

# GIFT: A Real-time and Scalable 3D Shape Search Engine

Song Bai<sup>1</sup> Xiang Bai<sup>1\*</sup> Zhichao Zhou<sup>1</sup> Zhaoxiang Zhang<sup>2</sup> Longin Jan Latecki<sup>3</sup>

<sup>1</sup>Huazhong University of Science and Technology, <sup>3</sup>Temple University

<sup>2</sup>CAS Center for Excellence in Brain Science and Intelligence Technology, CASIA

{songbai,xbai,zzc}@hust.edu.cn zhaoxiang.zhang@ia.ac.cn latecki@temple.edu

## Abstract

*Projective analysis is an important solution for 3D shape retrieval, since human visual perceptions of 3D shapes rely on various 2D observations from different view points. Although multiple informative and discriminative views are utilized, most projection-based retrieval systems suffer from heavy computational cost, thus cannot satisfy the basic requirement of scalability for search engines.*

*In this paper, we present a real-time 3D shape search engine based on the projective images of 3D shapes. The real-time property of our search engine results from the following aspects: (1) efficient projection and view feature extraction using GPU acceleration; (2) the first inverted file, referred as F-IF, is utilized to speed up the procedure of multi-view matching; (3) the second inverted file (S-IF), which captures a local distribution of 3D shapes in the feature manifold, is adopted for efficient context-based re-ranking. As a result, for each query the retrieval task can be finished within one second despite the necessary cost of IO overhead. We name the proposed 3D shape search engine, which combines GPU acceleration and Inverted File (Twice), as GIFT. Besides its high efficiency, GIFT also outperforms the state-of-the-art methods significantly in retrieval accuracy on various shape benchmarks and competitions.*

## 1. Introduction

3D shape retrieval is a fundamental issue in computer vision and pattern recognition. With the rapid development of large scale public 3D repositories, *e.g.*, Google 3D Warehouse or TurboSquid, and large scale shape benchmarks, *e.g.*, ModelNet [39], SHape REtrieval Contest (SHREC) [14, 31], the scalability of 3D shape retrieval algorithms becomes increasingly important for practical applications. However, *efficiency* issue has been more or less ignored by previous works, though enormous efforts have been devoted to retrieval *effectiveness*, that is to say, to design informative and discriminative features [12, 2, 17, 6,

40, 15, 18] to boost the retrieval accuracy. As suggested in [14], plenty of these algorithms do not scale up to large 3D shape databases due to their high time complexity.

Meanwhile, owing to the fact that human visual perception of 3D shapes depends upon 2D observations, projective analysis has become a basic and inherent tool in 3D shape domain for a long time, with applications to segmentation [38], matching [24], reconstruction, *etc.*. Specifically in 3D shape retrieval, projection-based methods demonstrate impressive performances. Especially in recent years, the success of planar image representation [7, 35, 43], makes it easier to describe 3D models using depth or silhouette projections.

Generally, a typical 3D shape search engine is comprised of the following four components (see also Fig. 1):

1. *Projection rendering.* With a 3D model as input, the output of this component is a collection of projections. Most methods set an array of virtual cameras at pre-defined view points to capture views. These view points can be the vertices of a dodecahedron [4], located on the unit sphere [35], or around the lateral surface of a cylinder [24]. In most cases, pose normalization [22] is needed for the sake of invariance to translation, rotation and scale changes.
2. *View feature extraction.* The role of this component is to obtain multiple view representations, which affects the retrieval quality largely. A widely-used paradigm is Bag-of-Words (BoW) [7] model, since it has shown its superiority as natural image descriptors. However, in order to get better performances, many features [14] are of extremely high dimension. As a consequence, raw descriptor extraction (*e.g.*, SIFT [20]), quantization and distance calculation are all time-consuming.
3. *Multi-view matching.* This component establishes the correspondence between two sets of view features, and returns a matching cost between two 3D models. Since at least a set-to-set matching strategy [25, 26, 27, 16, 9] is required, this stage suffers from high time complexity even when using the simplest Hausdorff matching.

\*Corresponding author

Hence, the usage of algorithms incorporated with some more sophisticated matching strategies on large scale 3D datasets is limited due to their heavy computational cost.

4. *Re-ranking*. It aims at refining the initial ranking list by using some extra information. For retrieval problems, since no prior or supervised information is available, contextual similarity measure is usually utilized. A classic context-based re-ranking methodology for shape retrieval is diffusion process [5], which exhibits outstanding performance on various datasets. However, as graph-based and iterative algorithms, many variants of diffusion process (*e.g.*, locally constrained diffusion process [41]), generally require the computational complexity of  $O(TN^3)$ , where  $N$  is the total number of shapes in the database and  $T$  is the number of iterations. In this sense, diffusion process does not seem to be applicable for real-time analysis.

In this paper, we present a real-time 3D shape search engine using projections that includes all the aforementioned components. It combines Graphics Processing Unit (GPU) acceleration and Inverted File (Twice), hence we name it **GIFT**. In on-line processing, once a user submits a query shape, GIFT can react and present the retrieved shapes within one second (the off-line preprocessing operations, such as CNN model training and inverted file establishment, are excluded). GIFT is evaluated on several popular 3D benchmarks datasets, especially on one track of SHape REtrieval Contest (SHREC) which focuses on scalable 3D retrieval. The experimental results on retrieval accuracy and query time demonstrate the capability of GIFT in handling large scale data.

In summary, our main contributions are as follows. Firstly, GPU is used to speed up the procedure of projection rendering and feature extraction. Secondly, in multi-view matching procedure, a robust version of Hausdorff distance for noise data is approximated with an inverted file, which allows for extremely efficient matching between two view sets without impairing the retrieval performances too much. Thirdly, in the re-ranking component, a new context-based algorithm based on fuzzy set theory is proposed. Different from diffusion processes of high time complexity, our re-ranking here is ultra time efficient on the account of using inverted file again.

## 2. Proposed Search Engine

### 2.1. Projection Rendering

Prior to projection rendering, pose normalization for each 3D shape is needed in order to attain invariance to some common geometrical transformations. However, apart from many pervious algorithms [23, 24, 22] that re-

quire rotation normalization using some Principal Component Analysis (PCA) techniques, we only normalize the scale and the translation in our system. Our concerns are two-fold: 1) PCA techniques are not always stable, especially when dealing with some specific geometrical characteristics such as symmetries, large planar or bumpy surfaces; 2) the view feature used in our system can tolerate the rotation issue to a certain extent, though cannot be completely invariant to such changes. In fact, we observe that if enough projections (more than 25 in our experiments) are used, one can achieve reliable performances.

