# Multi-feature Fusion for Image Retrieval Using Constrained Dominant Sets

Alemu Leulseged Tesfaye[a,*], Marcello Pelillo[a,b]

[a]*DAIS, Ca' foscari University of Venice, Via Torino 155, Mestre, Venezia, Italy*
[b]*ECLT, European Center for Living Technology, S. Marco 2940, 30124 Venezia, Italy*

## Abstract

Aggregating different image features for image retrieval has recently shown its effectiveness. While highly effective, though, the question of how to uplift the impact of the best features for a specific query image persists as an open computer vision problem. In this paper, we propose a computationally efficient approach to fuse several hand-crafted and deep features, based on the probabilistic distribution of a given membership score of a constrained cluster in an unsupervised manner. First, we introduce an incremental nearest neighbor (NN) selection method, whereby we dynamically select k-NN to the query. We then build several graphs from the obtained NN sets and employ constrained dominant sets (CDS) on each graph G to assign edge weights which consider the intrinsic manifold structure of the graph, and detect false matches to the query. Finally, we elaborate the computation of feature positive-impact weight (PIW) based on the dispersive degree of the characteristics vector. To this end, we exploit the entropy of a cluster membership-score distribution. In addition, the final NN set bypasses a heuristic voting scheme. Experiments on several retrieval benchmark datasets show that our method can improve the state-of-the-art result.

*Keywords:* Image retrieval, multi-feature fusion, diffusion process.

## 1. Introduction

The goal of semantic image search, or content-based image retrieval (CBIR), is to search for a query image from a given image dataset. This is done by computing image similarities based on low-level image features, such as color, texture, shape and spatial relationship of images. Variation of images in illumination, rotation, and orientation has remained a major challenge for CBIR. Scale-invariant feature transform (SIFT) [1] based local feature such as Bag of words (BOW) [2], [3], [4], has served as a backbone for most image retrieval processes. Nonetheless, due to the inefficiency of using only a local feature to describe the content of an image, local-global feature fusion has recently been introduced.

Multi-feature based CBIR attacks the CBIR problem by introducing an approach which utilizes multiple low-level visual features of an image. Intuitively, if the to-be-fused feature works well by itself, it is expected that its aggregation with other features will improve the accuracy of the retrieval. Nevertheless, it is quite hard to learn in advance the effectiveness of the to-be-fused features for a specific query image. Different methods have recently been proposed to tackle this problem [5], [6], [7]. Zhang et al. [6] developed a graph-based query specific fusion method, whereby local and global rank lists are merged with equal weight by conducting a link analysis on a fused graph.

Zheng et al. [7] proposed a score level fusion model called Query Adaptive Late Fusion (QALF) [7], in which, by approximating a score curve tail with a reference collected on irrelevant data, they able to estimate the effectiveness of a feature as negatively related to the area under the normalized curve. Yang *et al.* [5] used a mixture Markov model to combine given graphs into one. Unlike [6] where graphs are equally weighted, [5] proposed a method to compute a weight which quantifies the usefulness of the given graph based on a naive Bayesian formulation, which depends only on the statistics of image similarity scores.

However, existing multi-feature fusion methods have different drawbacks. For instance, [7], [6], [8], [9] heavily rely on a pre-calculated and offline stored data, which turns out to be computationally expensive when new images are constantly added to the dataset. On the other hand, Ensemble Diffusion (ED)[10] requires $O(n^3)$ to perform a similarity diffusion. In addition to that, its feature-weight computation approach is not a query specific.

Inspired by [7], in this work we present a novel and simple CBIR method based on a recently introduced constrained cluster notion. Our approach presents two main advantages. Firstly, compared to the state of the art methods, it can robustly quantify the effectiveness of features for a specific query, without any supervision. Secondly, by diffusing the pairwise similarity between the nearest neighbors, our model can easily avoid the inclusion of false positive matches in the final shortlist. Towards this end, we first dynamically collect the nearest neighbors to the query, therefore, for each feature, we will have a different number

---

*Corresponding author
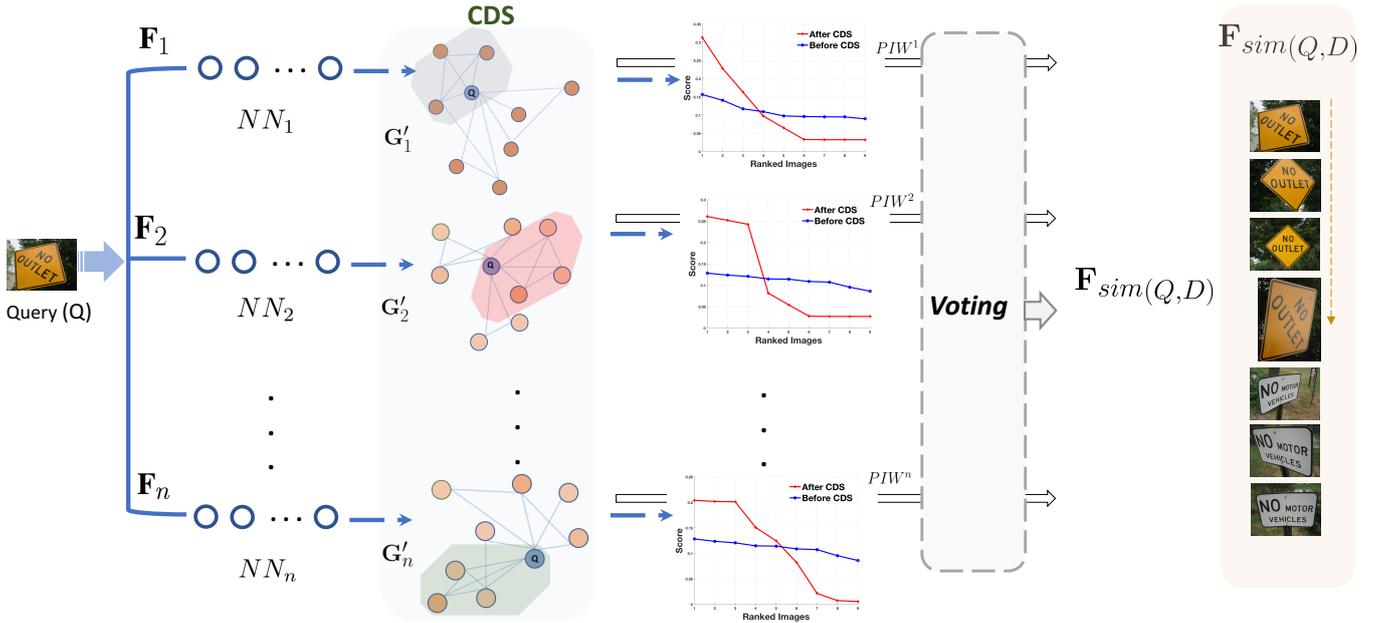*Email address:* `leulseged.alemu@unive.it` (Alemu Leulseged Tesfaye)

