

3-D Object Retrieval and Recognition With Hypergraph Analysis

Yue Gao, Meng Wang, *Member, IEEE*, Dacheng Tao, *Senior Member, IEEE*,
Rongrong Ji, *Member, IEEE*, and Qionghai Dai, *Senior Member, IEEE*

Abstract—View-based 3-D object retrieval and recognition has become popular in practice, e.g., in computer aided design. It is difficult to precisely estimate the distance between two objects represented by multiple views. Thus, current view-based 3-D object retrieval and recognition methods may not perform well. In this paper, we propose a hypergraph analysis approach to address this problem by avoiding the estimation of the distance between objects. In particular, we construct multiple hypergraphs for a set of 3-D objects based on their 2-D views. In these hypergraphs, each vertex is an object, and each edge is a cluster of views. Therefore, an edge connects multiple vertices. We define the weight of each edge based on the similarities between any two views within the cluster. Retrieval and recognition are performed based on the hypergraphs. Therefore, our method can explore the higher order relationship among objects and does not use the distance between objects. We conduct experiments on the National Taiwan University 3-D model dataset and the ETH 3-D object collection. Experimental results demonstrate the effectiveness of the proposed method by comparing with the state-of-the-art methods.

Index Terms—3-D object recognition, 3-D object retrieval, hypergraph learning, view-based.

I. INTRODUCTION

THE RAPID advances in computer techniques, graphics hardware, and networks have led to the wide application of 3-D models in various domains [1], [2], such as 3-D graphics, the medical industry, movie production, architecture design, and the engineering community. Therefore, efficient 3-D object retrieval and recognition technologies [3]–[6] are of great importance for many applications, including computer-aided design (CAD) [7] and molecular biology [8]. This

paper refers to 3-D object retrieval [1] and recognition [9] as the tasks of finding relevant objects for a given query and categorizing an object into one of a set of classes, respectively.

Existing 3-D object retrieval and recognition approaches can be divided into two paradigms, namely, model-based and view-based. Early methods [10]–[13] are mainly model-based and require 3-D models. When 3-D models are not available in many practical applications, an alternative approach is to reconstruct 3-D models based on a carefully collected set of 2-D images. The 3-D model reconstruction is computationally expensive and the sampling of the images needs to be sufficiently fine-grained to reconstruct a reasonably good 3-D model. These difficulties severely limit the practical applications of model-based 3-D object analysis methods.

Recently, extensive research efforts [14]–[16] have been dedicated to view-based 3-D object retrieval and recognition because of its flexibility in representing 3-D objects by multiple views. These methods describe a 3-D object by a set of views that are captured by a group of cameras. The retrieval and recognition of 3-D objects are typically accomplished by matching the views of a query object to those of reference objects in the database. View-based 3-D object retrieval and recognition are highly flexible and useful because the 3-D model information is not required. Visual analyses [17]–[22] have shown its superiority in multimedia applications. According to [23], “view-based methods have the advantages of being highly discriminative, can work for articulated objects, can be effective for partial matching and can also be beneficial for 2-D sketch-based and 2-D image-based queries.” In addition, experimental results reported in [2] and [24] demonstrate that view-based approaches can exhibit better overall retrieval performance than model-based methods.

Most classification and ranking algorithms are built upon the distance estimation of samples. Since each 3-D object is represented by a set of views, as illustrated in Fig. 1, it is difficult to precisely estimate the distance between two objects. Most existing methods employ view-matching to measure the distance between two objects [23], [25]–[27], such as the Hausdorff distance [28], [29] or the Earth Mover’s distance [30] and by using bipartite graph-matching [31]. In [25] and [32], probabilistic methods have been employed for matching multiple views. These methods can be effective if two objects from the same category have several close views. It is, however, inconvenient to require objects to have close views in practice. In addition, they ignore the higher order relationship among objects. For example, we cannot use the Hausdorff distance to indicate which three or more objects share close

Manuscript received November 14, 2010; revised March 22, 2012; accepted May 2, 2012. Date of publication May 15, 2012; date of current version August 22, 2012. This work was supported in part by the National Basic Research Project, China, under Grant 2010CB731800, the Key Project of the NSFC, China, under Grant 61035002 and Grant 61021063, and the Australia Research Council Discovery under Project DP-120103730. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Sharath Pankanti.

Y. Gao and Q. Dai are with the Department of Automation, Tsinghua National Laboratory for Information Science and Technology, Tsinghua University, Beijing 100084, China.

M. Wang is with the School of Computer and Information, Hefei University of Technology, Hefei 230009, China (e-mail: eric.mengwang@gmail.com).

D. Tao is with the Centre for Quantum Computation & Intelligent Systems and the Faculty of Engineering & Information Technology, University of Technology, Sydney NSW 2007, Australia.

R. Ji was with the Department of Computer Science, Harbin Institute of Technology, Harbin 150001, China. He is now with Columbia University, New York, NY 10027 USA.

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TIP.2012.2199502



Fig. 1. Multiple views of two example 3-D objects.

views because the distance only reflects the information of the closest views of a pair of objects.

We propose a 3-D object retrieval and recognition approach by exploring the higher order relationship among 3-D objects via hypergraphs. In our framework, multiple hypergraphs are constructed and the learning on the fused hypergraph is conducted for 3-D object retrieval and recognition. We conduct experiments on the National Taiwan University (NTU) [33] and ETH [34] 3-D object datasets to demonstrate the effectiveness of the proposed method by comparing with several existing retrieval and recognition algorithms. In contrast to the image retrieval approaches in [35] and [36], the proposed method has the following differences and advantages. First, multiple hypergraphs are employed in our method. Considering the complex relationship among the views of 3-D objects, it is difficult to construct a single optimal hypergraph. Using multiple hypergraphs can reduce the risk from using a single hypergraph. Multiple hypergraphs are fused by equivalent weights in retrieval, while the optimal weights for combining multiple hypergraphs are learned together with the categories of objects in recognition. Second, hypergraphs are constructed by different methods. In [35] and [36], the hypergraphs are constructed by using different visual features, while in our approach the hypergraphs are built by using view clustering.

The rest of this paper is organized as follows. In Section II, we briefly introduce hypergraph analysis and the involved notations. Section III presents our hypergraph construction approach and then introduces the 3-D object retrieval and recognition algorithms. Experimental settings are provided in Section IV. Experimental results on the NTU and ETH 3-D object datasets are provided in Section V. We conclude this paper in Section VI.

II. BRIEF INTRODUCTION TO HYPERGRAPH ANALYSIS

This section briefly reviews the hypergraph analysis theory. In a conventional (simple) graph, samples are represented by vertices and two related vertices are linked by an edge. The simple graph can be undirected or directed, depending on whether the pairwise relationships between objects are symmetric or not [37]. Learning tasks can be performed on the graph. For example, when categorizing samples that are represented by vectors in a feature space, we can construct a undirected graph based on their pairwise distances and, consequently, many graph-based learning methods can be performed [37]–[40]. However, the graph does not reflect the

higher order information. Different from a simple graph, an edge in a hypergraph can connect three or more vertices. Table I summarizes important notations and the corresponding definitions throughout this paper.

A hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w)$ is composed of a vertex set \mathcal{V} , an edge set \mathcal{E} , and the weights of the edges w . Each edge e is assigned a weight $w(e)$. The hypergraph \mathcal{G} can be denoted by a $|\mathcal{V}| \times |\mathcal{E}|$ incidence matrix \mathbf{H} , in which each entry is defined by

$$h(v, e) = \begin{cases} 1, & \text{if } v \in e \\ 0, & \text{if } v \notin e. \end{cases} \quad (1)$$

For a vertex $v \in \mathcal{V}$, its degree is defined by

$$d(v) = \sum_{e \in \mathcal{E}} w(e) h(v, e). \quad (2)$$

For an edge $e \in \mathcal{E}$, its degree is defined by

$$\delta(e) = \sum_{v \in \mathcal{V}} h(v, e). \quad (3)$$

We let \mathbf{D}_v and \mathbf{D}_e denote the diagonal matrices of the vertex degrees and the edge degrees, respectively. Let \mathbf{W} denote the diagonal matrix of the edge weights.

Different machine learning tasks can be performed on hypergraphs, such as classification, clustering, ranking, and embedding. Here we take binary classification as an example. Zhou *et al.* [37] presented a regularization framework

$$\arg \min_f \{ \lambda R_{\text{emp}}(f) + \Omega(f) \} \quad (4)$$

where f is the classification function, $\Omega(f)$ is a regularizer on the hypergraph, $R_{\text{emp}}(f)$ is an empirical loss, and $\lambda > 0$ is the tradeoff parameter. The regularizer on the hypergraph is defined by

$$\Omega(f) = \frac{1}{2} \sum_{e \in \mathcal{E}} \sum_{u, v \in \mathcal{V}} \frac{w(e) h(u, e) h(v, e)}{\delta(e)} \left(\frac{f(u)}{\sqrt{d(u)}} - \frac{f(v)}{\sqrt{d(v)}} \right)^2. \quad (5)$$

Let $\Theta = \mathbf{D}_v^{-(1/2)} \mathbf{H} \mathbf{W} \mathbf{D}_e^{-1} \mathbf{H}^T \mathbf{D}_v^{-(1/2)}$, and $\Delta = \mathbf{I} - \Theta$. The normalized cost function can be written as

$$\Omega(f) = f^T \Delta f \quad (6)$$

where Δ is a positive semidefinite matrix, and it is usually called the hypergraph Laplacian.