The projection procedure is as follows. Firstly, we place the centroid of each 3D shape at the origin of a spherical coordinate system, and resize the maximum polar distance of the points on the surface of the shape to unit length. Then  $N_v$  virtual cameras are set on the unit sphere evenly, and they are located by the azimuth  $\theta_{az}$  and the elevation  $\theta_{el}$  angles. At last, we render one projected view in depth buffer at each combination of  $\theta_{az}$  and  $\theta_{el}$ . For the sake of speed, GPU is utilized here such that for each 3D shape, the average time cost of rendering 64 projections is only 30ms.

### 2.2. Feature Extraction via GPU Acceleration

Feature design has been a crucial problem in 3D shape retrieval for a long time owing to its great influence on the retrieval accuracy. Though extensively studied, almost all the existing algorithms ignore the efficiency of the feature extraction.

To this end, our search engine adopts GPU to accelerate the procedure of feature extraction. Impressed by the superior performance of deep learning approaches in various visual tasks, we propose to use the activation of a Convolutional Neural Network (CNN). The CNN used here takes depth images as input, and the loss function is exerted on the classification error for projections. The network architecture consists of five successive convolutional layers and three fully connected layers as in [3]. We normalize each activation in its Euclidean norm to avoid scale changes. It only takes 56ms on average to extract the view features for a 3D model.

Since no prior information is available to judge the discriminative power of activations of different layers, we propose a robust re-ranking algorithm described in Sec. 2.4. It can fuse those homogenous features efficiently based on fuzzy set theory.

### 2.3. Inverted File for Multi-view Matching

Consider a query shape  $x_q$  and a shape  $x_p$  from the database  $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ . Let  $\mathcal{V}$  denote a mapping function from 3D shapes to their feature sets. We can obtain two sets  $\mathcal{V}(x_q) = \{q_1, q_2, \dots, q_{N_v}\}$  and  $\mathcal{V}(x_p) = \{p_1, p_2, \dots, p_{N_v}\}$  respectively, where  $N_v$  is the number of views.  $q_i$  (or  $p_i$ ) denotes the view feature assigned to the

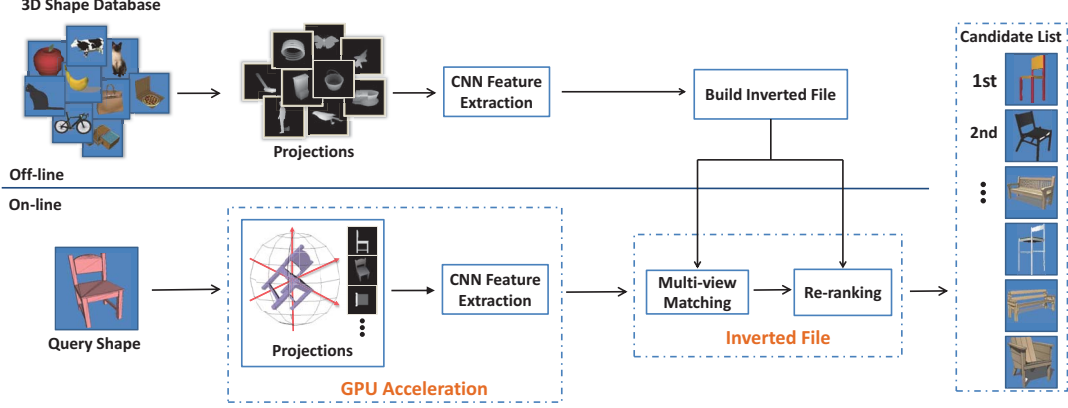


Figure 1. The structure of the proposed 3D shape search engine **GIFT**.

$i$ -th view of shape  $x_q$  (or  $x_p$ ).

A 3D shape search engine requires a multi-view matching component to establish a correspondence between two sets of view features. These matching strategies are usually metrics defined on sets (*e.g.*, Hausdorff distance) or graph matching algorithms (*e.g.*, Hungarian method, Dynamic Programming, clock-matching). However, these pairwise strategies are time-consuming for a real-time search engine. Among them, Hausdorff distance may be the most efficient one, since it only requires some simple algebraic operations without sophisticated optimizations.

Recall that the standard Hausdorff distance measures the difference between two sets, and it is defined as

$$D(x_q, x_p) = \max_{q_i \in \mathcal{V}(x_q)} \min_{p_j \in \mathcal{V}(x_p)} d(q_i, p_j), \quad (1)$$

where function  $d(\cdot)$  measures the distance between two input vectors. In order to eliminate the disturbance of isolated views in the query view set, a more robust version of Hausdorff distance is given by

$$D(x_q, x_p) = \frac{1}{N_v} \sum_{q_i \in \mathcal{V}(x_q)} \min_{p_j \in \mathcal{V}(x_p)} d(q_i, p_j). \quad (2)$$

For the convenience of analysis, we consider its dual form in the similarity space as

$$S(x_q, x_p) = \frac{1}{N_v} \sum_{q_i \in \mathcal{V}(x_q)} \max_{p_j \in \mathcal{V}(x_p)} s(q_i, p_j), \quad (3)$$

where  $s(\cdot)$  measures the similarity between the two input vectors. In this paper, we adopt the cosine similarity.

As can be seen from Eq. (2) and Eq. (3), Hausdorff matching requires the time complexity  $O(N \times N_v^2)$  for retrieving a given query (assuming that there are  $N$  shapes in the database). Though the complexity grows linearly with respect to the database size, it is still intolerable when  $N$  gets larger. However, by analyzing Eq. (3), we can make several observations: (1) let

$s^*(q_i) = \max_{1 \leq j \leq N_v} s(q_i, p_j)$ , the similarity calculations of  $s(q_i, p_j)$  are unnecessary when  $s(q_i, p_j) < s^*(q_i)$ , since these similarity values are unused due to the *max* operation, *i.e.*, only  $s^*(q_i)$  is kept; (2) when considering from the query side, we can find that  $s^*(q_i)$  counts little to the final matching cost if  $s^*(q_i) < \xi$  and  $\xi$  is a small threshold. Those observations suggest that although the matching function in Eq. (3) requires the calculation of all the pairwise similarities between two view sets, some similarity calculations, which generate small values, can be eliminated without impairing the retrieval performance too much.

