

Sparse Contextual Activation for Efficient Visual Re-ranking

Song Bai, *Student Member, IEEE*, and Xiang Bai, *Senior Member, IEEE*

Abstract—In this paper, we propose an extremely efficient algorithm for visual re-ranking. By considering the original pairwise distance in the contextual space, we develop a feature vector called Sparse Contextual Activation (SCA) that encodes the local distribution of an image. Hence, re-ranking task can be simply accomplished by vector comparison under the generalized Jaccard metric, which has its theoretical meaning in fuzzy set theory. In order to improve the time efficiency of re-ranking procedure, inverted index is introduced successfully to speed up the computation of generalized Jaccard metric. As a result, the average time cost of re-ranking for a certain query can be controlled within one millisecond. Furthermore, inspired by Query Expansion, we also develop an additional method called Local Consistency Enhancement (LCE) on the proposed sparse contextual activation to improve the retrieval performance in an unsupervised manner. On the other hand, the retrieval performance using a single feature may not be satisfactory enough, which inspires us to fuse multiple complementary features for accurate retrieval. Based on sparse contextual activation, a robust feature fusion algorithm is also exploited that also preserves the characteristic of high time efficiency. We assess our proposed method in various visual re-ranking tasks. Experimental results on Ukbench dataset (image), YAEL dataset B (face), Princeton Shape Benchmark (3D object), WM-SRHEC07 (3D competition) and MPEG-7 dataset (shape) manifest the effectiveness and efficiency of SCA.

Index Terms—Jaccard distance, Feature Fusion, Re-ranking, Retrieval, Inverted index.

I. INTRODUCTION

CONTEXTUAL similarity/dissimilarity [1], [2], [3] has been extensively exploited recently due to its effectiveness in various visual retrieval tasks, such as natural image search, shape retrieval, biological information retrieval, analysis of time series, *etc.* Unlike traditional Content-based Image Retrieval (CBIR) systems that consider only pairwise dissimilarity measure for ranking and indexing, the approaches on contextual dissimilarity measure are proposed to explore the contextual information from the database instances, and hence enhance and refine the dissimilarity measure for improving the retrieval performance, which is usually considered as an unsupervised re-ranking procedure based on the given distance measure.

In general, the re-ranking procedure is often performed as a post-processing of ranking initialization, which creates a ranking list for a given query. For image retrieval, the dissimilarity measure between a pair of images is often obtained

by calculating the distance of their corresponding features under a certain metric. Given a query image, all the database images are sorted in an ascending/descending order according to their dissimilarities/similarities to the query. The ranking list for the query image can be finally initialized, where the most similar images occupy its top positions. A key issue in ranking initialization is to design proper features with enough discriminative power to represent an image, for which the Bag-of-Features (BoF) image representation [4] is often suggested.

Instead of ranking with pairwise dissimilarity measure, the contextual re-ranking algorithms have been proposed and proven their effectiveness by considering the relationships among all database instance [2], [3], [1], [5], [6], [7]. In these re-ranking algorithms, the dissimilarity measure between two instances is iteratively updated and refined by taking into account their local distributions (neighborhood structure). As the neighborhood of each instance can be directly obtained from the ranking list, the key advantage of these re-ranking approaches consists in the fact that training/labeled data is not required, operating in an unsupervised manner.

Though extensively studied, almost all the existing contextual re-ranking algorithms only pay much attention to the *effectiveness*, which refers to the level of retrieval accuracy. The *efficiency*, which refers to the time cost for the procedure of re-ranking, has been more or less neglected. However, both effectiveness and efficiency are quite important for a real-time retrieval system, and the tradeoff between them is badly required at present. The contextual re-ranking approaches often consider all the distances among instances of a given dataset, and the contextual dissimilarities are often achieved by operating on such a distance matrix. Therefore, a large computational effort is essential (typically, between $O(N^2)$ and $O(N^3)$), which seriously hinders their use in large-scale retrieval services.

In this paper, we address the contextual re-ranking problem in an alternative and simpler manner. The contextual dissimilarity measure between two images is calculated by comparing two neighborhood sets, *e.g.* the neighborhood sets of a query and a target, respectively. It shares a similar intuition of the traditional approaches that the relation between two images should not be determined by only the distance between them, but also influenced by the relation among their neighbors on the distance manifold. Our main contribution is to propose an extremely fast re-ranking algorithm called **Sparse Contextual Activation** (SCA) for computing the dissimilarity between such two neighborhood sets. Given a certain image, the basic idea of SCA is encoding its local distribution according to the original pairwise similarity into a single vector. Consequently,

S. Bai and X. Bai are with the school of Electronics Information and Communications, Huazhong University of Science and Technology, Wuhan, China, 430074

Manuscript received April 19, 2005; revised December 27, 2012.

the contextual dissimilarity measure (set-to-set distance) between two images can be simply obtained by comparing two encoded vectors. Due to the sparsity property of such encoded vectors, the inverted index [8] can be further used to speed up the computation of the vector comparison for re-ranking. Besides the efficiency of SCA, it also achieves state-of-the-arts retrieval performance on several standard benchmarks. In addition, we extend the proposed SCA for efficiently fusing multiple kinds of distance metrics for highly effective re-ranking while the inverted index can be also incorporated.

The rest of the paper is organized as follows. We first review some related work in Section II. The details of Sparse Contextual Activation (SCA) are introduced in Section III, and feature fusion based on SCA is described in Section IV. Experiments are carried out in Section V. Conclusions and future work are summarized in Section VI.

II. RELATED WORK

Since there exist a large amount of works on re-ranking algorithms, we only review the unsupervised re-ranking algorithms in this section.

A. Re-ranking with single distance measure

In recent years, contextual information has been successfully explored to improve the retrieval accuracy by replacing a given pairwise similarity with a more faithful one, which is learned by considering the relation among the database objects [3], [9], [10], [11], [12], [13]. The contextual re-ranking has a very diverse taxonomy (graph transduction [3], [11], diffusion process [12], [13], [2], rank aggregation [14], [15], contextual similarity/dissimilarities measure [1], [10], [9], query expansion [16], [17]). These post-processing approaches share the common spirit that the effectiveness of retrieval tasks is improved by relationships among dataset objects in an unsupervised manner, without labeled data.

One of the most classical algorithms is graph transduction (GT) [3]. As a semi-supervised method, GT spreads the information from the labeled data to unlabeled data by regarding the query itself as the only labeled data.

A popular branch for re-ranking is diffusion process, which is summarized as a generic framework in [2]. Most variants of diffusion process share the same perspective that the pairwise similarity is context-sensitive, and the geometric structure of data manifold should be considered. In [12], Locally Constrained Diffusion Process (LCDP) is proposed to apply the affinity propagation with the constraint of locality. Tensor Product Graph (TPG) [13] diffuses the similarity information in the tensor product graph achieved by the tensor product of the original graph with itself.

Based on the observation that a good ranking is usually asymmetric, Contextual Dissimilarity Measure (CDM) [1] improves the retrieval performance of BoF vectors by modifying the neighborhood structure using Sinkhorn's scaling algorithm. Query Expansion [16], [17] can substantially improve the retrieval performance by using relevant images as extra queries. Spatial verification [18] is proposed for re-ranking by considering the spatial constraints.

These aforementioned algorithms, as well as Self Diffusion (SD) [19], diffusion maps [20], aims at improving the retrieval accuracy, however the re-ranking efficiency is more or less neglected. By contrast, our proposed SCA is a highly efficient re-ranking algorithm, while also performs better in the retrieval accuracy.

B. Re-ranking with multiple distance measures

Considering that one distance measure only focuses one aspect of images, some re-ranking algorithm also deals with multiple complementary features.

Zhang *et al.* [21] fuse BoW feature and holistic feature by a graph-based query specific fusion, and re-ranking is performed by using the local PageRank algorithm or finding the weighted maximum density subgraph. Co-transduction [22] adopts a semi-supervised framework based on co-training [23], [24] to combine complementary features for image and shape retrieval.

Based on sparse contextual activation, we also propose a re-ranking version that deals with multiple distance measures, which maintains the characteristic of efficiency, but improves the performance significantly.

III. SPARSE CONTEXTUAL ACTIVATION

Let $X = \{x_1, x_2, \dots, x_N\}$ denote a collection of images. We define two functions listed as follows:

- Function $f : x \rightarrow \mathbb{R}^n$: it extracts a n -dimensional feature to represent the input image x .
- Function $d : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$: it computes the distance for the two input feature vectors $f(x_q)$ and $f(x_p)$ under a certain metric, where x_q and x_p represent two images in the database.

