, z_2^t, \dots, z_C^t]$, where $z_i^t = [z_{i,1}^t, z_{i,2}^t, \dots, z_{i,O}^t]$, $z_{i,j}^t$ is the hidden layer feature of the j^{th} sample from class i , for $i = 1, 2, \dots, C$ and $j = 1, 2, \dots, O$. Using the Fisher discrimination criterion, discrimination is achieved by minimizing the within-class scatter of z^t , denoted by $S_w(z^t)$, and maximizing the between-class scatter of z^t , denoted by $S_b(z^t)$. $S_w(z^t)$ and $S_b(z^t)$ are defined as:

$$S_w(z^t) = \sum_{i=1}^C \sum_{z_{i,j}^t \in i} (z_{i,j}^t - m_i^t)(z_{i,j}^t - m_i^t)^T \quad (10)$$

$$S_b(z^t) = \sum_{i=1}^C n_i (m_i^t - m^t)(m_i^t - m^t)^T$$

where m_i^t and m^t are the mean vectors of z_i^t and z^t , respectively, and n_i is the number of samples from class i . We define the discriminative regularization term $tr(S_w(z^t)) - tr(S_b(z^t))$ and incorporate it into the objective function of the discriminative auto-encoder:

$$J(W^t, b^t) = \underset{W^t, b^t}{\operatorname{argmin}} \sum_{i=1}^C \frac{1}{2} \|x_i^t - G(F(x_i^t))\|_F^2 \quad (11)$$

$$+ \frac{1}{2} \lambda \|W^t\|_F^2 + \frac{1}{2} \gamma (tr(S_w(z^t)) - tr(S_b(z^t))).$$

For shape distribution $h_{i,j}^t$, we define the following functions:

$$J_0(W^t, b^t, h_{i,j}^t) = \frac{1}{2} \|h_{i,j}^t - G(F(h_{i,j}^t))\|_2^2 \quad (12)$$

$$L_0(z_{i,j}^t) = \frac{1}{2} tr((z_{i,j}^t - m_i^t)(z_{i,j}^t - m_i^t)^T) \quad (13)$$

$$- \frac{1}{2} tr((m_i^t - m^t)(m_i^t - m^t)^T).$$

To optimize the objective function, we adopt the back-propagation method. We denote by $W_{p,q}^{l,t}$ the weight associated with the connection between unit p in layer l and unit q in layer $l+1$. Also, $b_p^{l,t}$ is the bias associated with the connection

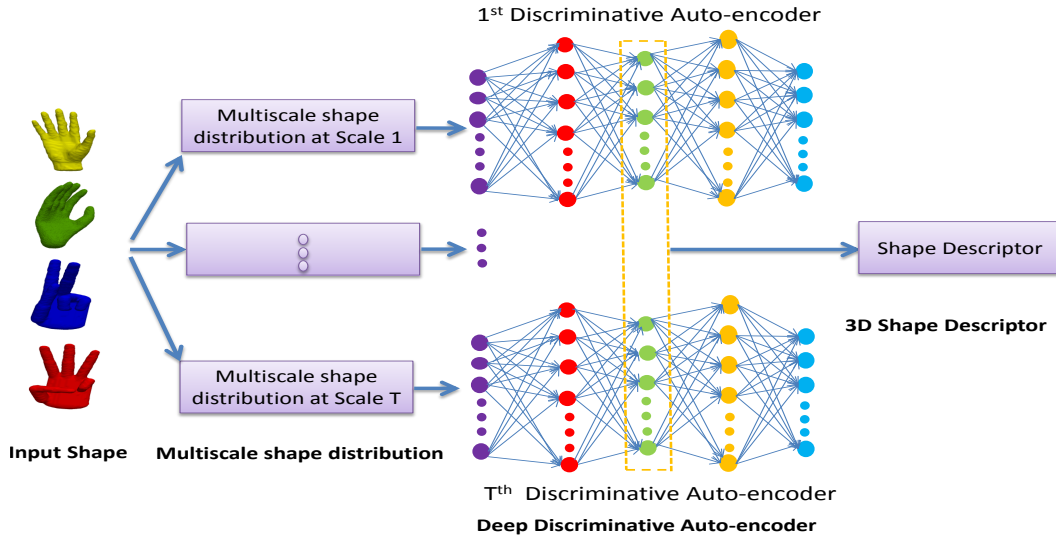


Fig. 1. The framework of the proposed discriminative auto-encoder based shape descriptor.

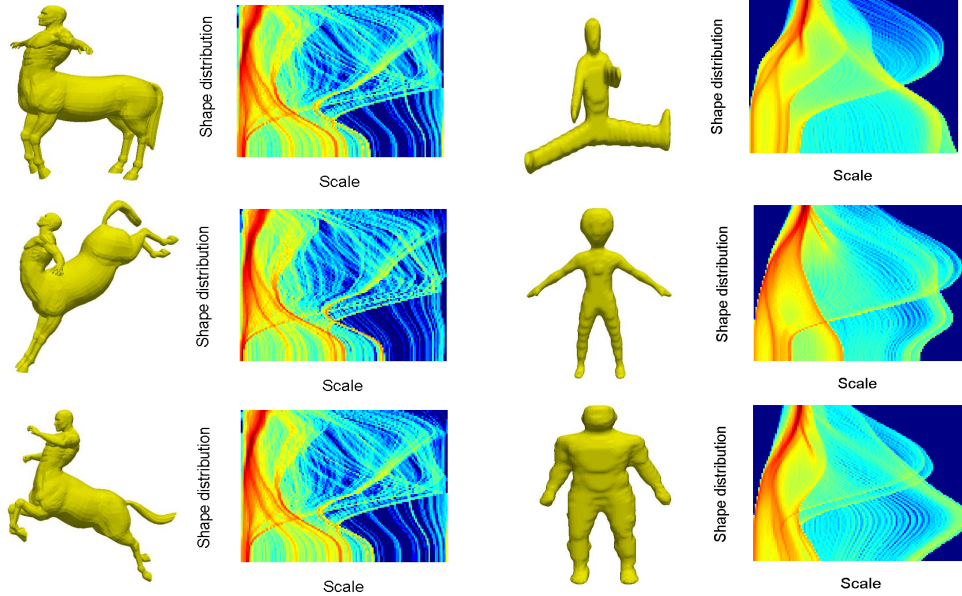


Fig. 2. The multiscale shape distributions of the Centaur and Human models. The left two columns show the Centaur models with isometric transformations and the corresponding multiscale shape distributions, respectively. The right two columns show the Human models with non-isometric structural variations and the corresponding multiscale shape distributions, respectively.