The hypergraph has been investigated in many data mining and information retrieval tasks [35], [37], [38], [41], [42] because of its capability of capturing the higher order relationship of samples. A probabilistic hypergraph matching method is proposed in [43] to match two feature sets. In the transductive learning framework for image retrieval [36], [44], each image is represented by a vertex in a probabilistic hypergraph and the image retrieval task is formulated as a hypergraph ranking problem. In the hypergraph, each object view is represented by visual features, and probabilistic hypergraph ranking is performed to explore the relationship between images in the visual feature space. A hypergraph-based social image search method is presented in [45], in which both the visual content of images and tags have been investigated in the constructed hypergraph structure. A hypergraph-based

TABLE I
NOTATIONS AND DEFINITIONS

Notation	Definition
$\mathcal{G} = (\mathcal{V}, \mathcal{E}, w)$	\mathcal{G} indicates a hypergraph, and \mathcal{V} , \mathcal{E} , and w indicate the set of vertices, the set of edges, and the weights of the edges, respectively.
\mathcal{V}	The set of vertices of the hypergraph, which contains n elements.
\mathcal{E}	The set of edges of the hypergraph that contains n_e elements.
x_i	The i th view of a 3-D object.
n	The number of objects in the database, i.e., $ \mathcal{V} $.
n_e	The number of edges in the hypergraph, i.e., $ \mathcal{E} $.
w	The $n_e \times 1$ weight vector of the edges in the hypergraph.
\mathbf{W}	The diagonal matrix of the edge weights.
$d(v)$	The degree of the vertex v .
$\delta(e)$	The degree of the edge e .
\mathbf{D}_v	The diagonal matrix of the vertex degrees. The dimension is $n \times n$.
\mathbf{D}_e	The diagonal matrix of the edge degrees. The dimension is $n_e \times n_e$.
\mathbf{H}_i	The incidence matrix for the i th hypergraph. The dimension is $n \times n_e$.
n_g	The number of hypergraphs.
$\alpha = [\alpha_1, \alpha_2, \dots, \alpha_{n_g}]^T$	The weights used for combining multiple hypergraphs.
n_c	The number of object classes in the database.
y_i	The label vector for the i th class. Its j th element is 1 if the j th object belongs to the i th class, and otherwise it is 0. The dimension is $n \times 1$.
$\mathbf{Y} = [y_1, y_2, \dots, y_{n_c}]$	The label matrix for all samples. The dimension is $n \times n_c$.
f_i	The to-be-learned relevance score vector for the i th class. The dimension is $n \times 1$.
$\mathbf{F} = [f_1, f_2, \dots, f_{n_c}]$	The to-be-learned relevance score matrix for all samples. The dimension is $n \times n_c$.

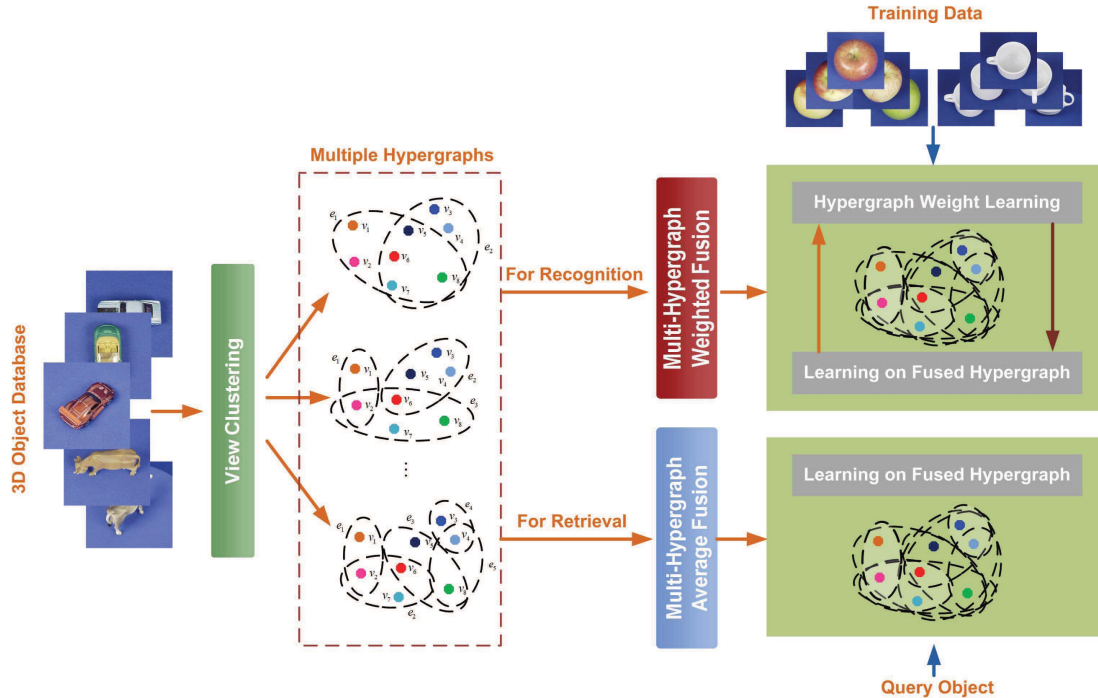


Fig. 2. Schematic illustration of the proposed method. Given a 3-D database, in which each 3-D object is represented by several views, the first step is grouping all views into clusters. With different clustering settings, multiple hypergraphs are constructed based on these clustering results. Each dotted circle in the figure represents an edge in the corresponding hypergraph. Retrieval and recognition are then performed based on these multiple hypergraphs.

3-D object representation method is presented in [46], which constructs the hypergraph by using the correlation among different surface boundary segments of an object in the CAD system. Generally, each 3-D model is composed of a group of surface patches. Here, each distinct surface patch is regarded

as a surface boundary segment. In the constructed hypergraph, each boundary segment is regarded as an attribute vertex. For each boundary segment, its associated attributes include the length of the segment and the angle between two connected segments. For two vertices, the edge represents the

connection between two corresponding segments, where the angle between the two connected segments is regarded as the edge attribute. These segments are grouped by different methods, including jump edges, where the depths of the two segments are largely different, concave angular-relation-based edges, edges adjacent to the background, and a combination of the aforementioned three types of edges. By using these segment groups, segments in one group are connected in the hypergraph. A class-specific hypergraph (CSHG) [47] is proposed to integrate local scale-invariant feature transform (SIFT) and global geometric constraints for object recognition, where the hypergraph is employed to model a specific category of objects with multiple appearance instances. The vertices of the hypergraph are the images that belong to the objects, and the selected SIFT points are employed as the features of vertices. It is further extended in [48] to learn a large-scale CSHG model for 3-D object recognition.

III. HYPERGRAPH-BASED 3-D OBJECT RETRIEVAL AND RECOGNITION

This section introduces the hypergraph-based 3-D object retrieval and recognition approaches. Fig. 2 shows the scheme of our approach. We first group the views of all objects into clusters. Each cluster is then regarded as an edge for connecting objects that have views in this cluster (note that an edge can connect multiple vertices in a hypergraph). A hypergraph is thus constructed, in which vertices denote objects in the database (this means the number of vertices is equivalent to the number of objects). We define the weight of an edge on the basis of the visual similarities between any two views in the cluster. By varying the number of clusters, multiple hypergraphs can be generated. These hypergraphs actually encode the relationships among objects with different granularities. By performing retrieval and recognition on these hypergraphs, we can avoid the object distance estimation problem because the hypergraphs already comprehensively describe the relationship of the objects. For example, two objects that are connected by more edges should be closer. Moreover, since each edge connects multiple objects, the higher order information, such as whether three or more objects share close views, can be explored. For retrieval, we fuse the hypergraphs by using equivalent weights. However, we learn the optimal combination coefficients for combining multiple hypergraphs by using the training data for recognition.

A. Hypergraph Construction

Hypergraph construction plays an important role in hypergraph-based 3-D object modeling. We regard each 3-D object in the database as a vertex in the hypergraph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w)$. For example, assuming there are in total n objects in the database, the generated hypergraph has n vertices.

We group all the views of the 3-D objects in the database into clusters by using the K -means algorithm [49]. The objects that have views in a cluster are connected by the corresponding edge. We let \mathbf{D}_v and \mathbf{D}_e denote the diagonal matrices of the vertex degrees and the edge degrees, respectively, and the

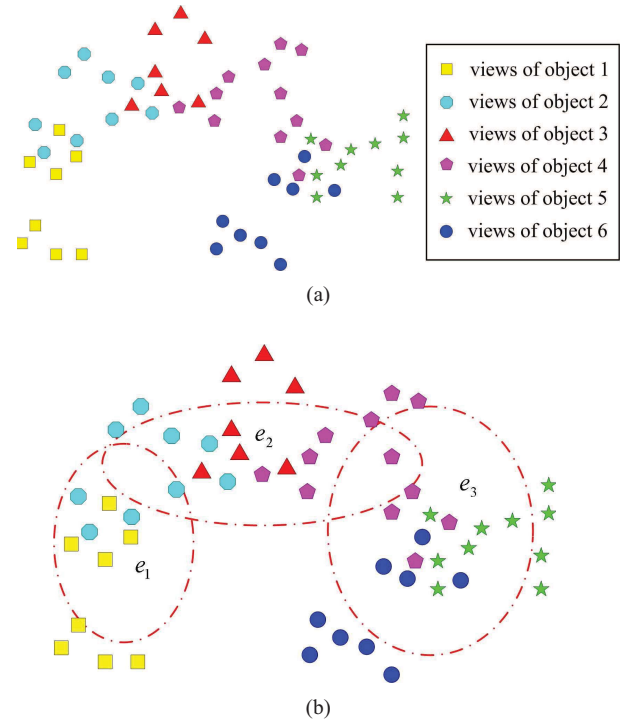


Fig. 3. Schematic illustration of the proposed hypergraph construction method. (a) Views of six 3-D objects, with each object's view indicated by a specific mark. (b) View clusters employed to construct hypergraph edges.

incidence matrix \mathbf{H} is constructed by using (1). The weight w of an edge e is estimated by the sum of the similarities between two views in the cluster

$$w(e) = \sum_{x_a, x_b \in e} \exp\left(-\frac{d(x_a, x_b)^2}{\sigma^2}\right) \quad (7)$$

where x_a and x_b are two views in the same view cluster and $d(x_a, x_b)$ is the distance between them. Euclidean distance is used for calculating $d(x_a, x_b)$. The parameter σ is empirically set to the median value of the distances of all view pairs.

The aforementioned process discovers the rationality of the proposed method: two objects tend to be connected by more edges if they have many close views, and an edge will be assigned a higher weight if the views in the cluster are closer. By varying the number of clusters, i.e., the parameter K in the K -means clustering, different hypergraphs can be constructed. These hypergraphs capture the relationship of the 3-D objects with different granularities. Fig. 3 uses an example to explain the proposed hypergraph construction method.

B. 3-D Object Retrieval by Hypergraph Modeling

Let $\mathcal{G}_1 = (\mathcal{V}_1, \mathcal{E}_1, w_1)$, $\mathcal{G}_2 = (\mathcal{V}_2, \mathcal{E}_2, w_2), \dots$, and $\mathcal{G}_{n_g} = (\mathcal{V}_{n_g}, \mathcal{E}_{n_g}, w_{n_g})$ denote the n_g hypergraphs, and $\{\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_{n_g}\}$, $\{\mathbf{D}_{v1}, \mathbf{D}_{v2}, \dots, \mathbf{D}_{vn_g}\}$, and $\{\mathbf{D}_{e1}, \mathbf{D}_{e2}, \dots, \mathbf{D}_{en_g}\}$ be the corresponding incidence matrices, vertex degree matrices, and hyperedge degree matrices, respectively. The retrieval needs to be performed based on the fusion of these hypergraphs. We denote the weight of i th hypergraph as α_i , where $\sum_{i=1}^{n_g} \alpha_i = 1$, and $\alpha_i \geq 0$.