Figure 1: Overview of the proposed image retrieval framework. Based on the given features, $F_1, F_2, ...F_n$, we first incrementally collect the $NN's$ to the query $Q$, denoted as $NN_1, NN_2, ...NN_n$. Next, for each $NN$ we build the corresponding graph $G'_1, G'_2, ...G'_n$, and then, we apply $CDS$ on each graph to learn the $PIW$ of each feature, $PIW_1, PIW_2, ...PIW_n$, in the subsequent plot, the blue and red curves depict the ranked score of NN's before and after the application of CDS, respectively. Following, the final candidates, which come from each feature, pass through a voting scheme. Finally, using the obtained votes and PIW's we compute the final similarity, $F_{sim}(Q, D)$, between the query and the dataset images by equ. 10 .

of NNs. Subsequently, we set up the problem as finding a cluster from the obtained NNs, which is constrained to contain the given query image. To this end, we employ a graph-theoretic method called constrained dominant sets [11]. Here is our assumption: if the nearest neighbor to the query image is a false match, after the application of CDS its membership score to the resulting constrained cluster should be less than the fixed threshold $\zeta$, which leads us to detect and exclude the outliers. Furthermore, we introduce the application of entropy to quantify the effectiveness of the given features based on the obtained membership score. In contrast to [7], our method does not need any reference or external information to learn a query specific feature-weight. Fig. 1 shows the pipline of the proposed method.

In particular, we make the following contributions. 1) Compared to the previous work [6], [7], we propose a simple but efficient entropy-based feature effectiveness weighting system; in addition to that, we demonstrate an effective way of outlier or false nearest neighbor detection method. 2) Most importantly, our proposed model is a generic approach, which can be adapted to distinct computer vision problems, such as object detection and person re-identification. 3) We show that our unsupervised graph fusion model easily alleviates the asymmetry neighborhood problem.

This paper is structured as follows. In section 2 we briefly survey literature relevant to our problem, followed by technical details of the proposed approach in Sec. 3.

And, in Sec. 4 we show the performance of our framework on different benchmark datasets.

## 2. Related Work

CBIR has become a well-established research topic in the computer vision community. The introduction of SIFT feature plays a vital role in the application of BOW model on the image retrieval problem. Particularly, its robustness in dealing with the variation of images in scale, translation, and rotation provide a significant improvement in the accuracy of similar image search. Sivic et al. [2] first proposed BOW-based image retrieval method by using SIFT, in that, local features of an image are quantized to visual words. Since then, CBIR has made a remarkable progress by incorporating k-reciprocal neighbor [12], query expansion [13], [12], [14], large visual codebook [15], [16], diffusion process [5] [17] and spacial verification [15]. Furthermore, several methods, which consider a compact representation of an image to decrease the memory requirement and boost the search efficiency have been proposed. Jegou et al. [18] developed a Vector of Locally Aggregated Descriptor(VLAD), whereby the residuals belonging to each of the codewords are accumulated.

While SIFT-based local features have considerably improved the result of image search, it does not leverage the discriminative information encoded in the global feature of an image, for instance, the color feature yields a better representation for smooth images. This motivates the

2

introduction of multiple feature fusion for image retrieval. In [6], a graph-based query specific fusion model has been proposed, in which multiple graphs are combined and re-ranked by conducting a link analysis on a fused graph. Following, [5] developed a re-ranking algorithm by fusing multi-feature information, whereby they apply a locally constrained diffusion process (LCDP) on the localized NNs to obtain a consistent similarity score.

Although the aggregation of handcrafted local and global features has shown promising results, the advent of a seminal work by A.Krizhevsky *et al.* [19] in 2012 changed the focus of the computer vision community. Since then, convolutional neural network (CNN) feature has been used as a main holistic cue in different computer vision problems, including CBIR. Despite its significant improvement on the result of image retrieval, CNN feature still can not endow the demanded accuracy on different benchmark retrieval datasets, especially without the use of fine-tuning. Thus, aggregating graphs which are built from a hand-engineered and CNN-based image features has shown improvement in the accuracy of the retrieval [2], [18], [20], [21], [22], [23].

In addition to that, Yang *et al.* [5] applied a diffusion process to understand the intrinsic manifold structure of the fused graph. Despite a significant improvement on the result, employing the diffusion process on the final (fused) graph restricts the use of the information which is encoded in the pairwise similarity of the individual graph. Instead, our proposed framework applies CDS on each graph which is built from the corresponding feature. Thus, we are able to propagate the pairwise similarity score throughout the graph. Thereby, we exploit the underutilized pairwise similarity information of each feature and alleviate the negative impact of the inherent asymmetry of a neighborhood.