The distance $d(f(x_q), f(x_p))$ is taken as the dissimilarity of the images x_q and x_p . To simplify the notation, we use $d(x_q, x_p)$ to replace $d(f(x_q), f(x_p))$ below where possible.

After all the pairwise dissimilarity values related to the given query x_q are achieved, we could initialize the ranking list by sorting the dissimilarity with an increasing order. The images with smaller dissimilarity values are ranked higher in the retrieval list, and vice versa.

Next, we introduce our proposed re-ranking algorithm to refine the original distance measure with high time efficiency.

A. Preliminary

We deem that if x_q and x_p are similar, their retrieval results, especially the top-ranked results, are exactly or approximately the same in the expected situation. Let $\mathcal{N}_k(x_q)$ represents the *neighborhood set* of x_q achieved by the k -nearest neighbors rule. $\mathcal{N}_k(x_q)$ is a mathematical set that contains the top- k candidates in the ranking list of x_q . Then, the distance of two neighborhood sets $\mathcal{N}_k(x_q)$ and $\mathcal{N}_k(x_p)$ is measured by Jaccard distance as

$$d_J(x_q, x_p) = 1 - \frac{|\mathcal{N}_k(x_q) \cap \mathcal{N}_k(x_p)|}{|\mathcal{N}_k(x_q) \cup \mathcal{N}_k(x_p)|}, \quad (1)$$

where $|\cdot|$ calculates the cardinality of the input set, and $|\mathcal{N}_k(x_q)| = k$. With the distance measure defined by Equation 1, the original ranking list of a given query x_q can be re-ranked.

The re-ranking distance measure in Equation 1 is expected to achieve better performance than the original one, for it utilizes the additional contextual information as diffusion process does. However it also has many shortcomings.

- 1) The neighbors in the neighborhood set contribute equally. It is not a proper behavior, since the top-ranked neighbors are more likely to be true positive patterns. Assigning larger weights to the top-ranked neighbors, and increasing their effects on the re-ranking distance measure is more reasonable.
- 2) The re-ranking distance measure is defined between two sets. In the specific scenario of re-ranking, it is more convenient to define the distance measure on two vectorial features.
- 3) The neighborhood set is simply defined as the k-nearest neighbors, which cannot guarantee that the images from the same category could own similar neighborhood sets, especially when a certain amount of outliers also occupy top positions in the ranking list.

In the next section, Sparse Contextual Activation (SCA) is proposed to address these problems. The Jaccard distance defined in Equation 1 will serve as the baseline method, and we will compare it with our proposed SCA in terms of retrieval performance and running time. For notation clarity, we refer to the baseline method as Jaccard re-ranking below.

B. The Proposed Sparse Contextual Activation

The neighborhood set $\mathcal{N}_k(x_q)$ is converted to a vector representation by defining an binary *indicator function* $F_q = [F_{q,1}, F_{q,2}, \dots, F_{q,N}]$ as

$$F_{q,p} = \begin{cases} 1 & \text{if } x_p \in \mathcal{N}_k(x_q) \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

As we can see, the binary vector F_q shows whether a certain image x_p appears in the neighborhood set of x_q .

Based on the definition of the indicator function, the intersection and union set of $\mathcal{N}_k(x_q)$ and $\mathcal{N}_k(x_p)$ are interpreted as

$$\mathcal{N}_k(x_q) \cap \mathcal{N}_k(x_p) \iff MIN(F_q, F_p), \quad (3)$$

$$\mathcal{N}_k(x_q) \cup \mathcal{N}_k(x_p) \iff MAX(F_q, F_p), \quad (4)$$

where MIN (or MAX) calculates the element-wise minimum (or maximum) value for two input vectors of the same length. Thus we attain the cardinality for the intersection and the union set by computing the L_1 norm of their corresponding indicators as

$$|\mathcal{N}_k(x_q) \cap \mathcal{N}_k(x_p)| = \|MIN(F_q, F_p)\|_1, \quad (5)$$

$$|\mathcal{N}_k(x_q) \cup \mathcal{N}_k(x_p)| = \|MAX(F_q, F_p)\|_1. \quad (6)$$

Based on Equation 5 and Equation 6, we can rewrite the definition of Jaccard distance in Equation 1 as

$$\hat{d}_J(x_q, x_p) = 1 - \frac{\sum_{i=1}^N \min(F_{q,i}, F_{p,i})}{\sum_{i=1}^N \max(F_{q,i}, F_{p,i})}. \quad (7)$$

Now, the definition of Jaccard distance in Equation 1 has been successfully defined in vector space. That is to say, we do not use a set, but use an indicator function to represent the neighbors of an image. As a result, the Jaccard distance of two neighborhood sets can be easily achieved by vector comparison through Equation 7.

Note that the indicator function in Equation 2 also considers the neighbors equally, but it is easy to implement different weights by restricting the value of the original binary indicator vector in the unit interval $[0, 1]$. However, it may be misleading that $F_{q,p}$ is assigned to a certain constant between 0 and 1, since the role of $F_{q,p}$ is to indicate the membership of the image x_p in the neighborhood set $\mathcal{N}_k(x_q)$. In classical set theory, the membership of x_p in the set $\mathcal{N}_k(x_q)$ is exact (x_q either belongs or does not belong to the set). It seems difficult to generalize the binary indicator function in the unit interval $[0, 1]$ with rational explanations.

To tackle with the problem, we introduce the *fuzzy set theory*. In mathematics, fuzzy set is a set whose elements have degrees of membership determined by a *membership function*. Compared with the classical set theory, fuzzy set allows the gradual assessment of the membership of elements in a set. At last, we define the neighborhood set as a fuzzy set in fuzzy set theory, and the membership grade of x_p in the neighborhood set $\mathcal{N}_k(x_q)$ is determined by the corresponding membership function $F_{q,p} \in [0, 1]$.

The problem we face now is how to determine the membership grade of neighbors. A natural solution is to use the elements in $\mathcal{N}_k(x_q)$ to reconstruct x_q in the feature space with non-negative constraint, and the weights for reconstruction can be used as the membership grades. It can be formulated as

$$\begin{aligned} \min & \sum_q \|f(x_q) - \sum_{i|x_i \in \mathcal{N}_k(x_q)} F_{q,i} f(x_i)\|^2, \\ \text{s.t.} & \quad 1^T F_q = 1, F_q \succeq 0. \end{aligned} \quad (8)$$

Except for the non-negative constraint, this formulation almost shares the same perspective with Locally Linear Embedding (LLE) [25], a classical non-linear dimension reduction algorithm. LLE expects each data point and its neighbors lie on or close to a locally linear patch of the manifold, and the reconstruction weights are used for dimension reduction by a neighborhood preserving mapping.

However, in the specific scenario, the above solution may be not fit enough for visual re-ranking for three reasons. First, LLE usually presumes that there is sufficient data so that the data manifold is well-sampled, but retrieval task may also be needed in small datasets. Second, the requirement for real-time retrieval is usually declared, and it is time-consuming to solve the least square optimization problem for each query presented in Equation 8. On the other hand, the image is represented by a set of vectors instead of a single vector in some cases (e.g. shape analysis in [26], [27], [28]), which makes the above optimization problem more difficult to solve. Hence, we just apply the (truncated) Gaussian kernel to the pairwise distances with the given query x_q , and the membership function is

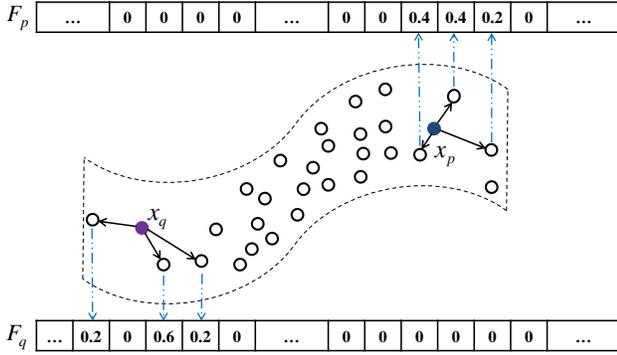


Fig. 1. The illustration of sparse contextual activation.

defined as

$$F_{q,p} = \begin{cases} \exp(-d(x_q, x_p)) & \text{if } x_p \in \mathcal{N}_k(x_q) \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

As a result, the top-ranked neighbors are assigned larger membership grades. F_q is subsequently L_1 normalized according to the sum-to-one constraint, and used for re-ranking through Jaccard distance Equation 7. In this way, the Jaccard distance is generalized to non-binary vectors.