In order to avoid these unnecessary operations and improve the efficiency of multi-view matching procedure, we adopt **inverted file** to approximate Eq. (3) by adding the Kronecker delta response as

$$S(x_q, x_p) = \frac{1}{N_v} \sum_{q_i \in \mathcal{V}(x_q)} \max_{p_j \in \mathcal{V}(x_p)} s(q_i, p_j) \cdot \delta_{c(q_i), c(p_j)}, \quad (4)$$

where  $\delta_{x,y} = 1$  if  $x = y$ , and  $\delta_{x,y} = 0$  if  $x \neq y$ . The quantizer  $c(x) = \arg \min_{1 \leq i \leq K} \|x - b_i\|^2$  maps the input feature into an integer index that corresponds to the nearest codeword of the given vocabulary  $B = \{b_1, b_2, \dots, b_K\}$ . As a result, the contribution of  $p_j$ , which satisfies  $c(q_i) \neq c(p_j)$ , to the similarity measure can be directly set to zero, without estimating  $s(q_i, p_j)$  explicitly.

In conclusion, our inverted file for multi-view matching is built as illustrated in Fig. 2. For each view feature, we store it and its corresponding shape ID in the nearest codeword. It should be mentioned that we can also use Multiple Assignment (MA), *i.e.*, assign each view to multiple codewords, to improve the matching precision at the sacrifice of memory cost and on-line query time.

## 2.4. Inverted File for Re-ranking

A typical search engine usually involves a re-ranking component [21], aiming at refining the initial candidate list by using some contextual information. In GIFT, we propose a new contextual similarity measure called Aggregated Contextual Activation (ACA), which follows the same

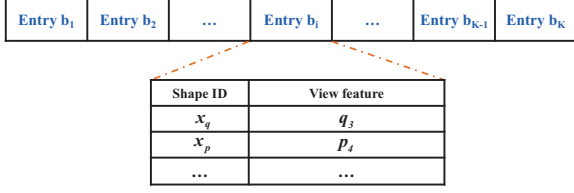


Figure 2. The structure of the first inverted file.

principles as diffusion process [5], *i.e.*, the similarity between two shapes should go beyond their pairwise formulation and is influenced by their contextual distributions along the underlying data manifold. However, apart from diffusion process which has high time complexity, ACA enables real-time re-ranking, which can be applied to large scale data.

Let  $\mathcal{N}_k(x_q)$  denote the neighbor set of  $x_q$ , which contains its top- $k$  neighbors. Similar to [42], our basic idea is that the similarity between two shapes can be more reliably measured by comparing their neighbors using Jaccard similarity as

$$S'(x_q, x_p) = \frac{|\mathcal{N}_k(x_q) \cap \mathcal{N}_k(x_p)|}{|\mathcal{N}_k(x_q) \cup \mathcal{N}_k(x_p)|}. \quad (5)$$

One can find that the neighbors are treated equally in Eq. (5). However the top-ranked neighbors are more likely to be true positives. So a more proper behavior is increasing the weights of top-ranked neighbors.

To achieve this, we propose to define the neighbor set using fuzzy set theory<sup>1</sup>. Different from classical (crisp) set theory where each element either belongs or does not belong to the set, fuzzy set theory allows a gradual assessment of the membership of elements in a set. We utilize  $S(x_q, x_i)$  to measure the membership grade of  $x_i$  in the neighbor set of  $x_q$ . Accordingly, Eq. (5) is re-written as

$$S'(x_q, x_p) = \frac{\sum_{x_i \in \mathcal{N}_k(x_q) \cap \mathcal{N}_k(x_p)} \min(S(x_q, x_i), S(x_p, x_i))}{\sum_{x_i \in \mathcal{N}_k(x_q) \cup \mathcal{N}_k(x_p)} \max(S(x_q, x_i), S(x_p, x_i))}. \quad (6)$$

Since considering equal-sized vector comparison is more convenient in real computational applications, we use  $F \in \mathbb{R}^N$  to encode the membership values. The  $i$ -th element in  $F_q$  is given as

$$F_q[i] = \begin{cases} S(x_q, x_i) & \text{if } x_i \in \mathcal{N}_k(x_q) \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

Based on this definition we replace Eq. 6 with

$$S'(x_q, x_p) = \frac{\sum_{i=1}^N \min(F_q[i], F_p[i])}{\sum_{i=1}^N \max(F_q[i], F_p[i])}. \quad (8)$$

<sup>1</sup>Please refer to **supplementary material** for the formal definitions of all the fuzzy operations used in this paper.

Considering vector  $F_q$  is sparse, we can view it as sparse activation of shape  $x_q$ , where the activation at coordinate  $i$  is the membership grade of  $x_i$  in the neighbor set  $\mathcal{N}_k(x_q)$ . Eq. (8) utilizes the sparse activations  $F_q$  and  $F_p$  to define the new contextual shape similarity measure.

Note that all the above analysis is carried out for only one similarity measure. However, in our specific scenario, the outputs of different layers of CNN are usually at different abstraction resolutions.

For example, two different layers of CNN lead to two different similarities  $S^{(1)}$  and  $S^{(2)}$  by Eq. (3), which in turn yield two different sparse activations  $F_q^{(1)}$  and  $F_q^{(2)}$  by Eq. (7). Since no prior information is available to assess their discriminative power, our goal now is to fuse them in a unsupervised way. For this we utilize the aggregation operation in fuzzy set theory, by which several fuzzy sets are combined in a desirable way to produce a single fuzzy set. We consider two fuzzy sets represented by the sparse activations  $F_q^{(1)}$  and  $F_q^{(2)}$  (the extension to more than two activations is similar). Their aggregation is then defined as

$$F_q = \left( \frac{(F_q^{(1)})^\alpha + (F_q^{(2)})^\alpha}{2} \right)^{\frac{1}{\alpha}}, \quad (9)$$

which computes the element-wise generalized means with exponent  $\alpha$  of  $F_q^{(1)}$  and  $F_q^{(2)}$ . Instead of using arithmetic mean, we use this generalized means ( $\alpha$  is set to 0.5 throughout our experiments). Our concern for this is to avoid the problem that some artificially large elements in  $F_q$  dominate the similarity measure. This motivation is very similar to handling bursty visual elements in Bag-of-Words (BoW) model (see [10] for examples).