## 3. Proposed Method

### 3.1. Incremental NN Selection

In this subsection, we show an incremental nearest neighbor collection method to the given query image. We start with an intuitive clustering concept that similar nodes with common features should have an approximate score distribution, while outliers, or nodes which do not belong to a similar semantic class, have different score values. Accordingly, we propose a technique to search for the transition point where our algorithm starts including the outlier nodes. To this end, we examine how distinctive two subsequent nodes are in a ranked list of neighbors. Thus, we define a criterion called neighbors proximity coefficient($NPC$), which is defined as the ratio of two consecutive NNs in the given ranked list. Therefore, images are added only if the specified criterion is met, which is designed in such a way that only images that are very likely to be similar to the query image are added. Thereby, we are able to decrease the number of false matches to the query in the k-nearest neighbors set.

Given an initial ranked list $R$. And then, we define top-k nearest neighbors (kNN) to query $Q$ as

$$kNN(q,k) = \begin{cases} \text{Add } n_i & if \quad \frac{Sim(q,n_{i+1})}{Sim(q,n_i)} > \text{NPC} \\ 0 & otherwize \end{cases} \quad (1)$$

where $|kNN(q,k)| = k$, and $|.|$ represents the cardinality of a set.

$$kNN(q,k) = \{n_1, n_2, ...n_k\}, \quad where \ kNN(q,k) \subseteq R \quad (2)$$

### 3.2. Graph Construction

Different features, $F = F_1, F_2...F_n$, are extracted from images in the dataset D and the query image $Q$, where each feature encodes discriminative information of the given image in different aspects. We then compute the distance between the given images based on a distance metric function $d'(I_i, I_j)$, where $I_i$ and $I_j$ denote the given feature vector extracted from image $i$ and $j$ respectively. Following, we compute symmetric affinity matrices $A'_1$, $A'_2$, . . . $A'_n$ from each distance matrix $D_i$ using a similarity function $S(D_i)$. We then apply minimax normalization on each similarity matrix as: $A_i = \frac{V_\alpha^{ij} - min(V_\alpha)}{max(V_\alpha) - min(V_\alpha)}$, where $V_\alpha$ is a column vector taken from matrix $A'_i$, which comprises the pairwise similarity score between a given image $V_\alpha^i$ and images in the dataset $V^j$, which is denoted as $V_\alpha^{ij}$. Next, we build undirected edge-weighted graphs with no self-loops $G_1, G_2...G_n$ from the affinity matrices $A_1, A_2, ...A_n$, respectively. Each graph $G_n$ is defined as $Gn = (V_n, E_n, w_n)$, where $V_n = 1, ..., n$ is vertex set, $E_n \subseteq V_n \times V_n$ is the edge set, and $w_n : E \longrightarrow \mathbb{R}_+^*$ is the (positive) weight function. Vertices in G correspond to the given images, edges represent neighborhood relationships, and edge-weights reflect similarity between pairs of linked vertices.

### 3.3. PIW Using Entropy of CDS

Since the nearest neighbor selection method heavily relies on the initial pairwise similarity, it is possible that the NN set easily includes false matches to the given query. This usually happens due to the lack of technics which consider the underlying structure of the data manifold, especially the inherent asymmetry of a neighborhood is a major shortcoming of such systems. For instance, although $Sim(n_i, q) = Sim(q, n_i)$, the nearest neighbor relationship between query $Q$ and image $n_i$ may not be symmetric, which implies that $m_i \in kNN(q, k)$ but $m_i \notin kNN(n_i, k)$. As demonstrated in the past retrieval works, the k-reciprocal neighbors [12] and similarity diffusion process [24] have been vastly taken as the optimal options to tackle this issue. However, the existing methods are not computationally efficient. In this work, we remedy the existing limitations using an unsupervised constrained clustering algorithm whereby we exploit the pairwise similarity

to find a cohesive cluster which incorporates the specified query.

### 3.3.1. Constrained Clustering for Coherent Neighbor Selection

Towards collecting true matches to the query image, we employ an unsupervised clustering algorithm on the top of the previous steps. Our hypothesis to tackle the asymmetry problem between the given query and its nearest neighbors is that images which are semantically similar to each other tend to be clustered in some feature space. As can be seen in the synthetic example (See Fig. 2), retrieved image $i_4$ and $i_6$ are outliers or false positives to the query image $Q$. We can confirm this by observing the common neighbors of $Q$ with $i_4$ and $i_6$. But due to the lack of contextual information, the system considers them as a true match (neighbor) to the query. In our proposed model, to attack this issue, we represent the set of $kNN$ as a graph $G'$ accordingly to subsection 3.2. Then, we treat outliers finding problem as an unsupervised clustering problem. We first convert graph $G'$ into a symmetric affinity matrix $A$, where the diagonal corresponding to each node is set to 0, and the $ij-th$ entry denotes the edge-weight $w_{ij}$ of the graph so that $A_{ij} \equiv A_{ji}$. Accordingly, given graph $G'$ and query $Q$, we cast detecting outliers from a given $NN$ set as finding the most compact and coherent cluster from graph $G'$, which is constrained to contain the query image $Q$. To this end, we adopt constrained dominant sets [11], [25], which is a generalization of a well known graph-theoretic notion of a cluster. We are given a symmetric affinity matrix $A$ and parameter $\mu > 0$, and then we define the following parametrized quadratic program

$$
\begin{aligned}
\text{maximize} \quad & f_Q^\mu(X) = X'(A - \mu\hat{\Gamma}_Q)X \\
& f_Q^\mu(X) = X'\hat{A}X \\
\text{subject to} \quad & X \in \Delta
\end{aligned}
\tag{3}
$$

where a prime denotes transposition and

$$
\Delta = \left\{ X \in R^n \ : \ \sum_{i=1}^n X_i = 1, \text{ and } X_i \geq 0 \text{ for all } i = 1 \ldots n \right\}
$$

$\Delta$ is the standard simplex of $R^n$. $\hat{\Gamma}_Q$ represents $n \times n$ diagonal matrix whose diagonal elements are set to zero in correspondence to the query $Q$ and to 1 otherwise. And $\hat{A}$ is defined as,

$$
\hat{A} = A - \mu\hat{\Gamma}_Q = \begin{pmatrix} 0 & . & . & . \\ . & -\mu & . & . \\ . & . & -\mu & . \\ . & . & . & -\mu \end{pmatrix}
$$

where the dots denote the $ij$ th entry of matrix $A$. Note that matrix $\hat{A}$ is scaled properly to avoid negative values.