with unit p in layer l . The partial derivatives of the overall cost function $J(\mathbf{W}^t, \mathbf{b}^t)$ can be computed as:

$$\frac{\partial J(\mathbf{W}^t, \mathbf{b}^t)}{\partial \mathbf{W}^{l,t}} = \sum_{i=1}^C \sum_{\mathbf{h}_{i,j}^t \in i} \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} + \lambda \mathbf{W}^{l,t} + \gamma \sum_{i=1}^C \sum_{\mathbf{z}_{i,j}^t \in i} \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} \quad (14)$$

$$\frac{\partial J(\mathbf{W}^t, \mathbf{b}^t)}{\partial \mathbf{b}^{l,t}} = \sum_{i=1}^C \sum_{\mathbf{h}_{i,j}^t \in i} \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} + \gamma \sum_{i=1}^C \sum_{\mathbf{z}_{i,j}^t \in i} \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}}. \quad (15)$$

We denote by $\delta^{L,t}$ the error of the output layer L in the auto-encoder. For the output layer L , we have:

$$\delta^{L,t} = -(\mathbf{h}_{i,j}^t - \mathbf{a}^{L,t}) \bullet \sigma'(\mathbf{u}^{L,t}) \quad (16)$$

where $\mathbf{a}^{L,t}$ is the activation of the output layer, $\mathbf{u}^{L,t}$ is the weighted sum of the outputs from layer $L-1$ to the output layer, $\sigma'(\mathbf{u}^{L,t})$ is the derivative of the activation function in the output layer and \bullet denotes the element-wise multiplication. For layers $l = L-1, L-2, \dots, 2$, the error $\delta^{l,t}$ can be recursively obtained by the back-propagation method in [28] using the following equation:

$$\delta^{l,t} = ((\mathbf{W}^{l,t})^T \delta^{l+1,t}) \bullet \sigma'(\mathbf{u}^{l,t}). \quad (17)$$

The partial derivatives of the function $J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)$, can

be computed as :

$$\begin{aligned} \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} &= \delta^{l+1,t} (\mathbf{a}^{l,t})^T \\ \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} &= \delta^{l+1,t}. \end{aligned} \quad (18)$$

Since $\mathbf{z}_{i,j}^t = \sigma(\mathbf{u}^{k,t}) = \sigma(\mathbf{W}^{k-1,t} \mathbf{a}^{k-1,t} + \mathbf{b}^{k-1,t})$, for $l > k-1$, $\frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} = 0$ and $\frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} = 0$. For $l \leq k-1$, $\frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}}$ and $\frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}}$ can be computed as:

$$\begin{aligned} \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{k-1,t}} &= \frac{\partial \mathbf{z}_{i,j,p}^t}{\partial \mathbf{W}^{k-1,t}} \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{z}_{i,j,p}^t} = a_q^{k-1,t} \sigma'(\mathbf{u}^{k,t})_p \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{z}_{i,j,p}^t} \\ \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{k-1,t}} &= \sigma'(\mathbf{u}^{k,t})_p \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{z}_{i,j,p}^t} \end{aligned} \quad (19)$$

where $\mathbf{z}_{i,j,p}^t$ is the p^{th} component of $\mathbf{z}_{i,j}^t$. The partial derivative of $L_0(\mathbf{z}_{i,j}^t)$ with respect to $\mathbf{z}_{i,j,p}^t$ can be obtained as:

$$\begin{aligned} \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{z}_{i,j,p}^t} &= (1 - \frac{1}{n_i})(\mathbf{z}_{i,j,p}^t - m_{i,p}^t) \\ &\quad - (\frac{1}{n_i} - \frac{1}{\sum n_i})(m_{i,p}^t - m_p^t) \end{aligned} \quad (20)$$

where $m_{i,p}^t$ and m_p^t are the p^{th} components of \mathbf{m}_i^t and \mathbf{m}^t , respectively.

Therefore, based on Eqs. (18), (19) and (20), for $l > k-1$, $\frac{\partial J(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}}$ and $\frac{\partial J(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}}$ can be obtained by Eq. (18). For $l \leq k-1$, $\frac{\partial J(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}}$ and $\frac{\partial J(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}}$ can be computed as:

$$\begin{aligned} \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} &= (\delta^{l+1,t} + \gamma((1 - \frac{1}{n_i}) \\ (\mathbf{z}_{i,j}^t - \mathbf{m}_i^t) - (\frac{1}{n_i} - \frac{1}{\sum n_i})(\mathbf{m}_i^t - \mathbf{m}^t)) \bullet \sigma'(\mathbf{u}^{l+1,t})) (\mathbf{a}^{l,t})^T \\ \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} &= \delta^{l+1,t} + \gamma((1 - \frac{1}{n_i}) \\ (\mathbf{z}_{i,j}^t - \mathbf{m}_i^t) - (\frac{1}{n_i} - \frac{1}{\sum n_i})(\mathbf{m}_i^t - \mathbf{m}^t)) \bullet \sigma'(\mathbf{u}^{l+1,t}). \end{aligned} \quad (21)$$

Once the partial derivatives with respect to \mathbf{W}^t and \mathbf{b}^t are computed, we can employ the conjugate gradient method to obtain \mathbf{W}^t and \mathbf{b}^t . The training algorithm of our proposed discriminative auto-encoder is summarized in Algorithm 1 above.

C. 3D Shape Descriptor

We use the outputs from the hidden layer of the discriminative auto-encoder to form the shape descriptor. In order to characterize the intrinsic structure of the shape more effectively, we train a discriminative auto-encoder at each scale by using a set of training shape distributions, $\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_C^t$, $t = 1, 2, \dots, T$. After training, we concatenate the outputs from the hidden layers of all scales to form the shape descriptor.

Denote the t^{th} encoder of the multiple discriminative auto-encoders by F^t . The shape descriptor of the j^{th} shape from

Algorithm 1 Training algorithm of our discriminative auto-encoder.

Input: training set \mathbf{x}_i^t ; the layer size of the auto-encoder; λ ; γ .
Output: \mathbf{W}^t and \mathbf{b}^t .

Initialize $\Delta \mathbf{W}^{l,t}$ and $\Delta \mathbf{b}^{l,t}$ with the restricted Boltzman machine for all l .