In summary, we call the feature in Eq. (9) Aggregated Contextual Activation (ACA). Next, we will introduce some improvements of Eq. (9) concerning its retrieval accuracy and computational efficiency.

#### 2.4.1 Improving Accuracy

Similar to diffusion process, the proposed ACA requires an accurate estimation of the context in the data manifold. Here we provide two alternative ways to improve the retrieval performance of ACA without depriving its efficiency.

**Neighbor Augmentation.** The first one is to augment  $F_q$  using the neighbors of second order, *i.e.*, the neighbors of the neighbors of  $x_q$ . Inspired by query expansion [24], the second order neighbors are added as

$$F_q^{(l)} := \frac{1}{|\mathcal{N}_k^{(l)}(x_q)|} \sum_{x_i \in \mathcal{N}_k^{(l)}(x_q)} F_i^{(l)}. \quad (10)$$

**Neighbor Co-augmentation.** Our second improvement is to use a so-called “neighbor co-augmentation”. Specific-

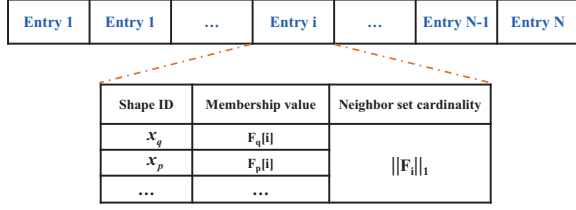


Figure 3. The structure of the second inverted file.

ly, the neighbors generated by one similarity measure are used to augment contextual activations of the other similarity measure, formally defined as

$$\begin{aligned}
 F_q^{(1)} &:= \frac{1}{|\mathcal{N}_k^{(2)}(x_q)|} \sum_{x_i \in \mathcal{N}_k^{(2)}(x_q)} F_i^{(1)}, \\
 F_q^{(2)} &:= \frac{1}{|\mathcal{N}_k^{(1)}(x_q)|} \sum_{x_i \in \mathcal{N}_k^{(1)}(x_q)} F_i^{(2)}.
 \end{aligned} \tag{11}$$

This formula is inspired by ‘‘co-training’’ [44]. Essentially, one similarity measure tells the other one that ‘‘I think these neighbors to be true positives, and lend them to you such that you can improve your own discriminative power’’.

Note that the size of neighbor set used here may be different from that used in Eq. (7). In order to distinguish them, we denote the size of neighbor set in Eq. (7) as  $k_1$ , while that used in Eq. (10) and Eq. (11) as  $k_2$ .

### 2.4.2 Improving Efficiency

Considering that the length of  $F_q$  is  $N$ , one may doubt the efficiency of similarity computation in Eq. (8), especially when the database size  $N$  is large. In fact,  $F_q$  is a sparse vector, since  $F_q$  only encodes the neighborhood structure of  $x_q$ , and the number of non-zero values is only determined by the size of  $\mathcal{N}_k(x_q)$ . This observation motivate us to utilize an **inverted file** again to leverage the sparsity of  $F_q$ . Now we derive the feasibility of applying inverted file in Jaccard similarity theoretically.

The numerator in Eq. (8) is computed as

$$\begin{aligned}
 \sum_i \min(F_q[i], F_p[i]) &= \sum_{i|F_q[i] \neq 0, F_p[i] \neq 0} \min(F_q[i], F_p[i]) \\
 + \sum_{i|F_q[i] = 0} \min(F_q[i], F_p[i]) &+ \sum_{i|F_p[i] = 0} \min(F_q[i], F_p[i]).
 \end{aligned} \tag{12}$$

Since all values of the aggregated contextual activation are non-negative, the last two items in Eq. (12) are equal to zero. Consequently, Eq. (12) can be simplified as

$$\sum_i \min(F_q[i], F_p[i]) = \sum_{i|F_q[i] \neq 0, F_p[i] \neq 0} \min(F_q[i], F_p[i]), \tag{13}$$

which only requires accessing non-zero entries of the query, and hence can be computed efficiently on-the-fly.

Although the calculation of the denominator in Eq. (8) seems sophisticated, it can be expressed as

$$\begin{aligned}
 &\sum_i \max(F_q[i], F_p[i]) \\
 &= \|F_q\|_1 + \|F_p\|_1 - \sum_i \min(F_q[i], F_p[i]) \\
 &= \|F_q\|_1 + \|F_p\|_1 - \sum_{i|F_q[i] \neq 0, F_p[i] \neq 0} \min(F_q[i], F_p[i]).
 \end{aligned} \tag{14}$$

Besides the query-dependent operations (the first and the last items), Eq. (14) only involves an operation of  $L_1$  norm calculation of  $F_p$ , which is simply equal to the cardinality of the fuzzy set  $\mathcal{N}_k(x_p)$  and can be pre-computed off-line.

Our inverted file for re-ranking is built as illustrated in Fig. 3. It has exactly  $N$  entries, and each entry corresponds to one shape in the database. For each entry, we first store the cardinality of its fuzzy neighbor set. Then, we find those shapes which have non-negative membership values in this entry. Those shape IDs and the membership values are stored in this entry.

## 3. Experiments

In this section, we evaluate the performance of GIFT on different kinds of 3D shape retrieval tasks. The evaluation metrics used in this paper include mean average precision (MAP), area under curve (AUC), Nearest Neighbor (NN), First Tier (FT) and Second Tier (ST). Refer to [39, 29] for their detailed definitions.

If not specified, we adopt the following setup throughout our experiments. The projection rendered for each shape is  $N_v = 64$ . For multi-view matching procedure, the approximate Hausdorff matching defined in Eq. (4) with an inverted file of 256 entries is used. Multiple Assignment is set to 2. We use two pairwise similarity measures, which are calculated using features from convolutional layer  $L_5$  and fully-connected layer  $L_7$  respectively. In re-ranking component, each similarity measure generates one sparse activation  $F_q$  to capture the contextual information for the 3D shape  $x_q$ , and neighbor co-augmentation in Eq. (11) is used to produce  $F_q^{(1)}$  and  $F_q^{(2)}$ . Finally, both  $F_q^{(1)}$  and  $F_q^{(2)}$  are integrated by (9) with exponent  $\alpha = 0.5$ .