In summary, the neighbors of x_q actually act as *Local Coordinate System*, and we can get a *Sparse Contextual Activation* denoted by F_q for x_q through Equation 9. The ‘‘Contextual’’ here indicates that F_q has non-zero values only in the index where the neighbors of x_q are located. Usually the cardinality of $\mathcal{N}_k(x_q)$ is much smaller than the size of the entire dataset, so the contextual activation is also a sparse vector. In Figure 1, we give an illustration of our proposed Sparse Contextual Activation (SCA).

C. Inverted Index Embedding

The proposed Sparse Contextual Activation (SCA) already manages assigning different weights to the neighbors at different positions in the ranking list, and gives a vector representation used for re-ranking. However, distance computation between a pair of SCAs is also waste of time, especially when the size of database becomes larger. Although the length of SCA is equal to the size of image database N , but the number of non-negative values in SCA is independent, only determined by the cardinality of neighborhood set. Considering the sparsity property of SCA, we introduce the inverted index [8] to reduce the computation complexity significantly.

Inverted index is a scalable indexing structure to store a large collection of images with their features. Although it has been applied to image retrieval successfully (e.g. [29], [30], [31]), it is the first work that introduces inverted index to visual re-ranking to my best knowledge now. Moreover, different from the usage of inverted index in Minkowski metric (Euclidean distance, Manhattan distance, etc.), we prove the feasibility of applying it in the metric of Jaccard distance theoretically.

Given two sparse contextual activations F_q and F_p , the

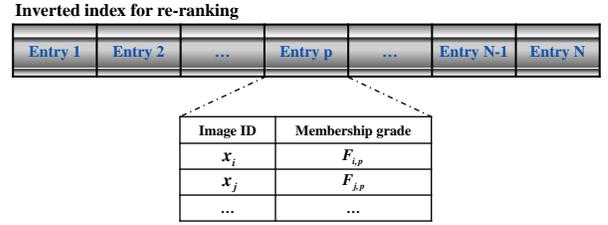


Fig. 2. The inverted index for re-ranking.

MIN operation can be computed as

$$\|MIN(F_q, F_p)\|_1 = \sum_{i|F_{q,i} \neq 0, F_{p,i} \neq 0} \min(F_{q,i}, F_{p,i}) + \sum_{i|F_{q,i} = 0} \min(F_{q,i}, F_{p,i}) + \sum_{i|F_{p,i} = 0} \min(F_{q,i}, F_{p,i}). \quad (10)$$

Since the sparse contextual activation only contains non-negative values, it is easy to find that last two items in Equation 10 are equal to zero. So the MIN operation can be achieved much more efficiently by

$$\|MIN(F_q, F_p)\|_1 = \sum_{i|F_{q,i} \neq 0, F_{p,i} \neq 0} \min(F_{q,i}, F_{p,i}). \quad (11)$$

By contrast, the computation of the MAX operation seems to be a bit complicated, since the item $\sum_{i|F_{q,i} = 0} \max(F_{q,i}, F_{p,i}) = \sum_{i|F_{q,i} = 0} F_{p,i}$ is not only determined by the query side. However, we also offer an efficient way to calculate the MAX operation. Note that

$$\|F_q\|_1 + \|F_p\|_1 = \|MIN(F_q, F_p)\|_1 + \|MAX(F_q, F_p)\|_1.$$

For two L_1 normalized sparse contextual activation, we can get

$$\begin{aligned} \|MAX(F_q, F_p)\|_1 &= 2 - \|MIN(F_q, F_p)\|_1, \\ &= 2 - \sum_{i|F_{q,i} \neq 0, F_{p,i} \neq 0} \min(F_{q,i}, F_{p,i}). \end{aligned} \quad (12)$$

In summary, our structure of inverted index is built as follows: (1) It has N entries as Figure 2 shows, where N is the size of database. Each entry relates to an image that acts as a base for activation. (2) For each entry p , we store the IDs of images whose neighborhood sets contain x_p and the corresponding membership grades. In other words, x_p owns non-zero membership grades in these neighborhood sets. (3) When re-ranking for the query x_q , distance computation can be conducted in a much smaller space using Equation 11 and Equation 12 as inverted index usually does.

Different from the inverted index on local descriptors in the paradigm of Bag of Words (BoW) model [4], our inverted index for re-ranking can be assumed as the second-level index, which is used in the level of images.

D. Local Consistency Enhancement

The images from the same category are expected to own the same neighborhood set, so that the distances between their sparse contextual activations are small. However, the

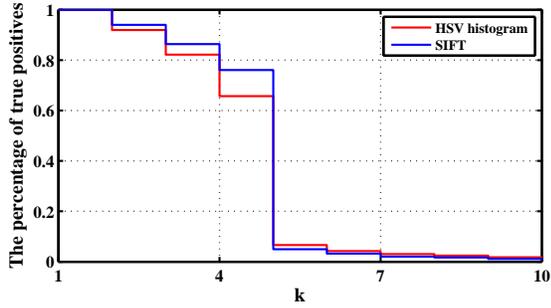


Fig. 3. The distribution of the percentage of true positives at the k -th position in the ranking list for all the queries.

neighborhood set is determined using k -nearest neighbors rule, which may be the simplest and the most efficient principle. Indeed, one can identify the neighborhood set by using more sophisticated rules (*e.g.* dominant neighbors [13]), but it will increase the computational time dramatically. Any extra time cost is not what we want in the proposed algorithm.

Inspired by the query expansion [16], [17] and local relevance feedback [32] in information retrieval, we propose to enhance the local consistency in generating SCAs of images from the same category using a similar way. In more detail, we define the Local Consistency Enhancement (LCE) on the sparse contextual activation as

$$F_q := \frac{1}{|\mathcal{N}_k(x_q)|} \sum_{i \in \mathcal{N}_k(x_q)} F_i, \quad (13)$$

by speculating that the images which are located very high in the initial ranking list of the query x_q are from the same category as the query.

In order to confirm our conjecture, an empirical analysis is performed on Ukbench dataset [8]. We plot the percentage of true positives at the k -th position in the ranking list for all the queries. It can be seen that the percentage of true positives decreases with the value of k increasing, and it is extremely high when k is small. Note that several false positives also exist at smaller value of k , which is known as *query drift*, but the small percentage of false positives will not impair the whole performance too much.

Based on the above analysis, the proposed LCE is an unsupervised method that does not need to know the labels. Compared with sparse contextual activation that considers the information from the direct neighbors of x_q , LCE is defined to consider the neighbors of the second order (*i.e.* the neighbors of the neighbors of x_q). As a result, LCE is more likely to include noise and outliers. Hence, we restrict the size of the neighborhood set used in LCE to a much smaller value in the experiments.

In order to distinguish the two neighborhood sets used in sparse contextual activation (Equation 9) and local consistency enhancement (Equation 13), we denote the size of the former one as k_1 and that of the latter one as k_2 below.

IV. SPARSE CONTEXTUAL ACTIVATION WITH MULTIPLE DISTANCE MEASURES

The proposed Sparse Contextual Activation presented above deals with only one input distance measure. In this section, we present how to introduce our method into rank aggregation for multiple input distance measures. We take two input distance measures as an example.

Let f^α and f^β denote two feature extractor functions for images, and d^α and d^β represent the corresponding distance functions respectively. Then for the image x_q with the feature extractor f^α (or f^β), we can also achieve the neighborhood set $\mathcal{N}_k^\alpha(x_q)$ (or $\mathcal{N}_k^\beta(x_q)$).

High Set and Low Set. Given a query image x_q and its two neighborhood sets $\mathcal{N}_k^\alpha(x_q)$ and $\mathcal{N}_k^\beta(x_q)$. We call the intersection set of the two neighborhood sets *High Set*, which can be described as

$$\mathcal{N}_k^{(H)}(x_q) = \mathcal{N}_k^\alpha(x_q) \cap \mathcal{N}_k^\beta(x_q), \quad (14)$$

and call the union set *Low Set* as

$$\mathcal{N}_k^{(L)}(x_q) = \mathcal{N}_k^\alpha(x_q) \cup \mathcal{N}_k^\beta(x_q). \quad (15)$$

The “high” indicates that the images in the high set are more likely to be true positives, since the images in the set occupy top positions in two ranking lists achieved by two distance functions d^α and d^β . By contrast, the low set may contains more false positives, for it is defined with a looser constraint, *i.e.* the images will be assigned to the low set as long as they are ranked high by either d^α or d^β .