For all $\mathbf{h}_{i,j}^t$:

- 1) Compute $\frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}}$ and $\frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}}$: $l > k-1$, compute them with Eq. (18); $l \leq k-1$, compute them with Eqs. (21) and (22).
- 2) Set $\Delta \mathbf{W}^{l,t}$ to $\Delta \mathbf{W}^{l,t} + \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{W}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{W}^{l,t}}$.
- 3) Set $\Delta \mathbf{b}^{l,t}$ to $\Delta \mathbf{b}^{l,t} + \frac{\partial J_0(\mathbf{W}^t, \mathbf{b}^t, \mathbf{h}_{i,j}^t)}{\partial \mathbf{b}^{l,t}} + \gamma \frac{\partial L_0(\mathbf{z}_{i,j}^t)}{\partial \mathbf{b}^{l,t}}$.

Update $\mathbf{W}^{l,t}$ and $\mathbf{b}^{l,t}$: $\mathbf{W}^{l,t} = \mathbf{W}^{l,t} - \beta(\Delta \mathbf{W}^{l,t} + \lambda \mathbf{W}^{l,t})$, $\mathbf{b}^{l,t} = \mathbf{b}^{l,t} - \beta \Delta \mathbf{b}^{l,t}$.

Output $\mathbf{W}^{l,t}$ and $\mathbf{b}^{l,t}$ until the values of $J(\mathbf{W}^t, \mathbf{b}^t, \mathbf{x}_i^t)$ in successive iterations are close enough or the maximum number of iterations is reached.

class i , i.e., outputs from the hidden layers, can be represented as:

$$\alpha_{i,j} = [F^1(\mathbf{h}_{i,j}^1); F^2(\mathbf{h}_{i,j}^2); \dots; F^T(\mathbf{h}_{i,j}^T)]. \quad (23)$$

IV. EXPERIMENTAL RESULTS

In this section, we evaluate our proposed shape descriptor and compare it to state-of-the-art methods on three benchmark datasets: McGill shape dataset [31], SHREC'10 ShapeGoogle dataset [20], SHREC'14 Human dataset [32] and SHREC'14 Large Scale Comprehensive Retrieval Track Benchmark (SHREC'14 LSCRTB) dataset [33].

A. Experimental Settings

We set time unit $\tau = 0.01$ and take 101 sampled time values (i.e., 101 scales) for the computation of HKS. For multiscale shape distribution, 128 bins are used to form the shape distribution at each scale, which results in a 128-dimensional input for the discriminative auto-encoder. We train an auto-encoder with an encoder of layer size 128-1000-500-250-30 and a symmetric decoder of the same layer size. In Eq. (11), λ and γ are set to 0.001, respectively.

In the McGill 3D shape dataset [31], there are 255 3D meshes with significant part articulations. They come from 10 classes: ant, crab, spectacle, hand, human, octopus, plier, snake, spider, and teddy-bear. Each class contains 3D shapes with large pose changes, which makes the McGill 3D shape dataset very challenging. Fig. 3 shows examples from the McGill shape dataset.

The SHREC'10 ShapeGoogle dataset [20] contains 1184 synthetic shapes, including 715 shapes from 13 classes generated by five simulated transformations: isometry, topology, isometry+topology, partiality and triangulation, and 456 unrelated shapes. Following the setting in [23], to make the dataset more challenging, all shapes are re-meshed to have the same vertices and samples having the same attribute are grouped into the same class. Fig. 4 shows examples from the ShapeGoogle dataset.

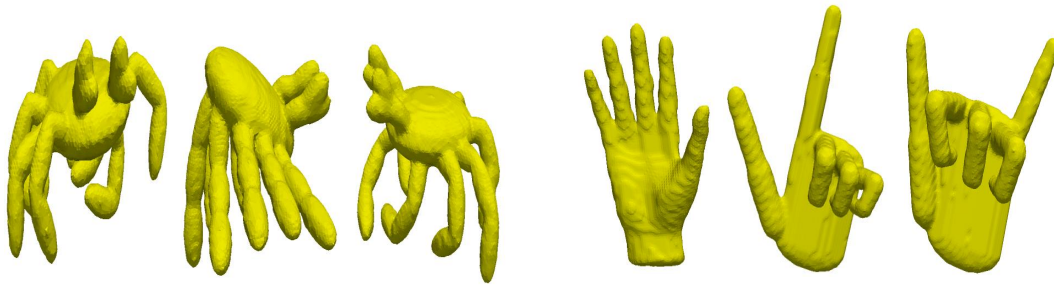


Fig. 3. Example shapes in the McGill dataset. The left three figures show the Crab shapes while the right three figures show the Hand shapes with nonrigid transformations.



Fig. 4. Example shapes with different transformations in the SHREC'10 ShapeGoogle dataset. From left to right, the Centaur shape with isometry, isometry+topology, topology, partiality and triangulation transformations are shown.

The SHREC'14 Human dataset [32] contains two subsets. The first sub-dataset contains 15 synthetic human models and each model has 20 different poses. The second subset consists of 40 scanned human models, with 10 different poses for each model. Following the setting in [23], all human shapes are re-scaled to 4500 triangles. The SHREC'14 Human dataset is an extremely challenging one because all human shapes share similar geometry information. Different poses and similar geometry structures will lead to large within-class variations and small inter-class variations. Fig. 5 shows two human shapes with different poses from the SHREC'14 Human dataset [32].

The SHREC'14 LSCRTB dataset [33] has 8987 3D shape models from 171 classes. It is a large-scale 3D shape dataset where the shapes are from 8 different 3D shape datasets, including the generic models, articulated models, architecture models and CAD models, etc. The average number of models in each class is 53. Most of the shapes are generic models such as bicycle, book, armchair. Moreover, the generic shape models of the same class are not deformed by a template. For example, for the armchair model, there are different kinds of armchairs in this dataset. Fig. 6 shows example shapes from the armchair class.

B. Evaluation of The Proposed 3D Shape Descriptor

We evaluate the effectiveness of our proposed shape descriptor on the McGill benchmark dataset [31]. To compare performance, we also generate test results using the multiscale shape distribution shape descriptor (without auto-encoder) and using auto-encoders but without the Fisher regularization term. We investigate the performance of the proposed shape descriptor in terms of robustness to deformations and noise.