### 3.1. ModelNet

ModelNet is a large-scale 3D CAD model dataset introduced by Wu *et al.* [39] recently, which contains 151, 128 3D CAD models divided into 660 object categories. Two subsets are used for evaluation, *i.e.*, ModelNet40 and ModelNet10. The former one contains 12, 311 models, and the latter one contains 4, 899 models. We evaluate the performance of GIFT on both subsets and adopt the same training and test split as in [39], namely randomly selecting 100

Methods	ModelNet40		ModelNet10	
	AUC	MAP	AUC	MAP
SPH [11]	34.47%	33.26%	45.97%	44.05%
LFD [4]	42.04%	40.91%	51.70%	49.82%
PANORAMA [24]	45.00%	46.13%	60.72%	60.32%
ShapeNets [39]	49.94%	49.23%	69.28%	68.26%
DeepPano [28]	77.63%	76.81%	85.45%	84.18%
MVCNN [32]	-	78.90%	-	-
$L_5$	63.70%	63.07%	78.19%	77.25%
$L_7$	77.28%	76.63%	89.03%	88.05%
GIFT	<b>83.10%</b>	<b>81.94%</b>	<b>92.35%</b>	<b>91.12%</b>

Table 1. The performance comparison with state-of-the-art on ModelNet40 and ModelNet10.

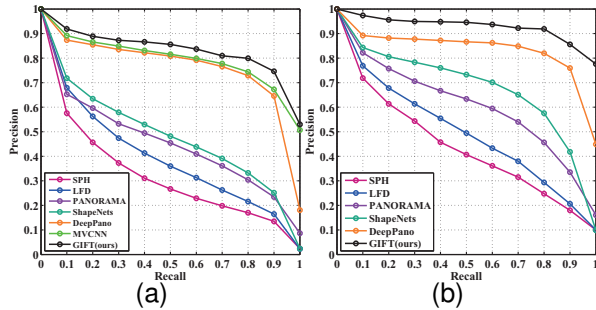


Figure 4. Precision-recall curves on ModelNet40 (a) and ModelNet10 (b).

unique models per category from the subset, in which 80 models are used for training the CNN model and the rest for testing the retrieval performance.

For comparison, we collected all the retrieval results publicly available<sup>2</sup>. The chosen methods are (Spherical Harmonic) SPH [11], (Light Field descriptor) LFD [4], PANORAMA [24], 3D ShapeNet [39], DeepPano [28] and MVCNN [32]. As Table 1 shows, GIFT outperforms all the state-of-the-art methods remarkably. We also present the performance of two baseline methods, *i.e.*, feature  $L_5$  or  $L_7$  with exact Hausdorff matching. As can be seen,  $L_7$  achieves a better performance than  $L_5$ , and GIFT leads to a significant improvement over  $L_7$  of 5.82% in AUC, 5.31% in MAP for ModelNet40 dataset, and 3.32% in AUC, 3.07% in MAP for ModelNet10 dataset.

Fig. 4 compares the precision-recall curves. It demonstrates again the discriminative power of the proposed search engine in 3D shape retrieval. Note that ModelNet also defines the 3D shape classification tasks. Considering GIFT is initially developed for real-time retrieval, its classification results are given in the supplementary material.

### 3.2. Large Scale Competition

As the most authoritative 3D retrieval competition held each year, SHape REtrieval Contest (SHREC) pays much attention to the development of scalable algorithms gradual-

ly. Especially in recent years, several large scale tracks [31], such as SHREC14LSGTB [14], are organized to test the scalability of algorithms. However, most algorithms that the participants submit are of high time complexity, and cannot be applied when the dataset becomes larger (millions or more). Here we choose SHREC14LSGTB dataset for a comprehensive evaluation. This dataset contains 8,987 3D models classified into 171 classes, and each 3D shape is taken in turn as the query. As for the feature extractor, we collected 54,728 unrelated models from ModelNet [39] divided into 461 categories to train a CNN model.

To keep the comparison fair, we choose two types of results from the survey paper [14] to present in Table 2. The first type consists of the top-5 best-performing methods on retrieval accuracy, including PANORAMA [24], DBSVC, MR-BF-DSIFT, MR-D1SIFT and LCDR-DBSVC. The second type is the most efficient one, *i.e.*, ZFDR [13].

As can be seen from the table, excluding GIFT, the best performance is achieved by LCDR-DBSVC. However, it requires 668.6s to return the retrieval results per query, which means that 69 days are needed to finish the query task on the whole dataset. The reason behind such a high complexity lies in two aspects: 1) its visual feature is 270K dimensional, which is time consuming to compute, store and compare; 2) it adopts locally constrained diffusion process (LCDP) [41] for re-ranking, while it is known that LCDP is an iterative graph-based algorithm of high time complexity. As for ZFDR, its average query time is shortened to 1.77s by computing parallel on 12 cores. Unfortunately, ZFDR achieves much less accurate retrieval performance, and its FT is 13% smaller than LCDR-DBSVC. In summary, a conclusion can be drawn that no method can achieve a good enough performance at a low time complexity.

By contrast, GIFT outperforms all these methods, including a very recent algorithm called Two Layer Coding (TLC) [1] which reports 0.585 in FT. What is more important that GIFT can provide the retrieval results within 63.14ms, which is 4 orders of magnitude faster than LCDR-DBSVC. Meanwhile, the two baseline methods  $L_5$  and  $L_7$  incur heavy query cost due to the usage of exact Hausdorff matching, which testifies the advantage of the proposed F-IF.

### 3.3. Generic 3D Retrieval

Following [34], we select three popular datasets for a generic evaluation, including Princeton Shape Benchmark (PSB) [29], Watertight Models track of SHape REtrieval Contest 2007 (WM-SHREC07) [8] and McGill dataset [30]. Among them, PSB dataset is probably the first widely-used generic shape benchmark, and it consists of 907 polygonal models divided into 92 categories. WM-SHREC07 contains 400 watertight models evenly distributed in 20 classes, and is a representative competition held by SHREC community.