We regard the retrieval task as a one-class classification problem [36]. Therefore, the transductive inference can be formulated as a regularization problem $\arg \min_f \{\lambda R_{\text{emp}}(f)\} + \Omega(f)$, and the regularizer term $\Omega(f)$ is defined by

$$\frac{1}{2} \sum_{i=1}^{n_g} \alpha_i \sum_{e \in \mathcal{E}_i} \sum_{u, v \in \mathcal{V}_i} \frac{w_i(e) h_i(u, e) h_i(v, e)}{\delta_i(e)} \times \left(\frac{f(u)}{\sqrt{d_i(u)}} - \frac{f(v)}{\sqrt{d_i(v)}} \right)^2 \quad (8)$$

where the vector f is the to-be-learned relevance score vector. Equation (8) further turns to

$$\begin{aligned} \Omega(f) &= \sum_{i=1}^{n_g} \alpha_i \sum_{e \in \mathcal{E}_i} \sum_{u, v \in \mathcal{V}_i} \frac{w_i(e) h_i(u, e) h_i(v, e)}{\delta_i(e)} \\ &\quad \times \left(\frac{f^2(u)}{d_i(u)} - \frac{f(u) f(v)}{\sqrt{d_i(u) d_i(v)}} \right) \\ &= \sum_{i=1}^{n_g} \alpha_i \left\{ \sum_{u \in \mathcal{V}_i} f^2(u) \sum_{e \in \mathcal{E}_i} \frac{w_i(e) h_i(u, e)}{d_i(u)} \sum_{v \in \mathcal{V}_i} \frac{h_i(v, e)}{\delta_i(e)} \right. \\ &\quad \left. - \sum_{e \in \mathcal{E}_i} \sum_{u, v \in \mathcal{V}_i} \frac{f(u) h_i(u, e) w_i(e) h_i(v, e) f(v)}{\sqrt{d_i(u) d_i(v)} \delta_i(e)} \right\} \\ &= \sum_{i=1}^{n_g} \alpha_i f^T (\mathbf{I} - \Theta_i) f \\ &= f^T \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} - \Theta_i) f \end{aligned} \quad (9)$$

where $\Theta_i = \mathbf{D}_{vi}^{-(1/2)} \mathbf{H}_i \mathbf{W}_i \mathbf{D}_{ei}^{-1} \mathbf{H}_i^T \mathbf{D}_{vi}^{-(1/2)}$. Let $\Delta = \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} - \Theta_i) = \mathbf{I} - \sum_{i=1}^{n_g} \alpha_i \Theta_i = \mathbf{I} - \Theta$, where $\Theta = \sum_{i=1}^{n_g} \alpha_i \Theta_i$. This Δ can be viewed as a fused hypergraph Laplacian. Thus we have

$$\Omega(f) = f^T \Delta f. \quad (10)$$

In this paper, all hypergraphs share the same notation \mathcal{V} . Thus, for all $i \in \{1, 2, \dots, n_g\}$, we let $\mathcal{V}_i = \mathcal{V}$. The loss function term is defined by

$$\|f - y\|^2 = \sum_{u \in \mathcal{V}} (f(u) - y(u))^2 \quad (11)$$

where y is the label vector. Let n denote the number of objects in the database and assume the i th object is selected as the query object. Let y denote an $n \times 1$ vector, where all the elements of y are 0 except the i th value which is 1. The learning task for 3-D object retrieval then becomes minimizing the sum of the two terms

$$\Phi(f) = f^T \Delta f + \lambda \|f - y\|^2 \quad (12)$$

where $\lambda > 0$ is the weighting parameter. By differentiating $\Phi(f)$ with respect to f , we obtain

$$f = \left(\mathbf{I} + \frac{1}{\lambda} \Delta \right)^{-1} y. \quad (13)$$

Algorithm 1 The Iterative Solution Method of (13)

Step 1: Initialize $f^{(t)}$, when $t = 0$.

Step 2: Update f by: $f^{(t+1)} = \frac{1}{1+\lambda} (\mathbf{I} - \Delta) f^{(t)} + \frac{\lambda}{1+\lambda} y$.

Step 3: Let $t = t + 1$, and then jump to Step 2 until convergence.

There are many existing nonnegative matrix factorization methods, such as in [50]–[52]. Similar to the approach developed in [53], the equation above can be efficiently solved by an iterative process. The process is illustrated in Algorithm 1.

Now we prove the convergence of this iterative process.

Theorem 1: The process in Algorithm 1 converges to (13).

Proof: To prove the theorem, we first prove that the eigenvalues of Θ are in $[-1, 1]$. Since $\Theta_i = \mathbf{D}_{vi}^{-(1/2)} \mathbf{H}_i \mathbf{W}_i \mathbf{D}_{ei}^{-1} \mathbf{H}_i^T \mathbf{D}_{vi}^{-(1/2)}$, we derive that its eigenvalues are in $[-1, 1]$. Therefore, $(\mathbf{I} \pm \Theta_i)$ are positive semidefinite.

Consequently, matrices $(\mathbf{I} \pm \Theta) = \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} \pm \Theta_i)$ are also positive semidefinite. This means that the eigenvalues of Θ are in $[-1, 1]$.

Now we prove the convergence of the iterative process. Without loss of generality, suppose $f^{(0)} = y$. From the iterative process, we have

$$\begin{aligned} f^{(t)} &= \left(\frac{\lambda}{1+\lambda} \right) \sum_{i=0}^{t-1} \left(\frac{1}{1+\lambda} \Theta \right)^i y + \left(\frac{1}{1+\lambda} \Theta \right)^t y \\ &= (1 - \zeta) \sum_{i=0}^{t-1} (\zeta \Theta)^i y + (\zeta \Theta)^t y \end{aligned} \quad (14)$$

where $\zeta = (1/(1+\lambda))$.

Since $0 < \zeta < 1$ and the eigenvalues of Θ are in $[-1, 1]$, we derive that

$$\lim_{t \rightarrow \infty} (\zeta \Theta)^t = 0 \quad (15)$$

and

$$\lim_{t \rightarrow \infty} \sum_{i=0}^{t-1} (\zeta \Theta)^i = (\mathbf{I} - \zeta \Theta)^{-1}. \quad (16)$$

There, we obtain

$$f = \lim_{t \rightarrow \infty} f^{(t)} = (1 - \zeta) (\mathbf{I} - \zeta \Theta)^{-1} y = \left(\mathbf{I} + \frac{1}{\lambda} \Delta \right)^{-1} y \quad (17)$$

which completes the proof. ■

We simply assign equal weights to all the n_g hypergraphs, i.e., $\alpha_i = (1/n_g)$. After obtaining f , the objects are ranked with decreasing relevance scores.

C. 3-D Object Recognition by Hypergraph Modeling

As a multiclass classification problem, 3-D object recognition differs from the retrieval task in terms of the following two aspects.

First, the vectors y and f need to be changed to matrices $\mathbf{Y} = [y_1, y_2, \dots, y_{n_c}]$ and $\mathbf{F} = [f_1, f_2, \dots, f_{n_c}]$, respectively. Both \mathbf{Y} and \mathbf{F} are $n \times n_c$ matrices. Second, since we have training data in the recognition task, we can learn the optimal combination coefficients to integrate all hypergraphs by using the training data.

We write the formulation of recognition as

$$\begin{aligned}
& \min_{\mathbf{F}, \alpha} \Phi(\mathbf{F}, \alpha) \\
& = \min_{\mathbf{F}, \alpha} \frac{1}{2} \sum_{k=1}^{n_c} \sum_{i=1}^{n_g} \alpha_i \sum_{e \in \mathcal{E}_{\mu}, v \in \mathcal{V}_i} \frac{w_i(e) h_i(u, e) h_i(v, e)}{\delta_i(e)} \\
& \quad \times \left(\frac{\mathbf{F}_{uk}}{\sqrt{d_i(u)}} - \frac{\mathbf{F}_{vk}}{\sqrt{d_i(v)}} \right)^2 + \mu \sum_{i=1}^{n_g} \alpha_i^2 + \lambda \sum_{k=1}^{n_c} \|f_k - y_k\|^2 \\
& = \min_{\mathbf{F}, \alpha} \sum_{k=1}^{n_c} f_k^T \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} - \Theta_i) f_k + \mu \sum_{i=1}^{n_g} \alpha_i^2 + \lambda \sum_{k=1}^{n_c} \|f_k - y_k\|^2 \\
& \text{s.t.} \quad \sum_{i=1}^{n_g} \alpha_i = 1, \quad \mu > 0, \quad \lambda > 0.
\end{aligned} \tag{18}$$

Here, $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_{n_g}]^T$, and \mathbf{Y} is the label matrix ($\mathbf{Y}_{ij} = 1$ if the i th object belongs to the j th class and $\mathbf{Y}_{ij} = 0$ otherwise). The matrix \mathbf{F} indicates the confidence scores of the classification. More specifically, \mathbf{F}_{ij} is the confidence score of categorizing the i th object into the j th class. Therefore, the final decision can be made by choosing the highest \mathbf{F}_{ij} for each object O_i . In the equation above, we have added a 2-norm regularizer on the weighting parameters of the hypergraphs. This enables us to learn a combination of the hypergraphs.

We adopt the alternating optimization to solve the above problem. First, we fix α and optimize \mathbf{F} , and (18) reduces to

$$\begin{aligned}
& \min_{\mathbf{F}} \sum_{k=1}^{n_c} f_k^T \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} - \Theta_i) f_k + \lambda \sum_{k=1}^{n_c} \|f_k - y_k\|^2 \\
& \text{s.t.} \quad \lambda > 0.
\end{aligned} \tag{19}$$

Similar to Section III-B, we obtain

$$\mathbf{F} = \left(\mathbf{I} + \frac{1}{\lambda} \Delta \right)^{-1} \mathbf{Y}. \tag{20}$$

Analogous to (13), it can also be solved with an iterative process. Actually, we only need to replace y and f with \mathbf{Y} and \mathbf{F} , respectively, in Algorithm 1.