Let $Q \subseteq V$, with $Q \neq \emptyset$ and let $\mu > \lambda_{max}(A_{V \setminus Q})$, where $\lambda_{max}(A_{V \setminus q})$ is the largest eigenvalue of the principal submatrix of $A$ indexed by the element of $V \setminus q$. If $X$ is a local maximizer of $f_Q^\mu(X)$ in $\Delta$, then $\delta(X) \cap Q \neq \emptyset$, where, $\delta(X) = i \in V : X_i > 0$. We refer the reader to [11] for the proof.

The above result provides us with a simple technique to determine a constrained dominant set which contains the query vertex $Q$. Indeed, if $Q$ is the vertex corresponding the query image, by setting

$$
\mu > \lambda(A_{V \setminus Q}) \tag{4}
$$

we are guaranted that all local solutions of eq (3) will have a support that necessarily contains the query element. The established correspondence between dominant set (coherent cluster) and local extrema of a quadratic form over the standard simplex allow us to find a dominant set using straightforward continuous optimization techniques known as replicator dynamics, a class of dynamical systems arising in evolutionary game theory [21]. The obtained solution provides a principled measure of a cluster cohesiveness as well as a measure of vertex participation. Hence, we show that by fixing an appropriate threshold $\zeta$ on the membership score of vertices, to extract the coherent cluster, we could easily be able to detect the outlier nodes from the k-nearest neighbors set. For each $X^i$, $\zeta^i$ is dynamically computed as

$$
\zeta^i = \Lambda(1 - max(X^i) + min(X^i))/L \tag{5}
$$

where $max(X)$ and $min(X)$ denote the maximum and minimum membership score of $X^i$, respectively. $\Lambda$ is a scaling parameter and $L$ stands for length of $X^i$. Moreover, we show an effective technique to quantify the usefulness of the given features based on the dispersive degree of the obtained characteristics vector $X$.

### 3.3.2. PIW Using Entropy of Constrained Cluster.

Entropy has been successfully utilized in a variety of computer vision applications, including object detection [26], image retrieval [27] and visual tracking [28]. In this paper, we exploit the entropy of a membership-score of nodes in the constrained dominant set to quantify the usefulness of the given features. To this end, we borrowed the concept of entropy in the sense of information theory (Shannon entropy). We claim that the discriminative power of a given feature is inversely proportional to the entropy of the score distribution, where the score distribution is a stochastic vector. Let us say we are given a random variable $C$ with possible values $c_1, c_2, ...c_n$, according to statistical point of view the information of the event $(C = c_i)$ is inversely proportional to its likelihood, which is denoted by $I(C_i)$ and defined as

$$
I(C_i) = log\left(\frac{1}{P(c_i)}\right) = -log(p(c_i)). \tag{6}
$$

Thus, as stated by [29], the entropy of $C$ is the expected value of I, which is given as

$$
H(C) = -\sum_{i=1}^N P(c_i)log(P(c_i)). \tag{7}
$$

4

For each characteristic vector $X^i, X^{i+1}...X^z$, where $X^i = \{X^i_\mu, X^i_{\mu+1}...X^i_n\}$, we compute the entropy $H(exp(X^i))$. Each $X^i$ corresponds to the membership score of nodes in the CDS, which is obtained from the given feature $F^i$. Assume that the top NNs obtained from feature x are irrelevant to the query Q, thus the resulting CDS will only contain the constraint element Q. Based on our previous claim, since the entropy of a singleton set is 0, we can infer that the feature is highly discriminative. Although this conclusion is right, assigning a large weight to feature with irrelevant NNs will have a negative impact on the final similarity. To avoid such unintended impact, we consider the extreme case where the entropy is 0. Following, we introduce a new term $C_a$, which is obtained from the cardinality of a given cluster, $K_c$, as $Ca^i = \frac{K^i_c}{\sum_{i=1}^z K^i_c}$. As a result, we formulate the PIW computation from the additive inverse of the entropy $\varepsilon^i = 1 - H(X^i)$, and $C^i_a$, as:

$$PIW^i = \frac{\vartheta^i}{\sum_{i=1}^z \vartheta^i} \quad Thus, \sum_{i=1}^z PIW^i = 1 \qquad (8)$$

where $\vartheta^i = \varepsilon^i + C^i_a$, and $i$ represents the corresponding feature.