<sup>2</sup><http://modelnet.cs.princeton.edu/>

Methods	Accuracy			Query time
	NN	FT	ST	
ZFDR	0.879	0.398	0.535	1.77s
PANORAMA	0.859	0.436	0.560	370.2s
DBSVC	0.868	0.438	0.563	62.66s
MR-BF-DSIFT	0.845	0.455	0.567	65.17s
MR-D1SIFT	0.856	0.465	0.578	131.04s
LCDR-DBSVC	0.864	0.528	0.661	668.6s
$L_5$	0.879	0.460	0.592	22.73s
$L_7$	0.884	0.507	0.642	4.82s
GIFT	<b>0.889</b>	<b>0.567</b>	<b>0.689</b>	<b>63.14ms</b>

Table 2. The performance comparison on SHREC14LSGTB.

Datasets	Off-line	On-line Indexing
ModelNet40	$\approx 0.7h$	27.02ms
ModelNet10	$\approx 0.3h$	10.25ms
SHREC14LSGTB	$\approx 8.5h$	63.14ms
PSB		16.25ms
WMSHREC07	$\approx 1.8h$	16.05ms
McGill		9.38ms

Table 4. The time cost analysis of GIFT.

McGill dataset focuses on non-rigid analysis, and contains 255 articulated objects classified into 10 classes. We train CNN on an independent TSB dataset [36], and then use the trained CNN to extract view features for the shapes in all the three testing datasets.

In Table 3, a comprehensive comparison between GIFT and various state-of-the-art methods is presented, including LFD [4], the curve-based method of Tabia *et al.* [33], DESIRE descriptor [37], total Bregman Divergences (tB-D) [19], Covariance descriptor [34], the Hybrid of 2D and 3D descriptor [23], Two Layer Coding (TLC) [1] and PANORAMA [24]. As can be seen, GIFT exhibits encouraging discriminative ability in retrieval accuracy and achieves state-of-the-art performances consistently on all the three evaluation metrics.

### 3.4. Execution Time

In addition to state-of-the-art performances on several datasets and competitions, the most important property of GIFT is the “real-time” performance with the potential of handling large scale shape corpora. In Table 4, we give a deeper analysis of the time cost. The off-line operations mainly include projection rendering and feature extraction for database shapes, training CNN, and building two inverted files. As the table shows, the time cost of off-line operations varies significantly for different datasets. Among them, the most time-consuming operation is training CNN, followed by building the first inverted file with k-means. However, the average query time for different datasets can be controlled within one second, even for the biggest SHREC14LSGTB dataset.

Feature	Hausdorff	Re-ranking		First Tier
		$\alpha$	NA	
$L_5$	$\times$			0.588
$L_7$	$\times$			0.653
$L_5 + L_7$	$\times$	1		0.688
$L_5 + L_7$	$\times$	0.5		0.692
$L_5 + L_7$	$\times$	0.5	$\times$	0.710
$L_5 + L_7$	$\times$	0.5	$\checkmark$	0.717
$L_5 + L_7$	$\checkmark$	0.5	$\checkmark$	0.712

Table 5. The performance improvements brought by various components in GIFT over baseline. In column “Hausdorff”,  $\checkmark$  denotes approximate Hausdorff matching in Eq. (4), while  $\times$  denotes exact matching in Eq. (3). Column “ $\alpha$ ” present the value of exponent in Eq. (9). Column “NA” describes the procedure of neighbor augmentation in Sec. 2.4.1:  $\checkmark$  is associated with Eq. (11) and  $\times$  is associated with Eq. (10). The blanks mean that this improvement is not used.

### 3.5. Parameter Discussion

Due to the space limitation, the discussion is conducted only on PSB dataset.

**Improvements Over Baseline.** In Table 5, a thorough discussion is given about the influence of various components of GIFT. We can observe a consistent performance boost by those improvements. The performance jumps a lot especially when re-ranking component is embedded. One should note a slight performance decrease when approximate Hausdorff matching with F-IF is used as compared with its exact version. However, as discussed below, the embedding with inverted file does not necessarily result in a poorer performance, but shortens the query time significantly.

**Discussion on F-IF.** In Fig. 5, we plot the retrieval performance and the average query time using feature  $L_7$ , as the number of entries used in the first inverted file changes. As Fig. 5(a) shows, the retrieval performance generally decreases with more entries, and multiple assignment can boost the retrieval performance significantly. However, it should be addressed that a better approximation to Eq. (3) using fewer entries (decreasing  $K$ ) or larger multiple assignments (increasing MA) does not necessarily imply a better retrieval performance. For example, when  $K = 256$  and MA = 2, the performance of approximate Hausdorff matching using inverted file surpasses the baseline using exact Hausdorff matching. The reason for this “abnormal” observation is that the principle of inverted file here is to reject those view matching operations that lead to smaller similarities, and sometimes they are noisy and false matching pairs which can be harmful to retrieval performance.

As can be seen from Fig. 5(b), the average query time is higher at smaller  $K$  and larger MA, since the two cases both increase the number of candidate matchings in each entry. The baseline query time using exact Hausdorff matching is

Methods	PSB dataset			WM-SHREC07 competition			McGill dataset		
	NN	FT	ST	NN	FT	ST	NN	FT	ST
LFD [4]	0.657	0.380	0.487	0.923	0.526	0.662	-	-	-
Tabia <i>et al.</i> [33]	-	-	-	0.853	0.527	0.639	-	-	-
DESIRE [37]	0.665	0.403	0.512	0.917	0.535	0.673	-	-	-
tBD [19]	0.723	-	-	-	-	-	-	-	-
Covariance [34]	-	-	-	0.930	0.623	0.737	0.977	0.732	0.818
2D/3D Hybrid [23]	0.742	0.473	0.606	0.955	0.642	0.773	0.925	0.557	0.698
PANORAMA [24]	0.753	0.479	0.603	0.957	0.673	0.784	0.929	0.589	0.732
PANORAMA + LRF [24]	0.752	0.531	0.659	0.957	0.743	0.839	0.910	0.693	0.812
TLC [1]	0.763	0.562	0.705	0.988	0.831	0.935	0.980	0.807	0.933
$L_5$	0.849	0.588	0.721	0.980	0.777	0.877	0.984	0.747	0.881
$L_7$	0.837	0.653	0.784	0.980	0.805	0.898	0.980	0.763	0.897
GIFT	<b>0.849</b>	<b>0.712</b>	<b>0.830</b>	<b>0.990</b>	<b>0.949</b>	<b>0.990</b>	<b>0.984</b>	<b>0.905</b>	<b>0.973</b>

Table 3. The performance comparison with other state-of-the-art algorithms on PSB dataset, WM-SHREC07 dataset and McGill dataset.