Next, we fix \mathbf{F} and estimate α , and then we have

$$\begin{aligned}
& \min_{\alpha} \sum_{k=1}^{n_c} f_k^T \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} - \Theta_i) f_k + \mu \sum_{i=1}^{n_g} \alpha_i^2 \\
& \text{s.t.} \quad \sum_{i=1}^{n_g} \alpha_i = 1, \quad \mu > 0.
\end{aligned} \tag{21}$$

Here, the Lagrangian is employed and the optimization problem turns to

$$\min_{\alpha, \eta} \sum_{k=1}^{n_c} f_k^T \sum_{i=1}^{n_g} \alpha_i (\mathbf{I} - \Theta_i) f_k + \mu \sum_{i=1}^{n_g} \alpha_i^2 + \eta \left(\sum_{i=1}^{n_g} \alpha_i - 1 \right). \tag{22}$$

We can derive

$$\eta = \frac{-\sum_{k=1}^{n_c} f_k^T \sum_{i=1}^{n_g} (\mathbf{I} - \Theta_i) f_k - 2\mu}{n_g} \tag{23}$$

and

$$\alpha_i = \frac{1}{n_g} + \frac{\sum_{k=1}^{n_c} f_k^T \sum_{i=1}^{n_g} (\mathbf{I} - \Theta_i) f_k}{2n_g \mu} - \frac{\sum_{k=1}^{n_c} f_k^T (\mathbf{I} - \Theta_i) f_k}{2\mu}. \tag{24}$$

Since each of the steps above decreases the objective function $\Phi(\mathbf{F}, \alpha)$ which has a lower bound 0, the convergence of the alternating optimization is guaranteed.

IV. EXPERIMENTAL SETUPS

In this section, we present the experimental setups, including the databases, baseline methods, and the evaluation criteria. Although there are several publicly available datasets, the standard settings are not available, i.e., different articles use different experimental settings, such as the number of views, query objects, and the splitting method of the training and test datasets. Therefore, it is difficult to directly compare our method to the baseline methods by using the previously reported results. In this paper, we implement the baseline methods and evaluate them under the same setting with the optimized model parameters.

A. Experimental Settings

We conducted experiments on the NTU 3-D model dataset [33] and the ETH 3-D object dataset (ETH) [34]. In the NTU dataset, 500 objects were selected. We first selected 50 categories of 3-D objects, including *Aqua, Ball, Bed, Bike, Bird, Boat, Bomb, Book, Bottle, Car, Chair, Chip, Cube, Cup, Door, Driver, Drum, Facemask, Finger, Flower, Glasses, Guitar, Gun, Hat, Head, Helicopter, House, Knife, Lamp, Man, Map, Motorcycle, Ocd, Pen, Phone, Plane, Plant, Pot, Starship, Stick, Submarine, Sword, Table, Table-oneleg, Tank, Train, Tree, Truck, Weed, and Wheel*. The 3-D objects in the NTU database contain the model information. To compare the proposed method to baselines, we captured views from these objects to generate a view-based database. There were 20 cameras set at the vertices of a regular dodecahedron to produce 20 views. This view generation procedure increases the cost for processing. In the following experiments, we assume that these views have been obtained and the proposed method is conducted on the view-based database. The ETH database contains 80 objects that belong to eight categories. Each object has 41 views spaced evenly over the upper viewing hemisphere, and all the positions for cameras are determined by subdividing the faces of an octahedron to the third recursion level. Fig. 4 shows example views of the 3-D objects in the NTU and ETH databases.

Based on the common assumption that visually close objects have a high probability of being relevant, we use the Zernike moments [54] as the visual descriptors in our experiments. This feature has been widely applied to 3-D object retrieval and recognition in existing work [25], [33], [55] because of its robustness to image translation, scaling, and rotation. The Euclidean distance between two views x_i and x_j is $d(x_i, x_j)$. In 3-D object recognition experiments, the parameter ζ was empirically set to 0.9 (17) and μ was simply set to the number of views in the database.

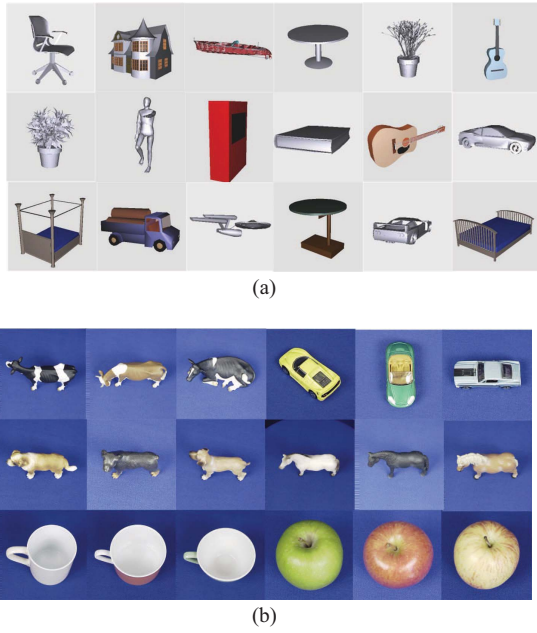


Fig. 4. Example views of 3-D objects in (a) NTU and (b) ETH datasets.

B. Evaluation Criteria

For retrieval, we adopt several performance evaluation measures that have been widely used in many articles [23], [25], [27], [33].

- 1) Precision–recall curve [25]: The precision–recall curve comprehensively demonstrates retrieval performance; it is assessed in terms of average recall and average precision, and has been widely used in multimedia applications [56], [57]. Here precision and recall are defined as follows:

Precision =

$$\frac{\text{No. of } \{(\text{relevant objects}) \cap (\text{retrieved objects})\}}{\text{No. of (retrieved objects)}}$$

Recall =

$$\frac{\text{No. of } \{(\text{relevant objects}) \cap (\text{retrieved objects})\}}{\text{No. of (relevant objects)}}$$

- 2) The precision of the first retrieved result (P1): This metric evaluates the retrieval accuracy of the first returned result.
- 3) First tier (FT): It is defined as the recall of the top τ results, i.e., $FT = (n_\tau/\tau)$, where τ is the number of relevant samples in the whole dataset, and n_τ is the number of relevant objects in the top τ retrieved results.
- 4) Second tier (ST): It is defined as the recall of the top 2τ results, i.e., $ST = (n_{2\tau}/\tau)$, where τ is the number of relevant samples in the whole dataset, and $n_{2\tau}$ is the number of relevant objects in the top 2τ retrieved results.
- 5) F-measure (F): In our experiments, we consider the top 20 retrieved results. Therefore, the F-measure is defined as $F = (2 \times P_{20} \times R_{20}/P_{20} + R_{20})$, where P_{20} and R_{20}

are the precision and the recall of the top 20 retrieval results.

- 6) Discounted cumulative gain (DCG) [58]: DCG is a statistic that assigns relevant results at the top ranking positions with higher weights under the assumption that a user is less likely to consider lower results.
- 7) Average normalized modified retrieval rank (ANMRR) [59]: ANMRR is a rank-based measure, and it considers the ranking information of relevant objects among the retrieved objects. A lower ANMRR value indicates a better performance, i.e., relevant objects rank at top positions.

Precision, recall, P1, FT, ST, F , DCG, and ANMRR, all range between 0 and 1.

For recognition, we use the recognition rate as our performance evaluation metric.

C. Compared Methods and Settings in 3-D Object Retrieval

We selected one object as query in the 3-D object retrieval experiments. We alternated the query objects until all the objects in the database were used once, and then we tested the average performance. The proposed method generated hypergraphs by setting K to different values which varied in a wide range 50, 100, 200, 400, 600, 1000, 1500, 2000, and 3000. We compared our hypergraph-based approach (denoted as “Hypergraph”) with the following methods that adopted different object distance measures.

- 1) HAUS employs the Hausdorff distance to estimate the distance between the query object and the objects in the database and then ranks the objects accordingly. The distance between two objects O_1 and O_2 can be written as

$$\begin{aligned} D(O_1, O_2) &= \max_{x_i \in O_1} \left\{ \min_{x_j \in O_2} d(x_i, x_j) \right\} \\ D(O_2, O_1) &= \max_{x_i \in O_2} \left\{ \min_{x_j \in O_1} d(x_i, x_j) \right\} \\ D_{\text{HAUS}} &= \max \{D(O_1, O_2), D(O_2, O_1)\} \end{aligned}$$

- 2) MEAN averages the distances of all view pairs across two objects. The objects in the database are ranked according to their distances to the query object. The distance can be written as

$$D_{\text{MEAN}} = \frac{1}{n_1 n_2} \sum_{x_i \in O_1} \sum_{x_j \in O_2} d(x_i, x_j) \quad (25)$$

where n_1 and n_2 are the numbers of views for O_1 and O_2 , respectively.

- 3) SumMin estimates the distance of a query view and the views of an object with minimum principle and then the results of all query views are accumulated. The distance can be written as

$$D_{\text{SumMin}} = \sum_{x_i \in O_1} \min_{x_j \in O_2} d(x_i, x_j). \quad (26)$$

Note that this distance is asymmetric, and here O_1 is the query object.

We implemented several existing methods for performance comparison as follows.

- 1) Adaptive views clustering (AVC) [25]: In AVC, 320 initial views are captured and representative views are optimally selected by AVC with Bayesian information criteria. A probabilistic method is then employed to calculate the similarity between two 3-D models, and those objects with high probability are selected as the retrieval results. There are two parameters in the method that are used to modulate the probabilities of objects and views, respectively. We tune these two parameters to their optimal values by a grid search in $[0, 1]$ with a granularity of 0.05.
- 2) The “Bag of visual features” (BoVF) method [60]: In BoVF, local SIFT features are extracted from all views of a 3-D object and they are quantized into visual words based on a pretrained visual vocabulary. According to BoVF, the visual vocabulary is trained by grouping all SIFT features in clusters by using the K -means clustering method. These extracted local features are further accumulated into a histogram as the representation of the 3-D object. The distance measurement between two 3-D objects is conducted by Kullback–Leibler divergence of the BoVF features. There are two parameters in the method, i.e., the size of the visual vocabulary, and the number of views used per object. We tune these two parameters to their optimal values by a grid search in $[1000, 2000]$ with a granularity of 200 and in $[6, 42]$ with a granularity of 4, respectively.
- 3) Compact multiview descriptor (CMVD) [23]: This method first selects 18 characteristic views of each 3-D object through 18 vertices of the corresponding bounding 32-hedron. Both the binary images and the depth images are taken to represent the views. The comparison between 3-D models is then accomplished by matching the selected views.
- 4) Elevation descriptor (ED) [27]: In ED, six range views are captured to describe the original 3-D object. These views contain the altitude information of the 3-D model from six directions. 3-D models are compared based on the matching of EDs.
- 5) Extension ray-based descriptor (ERD) [61]: ERD employs concentric spheres to extract the surface information of the 3-D model. Each sampling surface point provides the nearest sphere surface with a corresponding value, and the descriptor is the vector of the points in these concentric spheres’ surfaces. This method can describe the inside surface of 3-D models. The matching of 3-D models is performed by the distance of feature vectors.