1) Comparison to Multiscale Shape Distribution Descriptor: In our proposed method, we use the multiscale shape distribution as input to the discriminative auto-encoder. Learning deep features from the multiscale shape distribution with the discriminative auto-encoders can be viewed as extracting high-level features from the multiscale shape distribution. To reduce the dimension of the multiscale shape distribution from 101 scales, we concatenate the shape distributions from 26 diffusion time samples to form a 3328-dimensional feature vector. For a fair comparison, we use 26 discriminative auto-encoders to form the shape descriptor. With a hidden layer size of 30 for each auto-encoder, a 780-dimensional shape descriptor is formed to represent the shape. Fig. 7 shows the precision-recall curves for the multiscale shape distribution descriptor and the proposed shape descriptor. As can be seen in this figure, the proposed shape descriptor with a lower dimension has significantly better retrieval performance.

2) Comparison to Standard Auto-encoder: In order to demonstrate effectiveness of the proposed discriminative auto-encoder, we also compare the proposed shape descriptor to the shape descriptor obtained by employing a regular auto-encoder without the Fisher discrimination term. For both shape descriptors, we concatenate 101 auto-encoders to form a descriptor of the same dimension. Fig. 8 shows the precision-recall curves for the two shape descriptors. One can see that the proposed shape descriptor with the discriminative auto-encoder performs better than the shape descriptor obtained by using standard auto-encoder without the Fisher discrimination term. It implies that by imposing the Fisher discrimination constraint on the hidden layers the learned shape descriptor can reduce within-class variations and increase between-class variations, therefore improving the retrieval performance.

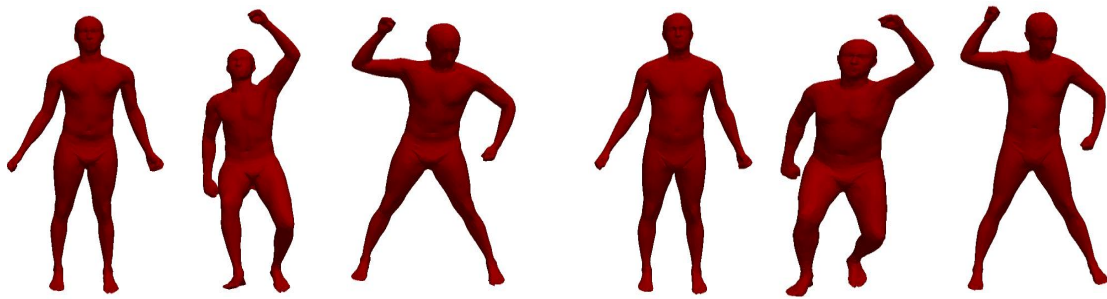


Fig. 5. Two human shapes with different poses in the SHREC'14 Human dataset. The left three figures show shapes with pose changes from one person while the right three figures show shapes with different poses from another person.



Fig. 6. Different kinds of armchair shapes in the SHREC'14 LSCRTB dataset.

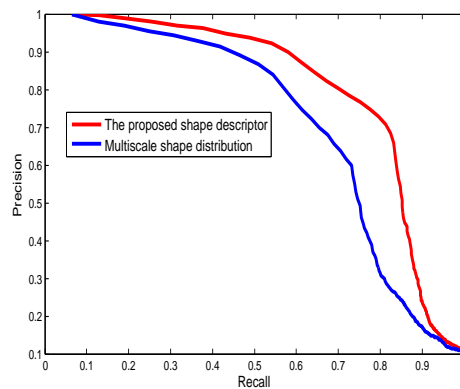


Fig. 7. The precision-recall curves for the multiscale shape distribution descriptor and the proposed shape descriptor on the McGill shape dataset.

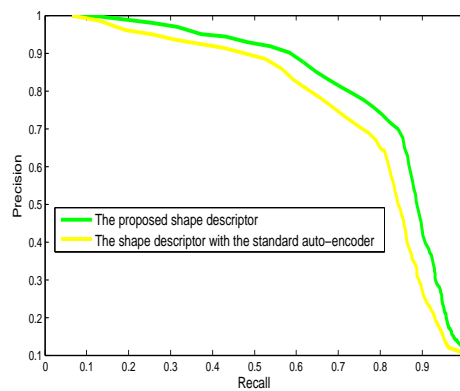


Fig. 8. The precision-recall curves for the shape descriptor using standard auto-encoder (without Fisher discrimination criterion) and the proposed shape descriptor on the McGill dataset.

3) *Robustness to Deformations and Noise:* A good shape descriptor should be robust to pose changes and noise corruptions. We evaluate robustness of the proposed shape descriptor against pose changes and noise. We chose Teddy-bear and

Human models with different poses from the McGill dataset [31] in our experiment. The shape descriptors of the deformed shapes are illustrated in Fig. 9. From the figure, we can see that the descriptors for the model with different poses are very similar. On the other hand, the shape descriptors for different models are distinctive. This shows that the hidden-layer features in the proposed discriminative auto-encoder have small within-class variations but large between-class variations.

By perturbing the vertices of the mesh with various levels of numerical noise, we also demonstrate that the proposed shape descriptor is robust to noise. The noise, represented as a 3-dimensional vector, is randomly generated from a multivariate normal distribution, $N_3(\mu, R \times \Sigma)$, where μ is the 3-dimensional mean vector of the coordinates of the vertices, Σ is the 3×3 covariance matrix of the vertices, and R denotes the ratio between the variance of noise and the variance of the coordinates of the vertices.

Fig. 10 shows the clean Crab and Hand models, and their noisy models. In (a) and (c), the green and red noisy models are generated by noise of $R = 0.01$ and $R = 0.04$, respectively. Particularly, in the noisy model with noise of $R = 0.04$, geometric structures of the mesh have moderately deteriorated. As shown in Fig. 10, the variations in the proposed shape descriptors for the clean and noisy models (plotted with yellow, green and red curves) are small. Since the level of noise for $R = 0.01$ is low, the yellow and green curves basically overlap. The experimental results demonstrate that the proposed shape descriptor is robust to noise.

C. Comparison with the State-of-the-art Methods

We tested our proposed shape descriptor on four benchmark datasets – McGill [31], SHREC’10 ShapeGoogle [20], SHREC’14 Human [32] and SHREC’14 LSCRTB [33]– and compare the results with several state-of-the-art methods. Each shape is represented by a compact 1D shape descriptor and L_2 norm is used to compute the distance between two shape descriptors in our retrieval experiments.