In order to perform extremely fast rank aggregation later, we also generate the sparse contextual activations for high set and low set similarly. In more detail, let F_q^α (or F_q^β) denote the sparse contextual activation of x_q computed using Equation 9 under the distance measure d^α (or d^β). The membership functions of high set and low set are defined respectively based on Equation 3 and Equation 4

$$\mathcal{N}_k^{(H)}(x_q) \iff F_q^{(H)} = \text{MIN}(F_q^\alpha, F_q^\beta), \quad (16)$$

$$\mathcal{N}_k^{(L)}(x_q) \iff F_q^{(L)} = \text{MAX}(F_q^\alpha, F_q^\beta). \quad (17)$$

As a result, two sparse contextual activations $F_q^{(H)}$ and $F_q^{(L)}$, corresponding with high set and low set respectively, are achieved for x_q .

One can find that we perform feature fusion naturally by defining the two sets, so the unsupervised estimation about the weights of two individual features during fusion procedure are avoided. As we all know, different features are often in different scales or measured by different statistics, which makes the weight learning in feature fusion difficult, especially using an unsupervised way. Our proposed sparse contextual activation for feature fusion is simple yet delicate. It does not require any complicated learning or optimization methods.

Distance Fusion. After getting two sparse contextual activations for each image in the database, we define the distance measure for rank aggregation as

$$\hat{d}_J(x_q, x_p) = 1 - \frac{1}{2} \sum_{I=H,L} \frac{\sum_{i=1}^N \min(F_{q,i}^{(I)}, F_{p,i}^{(I)})}{\sum_{i=1}^N \max(F_{q,i}^{(I)}, F_{p,i}^{(I)})}. \quad (18)$$

The contribution of high set and low set are treated equally, which avoids the parameter tuning at the stage. It seems that the individual performance of high set is better than low set, due to its high percentage of true positives. The true fact is that the definition of high set may be too strict in the specific scenario of re-ranking. In some cases (especially when parameter k_1 is not large enough), a much small number of images will be assigned to the high set, which deprives its discriminative ability. However, with k_1 increasing, the individual performance of high set will surpass low set after a certain threshold without doubt. We will experimentally analyse the performance difference between high set and low set in Section V-E, and prove that a simple linear combination of them is better than using either one only.

It should be mentioned that the inverted index (Section III-C) and local consistency enhancement (Section III-D) can also be introduced into this feature fusion paradigm. The two additional methods can improve the efficiency and accuracy remarkably.

V. EXPERIMENTS

In this section, we will evaluate the performance of the proposed Sparse Contextual Activation (SCA) for various visual re-ranking task. The datasets we use are Ukbench image dataset [8], YALE face dataset B [33], Princeton Shape Benchmark (PSB) [34], Watertight Models track of SHape REtrieval Contest 2007 dataset (WM-SHREC07) [35], and MPEG-7 shape dataset [36]. The discussion about the parameter and the analysis of the algorithm complexity are presented in Section V-E and Section V-F respectively.

A. 3D Object Retrieval

In this section, we show the application of the proposed SCA to 3D object retrieval on the well-known Princeton Shape Benchmark (PSB) [34] and Watertight Models track of SHape REtrieval Contest 2007 dataset (WM-SHREC07) [35].

The PSB benchmark contains 1,804 3D polygonal models, which are divided into training set and testing set with 907 models each. Following the common settings, only the testing set is used to evaluate the performance of 3D object retrieval. The testing set is spilt into 92 categories, and the number of models per category ranges from 4 to 50.

SHape REtrieval Contest (SHREC) is the most authoritative competition for evaluating the effectiveness of 3D object retrieval algorithms. It will be held each year, involving multiple tracks, such as sketch-based 3D retrieval, textured 3D retrieval, *etc.* In this paper, WM-SHREC07 is chosen, which consists of 400 watertight mesh models that are evenly distributed into 20 classes. The models exhibit sufficient and diverse variation, from pose change to shape variability in the same semantic category.

Four evaluation metrics are adopted to assess the retrieval performance, listed as follows:

- Nearest Neighbor (NN): the percentage of the closest matches that belongs to the same class as the query.
- First Tier (FT): the recall for the top $C - 1$ matches in the ranked list, where C is the number of shapes in the category that query belongs to.
- Second Tier (ST): the recall for the top $2(C - 1)$ matches in the ranked list, where C is the number of shapes in the category that query belongs to.
- Discounted Cumulative Gain (DCG): a statistic that attaches more importance to the correct results near the front of the ranked list than the correct results at the end of the ranked list, under the assumption that a user is more likely to consider the retrieved candidates in the front of the list.

Please refer to [34] for more details about the definition of NN, FT, ST and DCG if needed. The values of all the aforementioned metrics range from 0 to 1, and larger values indicates better performance.

In the 3D object retrieval task, we use two view-based baseline methods (one is Vector of Aggregated Local Descriptor (VLAD) [37], and the other one is deep feature learned with Convolutional Neural Network (CNN)). For VLAD, we use the same pipeline as [38]. The number of depth views is set to 64, and the codebook size is 2048. For the deep feature, we utilize the descriptor proposed in XXXXXX.

We first demonstrate the performance of different re-ranking algorithms using the same baseline methods (*i.e.*, VLAD and deep feature) in Table I, and FT is chosen as the evaluation metric. When computing SCA, we fix the size of neighborhood set k_1 to 10 for PSB dataset and 17 for WM-SHREC07 dataset, and the parameter k_2 for local enhancement to 4. We report the performances of some typical re-ranking algorithms (Self diffusion (SD) [19], Tensor Product Graph (TPG) [13], Locally Constrained Diffusion Process (LCDP) [12]) in the optimal parameter setup.

It can be observed that our proposed SCA achieves the best performance among all the compared re-ranking methods. SD obtains the worst results since it loses the constraint of “locality” that used in LCDP and TPG. The percent gain in performance of LCDP and TPG in PSB dataset is not as large as in WM-SHREC07. It can be explained in two aspects: the first one is that the baseline in PSB dataset is much lower than that in WM-SHREC07 dataset, which means that it will include more noise and outliers in the neighborhood set; the second one lies that the number of objects per category in PSB dataset is various. By contrast, the distribution of objects in WM-SHREC07 is exactly balanced, *i.e.* 20 objects are assigned to one category. The imbalanced distribution of objects in PSB dataset makes kNN rule difficult to generalize well. Nevertheless, our proposed SCA also performs stably and well in both datasets although kNN is used to define the neighborhood set.

The comparison with other state-of-the-art algorithms are presented in Table II. As we can see, our proposed SCA achieves the new state-of-the art performance among all other algorithms for all the four evaluation metrics by fusion VLAD and deep feature in both datasets. PANORAMA is one of the most representative 3D shape descriptors in recent years, and our proposed SCA outperforms it by 8.4% in NN, 22.1% in FT and 20.7% in ST in PSB dataset. In [32], the performance of PANORAMA is improved by large margins using Local Relevance Feedback (LRF). LRF also exploits the contextual contribution as SCA, and the superior results reveal the

Methods	PSB dataset				WM-SHREC07 competition			
	VLAD	Gain	Deep feature	Gain	VLAD	Gain	Deep feature	Gain
Baseline	0.575	-	0.572	-	0.708	-	0.783	-
SD [19]	0.576	+0.17%	0.582	+1.75%	0.711	+0.42%	0.760	-2.94%
TPG [13]	0.613	+6.61%	0.598	+4.55%	0.761	+7.49%	0.839	+7.15%
LCDP [12]	0.617	+7.30%	0.603	+5.42%	0.766	+8.19%	0.843	+7.66%
SCA	0.660	+14.78%	0.617	+7.87%	0.795	+12.23%	0.876	+11.88%

TABLE I
THE PERFORMANCE COMPARISON IN FT OF DIFFERENT RE-RANKING ALGORITHMS ON THE PSB DATASET.

Methods	PSB dataset				WM-SHREC07 competition			
	NN	FT	ST	DCG	NN	FT	ST	DCG
LFD [39]	0.657	0.380	0.487	0.643	0.923	0.526	0.662	-
Tabia <i>et al.</i> [40]	-	-	-	-	0.853	0.527	0.639	0.719
DESIRE [41]	0.665	0.403	0.512	0.663	0.917	0.535	0.673	-
tBD [42]	0.723	-	-	0.667	-	-	-	-
Covariance [43]	-	-	-	-	0.930	0.623	0.737	0.864
2D/3D Hybrid [44]	0.742	0.473	0.606	-	0.955	0.642	0.773	-
PANORAMA [32]	0.753	0.479	0.603	-	0.957	0.673	0.784	-
PANORAMA + LRF [32]	0.752	0.531	0.659	-	0.957	0.743	0.839	-
3DVFF [45]	-	-	-	0.841	-	-	-	-
SCA+Deep	0.773	0.617	0.740	0.800	0.952	0.876	0.951	0.957
SCA+VLAD	0.803	0.660	0.786	0.826	0.955	0.795	0.911	0.938
SCA+VLAD+Deep	0.837	0.700	0.810	0.850	0.990	0.900	0.956	0.972

TABLE II
THE PERFORMANCE COMPARISON WITH OTHER STATE-OF-THE-ART ALGORITHMS ON THE PSB DATASET AND WM-SHREC07 DATASET.

effectiveness of our proposed method.