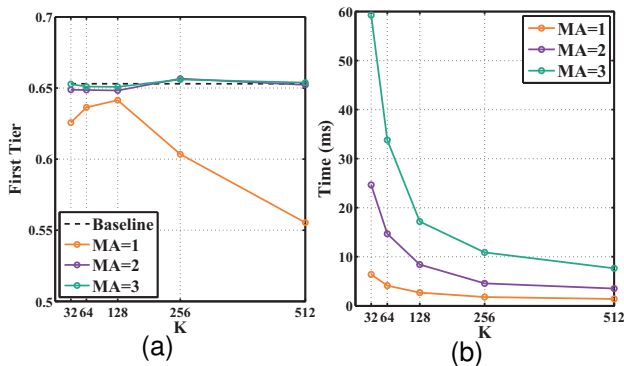


Figure 5. The performance difference between Hausdorff matching and its approximate version in terms of retrieval accuracy (a) and average query time (b).

0.69s, which is at least one order of magnitude larger than the approximate one.

**Discussion on S-IF.** Two parameters,  $k_1$  and  $k_2$ , are involved in the second inverted file, which are determined empirically. We plot the influence of them in Fig. 6. As can be drawn from the figure, when  $k_1$  increases, the retrieval performance increases at first. Since noise contextual information can be included at a larger  $k_1$ , we can observe the performance decreases after  $k_1 > 10$ . Meanwhile, neighbor augmentation can boost the performance further. For example, the best performance is achieved when  $k_2 = 4$ . However, when  $k_2 = 5$ , the performance tends to decrease. One may find that the optimal value of  $k_2$  is much smaller than that of  $k_1$ . The reason for this is that  $k_2$  defines the size of the second order neighbor, which is more likely to return noise context compared with the first order neighbor defined by  $k_1$ .

#### 4. Conclusions

In the past years, 3D shape retrieval was evaluated with only small numbers of shapes. In this sense, the problem of 3D shape retrieval has stagnated for a long time. Only

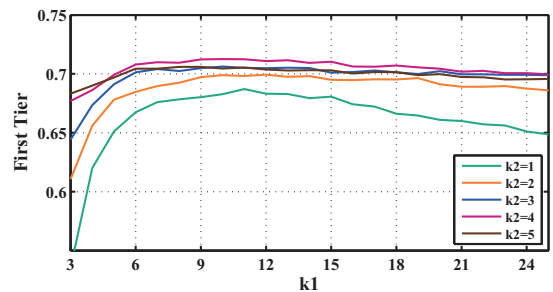


Figure 6. The influence of neighbor set sizes  $k_1$  and  $k_2$  used in the second inverted file.

recently, shape community started to pay more attention to the scalable retrieval issue gradually. However, as suggested in [14], most classical methods encounter severe obstacles when dealing with larger databases.

In this paper, we focus on the scalability of 3D shape retrieval algorithms, and build a well-designed 3D shape search engine called GIFT. In our retrieval system, GPU is utilized to accelerate the speed of projection rendering and view feature extraction, and two inverted files are embedded to enable real-time multi-view matching and re-ranking. As a result, the average query time is controlled within one second, which clearly demonstrates the potential of GIFT for large scale 3D shape retrieval. What is more impressive is that while preserving the high time efficiency, GIFT outperforms state-of-the-art methods in retrieval accuracy by a large margin. Therefore, we view the proposed search engine as a promising step towards larger 3D shape corpora.

We submitted a version of GIFT to the latest SHREC2016 large scale track (the results are available in <https://shapenet.cs.stanford.edu/shrec16/>), and won the first place on perturbed dataset.

**Acknowledgements.** This work was supported in part by NSFC 61573160, Program for New Century Excellent Talents in University under Grant NCET-12-0217, NSF No. IIS-1302164 and China Scholarship Council.