1) *McGill Shape Dataset*: For the McGill 3D shape dataset [31], we compare our method to the Hybrid BOW method [34], the PCA based VLAT method [35], the graph-based method [36], the hybrid 2D/3D approach [21], the covariance descriptor [37] and the CompactBoFHKS method (CBoFHKS) [22]. We denote our proposed discriminative auto-encoder-based shape descriptor by DASD. In the CompactBoFHKS method, 21 scales are chosen every 5 scales from 101 scales for the HKS. The size of the Bag-of-Words is set to 64. In our experiments, 10 shapes per class are randomly chosen to train the discriminative auto-encoder and the remaining shapes in each class are used for testing. We use different performance measures in our evaluation, namely, Nearest Neighbor (NN), the First Tier (1-Tier), the Second Tier (2-Tier) and the Discounted Cumulative Gain (DCG). The retrieval performance of our method and other state-of-the-art methods is illustrated in Table I. From this table, we see that the proposed method achieves the best performance with the NN, 1-Tier, and DCG measures. There are large nonrigid deformations with the

objects in the McGill shape dataset, which results in large within-class variations. Nonetheless, due to the discriminative feature representation of our method, as shown earlier in Fig. 9, DASD is robust to large nonrigid deformations.

TABLE I
RETRIEVAL RESULTS ON THE MCGILL DATASET.

Methods	NN	1-Tier	2-Tier	DCG
Covariance method [37]	0.977	0.732	0.818	0.937
Graph-based method [36]	0.976	0.741	0.911	0.933
PCA-based VLAT [35]	0.969	0.658	0.781	0.894
Hybrid BOW [34]	0.957	0.635	0.790	0.886
Hybrid 2D/3D [21]	0.925	0.557	0.698	0.850
CBoFHKS [22]	0.901	0.778	0.876	0.891
DASD	0.988	0.782	0.834	0.955

2) *SHREC’10 ShapeGoogle Dataset*: We also compared our proposed DASD method to the bag of feature (BOF) descriptor based on standard vector quantization (VQ) [20], sparse coding with unsupervised dictionary learning (UDL) [23], sparse coding with supervised dictionary learning (SDL) [23] and the CompactBoFHKS method (CBoFHKS) [22] on the SHREC’10 ShapeGoogle dataset [20]. We used the mean average precision criterion in our evaluations. Evaluation results are summarized in Table II. From this table, one can see that our proposed DASD is superior to the BOF descriptors based on standard VQ [20], UDL [23], SDL [23] and CBoFHKS [22] in the cases of isometry, isometry+topology, partiality and triangulation. Since deep auto-encoder has good ability to model nonlinearity, DASD can characterize the low-dimensional manifold embedded in the high-dimensional shape space better and therefore achieve better performance. For example, in the cases of isometry+topology and partiality, the supervised dictionary learning based shape descriptor can achieve accuracies of 0.956 and 0.951, while our proposed DASD can achieve accuracies of 0.982 and 0.973, respectively.

TABLE II
RETRIEVAL RESULTS (MEAN AVERAGE PRECISION) ON THE SHREC’10 SHAPEGOOGLE DATASET.

Transformation	VQ [20]	UDL [23]	SDL [23]	CBoFHKS [22]	DASD
Isometry	0.988	0.977	0.994	0.966	0.998
Topology	1.000	1.000	1.000	1.000	0.996
Isometry+Topology	0.933	0.934	0.956	0.915	0.982
Partiality	0.947	0.948	0.951	0.968	0.973
Triangulation	0.954	0.950	0.955	0.891	0.955

3) *SHREC’14 Human Dataset*: For the synthetic and scanned human sub-datasets, we compare the proposed DASD method to several recently proposed shape retrieval methods: Histogram of area projection transform (HAPT) [38], intrinsic pyramid matching (ISPM) [39], reduced Bi-harmonic distance matrix (RBiHDM) [40], deep belief network (DBN) [32], the bag of feature descriptor based on standard vector quantization (VQ) [20], and sparse coding with unsupervised dictionary learning (UDL) [23]. For the synthetic sub-dataset, 11 shapes in each class are used to train the discriminative auto-encoder and the remaining shapes in each class are used for testing. For the scanned sub-dataset, 6 shapes in each class are used for training and the remaining shapes are used for testing. The mean average precisions are reported in Table III. The scanned sub-dataset is an extremely challenging dataset, as can be seen in the table, our method can achieve better performance. For

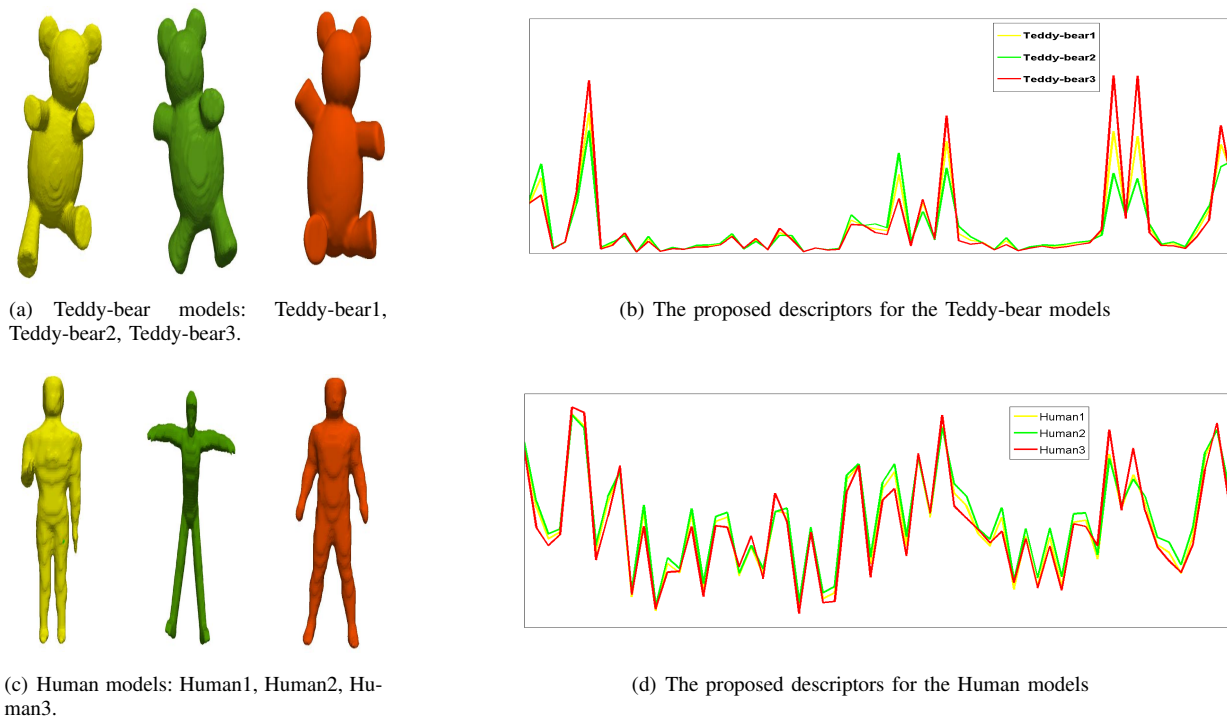


Fig. 9. The proposed descriptors for the Teddy-bear model and the Human model.