We implemented the five methods above according to [23], [25], [27], [60], and [61], and they were tested under the same experimental settings as that used in our approach. It is worth noting that both ED and ERD are model-based methods, whereas AVC, CMVD, and BoVF are view-based methods. In these view-based methods, views of AVC and CMVD are generated from 3-D models. We conducted all

the five methods on the NTU dataset [33]. Since the ETH database [34] does not contain the 3-D model information, we only conducted the BoVF method.

D. Baseline Methods and Settings in 3-D Object Recognition

In 3-D object recognition, we implemented two classification methods for comparison, which are K -nearest neighbor (KNN) (with $k = 1$) and graph-based semisupervised learning (or GRAPH for short) introduced in [53]. Note that, in the subsection above, three object distance estimation methods were introduced, i.e., HAUS, MEAN, and SumMin. Therefore, for both KNN and GRAPH, we applied these three distance estimators separately, and we compared the proposed method with the following six methods, which were: 1) KNN + HAUS; 2) KNN + MEAN; 3) KNN + SumMin; 4) GRAPH + HAUS; 5) GRAPH + MEAN; and 6) GRAPH + SumMin.

In addition, we also implemented the following methods that were used in other articles for comparison.

- 1) Weighted subspace distance (WSD) [62]: WSD first generates the subspaces of training data, and the distance between each object in the database and an object category is calculated by WSD. The nearest neighbor method is then employed to accomplish object recognition.
- 2) Feature sharing and view clustering (FSVC) [63]: FSVC employs a scalable 3-D object representation scheme by FSVC. An automatic learning method with automatic feature detection, appearance library generation, and view clustering is then performed to achieve accurate 3-D object representation. Objects are recognized based on a hypothesis and verification framework by using the shared feature-based view clustering.
- 3) Pyramid matching (PM) [64]: In PM, range images are employed to represent 3-D objects. Salient points are detected and each salient point is associated with a surface descriptor. The pyramid match kernel is employed to measure the distance between two sets of features. The number of the selected points for feature extraction is tuned to its optimal value by grid search in $[100, 200]$ with a granularity of 20.

In the recognition experiments, we randomly split the database into a training set and a test set for 10 trials. We let the number of training samples per class increase from 2 to 8 in steps of 2. All the demonstrated results are the average recognition rates of the 10 trials. For the graph-based method, the radius parameter for similarity estimation and the regularization weight have been tuned to their optimal values by grid search (see [53]). In the grid search process for parameter selection, the radius parameter for similarity estimation was set to $r\bar{d}$, where \bar{d} was the average pairwise distance of the objects in the database. We varied r and the regularization weight in $\{0.01, 0.02, 0.05, 0.1, 0.2, 0.5, 1, 2, 5, 10\}$ and $\{0, 0.2, 0.4, 0.6, 0.8, 1\}$, respectively. The optimal values of the two parameters were obtained by jointly tuning them.

V. EXPERIMENTAL RESULTS

This section compares the different methods for 3-D object retrieval and recognition. Experimental results on retrieval

TABLE II

PERFORMANCE OF DIFFERENT METHODS FOR SEVERAL MEASURES ON THE NTU DATASET, WHERE THE VALUES IN () ARE THE GAINS ACHIEVED BY THE PROPOSED METHOD COMPARED WITH EACH OF THE OTHER METHODS IN THE TABLE IN TERMS OF %

	Hypergraph	HAUS	MEAN	SumMin	AVC	BoVF	CMVD	ED	ERD
P1	0.5939	0.4212 (29.1)	0.2697 (54.6)	0.5333 (10.2)	0.4242 (28.6)	0.4443 (25.2)	0.5461 (8.0)	0.2641 (55.5)	0.3263 (45.1)
FT	0.3051	0.2185 (28.4)	0.1684 (44.8)	0.2420 (20.7)	0.2171 (28.8)	0.2367 (22.4)	0.2730 (10.5)	0.1563 (48.8)	0.1860 (39.0)
ST	0.4175	0.3195 (23.5)	0.2680 (35.8)	0.3495 (16.3)	0.3114 (25.4)	0.3281 (21.4)	0.3528 (15.5)	0.2192 (47.5)	0.2572 (38.4)
F	0.3226	0.2630 (18.5)	0.2331 (27.7)	0.2834 (12.2)	0.2579 (20.1)	0.2713 (15.9)	0.3021 (6.4)	0.1653 (48.8)	0.1942 (39.8)
DCG	0.6951	0.6318 (9.1)	0.5883 (15.4)	0.6580 (5.3)	0.6264 (9.9)	0.6490 (6.6)	0.6685 (3.8)	0.4358 (37.3)	0.5147 (26.0)
ANMRR	0.5870	0.6685 (13.9)	0.7187 (22.4)	0.6428 (9.5)	0.6690 (14.0)	0.6492 (10.6)	0.6189 (5.4)	0.7351 (25.2)	0.7021 (19.6)

TABLE III

PERFORMANCE OF DIFFERENT METHODS FOR SEVERAL MEASURES ON THE ETH DATASET, WHERE THE VALUES IN () ARE THE GAINS ACHIEVED BY THE PROPOSED METHOD COMPARED WITH EACH OF THE OTHER METHODS IN THE TABLE IN TERMS OF %

	Hypergraph	HAUS	MEAN	SumMin	BoVF
P1	0.9125	0.7750 (15.1)	0.7625 (16.4)	0.8875 (2.7)	0.7625 (16.4)
FT	0.7097	0.6333 (10.8)	0.6750 (4.9)	0.7015 (1.2)	0.6472 (8.8)
ST	0.9194	0.8792 (4.4)	0.9001 (2.1)	0.9139 (0.6)	0.7861 (14.5)
F	0.6292	0.6003 (4.6)	0.6150 (2.3)	0.6242 (0.8)	0.5458 (13.3)
DCG	0.9267	0.8848 (4.5)	0.8918 (3.8)	0.9222 (0.5)	0.8774 (5.3)
ANMRR	0.2341	0.3047 (30.2)	0.2765 (18.1)	0.2422 (3.5)	0.2935 (25.4)

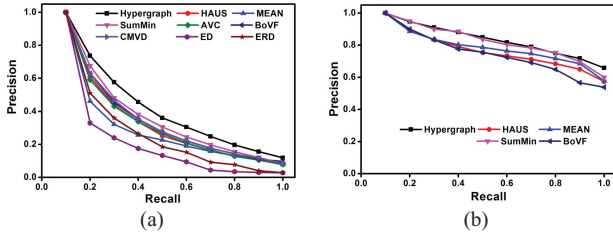


Fig. 5. Precision-recall curves of different methods on (a) NTU and (b) ETH datasets.

and recognition are demonstrated in Sections V-A and V-B, respectively. Comparison of the computational cost for different methods is provided in Section V-C. We implemented our method according to the settings used in other articles and directly compare our results with those reported in Section V-D. Due to the limited page length, we selected four methods: AVC [25], BoVF [60], CMVD [23], and WSD [62], as baselines.

A. Experiments on 3-D Object Retrieval

Fig. 5 shows the precision-recall curves on the NTU and ETH datasets. Tables II and III compare other performance measures on the two datasets, respectively. As shown in the figures and tables, the proposed approach outperforms all the other methods. The performance gains on the ETH dataset are less significant in comparison with those on the NTU dataset, which can be attributed to the fact that the retrieval performance on ETH is already fairly good (the P1 measurements are all above 0.7), and thus it leaves relatively less room for performance improvement.

We then observe the retrieval performance by using different individual hypergraphs. Figs. 6 and 7 illustrate the performance comparison between using the individual hypergraphs

generated with different K and using the fused hypergraph. We can see that the performance curves mainly exhibit a “wedge” shape when K varies. This is because the discriminative ability of a hypergraph will be weak when K is too small, and many objects will not be connected when K is too large. Although we simply employ an average fusion for the hypergraphs, we can see that the result can be better than using the best individual hypergraph in most cases. This makes our method robust and feasible because the retrieval performance is not highly sensitive to the parameter K .

B. Experiments on 3-D Object Recognition

We illustrate the comparison of the average object recognition performance of the proposed method and other compared methods on the NTU and ETH datasets in Fig. 8.

From the results we can see that more training samples can lead to better classification performance. The proposed hypergraph-based approach performs significantly better than other methods in most cases. The gains achieved by the proposed method compared with other methods in terms of % on the NTU dataset and the ETH dataset are provided in Tables IV and V, where n_t is the number of training sample per class.

On the ETH dataset, it performs closely to the GRAPH + SumMin method and is slightly worse when the number of training sample per class is set to 2. This is partially due to the fact that the improvement space is small for the ETH dataset in comparison with the NTU dataset. In addition, for the graph-based methods, the two parameters are set to their optimal values (this is infeasible in practical applications) and this also overestimates their performance.

C. Comparison of the Computational Cost

We first analyze the computational cost of the proposed method, and then empirically compare the computational

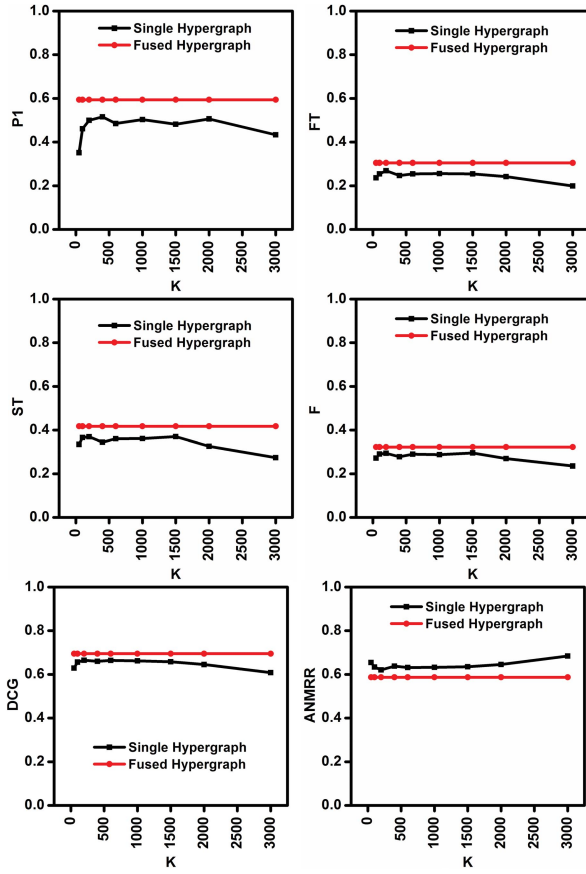


Fig. 6. Retrieval performance comparison using the individual hypergraphs generated with different K and using the fused hypergraph on the NTU dataset.