Among the all compared methods, 2D/3D Hybrid [44], PANORAMA [32] and 3DVFF [45] consider the fusion of multiple complementary features in the hope of more robust distance measure. 2D/3D Hybrid just simply concatenates 2D features based on depth buffers and 3D features based on spherical harmonics, and the dissatisfactory performance verifies the importance of designing more discriminative feature fusion methods. 3DVFF employs Multi-Feature Anchor Manifold that approximates multiple manifolds of heterogeneous features, which can be assumed as one variant of diffusion-based re-ranking methods. SCA also leads to a better retrieval accuracy than 3DVF in DCG (DCG is the only widely-accepted evaluation metric adopted in [45]).

B. Face Retrieval

We also evaluate the performance of SCA in face retrieval on YALE face dataset B [33]. YALE face dataset B is a standard benchmark widely used for face clustering, which is composed of face images sharing various poses and illumination conditions. In order to keep the comparison fair, we use the same subset and the same baseline method as generic diffusion framework [2]. Specifically, 15 subjects with 11 different conditions are gathered to generate a new dataset. Each image is normalized to 0-mean and 1-variance, and Euclidean distance between the vectorized representations is adopted to measure the pairwise dissimilarity directly. The evaluation metric is bull's eye score, which counts the recall before top-15 ranking list.

The baseline bull's eye score for the selected subset is 69.48%. Note that our goal is not achieving superior retrieval performance in this dataset, since better performance can be

Algorithm	Bull's eye score
Self Diffusion [19]	71.46%
Tensor Product Graph [13]	75.32%
Locally Constrained Diffusion Process [12]	75.59%
Generic diffusion [2]	77.30%
SCA	77.80%

TABLE III
THE PERFORMANCE COMPARISON WITH OTHER STATE-OF-THE-ART ALGORITHMS ON THE YALE FACE DATASET B.

obtained easily by using more discriminative descriptors, such as LBP [46], *etc.*. Instead, we focus on the performance improvement by our proposed re-ranking method. In [2], a classical branch of re-ranking methods called diffusion process is summarized, so the comparison with these methods will be more convincing. Similar to SCA, diffusion-based re-ranking methods capture the structure of the underlying data manifold by using the contextual information. The primary difference lies that they propagate the affinity values with random walks on a pre-defined graph in an iterative manner, while SCA does not need the iterative procedure that is often of great computational cost. We refer the readers to [2] for more details about diffusion process if necessary.

In Table III, the retrieval performances resulting from other state-of-the-art methods and our proposed method are reported. All the numerical results are borrowed from [2] or computed by using the public-available codes. When computing our sparse contextual activation, we manually set the size of neighborhood set k_1 to 4, and the parameter k_2 for local consistency enhancement to 5.

As we can see from the table, SCA outperforms Self Diffusion (SD) [19], Tensor Product Graph (TPG) [13], Locally

Descriptor	Re-ranking algorithm	Score
SC	-	86.80%
	GT [3]	92.91%
	SCA	95.21%
IDSC	-	85.40%
	CDM [1]	88.30%
	Indexing Re-Ranking [6]	91.56%
	GT [3]	91.61%
	Pairwise Recom. [47]	92.21%
	RL-Sim [14]	92.62%
	LCDP [12]	93.32%
	SSP [5]	93.35%
	SCA	93.44%
SC+IDSC	-	92.16%
	Co-T [22]	97.72%
	SCA	99.01%

TABLE IV

THE PERFORMANCE COMPARISON WITH OTHER STATE-OF-THE-ART ALGORITHMS ON MPEG-7 DATASET.

Constrained Diffusion Process (LCDP) [12] by 6.34%, 2.48% and 2.21%. As for the generic diffusion process, it enumerates all the variants of diffusion process by using four different types of initialization, six different types of transition and matrices and three different update schemes. The bull’s eye score of the generic version of diffusion process is chosen as the best performance for all the variants. However, SCA also achieves 77.80%, which is the best performance among all re-ranking methods, including the generic diffusion.

C. Shape Retrieval

Finally, we test our method for shape retrieval on MPEG-7 dataset [36]. It consists of 1,400 binary images divided into 70 categories evenly. Bull’s eye score is used to evaluate the performance, which counts the recall in the top-40 ranking list. We use the same baseline methods as Co-transduction [22], where Shape Context (SC) [26] and Inner Distance Shape Context (IDSC) [27] are used as the raw descriptors.

The comparison with other state-of-the-art algorithms is presented in Table IV. In order to improve the readability of the comparison, the table is divided into three parts with each part using the same baseline method. The baseline of the direct sum of SC and IDSC is 92.16%. As we can see, our proposed SCA leads to the best performance in each part, and improves the baseline method by a large margin.

D. Natural Image Retrieval

We first demonstrate the performance of SCA for natural image retrieval in Ukbench image dataset [8], which consists of 10,200 images. The whole dataset is divided into 2,550 categories with only 4 images per category. Each image is used in turn as a query. The performance is measured by the average recall of the top four ranked images, referred as N-S score (maximum is 4). The limited ground-truth images per category makes it challenging to achieve good performance for context-based re-ranking methods.

We implement two baseline methods using a local feature and a holistic feature. For local feature, SIFT [48] descriptors are extracted at Hessian-affine [49] interest points, and

Algorithms (single-feature)					
[37]	[51]	[52]	[53]	[13]	[29]
3.47	3.53	3.56	3.56	3.61	3.64
[54]	[1]	SCA+HSV		SCA+SIFT	
3.67	3.68	3.56		3.69	
Algorithms (multi-feature)					
[55]	[56] ^a	[57]	[22]	[58]	[50] ^a
3.60	3.60	3.62	3.66	3.68	3.71
[21] ^a	[21] ^b	[56] ^b	[50] ^b	SCA+SIFT+HSV	
3.76	3.77	3.79	3.85	3.86	

TABLE V

THE PERFORMANCE COMPARISON WITH STATE-OF-THE-ART METHODS IN UKBENCH DATASET. THE DIFFERENCE BETWEEN [50]^a, [56]^a AND [50]^b, [56]^b IS THE METHODS WITH SUPERSCRIP “b” USE THE ADDITIONAL GRAPH FUSION ALGORITHM PROPOSED IN [21]. [21]^a AND [21]^b DENOTE GRAPH PAGERANK AND GRAPH DENSITY RESPECTIVELY.

RootSIFT [17] variant is used. We train the codebook of size 1M using K-means on the independent Flickr60k data [30]. The Bag-of-Words (BoW) feature is computed with hard-assignment and TF-IDF weighting [31]. For holistic feature, we use the 1000-dimensional HSV color histogram (20×10×5 bins for H, S, V components) following [21], [50]. The HSV histogram is first normalized by its L_1 norm, and finally each element of the histogram is square rooted. The baselines of SIFT and HSV histogram are 3.56 and 3.40 respectively.

We compare our proposed SCA with several state-of-the-art algorithms in Table V. For computing the sparse contextual activation, the size of neighborhood set k_1 is to 4. The parameter for local enhancement enhancement k_2 is set to 2 empirically.

In order to keep the comparison fair, we first compare those methods which use only single feature, but different extra improvements. These selected methods are Vector of Aggregated Local Descriptors (VLAD) [37], Triangulation embedding [51], Spatially Constrained Similarity Measure (SCSM) [52], Spatial Contextual Weighting (SCW) [53], Tensor Product Graph (TPG) [13], Burstiness [29], RNN re-ranking [54] and Contextual Dissimilarity Measure [1].