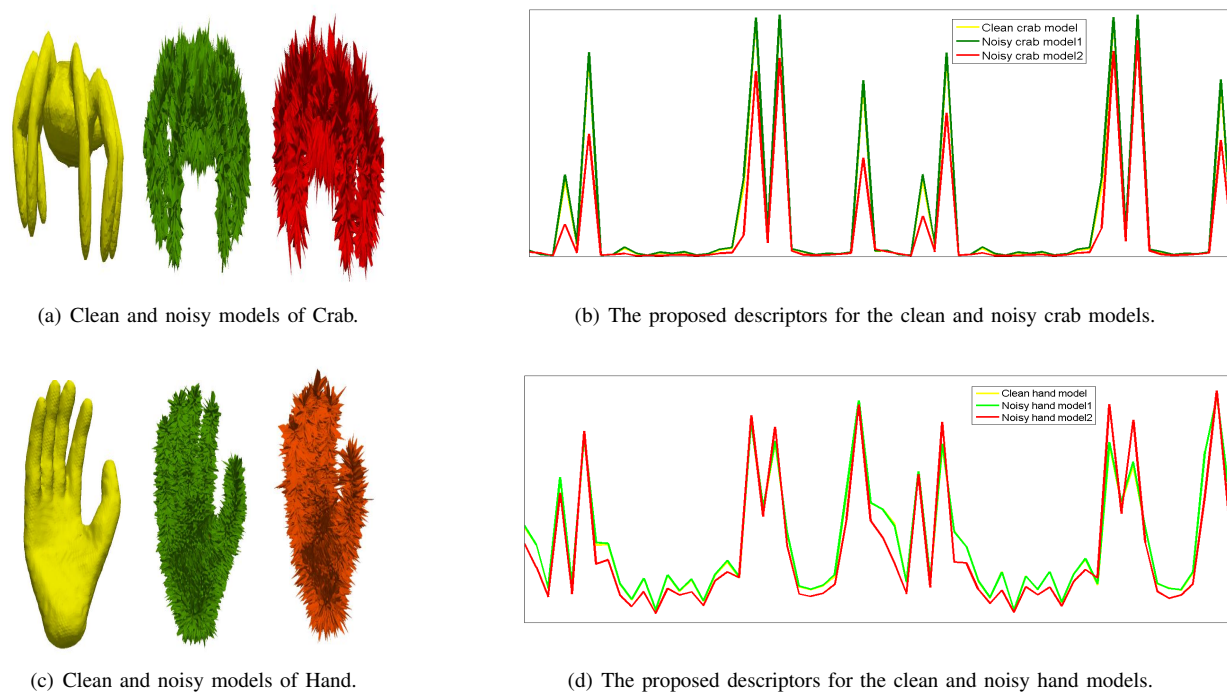


Fig. 10. The proposed descriptors for the clean and noisy models of Crab and Hand. In (a) and (c), the green and red shapes were generated with noise of $R = 0.01$ and $R = 0.04$, respectively. In (b) and (d), the descriptors for the shapes are represented by the yellow, green and red curves, corresponding to the clean model, the noisy model with noise of $R = 0.01$, and the noisy model with noise of $R = 0.04$, respectively.

the synthetic sub-dataset, the mean average precision of our method is lower than the ISPM method [39], and slightly lower than the DBN [32] and UDL [23] methods.

TABLE III
RETRIEVAL RESULTS (MEAN AVERAGE PRECISION) ON THE SHREC'14 HUMAN DATASET.

Method	Synthetic model	Scanned model
HAPT[38]	0.817	0.637
ISPM[39]	0.92	0.258
RBHDM[40]	0.642	0.640
DBN[32]	0.842	0.304
VQ [20]	0.813	0.514
UDL [23]	0.842	0.523
DASD	0.823	0.657

4) *SHREC'14 LSCRTB Dataset*: For the SHREC'14 LSCRTB dataset, we compared our proposed DASD method to the state-of-the-art methods [33]: CSLBP, HSR-DE, KVLAD, DBNAA_DERE, BF-DSIFT, VM-1SIFT, ZFDR and DBSVC. We used the NN, 1-Tier, 2-Tier, E-Measures (E) and DCG for performance evaluation. The evaluation results are listed in Table IV. From this table, one can see that the proposed DASD method can obtain the best performance with the NN and E measures. In terms of the 1-Tier and 2-Tier measures, our proposed DASD method is lower than the DBSVC method. Nonetheless, our proposed DASD method is comparable to the CSLBP, HSR-DE, DBNAA_DERE, BF-DSIFT, VM-1SIFT and ZFDR methods. Note that we did not compare our method to the LCDR-DBSRC method in [33]. This is because that we only calculated the Euclidean distance between the learned 3D descriptors while in the LCDR-DBSRC method the manifold ranking algorithm [41] is employed on the learned 3D shape features.

In the compared methods [33], the CSLBP, HSR-DE, DBNAA_DERE and ZFDR methods use a combination of hand-crafted features while the KVLAD, BF-DSIFT, VM-1SIFT and DBSVC methods use learned features. In all these methods, the features are extracted on the rendered images by projecting the 3D shape from different viewpoints. This type of approach suffers from the following two drawbacks: 1) the rendered images are sensitive to the deformation of the 3D model; 2) pre-alignment of 3D models (pose, transformation) in the same category is usually employed prior to 3D surface parameterization in order to normalize the models that are used for learning. However, normalization of 3D models may be difficult due to different types of structural variations.

Nonetheless, in our proposed DASD method, based on the HKS, the formed multiscale shape distribution is a statistic of the local geometrical structure information of the shape, which is robust to the deformations. Then, by imposing the Fisher discrimination criteria on the neurons in the hidden layer of the neural network, the proposed discriminative auto-encoder can minimize within-class variations and maximize between-class variations of the shape features. Furthermore, the developed discriminative auto-encoder can improve robustness of the shape features to large deformations. And it does not need shape alignment.

TABLE IV
RETRIEVAL RESULTS ON THE SHREC'14 LSCRTB DATASET.