## References

- [1] X. Bai, S. Bai, Z. Zhu, and L. J. Latecki. 3d shape matching via two layer coding. *TPAMI*, 37(12):2361–2373, 2015.
- [2] M. M. Bronstein and I. Kokkinos. Scale-invariant heat kernel signatures for non-rigid shape recognition. In *CVPR*, pages 1704–1711, 2010.
- [3] K. Chatfield, K. Simonyan, A. Vedaldi, and A. Zisserman. Return of the devil in the details: Delving deep into convolutional nets. *CoRR*, abs/1405.3531, 2014.
- [4] D. Y. Chen, X. P. Tian, Y. T. Shen, and M. Ouhyoung. On visual similarity based 3d model retrieval. *Comput. Graph. Forum*, 22(3):223–232, 2003.
- [5] M. Donoser and H. Bischof. Diffusion processes for retrieval revisited. In *CVPR*, pages 1320–1327, 2013.
- [6] Y. Fang, J. Xie, G. Dai, M. Wang, F. Zhu, T. Xu, and E. Wong. 3d deep shape descriptor. In *CVPR*, pages 2319–2328, 2015.
- [7] L. Fei-Fei and P. Perona. A bayesian hierarchical model for learning natural scene categories. In *CVPR*, pages 524–531, 2005.
- [8] D. Giorgi, S. Biasotti, and L. Paraboschi. Shape retrieval contest 2007: Watertight models track. *SHREC competition*, 8, 2007.
- [9] M. Havlena and K. Schindler. Vocmatch: Efficient multiview correspondence for structure from motion. In *ECCV*, pages 46–60, 2014.
- [10] H. Jégou, M. Douze, and C. Schmid. On the burstiness of visual elements. In *CVPR*, pages 1169–1176, 2009.
- [11] M. Kazhdan, T. Funkhouser, and S. Rusinkiewicz. Rotation invariant spherical harmonic representation of 3d shape descriptors. In *SGP*, pages 156–164, 2003.
- [12] I. Kokkinos, M. M. Bronstein, R. Litman, and A. M. Bronstein. Intrinsic shape context descriptors for deformable shapes. In *CVPR*, pages 159–166, 2012.
- [13] B. Li and H. Johan. 3d model retrieval using hybrid features and class information. *Multimedia tools and applications*, 62(3):821–846, 2013.
- [14] B. Li, Y. Lu, C. Li, A. Godil, T. Schreck, M. Aono, et al. A comparison of 3d shape retrieval methods based on a large-scale benchmark supporting multimodal queries. *CVI-U*, 131:1–27, 2015.
- [15] C. Li, M. Ovsjanikov, and F. Chazal. Persistence-based structural recognition. In *CVPR*, pages 2003–2010, 2014.
- [16] Z. Lian, A. Godil, X. Sun, and J. Xiao. Cm-bof: visual similarity-based 3d shape retrieval using clock matching and bag-of-features. *Mach. Vis. Appl.*, 24(8):1685–1704, 2013.
- [17] H. Ling and D. W. Jacobs. Shape classification using the inner-distance. *TPAMI*, 29(2):286–299, 2007.
- [18] R. Litman, A. Bronstein, M. Bronstein, and U. Castellani. Supervised learning of bag-of-features shape descriptors using sparse coding. *Computer Graphics Forum*, 33(5):127–136, 2014.
- [19] M. Liu, B. C. Vemuri, S. ichi Amari, and F. Nielsen. Shape retrieval using hierarchical total bregman soft clustering. *TPAMI*, 34(12):2407–2419, 2012.
- [20] D. G. Lowe. Distinctive image features from scale-invariant keypoints. *IJCV*, 60(2):91–110, 2004.
- [21] T. Mei, Y. Rui, S. Li, and Q. Tian. Multimedia search reranking: A literature survey. *ACM Comput. Surv.*, 46(3):38:1–38:38, 2014.
- [22] P. Papadakis, I. Pratikakis, S. J. Perantonis, and T. Theoharis. Efficient 3d shape matching and retrieval using a concrete radialized spherical projection representation. *Pattern Recognition*, 40(9):2437–2452, 2007.
- [23] P. Papadakis, I. Pratikakis, T. Theoharis, G. Passalis, and S. J. Perantonis. 3d object retrieval using an efficient and compact hybrid shape descriptor. In *3DOR*, pages 9–16, 2008.
- [24] P. Papadakis, I. Pratikakis, T. Theoharis, and S. J. Perantonis. Panorama: A 3d shape descriptor based on panoramic views for unsupervised 3d object retrieval. *IJCV*, 89(2-3):177–192, 2010.
- [25] E. Rodolà, T. Harada, Y. Kuniyoshi, and D. Cremers. Efficient shape matching using vector extrapolation. In *BMVC*, volume 1, 2013.
- [26] E. Rodola, S. Rota Bulò, T. Windheuser, M. Vestner, and D. Cremers. Dense non-rigid shape correspondence using random forests. In *CVPR*, pages 4177–4184, 2014.
- [27] E. Rodola, A. Torsello, T. Harada, Y. Kuniyoshi, and D. Cremers. Elastic net constraints for shape matching. In *ICCV*, pages 1169–1176, 2013.
- [28] B. Shi, S. Bai, Z. Zhou, and X. Bai. Deeppano: Deep panoramic representation for 3-d shape recognition. *IEEE Signal Processing Letters*, 22(12):2339–2343, 2015.
- [29] P. Shilane, P. Min, M. M. Kazhdan, and T. A. Funkhouser. The princeton shape benchmark. In *SMI*, 2004.
- [30] K. Siddiqi, J. Zhang, D. Macrini, A. Shokoufandeh, S. Bouix, and S. J. Dickinson. Retrieving articulated 3-d models using medial surfaces. *Mach. Vis. Appl.*, 19(4):261–275, 2008.
- [31] I. Sipiran, B. Bustos, T. Schreck, A. Bronstein, S. Choi, L. Lai, H. Li, R. Litman, and L. Sun. Scalability of non-rigid 3d shape retrieval. In *3DOR*, pages 121–128, 2015.
- [32] H. Su, S. Maji, E. Kalogerakis, and E. G. Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *ICCV*, 2015.
- [33] H. Tabia, M. Daoudi, J.-P. Vandeborre, and O. Colot. A new 3d-matching method of nonrigid and partially similar models using curve analysis. *TPAMI*, 33(4):852–858, 2011.
- [34] H. Tabia, H. Laga, D. Picard, and P.-H. Gosselin. Covariance descriptors for 3d shape matching and retrieval. In *CVPR*, pages 4185–4192, 2014.
- [35] H. Tabia, D. Picard, H. Laga, and P. H. Gosselin. Compact vectors of locally aggregated tensors for 3d shape retrieval. In *3DOR*, pages 17–24, 2013.
- [36] A. Tatsuma, H. Koyanagi, and M. Aono. A large-scale shape benchmark for 3d object retrieval: Toyohashi shape benchmark. In *APSIPA*, pages 1–10, 2012.
- [37] D. V. Vranic. Desire: a composite 3d-shape descriptor. In *ICME*, pages 962–965, 2005.
- [38] Y. Wang, M. Gong, T. Wang, D. Cohen-Or, H. Zhang, and B. Chen. Projective analysis for 3d shape segmentation. *ACM Trans. Graph.*, 32(6):192:1–192:12, 2013.
- [39] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shape modeling. In *CVPR*, 2015.

- [40] J. Xie, Y. Fang, F. Zhu, and E. Wong. Deepshape: Deep learned shape descriptor for 3d shape matching and retrieval. In *CVPR*, pages 1275–1283, 2015.
- [41] X. Yang, S. Koknar-Tezel, and L. J. Latecki. Locally constrained diffusion process on locally densified distance spaces with applications to shape retrieval. In *CVPR*, pages 357–364, 2009.
- [42] S. Zhang, M. Yang, T. Cour, K. Yu, and D. N. Metaxas. Query specific rank fusion for image retrieval. *TPAMI*, 37(4):803–815, 2015.
- [43] L. Zheng, S. Wang, L. Tian, F. He, Z. Liu, and Q. Tian. Query-adaptive late fusion for image search and person re-identification. In *CVPR*, pages 1741–1750, 2015.
- [44] Z.-H. Zhou and M. Li. Semi-supervised regression with co-training. In *IJCAI*, 2005.