As shown is Table V, our proposed SCA improves the baseline of SIFT from 3.56 to 3.69 and improves the baseline of HSV histogram from 3.40 to 3.56, which makes a significant gain in performance. Among those methods, SCSM, TPG, RNN re-ranking and CDM follow the similar principle with our method, they both attach importance to the local context of a given image in the manifold. In [52], the performance of SCSM is improved by a kNN-based re-ranking algorithms from 3.52 to 3.56. Its percent gain in N-S score (1.1%) is much lower than our method, whose percent gains are 4.7% and 3.65% for HSV histogram and SIFT respectively. Moreover, TPG, RNN re-ranking and CDM adopt an iterative strategy to update the similarity measure, and also achieve inferior performances compared with SCA, which does not need to iterate until convergence. It can be drawn that SCA can get

better performance by exploiting a more reliable distance measure with higher time efficiency.

The performance comparison of the algorithms using multi-feature is also conducted in Table V. We list nearly all the feature fusion algorithms which report their N-S scores to my best knowledge now, including Co-indexing [55], Coupled Binary Embedding [56], Bayes [57], Co-transduction [22], CrDP [58], c-MI [50] and Graph Fusion [21].

Graph Fusion is a representative re-ranking algorithm which deals with multiple input features. The performances of two variants of Graph Fusion, Graph PageRank and Graph Desnity, are 3.76 and 3.77 respectively. By contrast, our proposed feature fusion strategy with SIFT and HSV histogram achieves **3.86 N-S score**. Considering that c-MI [50]^b and Coupled Binary Embedding [56]^b actually fuse three features (fuse SIFT and Color descriptor at the indexing level, and HSV histogram descriptor is also combined at the post-processing level), it seems unfair to compare them with SCA, since we actually utilize two features only. However, as can be seen from the table, the performance of SCA is also higher than both of them. Compared with our baseline methods used here, we improve the performance from 3.40 (HSV histogram) and 3.56 (SIFT) to 3.86 significantly. The superior performance demonstrates the discriminative power of our SCA-based feature fusion algorithm.

E. Discussion

Two important parameters are involved in our method: the size of neighborhood set k_1 in sparse contextual activation and the size of neighborhood set k_2 for local consistency enhancement. In this section, we will discuss their effect on the retrieval performance, and compare with the Jaccard re-ranking to show the performance improvement brought by SCA simultaneously.

The parameter k_1 . The performance of standard Jaccard distance (Equation 1) and the proposed sparse contextual activation (Equation 7) with regard to the value of k_1 are reported in Figure 4. We utilize the N-S score for Ukbench dataset, bull’s eye score for YALE face dataset B, First Tier for PSB dataset and WM-SHREC07 dataset to evaluate the retrieval performance.

It can be observed that SCA outperforms the standard Jaccard distance consistently when the size of neighborhood set k_1 varies, which demonstrates firmly that it is beneficial to bring into the weights to increase the importance of top-ranked images. On the other hand, we can find that when k_1 increase, the performance first increases and drops later after k_1 reaches a certain threshold. It is easy to understand the decreasing of performance when k_1 is relatively larger, owing to the fact that the percentage of false positives will increase in that case.

The parameter k_2 . In Figure 5, we discuss the effect the parameter used in local consistency enhancement by fixing the size of neighborhood set to 4 for Ukbench dataset, 4 for YALE dataset, 10 for PSB dataset and 17 for WM-SHREC07 dataset following Section V-D, Section V-B, Section V-A respectively.

Metric	Without LCE	With LCE
L_0	0.572	0.550
L_1	0.574	0.610
L_2	0.565	0.583
L_∞	0.387	0.545
χ^2	0.573	0.612
Hellinger	0.573	0.612
ours	0.584	0.617

TABLE VI
THE PERFORMANCE COMPARISON OF SCA UNDER DIFFERENT METRICS.

Note that when k_2 is equal to 1, the local consistency enhancement is not applied at all. Hence, the performances of SCA and SCA with local consistency enhancement are identical at the starting point of the curves in Figure 5.

It is observed that when local consistency enhancement is applied by using a proper factor k_2 , the performance of sparse contextual activation can be improved further. Such a phenomenon demonstrates our previous analysis in Section III-D. It should also be mentioned that if a much larger k_2 is used, the performance will decrease without question. Since local consistency enhancement considers the contributions from the neighbors of the second order, it is more easy to include the negative effects of noise and outliers compared with only using the direct neighbors (*i.e.*, the neighbors of the first order). However, we find local consistency enhancement of great help to improve the retrieval performance when k_2 is small.

Metric. Note that our sparse contextual activation is computed using Equation 7, a generalized version of Jaccard distance. Indeed, many other metrics can be used here despite the fact that these metrics cannot be interpreted well using the set theory, such as L_r distance¹, χ^2 distance² and Hellinger distance³.

In Table VI, we list the performances of sparse contextual activation with LCE or without LCE under different metrics in PSB dataset. The baseline method is the deep feature and First Tier is chosen as the evaluation measure. We can find that the metric defined in Equation 7 outperforms all the other evaluation metrics. The performance of the widely-used Euclidean distance is not satisfactory enough.

Discussion on rank aggregation. In Figure 6, we compare the performances of using low set alone, using high set alone and using the linear combination of high set and low set.

First, as the figure shows, the performances of high set and low set reach the peak at different values of k_1 , and this value for low set is usually smaller than that for high set. Second, the two curves that represents the performance of high set and low set have a intersection point, before which the performance of low set is superior to high set. At last, the linear combination of high set and low set is always better than using either one alone.

¹The L_r distance for two SCAs F_q and F_p is $(\sum_i^N |F_{q,i} - F_{p,i}|^r)^{\frac{1}{r}}$

²The χ^2 distance for two SCAs F_q and F_p is $\frac{1}{2} \sum_i^N \frac{(F_{q,i} - F_{p,i})^2}{F_{q,i} + F_{p,i}}$

³The Hellinger distance for two SCAs F_q and F_p is $\sum_i^N (\sqrt{F_{q,i}} - \sqrt{F_{p,i}})^2$

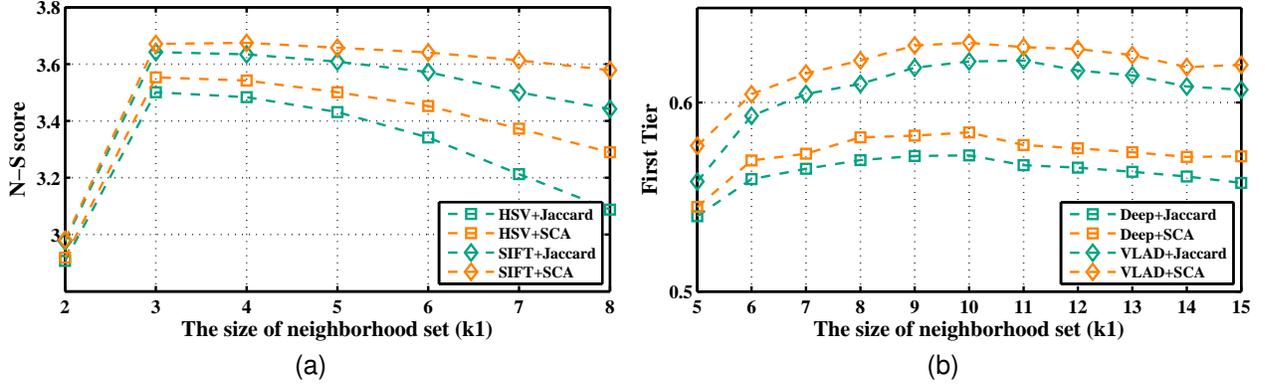


Fig. 4. The impact of the neighborhood set size on retrieval accuracy. N-S score for Ukbench (a), and First Tier for PSB (b) are presented.

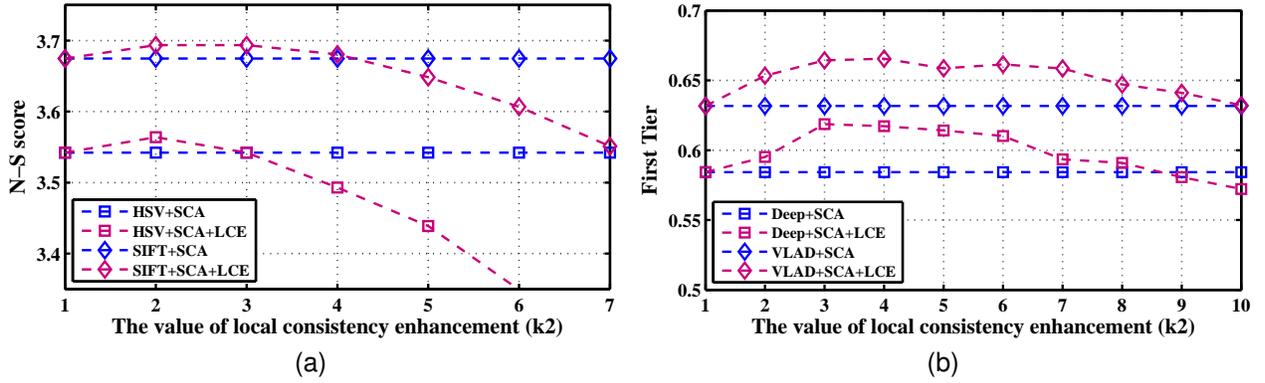


Fig. 5. The impact of the parameter k_2 in local consistency enhancement on retrieval accuracy. N-S score for Ukbench (a) and First Tier for PSB (b) are presented.