TABLE IV
GAINS ACHIEVED BY THE PROPOSED METHOD COMPARED WITH EACH OF THE OTHER METHODS IN THE TABLE IN TERMS OF % ON THE NTU DATASET

n_t	2	4	6	8
KNN + HAUS	50.0	52.0	38.8	36.4
KNN + MEAN	60.9	62.4	54.1	49.3
KNN + SumMin	36.2	43.2	28.2	22.7
GRAPH + HAUS	21.7	30.4	24.7	20.5
GRAPH + MEAN	47.8	48.0	43.5	43.2
GRAPH + SumMin	29.0	27.2	28.2	27.3
WSD	33.9	38.8	35.7	31.2
FSVC	26.3	25.9	22.7	21.4
PM	20.6	22.5	20.1	14.9

cost of the proposed method with those of the baseline methods.

According to the algorithm introduced in Section III, it can be obtained that the hypergraph construction process costs $O(n_g n_v \bar{K} T)$, where n_g is the number of the hypergraphs, n is the number of the objects in the database, n_v is the average number of views for all objects, \bar{K} is the average number of K in the K -means clustering, and T is the iteration number in the K -means clustering. The computational costs for retrieval and recognition are $O(n^2 T_b)$ and $O(T_a n^2 n_c T_b)$, respectively,

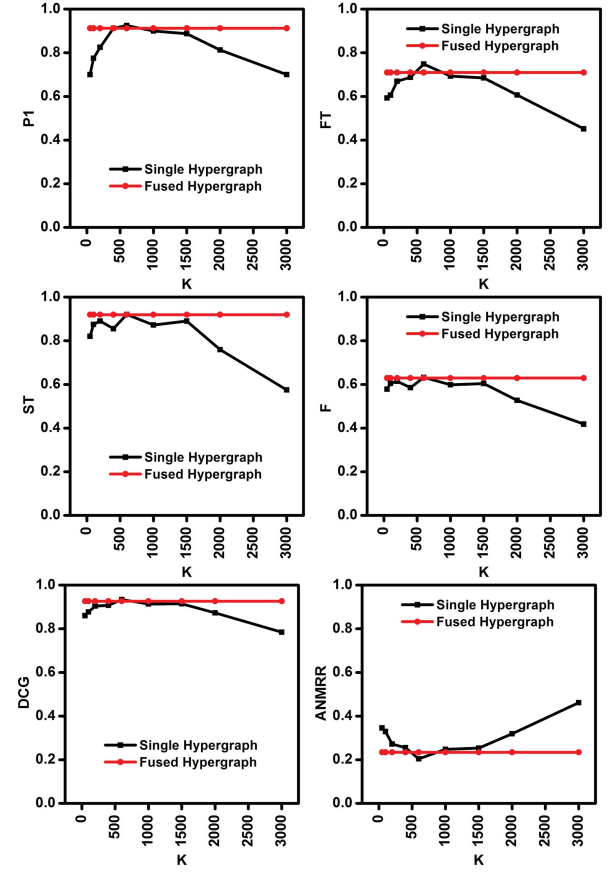


Fig. 7. Retrieval performance comparison using the individual hypergraphs generated with different K and using the fused hypergraph on the ETH dataset.

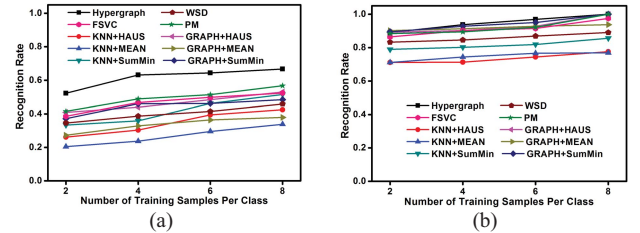


Fig. 8. Recognition performance of different methods on (a) NTU and (b) ETH datasets.

where T_a is the iteration number for the alternating optimization process, T_b is the iteration number of the iterative process in Algorithm 1, and n_c is the number of classes for classification. It is worth emphasizing that graph sparsification methods [37] can be adopted to reduce the computational costs.

We further compare different methods in terms of the retrieval and the recognition time. All experiments were conducted on a PC with Pentium 4 2.0-GHz CPU and 4-GB memory. The search time per query is shown in Fig. 9. As shown in this figure, the proposed method is slightly slower than most of the compared methods but its time is affordable in practice.

For 3-D recognition, the run time for recognition of all objects in the database with four training samples is shown in Fig. 10. As shown in this figure, the time of the proposed method is higher than that of KNN and GRAPH-

TABLE V
GAINS ACHIEVED BY THE PROPOSED METHOD COMPARED
WITH EACH OF THE OTHER METHODS IN THE TABLE IN
TERMS OF % ON THE ETH DATASET

n_t	2	4	6	8
KNN + HAUS	20.2	24.0	23.2	22.5
KNN + MEAN	20.2	20.7	21.0	23.1
KNN + SumMin	11.4	14.4	15.5	14.4
GRAPH + HAUS	-1.8	2.2	3.2	0.0
GRAPH + MEAN	-2.9	1.7	3.2	6.3
GRAPH + SumMin	-1.8	0.0	1.1	0.0
WSD	6.5	9.8	10.3	11.0
FSVC	3.0	3.8	5.6	2.5
PM	0.7	4.6	4.9	0.0

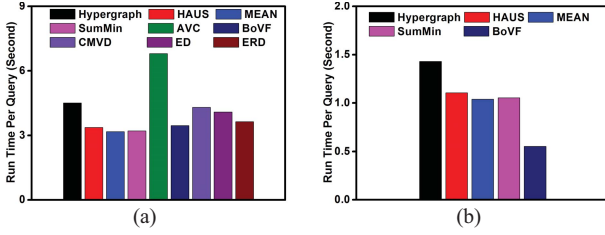


Fig. 9. Run time comparison of different methods in 3-D object retrieval on (a) NTU and (b) ETH datasets.

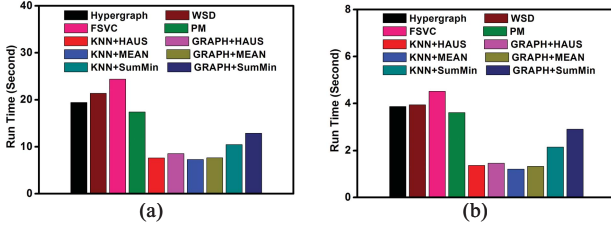


Fig. 10. Run time comparison of different methods in 3-D object recognition on (a) NTU and (b) ETH datasets.

based methods, but comparable to those of WSD, FSVC, and PM.

From the theoretical analysis, we can see that most of the computational cost lies in the hypergraph construction procedure. The proposed method can be further speeded up using preprocessing of the employed database. A hypergraph structure of the whole database can be constructed offline. When a query object is provided, the query object is further embedded in the hypergraph structure, and the retrieval process can be conducted directly without the complex hypergraph construction procedure. The online hypergraph construction procedure can be reduced to $O(n_g n_v \bar{K})$, which can be applied for large-scale search and recognition.

D. Comparison With the Reported Results of State-of-the-Art Methods

We implemented our approach under the settings of existing work and directly compared its performance with the reported results. For these experiments, the specific settings can be found in the related references.

For 3-D object retrieval, we first compare our approach with the AVC method on the Princeton shape database. Fig. 11

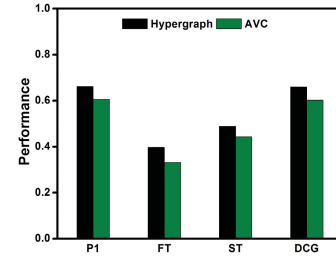


Fig. 11. Experimental results of the proposed method and the AVC method obtained under the settings of [25].

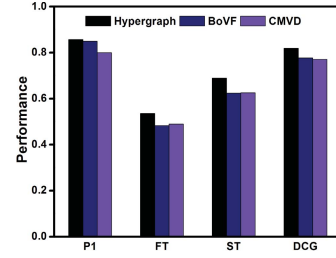


Fig. 12. Experimental results of BoVF, CMVD, and the proposed approach obtained under the settings of [65].

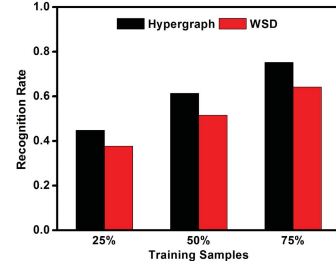


Fig. 13. Classification accuracies of our approach and WSD under the settings of [62].

compares the proposed method with AVC (probability) [25] on 3-D object retrieval under the same experimental setting used in [25] where the results obtained by AVC (probability) are directly taken from [25, Table I]. Experimental results suggest the proposed methods outperforms the baseline methods in terms of the four reported evaluation measures.

We then compare our approach with the BoVF and the CMVD methods on the SHREC09 database. Fig. 12 compares the proposed method with BoVF [60] and CMVD [23] on 3-D object retrieval under the same experimental setting used in [65], where the results obtained by BoVF and CMVD are directly taken from [65, Table II]. We also see the superiority of our approach.

For 3-D object recognition, we compare our method with WSD on the ALOI database [62]. We followed the settings described in [62]. In [62], the classification accuracies of WSD have been reported when the ratio of training data is 25%, 50%, and 75%. Fig. 13 compares the proposed method with WSD [62] on 3-D object recognition under the same experimental setting used in [62], where the results obtained by WSD are directly taken from [62, Table II]. From the comparison, we can observe the consistent superiority of the proposed method.

VI. CONCLUSION

In this paper, we proposed a hypergraph analysis method by integrating multiple hypergraphs for 3-D object retrieval and recognition. The proposed method groups the views of 3-D objects into clusters based on their visual descriptions. Hypergraphs were then constructed, where each vertex was an object and an edge was a cluster of views. Therefore, an edge connects multiple objects in the hypergraphs. By varying the number of view clusters, we could generate multiple hypergraphs that captured the higher order relationship of the 3-D objects at different granularities. Retrieval and recognition were formulated as learning tasks with different regularization frameworks on the hypergraphs.