Method	NN	1-Tier	2-Tier	E	DCG
CSLBP [33]	0.840	0.353	0.452	0.197	0.736
HSR-DE[33]	0.837	0.381	0.490	0.203	0.752
KVLAD[33]	0.605	0.413	0.546	0.214	0.746
DBNAA_DERE[33]	0.817	0.355	0.464	0.188	0.731
BF-DSIFT [33]	0.824	0.378	0.492	0.201	0.756
VM-1SIFT[33]	0.732	0.282	0.380	0.158	0.688
ZFDR [33]	0.838	0.386	0.501	0.209	0.757
DBSVC [33]	0.868	0.438	0.563	0.234	0.790
DASD	0.897	0.401	0.503	0.243	0.774

V. CONCLUSIONS

In this paper, we have proposed a novel deep shape descriptor for 3D shape retrieval. We first compute the multiscale shape distribution features and then train a set of discriminative auto-encoders to extract high-level shape features at different scales. By imposing the Fisher discrimination criterion on the hidden layers of the auto-encoders, our feature representation results in small within-class scatter and large between-class scatter. Our shape descriptor is formed by concatenating the high-level features from different scales. Experimental results demonstrated the superior performance of our proposed descriptor.

REFERENCES

- [1] R. M. Rustamov, "Laplace-beltrami eigenfunctions for deformation invariant shape representation," *Eurographics symposium on Geometry processing*, pp. 225–233, 2007.
- [2] R. Osada, T. Funkhouser, B. Chazelle, and D. Dokin, "Shape distributions," *ACM Transactions on Graphics*, vol. 33, pp. 133–154, 2002.
- [3] S. Katz, G. Leifman, and A. Tal, "Mesh segmentation using feature point and core extraction," *The Visual Computer*, vol. 21, pp. 649–658, 2005.
- [4] F. De Goes, S. Goldenstein, and L. Velho, "A hierarchical segmentation of articulated bodies," *Computer Graphics Forum*, vol. 27, pp. 1349–1356, 2008.
- [5] X. Chen, A. Golovinskiy, and T. Funkhouser, "A benchmark for 3D mesh segmentation," *ACM Transactions on Graphics*, 2009.
- [6] A. M. Bronstein, M. M. Bronstein, and R. Kimmel, "Efficient computation of isometry-invariant distances between surfaces," *SIAM Journal on Scientific Computing*, vol. 28, pp. 1812–1836, September 2006.
- [7] J. W. H. Tangelder and R. C. Veltkamp, "A survey of content based 3D shape retrieval methods," *In Shape Modeling International*, pp. 145–156, 2004.
- [8] N. Iyer, S. Jayanti, K. Lou, Y. Kalyanaraman, and K. Ramani, "Three-dimensional shape searching: state-of-the-art review and future trends," *Computer Aided Design*, vol. 37, no. 5, pp. 509 – 530, 2005.
- [9] M. Elad, A. Tal, and S. Ar, "Content based retrieval of VRML objects - an iterative and interactive approach," *Eurographics Workshop Multimedia*, pp. 97–108, 2001.
- [10] D. V. Vranic, D. Saupe, and J. Richter, "Tools for 3D-object retrieval: Karhunen-Loeve transform and spherical