### 3.4. Naive Voting Scheme and Similarity Fusion

In this section, we introduce a simple yet effective voting scheme, which is based on the member nodes of k-nearest neighbor sets and the constrained dominant sets, let $NN_1, NN_2...NN_z$ and $CDS_1, CDS_2...CDS_z$ represent the $NN$ and $CDS$ sets respectively, which are obtained from $G'_1, G'_2...G'_z$. Let us say $\xi = 2(z-1) - 1$, and then we build $\xi$ different combinations of $NN$ sets, $\varphi_1, \varphi_2...\varphi_\xi$. Each $\varphi$ represents an intersection between $z - 1$ unique combinations of $NN$ sets. We then form a super-set $\varpi$ which contains the union of $\varphi$ sets, with including repeated nodes. Assume that $NNs = \{NN_1, NN_2, NN_3\}, \xi = 3$, thus each $\varphi$ set contains the intersection of two $NN$ sets as $\varphi_1 = \{NN_1 \cap NN_2\}, \varphi_2 = \{NN_1 \cap NN_3\}$ and $\varphi_3 = \{NN_2 \cap NN_3\}$. Hence the resulting $\varpi$ is defined as $\varpi = (\varphi_1 \ominus \varphi_2 \ominus \varphi_3)$, where $(.\ominus.)$ is an operator which returns the union of given sets, including repeated nodes. We have also collected the union of $CDS$ sets as $\omega = (CDS_1 \ominus CDS_2 \ominus CDS_3)$. Following, we compute $\kappa$ as $(\kappa = \varphi_1 \cap \varphi_2 \cap ...\varphi_\xi)$. Thereby we find super-sets $\varpi, \omega$ and $\kappa$. Next, we design three different counters, which are formulated to increment when the NN node appears in the corresponding super-sets. Based on the value obtained from each counter, we finally compute the vote scores for each $NN$ node to the query as $v_1 = v_1/\eta, v_2 = v_2/\theta$ and $v_3 = v_3/\iota$, where $\eta, \theta$ and $\iota$ are parameters which are fixed empirically. Note that the outlier detecting capability of our framework is encoded in the voting process. Thus, if a NN node $n_i$ is contained in more than one cluster, its probability to be given a large weight is higher. This is due to the number of votes it gets from each cluster.

### 3.4.1. Final Similarity.

After obtaining the aforementioned terms, we compute the final similarity as follows: say we are given $n$ features, $Q$ is the query image and $D$ denotes image dataset, then the initial similarity of $D$ to $Q$, with respect to feature $F_i, i = 1...n,$ ,is given as $S^{(i)}_{D,Q}$. Let $PIW^{(i)}_Q, i = 1...n$, encode the weight of feature $F_i$ for query $Q$, and then the final similarity score, $F_{sim(Q,D)}$, between $Q$ and $D$ is given as

$$N_s = \prod_{i=1}^k (S^{(i)}_{D,Q})^{PIW^{(i)}_Q} \qquad (9)$$

$$F_{sim(Q,D)} = \lambda N_s + (1 - \lambda) \sum_{\Omega=1}^\Psi v_\Omega \qquad (10)$$

where $\Psi = 3$, is the total number of voter sets. And $\lambda \in [0, 1]$ defines the penalty factor which penalizes the similarity fusion, when $\lambda = 1$ only $F_s$ is considered, otherwise, if $\lambda = 0$, only $v$ is considered.

## 4. Experiments

In this section, we present the details about the features, datasets and evaluation methodology we used along with rigorous experimental analysis.

### 4.1. Datasets and Metrics

To provide a thorough evaluation and comparison, we evaluate our approach on INRIA Holiday, Ukbench, Oxford5k and Paris6k datasets.

**Ukbench Dataset [30].** Contains 10,200 images which are categorized into 2,550 groups, each group consists of three similar images to the query which undergo severe illumination and pose variations. Every image in this dataset is used as a query image in turn while the remaining images are considered as dataset images, in "leave-one-out" fashion. As customary, we used the N-S score to evaluate the performance of our method, which is based on the average recall of the top 4 ranked images.

**INRIA Holiday Dataset [31].** Comprises 1491 personal holiday pictures including 500 query images, where most of the queries have one or two relevant images. Mean average precision (MAP) is used as a performance evaluation metric.

**Oxford5k Dataset [15].** It is one of the most popular retrieval datasets, which contains 5062 images, collected from flicker-images by searching for landmark buildings in the Oxford campus. 55 queries corresponding to 11 buildings are used.

**Paris6k Dataset [32].** Consists of 6392 images of Paris landmark buildings with 55 query images that are manually annotated.
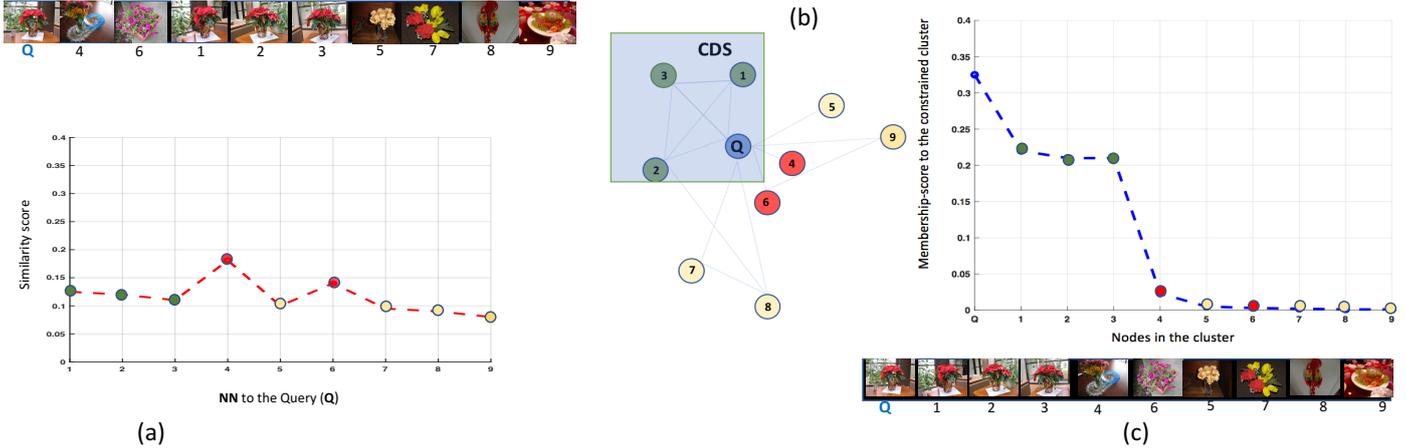
Figure 2: (a) Initial score distribution of the top k nearest neighbors to the query Q, green and red points denote the false-negative and false-posetive NNs. (b) Graph $G'$, built from the initial pairwise similarity of the k-nearest neighbor set. And the blue box contains the CDS nodes which are obtained by running CDS on graph $G'$. (c) The resulting constrained dominant set membership-score distribution.