To test the performance of the proposed approach, we conducted experiments on the NTU and ETH datasets by using several evaluation criteria for performance evaluation. Experimental results showed that our method achieved better results compared with many baseline methods, including multiple-view matching with common distance measures and several existing algorithms used by others, such as ED, ERD, AVC, CMVD, BoVF, WSD, FS, and VC, and PM. The superiority is remarkable on the NTU dataset. The proposed method only slightly outperformed existing top-level schemes for 3-D object retrieval and recognition on the ETH dataset, because the retrieval and recognition performance on the ETH dataset is already very high and the room for performance improvement is limited. We also compared our method, which integrates multiple hypergraphs, to that of using only individual hypergraphs. The results show that in most cases our approach is better than the best result obtained by using individual hypergraphs. This demonstrates that our method is robust and feasible because the retrieval performance is not highly sensitive to the view clustering results.

In summary, the proposed method has the following three advantages.

- 1) The proposed method avoids the distance estimation of 3-D objects because we only need to analyze the relationship of different view groups.
- 2) It explores the higher order relationship among 3-D objects. Here, the higher order relationship of objects is encoded in the hypergraph structure, i.e., the connection of each edge to multiple vertices.
- 3) The proposed learning on the hypergraph structure is essentially a particular implementation under the umbrella of the graph-based semisupervised learning. Thus, the proposed method enjoys all the advantages of graph-based semisupervised learning. The proposed method uses the unlabeled data to construct the hypergraph regularization, where the data distribution of the whole database revealed by the unlabeled data is taken into consideration. This hypergraph regularization makes the proposed method capable of making use of the unlabeled data to boost the performance of supervised learning. Therefore, the proposed method is robust even with very few labeled examples for 3-D object recognition.

Our method does not need any information from 3-D models. This makes our approach flexible and useful in

cases when 3-D models are not available. Although a 3-D model can be reconstructed from a set of 2-D views, it is computationally expensive. Therefore, directly using the 2-D view has specific benefits in reducing the computational cost in practical applications.

Our current method has some limitations. First, from the discussion in Section V-C, we can see that the computational cost scales as the square of the number of objects, and it makes our method difficult in handling large-scale datasets. Second, we need a strategy to handle newly added objects. When new objects are added into the database, we need to reconstruct hypergraphs to perform our method, which is time consuming. A method to efficiently update the constructed hypergraphs after adding new samples is required. Third, our approach is a transductive method for the 3-D object classification task. This means, it directly classifies a set of objects without yielding a classification model. Therefore, it cannot be applied to the cases when classification models are explicitly required for applications. We will further explore these problems and improve our method in the future.

ACKNOWLEDGMENT

The authors would like to thank the Handling Associate Editor S. Pankanti and the anonymous reviewers for their valuable and constructive comments. The authors would also like to thank B. Munday and S. Felix for proofreading this paper.

REFERENCES

- [1] A. Bimbo and P. Pala, "Content-based retrieval of 3D models," *ACM Trans. Multimedia Comput., Commun., Appl.*, vol. 2, no. 1, pp. 20–43, 2006.
- [2] B. Bustos, D. Keim, D. Saupe, T. Schreck, and D. Vranic, "Feature-based similarity search in 3D object databases," *ACM Comput. Surv.*, vol. 37, no. 4, pp. 345–387, 2005.
- [3] J. W. H. Tangelder and R. C. Veltkamp, "A survey of content based 3D shape retrieval methods," *Multimedia Tools Appl.*, vol. 39, no. 3, pp. 441–471, Sep. 2008.
- [4] Y. Yang, H. Lin, and Y. Zhang, "Content-based 3D model retrieval: A survey," *IEEE Trans. Syst., Man, Cybern., Part C, Appl. Rev.*, vol. 37, no. 6, pp. 1081–1035, Nov. 2007.
- [5] B. Leng and Z. Xiong, "Modelseek: An effective 3D model retrieval system," *Multimedia Tools Appl.*, vol. 51, no. 3, pp. 935–962, 2011.
- [6] Q. Xiao, H. Wang, F. Li, and Y. Gao, "3D object retrieval based on a graph model descriptor," *Neurocomputing*, vol. 74, no. 17, pp. 3486–3493, 2011.
- [7] S. Jayanti, K. Kalyanaraman, N. Iyer, and K. Ramani, "Developing an engineering shape benchmark for CAD models," *Comput.-Aided Design*, vol. 38, no. 9, pp. 939–953, 2006.
- [8] P. Daras, D. Zarpalas, A. Axenopoulos, D. Tzovaras, and M. Strintzis, "Three-dimensional shape-structure comparison method for protein classification," *IEEE/ACM Trans. Comput. Biol. Bioinf.*, vol. 3, no. 3, pp. 193–207, Jul. 2006.
- [9] R. J. Campbell and P. J. Flynn, "A survey of free-form object representation and recognition techniques," *Comput. Vis. Image Understand.*, vol. 81, no. 2, pp. 166–210, Feb. 2001.
- [10] A. E. Johnson and M. Hebert, "Using spin images for efficient object recognition in cluttered 3D scenes," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 5, pp. 433–449, May 1999.
- [11] R. Osada, T. Funkhouser, B. Chazelle, and D. Dobkin, "Shape distributions," *ACM Trans. Graphic*, vol. 21, no. 4, pp. 807–832, 2002.
- [12] E. Paquet and M. Rioux, "Nefertiti: A query by content system for three-dimensional model and image databases management," *Image Vis. Comput.*, vol. 17, no. 2, pp. 157–166, Feb. 1999.