- harmonics,” *Workshop on Multimedia Signal Processing*, pp. 293–298, 2001.
- [11] P. Shilane, P. Min, M. Kazhdan, and T. Funkhouser, “The Princeton shape benchmark,” *In Shape Modeling International*, pp. 167–178, 2004.
- [12] D. Saupe and D. V. Vranic, “3D model retrieval with spherical harmonics and moments,” *DAGM Symposium on Pattern Recognition*, pp. 392–397, 2001.
- [13] D.-Y. Chen, X.-P. Tian, Y.-T. Shen, and M. Ouhyoung, “On visual similarity based 3d model retrieval,” *Computer Graphics Forum*, vol. 22, no. 3, pp. 223–232, 2003.
- [14] V. Jain and H. Zhang, “A spectral approach to shape-based retrieval of articulated 3D models,” *Computer Aided Design*, vol. 39, no. 5, pp. 398–407, 2007.
- [15] J. Assfalg, M. Bertini, A. D. Bimbo, and P. Pala, “Content-based retrieval of 3D objects using spin image signatures,” *IEEE Transactions on Multimedia*, vol. 9, no. 3, pp. 589–599, 2007.
- [16] S. Belongie, J. Malik, and J. Puzicha, “Shape context: A new descriptor for shape matching and object recognition,” in *Advances in Neural Information Processing Systems 13, Denver, CO, USA, 2000*, pp. 831–837.
- [17] A. Zaharescu, E. Boyer, K. Varanasi, and R. Horaud, “Surface feature detection and description with applications to mesh matching,” in *IEEE Conference on Computer Vision and Pattern Recognition, Miami, Florida, USA, 2009*, pp. 373–380.
- [18] J. Sun, M. Ovsjanikov, and L. J. Guibas, “A concise and provably informative multi-scale signature based on heat diffusion,” *Computer Graphics Forum*, vol. 28, no. 5, pp. 1383–1392, 2009.
- [19] A. M. Bronstein, M. M. Bronstein, R. Kimmel, M. Mahmoudi, and G. Sapiro, “A Gromov-Hausdorff framework with diffusion geometry for topologically-robust non-rigid shape matching,” *International Journal of Computer Vision*, vol. 89, pp. 266–286, 2010.
- [20] A. M. Bronstein, M. M. Bronstein, L. J. Guibas, and M. Ovsjanikov, “Shape google: Geometric words and expressions for invariant shape retrieval,” *ACM Transactions on Graphics*, vol. 30, no. 1, p. 1, 2011.
- [21] G. Lavoué, “Combination of bag-of-words descriptors for robust partial shape retrieval,” *The Visual Computer*, vol. 28, no. 9, pp. 931–942, 2012.
- [22] Z. Lian and J. Zhang, et al., “SHREC’15 track: Non-rigid 3D shape retrieval,” in *Eurographics Workshop on 3D Object Retrieval, Zurich, Switzerland, 2015*, pp. 107–120.
- [23] R. Litman, A. M. Bronstein, M. M. Bronstein, and U. Castellani, “Supervised learning of bag-of-features shape descriptors using sparse coding,” *Computer Graphics Forum*, vol. 33, no. 5, pp. 127–136, 2014.
- [24] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, “3D shapenets: A deep representation for volumetric shapes,” in *IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 2015*, pp. 1912–1920.
- [25] D. Boscaini, J. Masci, S. Melzi, M. M. Bronstein, U. Castellani, and P. Vnderghenst, “Learning class-specific descriptors for deformable shapes using localized spectral convolutional networks,” *Computer Graphics Forum*, vol. 34, no. 5, pp. 13–23, 2015.
- [26] J. Masci, D. Boscaini, M. M. Bronstein, and P. Vnderghenst, “Intrinsic convolutional neural networks on riemannian manifolds,” in *IEEE Workshop on 3D Representation and Recognition (3dRR)*, 2015.
- [27] A. Dosovitskiy, J. T. Springenberg, and T. Brox, “Learning to generate chairs with convolutional neural networks,” in *IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 2015*, pp. 1538–1546.
- [28] G. Hinton and R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, vol. 313, no. 5786, pp. 504 – 507, 2006.
- [29] Y. Bengio, “Learning deep architectures for AI,” *Foundations and Trends in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
- [30] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification (2nd Ed)*. Wiley, 2001.
- [31] K. Siddiqi, J. Zhang, D. Macrini, A. Shokoufandeh, S. Bouix, and S. J. Dickinson, “Retrieving articulated 3D models using medial surfaces,” *Machine Vision Application*, vol. 19, no. 4, pp. 261–275, 2008.
- [32] D. Pickup and X. Sun, et al., “SHREC’14 track: Shape retrieval of non-rigid 3D human models,” in *Eurographics Workshop on 3D Object Retrieval, Strasbourg, France, 2014*.
- [33] B. Li and Y. Lu, et al., “SHREC’14 track: Large scale comprehensive retrieval track benchmark,” in *Eurographics Workshop on 3D Object Retrieval, Strasbourg, France, 2014*.
- [34] P. Papadakis, I. Pratikakis, T. Theoharis, G. Passalis, and S. J. Perantonis, “3D object retrieval using an efficient and compact hybrid shape descriptor,” in *Eurographics Workshop on 3D Object Retrieval, Crete, Greece, 2008*, pp. 9–16.
- [35] H. Tabia, D. Picard, H. Laga, and P. H. Gosselin, “Compact vectors of locally aggregated tensors for 3D shape retrieval,” in *Eurographics Workshop on 3D Object Retrieval, Girona, Spain, 2013*, pp. 17–24.
- [36] A. Agathos, I. Pratikakis, P. Papadakis, S. J. Perantonis, P. N. Azariadis, and N. S. Sapidis, “Retrieval of 3D articulated objects using a graph-based representation,” in *Eurographics Workshop on 3D Object Retrieval, Munich, Germany, 2009*, pp. 29–36.
- [37] H. Tabia, H. Laga, D. Picard, and P. H. Gosselin, “Covariance descriptors for 3D shape matching and retrieval,” in *IEEE Conference on Computer Vision and Pattern Recognition, Columbus, OH, USA, 2014*, pp. 4185–4192.
- [38] A. Giachetti and C. Lovato, “Radial symmetry detection and shape characterization with the multiscale area projection transform,” *Computer Graphics Forum*, vol. 31, no. 5, pp. 1669–1678, 2012.
- [39] C. Li and A. B. Hamza, “A multiresolution descriptor for deformable 3D shape retrieval,” *The Visual Computer*, vol. 29, no. 6-8, pp. 513–524, 2013.
- [40] J. Ye, Z. Yan, and Y. Yu, “Fast nonrigid 3D retrieval using

model space transform,” in *International Conference on Multimedia Retrieval, Dallas, TX, USA, 2013*, pp. 121–



Jin Xie received his Ph.D. degree from the Department of Computing, The Hong Kong Polytechnic University, with research on texture learning based texture classification. He is a postdoctoral associate at New York University Abu Dhabi. His research interests include image forensics, computer vision and machine learning. Currently he is focusing on 3D computer vision with the convex optimization and deep learning methods.



Guoxian Dai received his master degree from Fudan University, China. He is a Ph.D. candidate in the Department of Computer Science and Engineering at the NYU Tandon School of Engineering. His current research interests focus on 3D shape analysis such as 3D shape retrieval and cross-domain 3D model retrieval.



Fan Zhu received the MSc degree with distinction in Electrical Engineering and the Ph.D. degree at the Visual Information Engineering group from the Department of Electronic and Electrical Engineering, the University of Sheffield, Sheffield, U.K, in 2011 and 2015, respectively. He is currently a post-doctoral associate at New York University Abu Dhabi. His research interests include submodular optimization for computer vision, sparse coding, 3D feature learning, dictionary learning and transfer learning. He has authored/co-authored over 10 papers in well-known journals/conferences such as IJCV, IEEE TNNLS, CVPR, CIKM and BMVC, and two China patents. He has been awarded the National Distinguished Overseas Self-funded Student of China prize in 2014. He serves as a reviewer of IEEE Transactions on Cybernetics.

126.

- [41] D. Zhou, J. Weston, A. Gretton, O. Bousquet, and B. Schölkopf, “Ranking on data manifolds,” in *Advances in Neural Information Processing Systems, Vancouver and Whistler, British Columbia, Canada, 2003*, pp. 169–176.



Edward K. Wong received his Ph.D. degree in Electrical Engineering from Purdue University. He is currently an associate professor and the director of MS CS program in the Department of Computer Science and Engineering at the NYU Tandon School of Engineering. His research interests lie in the general areas of computer vision, pattern recognition, and machine learning. His current research focus is on developing novel machine-learning-based techniques for video surveillance applications. He had previously worked on funded projects in document image analysis and security, video scene segmentation and classification, fingerprint verification, morphological image processing, infrared target classification, three-dimensional object recognition, pavement image analysis, and optical character recognition, among others. He had published extensively in image processing and multimedia conferences and journals. Dr. Wong is currently an associate editor for two international journals in multimedia and security, and he had served on the organizing committee and technical program committee of several major IEEE and ACM technical conferences in image processing and multimedia.



Yi Fang received his Ph.D. degree from Purdue University with research focus on computer graphics and vision. Upon one year industry experience as a research intern in Siemens in Princeton, New Jersey and a senior research scientist in Riverain Technologies in Dayton, Ohio, and a half-year academic experience as a senior staff scientist at Department of Electrical Engineering and Computer science, Vanderbilt University, Nashville, he joined NYU Abu Dhabi as an Assistant Professor of Electrical and Computer Engineering. He is currently working on the development of state-of-the-art techniques in large-scale visual computing, deep visual learning, deep cross-domain and cross-modality model, and their applications in engineering, social science, medicine and biology.