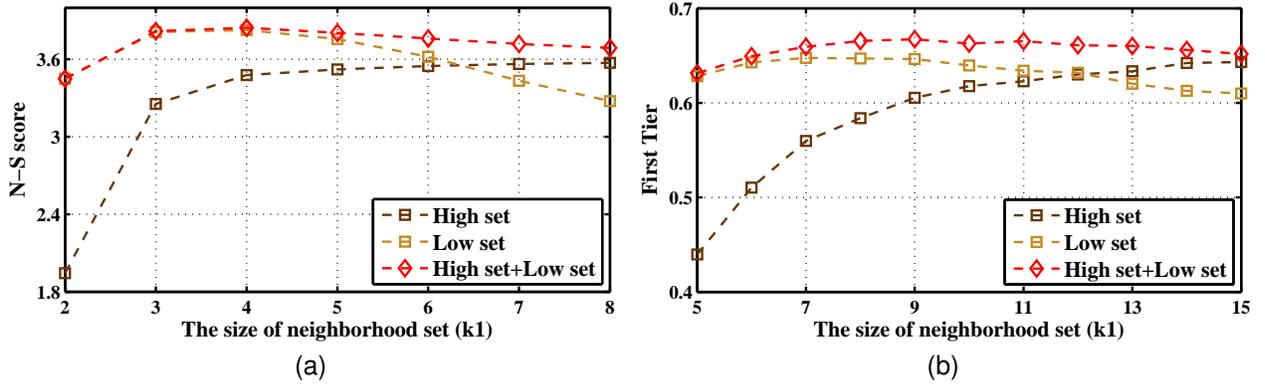


Fig. 6. The discussion about rank aggregation. The performances of high set, low set and their combination are presented. N-S score for Ukbench (a), First Tier for PSB (b) are presented.

F. Complexity Analysis

Considering that SCA is a post-processing procedure that focuses on re-ranking instead of ranking initialization, we only analyse the time efficiency and memory cost after the original ranking are given. The operation in the pre-processing, such as feature extraction, feature quantization, pairwise distance computation, *etc.*, are beyond the scope of our concern.

Time Complexity. Given an image collection with N elements and their pairwise distances known, the computation of the

sparse contextual activation for an given image x_q requires the computational complexity $O(k_1)$, where k_1 is the size of the neighborhood set. When local consistency enhancement determined by the parameter k_2 is applied, the computation time cost for SCA is equivalent to $O(k_2 \times k_1)$. We can find that the initiation of SCA is independent of the image collection size N , but related to the size of the local patch around x_q on the manifold. Usually, the values of k_1 and k_2 are much smaller than N , which indicates that the computation of SCA for a single image can be assumed to achieve in $O(1)$

Methods	Jaccard	TPG [13]	LCDP [12]	SCA
Ukbench	3.22s	3.81s	35.12ms	0.39ms
PSB	0.28s	75.27ms	12.22ms	0.17ms
MPEG-7	0.44s	0.26s	24.89ms	0.20ms

TABLE VII
THE COMPARISON OF AVERAGE RE-RANKING TIME.

efficiently.

As for re-ranking with SCA, the direct computation of the distance between x_q and all the other images in the database requires $O(2 \times N^2)$. However, the time complexity can be reduced significantly by using inverted index as presented in Section III-C. By embedded with inverted index, the *MIN* operation needs $O(M \times k_1)$ through Equation 11 in the worst case, where M denotes the number of images that have overlaps in neighborhood sets with x_q . *MAX* only needs $O(M)$ through Equation 12. In summary, the time complexity is reduced from $O(2 \times N^2)$ to $O(M \times (k_1 + 1))$ by introducing the inverted index. In fact, the value of M is also much smaller than N (e.g., for the HSV histogram in Ukbench dataset, the average value of M is merely 7.46).

All our experiments are carried out on a desktop machine with an Intel(R) Core(TM) i5 CPU (3.40 GHz) and 16 GB memory. In Ukbench dataset, the generation of SCA takes 0.14ms and local consistency enhancement takes 0.36ms per image. It takes 0.33s to build the inverted index for the whole dataset. For a given query, the average cost for re-ranking is only **0.39ms**.

Many previous re-ranking algorithms (e.g. [22], [12], [13]) are time-consuming due to their high time complexity ($O(N^3)$ in most cases). They usually need an iterative procedure to spread the affinity values on a huge graph which establishes the relationship among all the images. Table VII compares the average time for re-ranking of some typical algorithms, including Jaccard re-ranking defined in Equation 1. As the table shows, SCA reduces the time cost of Jaccard re-ranking by more than 8,000 times in Ukbench dataset when inverted index is applied. Compared with the two representative diffusion-based re-ranking methods, SCA is 90 times faster than LCDP, and 9,800 times faster than TPG. It indicates that our proposed re-ranking method has the potential for large scale re-ranking task.

Space Complexity. The scale of the inverted index is equivalent to the number of images in the database, which indicates the space complexity of SCA is $O(N)$.

For each entry in the inverted index used for re-ranking, we use 4 bytes to store one image index. 4 bytes (single format) are used to denote the activation for a certain image. On the Ukbench dataset, it only takes 65.2KB to save the inverted index in our implementation. Considering the big improvement on the computational efficiency, the minor extra memory cost can be tolerated.

VI. CONCLUSION

In this paper, we propose an extremely efficient algorithm called Sparse Contextual Activation (SCA) for visual re-

ranking. SCA is a merely single vector that encodes the contextual distribution for an image. The re-ranking procedure can be simply conducted by vector comparison using generalized Jaccard distance. For the first time, inverted index is introduced into re-ranking task, which makes it possible that re-ranking for a single query can be finished within one millisecond. Local Consistency Enhancement (LCE) is also developed to improve the performance of SCA further. The experimental results on Ukbench dataset (natural image), YALE dataset B (face), PSB dataset (3D object), WM-SHREC07 (3D competition) and MPEG-7 dataset (shape) demonstrate the effectiveness and efficiency of the proposed method.

Note that we design an inverted index for visual re-ranking, so an efficient method about the dynamical of insertion, deletion and modification should be taken into consideration. Meanwhile, since SCA is defined in the context, it requires for accurate contextual distribution if possible. So does it benefit from adding more extra images or artificial ghost points [12] to densify the feature space? We leave these issues for the future work.

ACKNOWLEDGMENT

This work was primarily supported by National Natural Science Foundation of China (NSFC) (No. 61222308), and in part by NSFC (No. 61173120), Program for New Century Excellent Talents in University (No. NCET-12-0217), Fundamental Research Funds for the Central Universities (No. HUST 2013TS115).