## 4.2. Image Features

**Object Level Deep Feature Pooling (OLDFP)[33].** OLDFP is a compact image representation, whereby images are represented as a vector of pooled CNN features describing the underlying objects. Principal Component Analysis (PCA) has been employed to reduce the dimensionality of the compact representation. We consider the top 512-dimensional vector in the case of the Holiday dataset while considering the top 1024-dimensional vector to describe images in the Ukbench dataset. As suggested in [33], we have applied power normalization (with exponent 0.5), and l2 normalization on the PCA projected image descriptor.

**BOW.** Following [34], [7], we adopt Hamming Embedding [31]. SIFT descriptor and Hessian-Affine detector are used in feature extraction, and we used 128-bit vector binary signatures of SIFT. The Hamming threshold and weighting parameters are set to 30 and 16 respectively, and three visual words are provided for each key-point. Flickr60k data [31] is used to train a codebook of size 20k. We also adopt root sift as in [35], average IDF as defined in [36] and the burstiness weighting [37].

**NetVLAD [38].** NetVLAD is an end-to-end trainable CNN architecture that incorporates the generalized VLAD layer.

**HSV Color Histogram.** Like [5], [7], for each image, we extract 1000-dimensional HSV color histograms where the number of bins for H, S, V are 20, 10, 5 respectively.

## 4.3. Experiment on Holiday and Ukbench Datasets

As it can be seen in Fig.3(a), the noticeable similarity between the query image and the irrelevant images, in the Holiday dataset, makes the retrieval process challenging. For instance, (See Fig.3(a)), at a glance all images seem similar to the query image while the relevant are only the first two ranked images. Moreover, we can observe that

the proposed scheme is invariant to image illumination and rotation change. Table 2 shows that our method significantly improves the MAP of the baseline method [33] on Holiday dataset by 7.3 % while improving the state-of-the-art method by 1.1 %. Likewise, it can be seen that our method considerably improves the N-S score of the baseline method [33] on the Ukbench dataset by 0.15 while improving the state-of-the-art method by 0.03.

Furthermore, to show how effective the proposed feature-weighting system is, we have experimented by fusing the given features with and without PIW. Naive fusion (NF) denotes our approach with a constant PIW for all features used, thus the final similarity $F_s$ defined as $F_s = \frac{1}{k}(\prod_{i=1}^{k}(S_{D,Q}^{(i)}))$. In Fig.6 we have demonstrated the remarkable impact of the proposed PIW. As can be observed, our scheme effectively uplifts the impact of a discriminative feature while downgrading the inferior one. Note that in the PIW computation we have normalized the minimum entropy (See eq.8), thus its values range between 0 and 1. Accordingly, one implies that the feature is highly discriminative, while zero shows that the feature is indiscriminate.

In order to demonstrate that our scheme is robust in handling outliers, we have conducted an experiment by fixing the number of NNs (disabling the incremental NNs selection) to different numbers. As is evident from Fig.6, the performance of our method is consistent regardless of the number of $kNN$. As elaborated in subsection 3.3.1, the robustness of our method to the number of $k$ comes from the proposed outlier detection method. Since the proposed outliers detector is formulated in a way that allows us to handle the outliers, we are easily able to alleviate the false matches which are incorporated in the nearest neighbors set. This results in finding a nearly constant number of nearest neighbors regardless of the choice of $k$.

Figure 3: five relevant images to the query where the green and red frame indicate the True and False posetives to the query, respectively. **Top-row (a) and (b):** show the top five relevant images of our proposed method. **Bottom row (a) and (b):** show the top five relevant images obtained from a Naive fusion of several features.
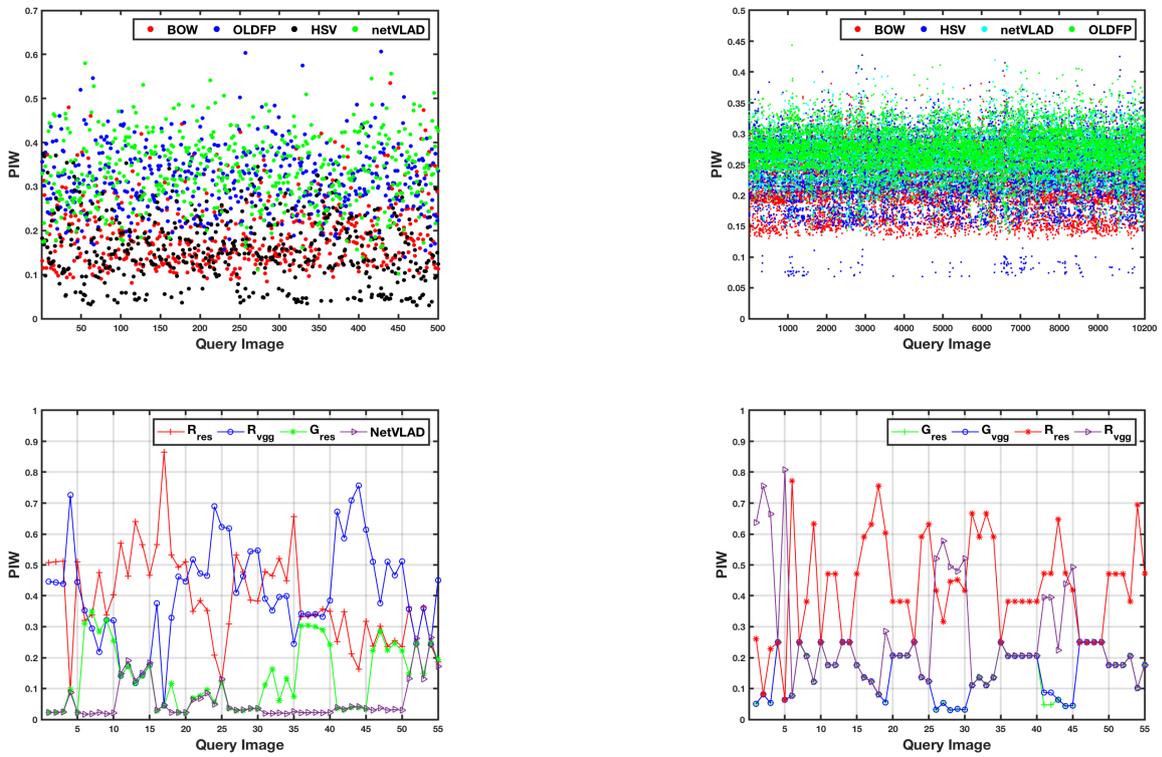