- [13] B. Leng and Z. Qin, "A powerful relevance feedback mechanism for content-based 3D model retrieval," *Multimedia Tools Appl.*, vol. 40, no. 1, pp. 135–150, 2008.
- [14] Y. Gao, Q. H. Dai, and N. Y. Zhang, "3D model comparison using spatial structure circular descriptor," *Pattern Recognit.*, vol. 43, no. 3, pp. 1142–1151, 2010.
- [15] W. Li, G. Bebis, and N. Bourbakis, "3-D object recognition using 2-D views," *IEEE Trans. Image Process.*, vol. 17, no. 11, pp. 2236–2255, Nov. 2008.
- [16] Y. Gao, J. Tang, H. Li, Q. Dai, and N. Zhang, "View-based 3D model retrieval with probabilistic graph model," *Neurocomputing*, vol. 73, nos. 10–12, pp. 1900–1905, 2010.
- [17] T. Chen, M.-M. Cheng, P. Tan, A. Shamir, and S.-M. Hu, "Sketch2photo: Internet image montage," *ACM Trans. Graph.*, vol. 28, no. 5, pp. 124–133, 2009.
- [18] M.-M. Cheng, G.-X. Zhang, N. J. Mitra, X. Huang, and S.-M. Hu, "Global contrast based salient region detection," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2011, pp. 409–416.
- [19] R. Ji, H. Yao, X. Sun, B. Zhong, P. Xu, and W. Gao, "Toward semantic embedding in visual vocabulary," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 918–925.
- [20] R. Ji, H. Yao, and X. Sun, "Actor-independent action search using spatiotemporal vocabulary with appearance hashing," *Pattern Recognit.*, vol. 44, no. 3, pp. 624–638, 2011.
- [21] N. Guan, D. Tao, Z. Luo, and B. Yuan, "Non-negative patch alignment framework," *IEEE Trans. Neural Netw.*, vol. 22, no. 8, pp. 1218–1230, Aug. 2011.
- [22] D. Tao, X. Tang, X. Li, and X. Wu, "Asymmetric bagging and random subspace for support vector machines-based relevance feedback in image retrieval," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 28, no. 7, pp. 1088–1099, Jul. 2006.
- [23] P. Daras and A. Axenopoulos, "A 3D shape retrieval framework supporting multimodal queries," *Int. J. Comput. Vis.*, vol. 89, no. 2, pp. 229–247, 2010.
- [24] P. Shilane, P. Min, M. Kazhdan, and T. Funkhouser, "The Princeton shape benchmark," in *Proc. Shape Model. Int.*, 2004, pp. 167–178.
- [25] T. F. Ansary, M. Daoudi, and J. P. Vandeborre, "A Bayesian 3-D search engine using adaptive views clustering," *IEEE Trans. Multimedia*, vol. 9, no. 1, pp. 78–88, Jan. 2007.
- [26] Y. Gao, M. Wang, Z. Zha, Q. Tian, Q. Dai, and N. Zhang, "Less is more: Efficient 3-D object retrieval with query view selection," *IEEE Trans. Multimedia*, vol. 11, no. 5, pp. 1007–1018, Oct. 2011.
- [27] J. L. Shih, C. H. Lee, and J. T. Wang, "A new 3D model retrieval approach based on the elevation descriptor," *Pattern Recognit.*, vol. 40, no. 1, pp. 283–295, Jan. 2007.
- [28] M. J. Atallah, "A linear time algorithm for the Hausdorff distance between convex polygons," *Inf. Process. Lett.*, vol. 17, no. 4, pp. 207–209, 1983.
- [29] M. P. Dubuisson and A. K. Jain, "Modified Hausdorff distance for object matching," in *Proc. IAPR Int. Conf. Pattern Recognit.*, 1994, pp. 566–568.
- [30] Y. Rubner, C. Tomasi, and L. J. Guibas, "The earth mover's distance as a metric for image retrieval," *Int. J. Comput. Vis.*, vol. 40, no. 2, pp. 99–121, 2000.
- [31] Y. Gao, Q. Dai, M. Wang, and N. Zhang, "3D model retrieval using weighted bipartite graph matching," *Signal Process.: Image Commun.*, vol. 26, no. 1, pp. 39–47, 2011.
- [32] Y. Gao, J. Tang, R. Hong, S. Yan, Q. Dai, N. Zhang, and T. Chua, "Camera constraint-free view-based 3-D object retrieval," *IEEE Trans. Image Process.*, vol. 21, no. 4, pp. 2269–2281, Apr. 2012.
- [33] D. Y. Chen, X. P. Tian, Y. T. Shen, and M. Ouhyoung, "On visual similarity based 3D model retrieval," *Comput. Graph. Forum*, vol. 22, no. 3, pp. 223–232, 2003.
- [34] B. Leibe and B. Schiele, "Analyzing appearance and contour based methods for object categorization," in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2003, pp. 409–415.
- [35] Y. Huang, Q. Liu, and D. Metaxas, "Video object segmentation by hypergraph cut," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 1738–1745.
- [36] Y. Huang, Q. Liu, S. Zhang, and D. Metaxas, "Image retrieval via probabilistic hypergraph ranking," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 3376–3383.
- [37] D. Zhou, J. Huang, and B. Schokopf, "Learning with hypergraphs: Clustering, classification, and embedding," in *Proc. Adv. Neural Inf. Process. Syst.*, 2007, pp. 1601–1608.
- [38] M. Wang, X.-S. Hua, R. Hong, J. Tang, G.-J. Qi, and Y. Song, "Unified video annotation via multigraph learning," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 19, no. 5, pp. 733–746, May 2009.
- [39] X. Zhu, "Semi-supervised learning with graphs," Ph.D. dissertation, School Comput. Sci., Carnegie Mellon Univ., Pittsburgh, PA, 2005.
- [40] M. Wang, K. Yang, X.-S. Hua, and H.-J. Zhang, "Toward a relevant and diverse search of social images," *IEEE Trans. Multimedia*, vol. 12, no. 8, pp. 829–842, Dec. 2010.
- [41] H. Chen and B. Bhanu, "Efficient recognition of highly similar 3D objects in range images," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 1, pp. 172–179, Jan. 2009.
- [42] S. Xia and E. Hancock, *Clustering Using Class Specific Hyper Graphs* (Lecture Notes in Computer Science), vol. 5342. New York: Springer-Verlag, 2008, pp. 318–328.
- [43] R. Zass and A. Shashua, "Probabilistic graph and hypergraph matching," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2008, pp. 1–8.
- [44] Y. Liu, X.-L. Wang, H.-Y. Wang, H. Zha, and H. Qin, "Learning robust similarity measures for 3D partial shape retrieval," *Int. J. Comput. Vis.*, vol. 89, no. 2, pp. 408–431, 2010.
- [45] Y. Gao, M. Wang, S. Yan, J. Shen, and D. Tao, "Tag-based social image search with visual-text joint hypergraph learning," in *Proc. Conf. Multimedia*, 2011, pp. 1517–1520.
- [46] A. K. C. Wong and S. W. Lu, "Recognition and shape synthesis of 3-D objects based on attributed hypergraphs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 11, no. 3, pp. 279–290, Mar. 1989.
- [47] S. Xia and E. Hancock, *3D Object Recognition Using Hyper-Graphs and Ranked Local Invariant Features* (Lecture Notes in Computer Science), vol. 5342. New York: Springer-Verlag, 2008, pp. 117–126.
- [48] S. Xia and E. Hancock, "Learning large scale class specific hyper graphs for object recognition," in *Proc. Int. Conf. Image Graph.*, 2008, pp. 366–371.
- [49] S. P. Lloyd, "Least square quantization in PCM," *IEEE Trans. Inf. Theory*, vol. 28, no. 2, pp. 129–137, Mar. 1957.
- [50] N. Guan, D. Tao, Z. Luo, and B. Yuan, "Online non-negative matrix factorization with robust stochastic approximation," *IEEE Trans. Neural Netw. Learn. Syst.*, 2012, DOI: 10.1109/TNNLS.2012.2197827.
- [51] N. Guan, D. Tao, Z. Luo, and B. Yuan, "NeNMF: An optimal gradient method for nonnegative matrix factorization," *IEEE Trans. Signal Process.*, vol. 60, no. 6, pp. 2882–2898, Jun. 2012.
- [52] N. Guan, D. Tao, Z. Luo, and B. Yuan, "Manifold regularized discriminative nonnegative matrix factorization with fast gradient descent," *IEEE Trans. Image Process.*, vol. 20, no. 7, pp. 2030–2048, Jul. 2011.
- [53] D. Zhou, O. Bousquet, T. N. Lal, J. Weston, and B. Schokopf, "Learning with local and global consistency," in *Proc. Adv. Neural Inf. Process. Syst.*, 2004, pp. 321–328.
- [54] A. Khotanzad and Y. H. Hong, "Invariant image recognition by Zernike moments," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 12, no. 5, pp. 489–497, May 1990.
- [55] W. Y. Kim and Y. S. Kim, "A region-based shape descriptor using Zernike moments," *Signal Process.: Image Commun.*, vol. 16, nos. 1–2, pp. 95–102, 2000.
- [56] R. Ji, L.-Y. Duan, J. Chen, H. Yao, J. Yuan, Y. Rui, and W. Gao, "Location discriminative vocabulary coding for mobile landmark search," *Int. J. Comput. Vis.*, vol. 96, no. 3, pp. 290–314, Feb. 2012.
- [57] M. Wang, X.-S. Hua, J. Tang, and R. Hong, "Beyond distance measurement: Constructing neighborhood similarity for video annotation," *IEEE Trans. Multimedia*, vol. 11, no. 3, pp. 465–476, Apr. 2009.
- [58] K. Jarvelin and J. Kekalainen, "Cumulated gain-based evaluation of IR techniques," *ACM Trans. Inf. Syst.*, vol. 20, no. 4, pp. 422–446, 2002.
- [59] *Description of Core Experiments for MPEG-7 Color/Texture Descriptors*, Standard ISO/MPEGJTC1/SC29/WG11 MPEG98/M2819, 1999.
- [60] R. Ohbuchi, K. Osada, T. Furuya, and T. Banno, "Salient local visual features for shape-based 3D model retrieval," in *Proc. IEEE Conf. Shape Model. Appl.*, Jun. 2008, pp. 93–102.
- [61] D. Vranic, "An improvement of rotation invariant 3D-shape based on functions on concentric spheres," in *Proc. IEEE Int. Conf. Image Process.*, Sep. 2003, pp. 757–760.

- [62] F. Li, Q. Dai, W. Xu, and G. Er, "Weighted subspace distance and its applications to object recognition and retrieval with image sets," *IEEE Signal Process. Lett.*, vol. 16, no. 3, pp. 227–230, Mar. 2009.
- [63] S. Kim and I. Kweon, "Scalable representation for 3D object recognition using feature sharing and view clustering," *Pattern Recognit.*, vol. 41, no. 2, pp. 754–773, 2008.
- [64] X. Li and I. Guskov, "3D object recognition from range images using pyramid matching," in *Proc. IEEE 11th Int. Conf. Comput. Vis.*, Oct. 2007, pp. 1–6.
- [65] A. Godil, H. Dutagaci, C. Akgul, A. Axenopoulos, B. Bustos, M. Chaouch, P. Daras, T. Furuya, S. Kreft, Z. Lian, T. Napoleon, A. Mademlis, R. Ohbuchi, P. L. Rosin, B. Sankur, T. Schreck, X. Sun, M. Tezuka, A. Verroust-Blondet, M. Walter, and Y. Yemez, "Shrec'09 track: Generic shape retrieval," in *Proc. Eurograph. Workshop 3D Object Retr.*, Munich, Germany, 2009, pp. 61–68.



Yue Gao received the B.S. degree from the Harbin Institute of Technology, Harbin, China, in 2005, and the M.E. and Ph.D. degrees from Tsinghua University, Beijing, China, in 2008 and 2012, respectively.

He is currently with the Department of Automation, Tsinghua National Laboratory for Information Science and Technology, Tsinghua University. He was a Visiting Scholar with Carnegie Mellon University, Pittsburgh, PA, where he worked with A. Hauptmann from October 2010 to March 2011. He was a Research Intern with the National University

of Singapore, Singapore, and with the Intel China Research Center, Beijing. His current research interests include multimedia information retrieval, 3-D object retrieval and recognition, and social media analysis.



Meng Wang (M'09) received the B.E. degree from the Department of Electronic Engineering and Information Science, and the Ph.D. degree in the Special Class for the Gifted Young, Department of Electronic Engineering and Information Science, University of Science and Technology of China, Hefei, China.

He is currently a Professor with the Hefei University of Technology, Hefei. He was an Associate Researcher with Microsoft Research Asia, Beijing, China, and was a Core Member in a startup in Silicon Valley. He was also with the National University of Singapore, Singapore, as a Senior Research Fellow. He has authored more than 100 articles published in books, journals, and conference proceedings in his areas of expertise. His current research interests include multimedia content analysis, search, mining, recommendation, and large-scale computing.

Prof. Wang was a recipient of Best Paper Awards at both the 17th and 18th ACM International Conference on Multimedia and the Best Paper Award at the 16th International Multimedia Modeling Conference. He is a member of the ACM.



Dacheng Tao (M'07–SM'12) is a Professor of computer science with the Centre for Quantum Computation and Information Systems, and the Faculty of Engineering and Information Technology, University of Technology, Sydney, Australia. He applies statistics and mathematics for data analysis problems in data mining, computer vision, machine learning, multimedia, and video surveillance. He has authored or co-authored more than 100 scientific articles presented at top venues, including IEEE T-PAMI, T-IP, AISTATS, ICDM, CVPR, and ECCV.

Prof. Tao was a recipient of the Best Theory and Algorithm Paper Runner Up Award at IEEE ICDM'07.



Rongrong Ji (M'12) received the Ph.D. degree in computer science from the Harbin Institute of Technology, Harbin, China.

He has been a Post-Doctoral Research Fellow with Columbia University, New York, NY, since 2011, working with S.-F. Chang. From March 2010 to May 2010, he was a Visiting Student with the University of Texas at San Antonio, San Antonio, working with Q. Tian. From April 2010 to November 2010, he was a Research Assistant with Peking University, Beijing, China, with W. Gao. He was a Research Intern

with Microsoft Research Asia, Beijing, working with X. Xie. From 2007 to 2010, he led the Multimedia Retrieval Group, Visual Intelligence Laboratory, Harbin Institute of Technology. He is the author of over 40 refereed papers published in journals and conference proceedings, such as IJCV, TIP, TMM, TOMCCAP, IEEE Multimedia, PR, CVPR, ACM Multimedia, IJCAI, and AAAI. His current research interests include image and video search, content understanding, mobile visual search and recognition, and interactive human-computer interfaces.

Dr. Ji was a recipient of the Best Paper Award at the ACM Multimedia in 2011 and the Microsoft Fellowship in 2007. He is an Associate Editor for the *International Journal of Computer Applications*, and a Guest Editor of the *International Journal of Advanced Computer Science and Applications*. He was a Session Chair of the ICME in 2008 and has been on the ACM Multimedia Program Committee since 2011.



Qionghai Dai (SM'05) received the B.S. degree in mathematics from Shanxi Normal University, Shanxi, China, in 1987, and the M.E. and Ph.D. degrees in computer science and automation from Northeastern University, Beijing, China, in 1994 and 1996, respectively.

He has been with the Faculty of Tsinghua University, Beijing, China, since 1997, and is currently a Professor and the Director of the Broadband Networks and Digital Media Laboratory. His current research interests include signal processing, broad-

band networks, video processing, and communication.