REFERENCES

- [1] H. Jegou, C. Schmid, H. Harzallah, and J. J. Verbeek, "Accurate image search using the contextual dissimilarity measure," *PAMI*, 2010. **1, 2, 8**
- [2] M. Donoser and H. Bischof, "Diffusion processes for retrieval revisited," in *CVPR*, 2013. **1, 2, 7**
- [3] X. Bai, X. Yang, L. J. Latecki, W. Liu, and Z. Tu, "Learning context-sensitive shape similarity by graph transduction," *PAMI*, 2010. **1, 2, 8**
- [4] J. Sivic and A. Zisserman, "Video google: A text retrieval approach to object matching in videos," in *ICCV*, 2003. **1, 4**
- [5] J. Wang, Y. Li, X. Bai, Y. Zhang, C. Wang, and N. Tang, "Learning context-sensitive similarity by shortest path propagation," *Pattern Recognition*, 2011. **1, 8**
- [6] D. Pedronette, J. Almeida, and R. Torres, "A scalable re-ranking method for content-based image retrieval," *Information Science*, 2014. **1, 8**
- [7] L. Luo, C. Shen, C. Zhang, and A. van den Hengel, "Shape similarity analysis by self-tuning locally constrained mixed-diffusion," *IEEE Trans. on Multimedia*, 2013. **1**
- [8] D. Nistér and H. Stewénus, "Scalable recognition with a vocabulary tree," in *CVPR*, 2006. **2, 4, 5, 6, 8**
- [9] P. Kotschieder, M. Donoser, and H. Bischof, "Beyond pairwise shape similarity analysis," in *ACCV*, 2009. **2**
- [10] A. Egozi, Y. Keller, and H. Guterman, "Improving shape retrieval by spectral matching and meta similarity," *TIP*, 2010. **2**
- [11] J. Wang and Y. Sun, "From one graph to many: Ensemble transduction for content-based database retrieval," *Knowledge-Based Systems*, 2014. **2**
- [12] 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*, 2009. **2, 6, 7, 8, 11**
- [13] X. Yang, L. Prasad, and L. J. Latecki, "Affinity learning with diffusion on tensor product graph," *PAMI*, 2013. **2, 5, 6, 7, 8, 11**
- [14] D. Pedronette and R. Torres, "Image re-ranking and rank aggregation based on similarity of ranked lists," *Pattern Recognition*, 2013. **2, 8**
- [15] D. Pedronette, O. Penatti, and R. Torres, "Unsupervised manifold learning using reciprocal knn graphs in image re-ranking and rank aggregation tasks," *Image Vision Computing*, 2014. **2**

- [16] O. Chum, A. Mikulik, M. Perdoch, and J. Matas, "Total recall ii: Query expansion revisited," in *CVPR*, 2011. [2](#), [5](#)
- [17] R. Arandjelovic and A. Zisserman, "Three things everyone should know to improve object retrieval," in *CVPR*, 2012. [2](#), [5](#), [8](#)
- [18] J. Philbin, O. Chum, M. Isard, J. Sivic, and A. Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *CVPR*, 2007. [2](#)
- [19] B. Wang and Z. Tu, "Affinity learning via self-diffusion for image segmentation and clustering," in *CVPR*, 2012. [2](#), [6](#), [7](#)
- [20] R. R. Coifman and S. Lafon, "Diffusion maps," *Applied and computational harmonic analysis*, 2006. [2](#)
- [21] S. Zhang, M. Yang, T. Cour, K. Yu, and D. N. Metaxas, "Query specific fusion for image retrieval," in *ECCV*, 2012. [2](#), [8](#), [9](#)
- [22] X. Bai, B. Wang, C. Yao, W. Liu, and Z. Tu, "Co-transduction for shape retrieval," *TIP*, 2012. [2](#), [8](#), [9](#), [11](#)
- [23] Z.-H. Zhou and M. Li, "Tri-training: Exploiting unlabeled data using three classifiers," *Knowledge and Data Engineering, IEEE Transactions on*, 2005. [2](#)
- [24] Z.-H. Zhou and M. Li, "Semi-supervised regression with co-training," in *IJCAI*, 2005. [2](#)
- [25] S. T. Roweis and L. K. Saul, "Nonlinear dimensionality reduction by locally linear embedding," *Science*, 2000. [3](#)
- [26] S. Belongie, J. Malik, and J. Puzicha, "Shape matching and object recognition using shape contexts," *PAMI*, 2002. [3](#), [8](#)
- [27] H. Ling and D. W. Jacobs, "Shape classification using the inner-distance," *PAMI*, 2007. [3](#), [8](#)
- [28] J. Wang, X. Bai, X. You, W. Liu, and L. J. Latecki, "Shape matching and classification using height functions," *PRL*, 2012. [3](#)
- [29] H. Jegou, M. Douze, and C. Schmid, "On the burstiness of visual elements," in *CVPR*, 2009. [4](#), [8](#)
- [30] H. Jégou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *ECCV*, 2008. [4](#), [8](#)
- [31] L. Zheng, S. Wang, Z. Liu, and Q. Tian, "Lp-norm IDF for large scale image search," in *CVPR*, 2013. [4](#), [8](#)
- [32] 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*, 2010. [5](#), [6](#), [7](#)
- [33] A. S. Georghiadis, P. N. Belhumeur, and D. J. Kriegman, "From few to many: Illumination cone models for face recognition under variable lighting and pose," *PAMI*, 2001. [6](#), [7](#)
- [34] P. Shilane, P. Min, M. M. Kazhdan, and T. A. Funkhouser, "The princeton shape benchmark," in *SMI*, 2004. [6](#)
- [35] D. Giorgi, S. Biasotti, and L. Paraboschi, "Shape retrieval contest 2007: Watertight models track," *SHREC competition*, 2007. [6](#)
- [36] L. J. Latecki, R. Lakamper, and T. Eckhardt, "Shape descriptors for non-rigid shapes with a single closed contour," in *CVPR*, 2000. [6](#), [8](#)
- [37] H. Jégou, F. Perronnin, M. Douze, J. Sánchez, P. Pérez, and C. Schmid, "Aggregating local image descriptors into compact codes," *PAMI*, 2012. [6](#), [8](#)
- [38] H. Tabia, D. Picard, H. Laga, and P. H. Gosselin, "Compact vectors of locally aggregated tensors for 3d shape retrieval," in *3DOR*, 2013. [6](#)
- [39] D. Y. Chen, X. P. Tian, Y. T. Shen, and M. Ouhyoung, "On visual similarity based 3d model retrieval," *Comput. Graph. Forum*, 2003. [7](#)
- [40] H. Tabia, M. Daoudi, J.-P. Vandeborre, and O. Colot, "A new 3d-matching method of nonrigid and partially similar models using curve analysis," *PAMI*, 2011. [7](#)
- [41] D. V. Vranic, "Desire: a composite 3d-shape descriptor," in *IEEE International Conference on Multimedia and Expo*, 2005. [7](#)
- [42] M. Liu, B. C. Vemuri, S.-I. Amari, and F. Nielsen, "Shape retrieval using hierarchical total bregman soft clustering," *PAMI*, 2012. [7](#)
- [43] H. Tabia, H. Laga, D. Picard, and P.-H. Gosselin, "Covariance descriptors for 3d shape matching and retrieval," in *CVPR*, 2014. [7](#)
- [44] 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*, 2008. [7](#)
- [45] T. Furuya and R. Ohbuchi, "Fusing multiple features for shape-based 3d model retrieval," in *BMVC*, 2014. [7](#)
- [46] T. Ojala, M. Pietikäinen, and T. Mäenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *PAMI*, 2002. [7](#)
- [47] D. C. G. e. Pedronette and R. da S Torres, "Exploiting pairwise recommendation and clustering strategies for image re-ranking," *Information Sciences*, 2012. [8](#)
- [48] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *IJCV*, 2004. [8](#)
- [49] K. Mikolajczyk and C. Schmid, "Scale and affine invariant interest point detectors," *IJCV*, 2004. [8](#)
- [50] L. Zheng, S. Wang, Z. Liu, and Q. Tian, "Packing and padding: Coupled multi-index for accurate image retrieval," in *CVPR*, 2014. [8](#), [9](#)
- [51] H. Jégou and A. Zisserman, "Triangulation embedding and democratic aggregation for image search," in *CVPR*, 2014. [8](#)
- [52] X. Shen, Z. Lin, J. Brandt, S. Avidan, and Y. Wu, "Object retrieval and localindex-based image re-ranking with spatially-constrained similarity measure and k-nn re-ranking," in *CVPR*, 2012. [8](#)
- [53] X. Wang, M. Yang, T. Cour, S. Zhu, K. Yu, and T. X. Han, "Contextual weighting for vocabulary tree based image retrieval," in *ICCV*, 2011. [8](#)
- [54] D. Qin, S. Gammeter, L. Bossard, T. Quack, and L. van Gool, "Hello neighbor: Accurate object retrieval with k-reciprocal nearest neighbors," in *CVPR*, 2011. [8](#)
- [55] S. Zhang, M. Yang, X. Wang, Y. Lin, and Q. Tian, "Semantic-aware co-indexing for image retrieval," in *ICCV*, 2013. [8](#), [9](#)
- [56] L. Zheng, S. Wang, and Q. Tian, "Coupled binary embedding for large-scale image retrieval," *TIP*, 2014. [8](#), [9](#)
- [57] L. Zheng, S. Wang, W. Zhou, and Q. Tian, "Bayes merging of multiple vocabularies for scalable image retrieval," in *CVPR*, 2014. [8](#), [9](#)
- [58] B. Wang, J. Jiang, W. Wang, Z.-H. Zhou, and Z. Tu, "Unsupervised metric fusion by cross diffusion," in *CVPR*, 2012. [8](#), [9](#)