Figure 4: Feature positive-impact weights (PIW's) learned by our algorithm. Top-left, top-right, bottom-left, and bottom-right: on Holiday, Ukbench, Oxford5k and Paris6k datasets, respectively.

Table 1: The performance of baseline features on Holidays, Ukbench, Oxford5k and Paris6k datasets.

| Datasets | Metrics | NetVLAD [38] | BOW | OLDFP | HSV | $R_{res}$[24] | $G_{res}$[24] | $R_{vgg}$ [24] | $G_{vgg}$[24] |
|---|---|---|---|---|---|---|---|---|---|
| **Holidays** | MAP | 84 | 80 | 87 | 65 | - | - | - | - |
| **Ukbench** | N-S score | 3.75 | 3.58 | 3.79 | 3.19 | | - | - | - |
| **Oxford5k** | MAP | 69 | - | - | - | 95.8 | 87.7 | 93 | - |
| **Paris6k** | MAP | - | - | - | - | 96.8 | 94.1 | 96.4 | 95.6 |



Holiday Dataset

Ukbench Dataset

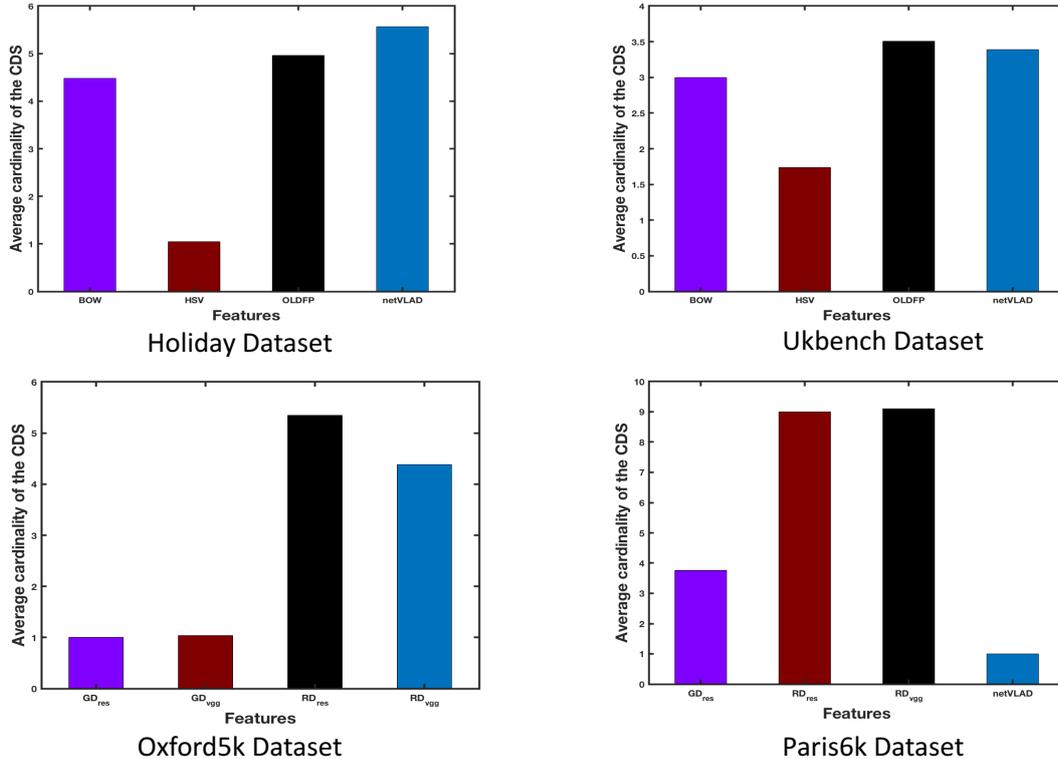Oxford5k Dataset

Paris6k Dataset

Figure 5: The cardinality of constrained dominant sets for the given features.

### 4.4. Experiment on Oxford5k and Paris6k Datasets

In the same fashion as the previous analysis, we have conducted extensive experiments on the widely used Oxford5k and Paris6k datasets. Unlike the Holiday and Ukbench datasets, we adapt affinity matrices which are obtained through a diffusion process on a regional *Resnet* and *VGG* representation [24], and they are denoted as $R_{res}$ and $R_{vgg}$ respectively, as well as affinity matrices $G_{res}$ and $G_{vgg}$ which are also obtained through a diffusion process on a global *Resnet* and *VGG* representation, respectively. Table 2 shows that the proposed method slightly improves the state-of-the-art result. Even if the performance gain is not significant, our scheme marginally achieves better MAP over the state-of-the-art methods. Furthermore, as shown in Fig 4, the proposed model learns the PIW of the given features effectively. Therefore, a smaller average weight is assigned to $G_{vgg}$ and $NetVLAD$ feature comparing to $R_{res}$ and $R_{vgg}$.

### 4.5. Robustness of Proposed PIW

As can be seen in Fig 4, for all datasets, our algorithm has efficiently learned the appropriate weights to the corresponding features. Fig. 4 shows how our algorithm assigns PIW in a query adaptive manner. In Holiday and Ukbench datasets, the average weight given to HSV feature is much smaller than all the other features used. Conversely, a large PIW is assigned to OLDFP and NetVLAD features. Nevertheless, it is evident that in some cases a large value of PIW is assigned to HSV and BOW features as well, which is appreciated considering its effectiveness on discriminating good and bad features in a query adaptive manner.

### 4.6. Impact of Parameters

To evaluate the robustness of our method we have performed different experiments by changing one parameter at a time. Thereby, we have observed that setting $\Lambda$ to a large value results in assigning insignificant PIW to indiscriminate features. The reason is that after the application of CDS, the cluster membership-score of the dissimilar

Table 2: Comparison among various retrieval methods with our method on benchmark datasets, where QALF is implemented with the same baseline similarities used in our experiments.

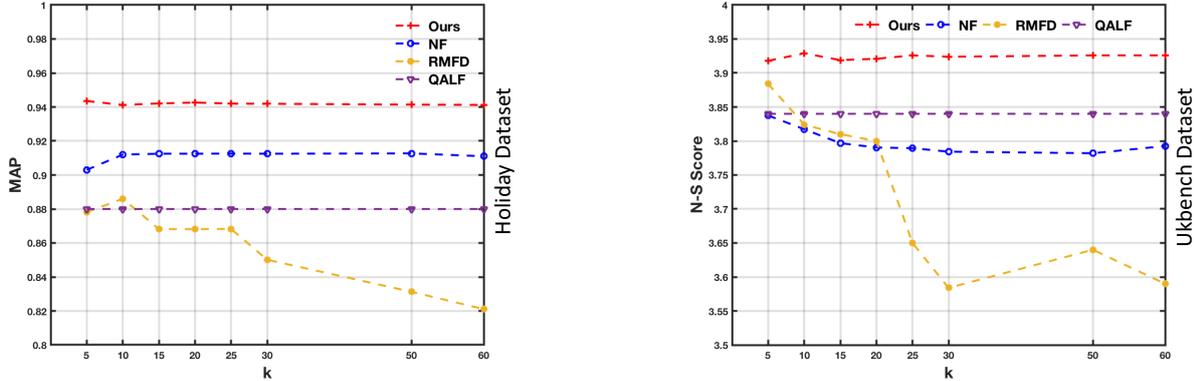| Datasets | Metrics | Baselines | QALF[7] | [5] | NF | ED[39] | [40] | [41] | [42] | [43] | Ours |
|----------|---------|-----------|---------|-----|-----|--------|------|------|------|------|------|
| **Ukbench** | N-S score | 3.79[33] | 3.84 | 3.86 | 3.86 | 3.93 | - | - | - | 3.76 | **3.94** |
| **Holiday** | MAP | 87[33] | 88 | 88 | 91 | 93 | 90 | 83 | 89 | 77 | **94** |
| **Oxford5k** | MAP | 95.8[24] | - | 76.2 | 94.4 | - | 89.1 | 79.7 | 81.4 | 67.6 | **96.2** |
| **Paris6k** | MAP | 96.8[24] | - | 83.3 | - | - | 91.2 | 83.8 | 88.9 | - | **97.4** |



Figure 6: Comparison with state-of-the-art fusion methods with respect to varying k. Naive Fusion (NF), Reranking by Multi-feature Fusion (RMFD) [5], and QALF [7].

images to the query will become smaller. Thus, since the threshold fixed to choose the true neighbors is tighter, the resulting constrained dominant set will be forced to yield a singleton cluster. As a result, we obtained a very small PIW due to the cardinality of the constrained-cluster. In addition to that, we observe that the MAP start to decline when $\lambda$ is set to a very large value (See. Fig 7, right).

### 4.7. Impact of Cluster Cardinality

On the Ukbench dataset, as can be observed in Fig. 5, the average cardinality of the constrained clusters which is obtained from HSV and BOW feature is 3 and 1.7, respectively. In contrast, for NetVLAD and OLDFP, the average cluster cardinality is 3.4 and 3.5, respectively . Similarly, in the case of the Holiday dataset, the cluster cardinality obtained from HSV feature is one while for BOW, NetVLAD and OLDFP is 4.5, 5 and 5.6, respectively. Thus, from this, we can draw our conclusion that the cardinality of a constrained dominant set, in a certain condition, has a direct relationship with the effectiveness of the given feature.

### 4.8. Computational Time

In Fig. 7 we depict the query time taken to search for each query image, red and blue lines represent our method and QALF, respectively. The vertical axis denotes the CPU time taken in seconds, and the horizontal axis shows the query images. As can be seen from the plot, the proposed framework is faster than the fastest state-of-the-art feature-fusion method [7]. As for time complexity, in our experiment we used a replicator dynamics to solve problem

(3), hence, for a graph with N nodes, the time complexity per step is $O(N^2)$, and the algorithm usually takes a few steps to converge, while that of [10] is $O(N^3)$. However, we note that by using the Infection-immunization algorithm [44] we can achieve even faster convergence as its per-step complexity would be linear in the number of nodes.

## 5. Conclusion

In this paper, we addressed a multi-feature fusion problem in CBIR. We developed a novel and computationally efficient CBIR method based on a constrained-clustering concept. In particular, we showed an efficient way of estimating a positive impact weight of features in a query-specific manner. Thus it can be readily used for feature combination. Furthermore, the proposed scheme is fully unsupervised, and can easily be able to detect false-positive NNs to the query, through the diffused similarity of the NNs. To demonstrate the validity of our method, we performed extensive experiments on benchmark datasets. Besides the improvements achieved on the state-of-the-art results, our method shows its effectiveness in quantifying the discriminative power of given features. Moreover, its effectiveness on feature-weighting can also be exploited in other computer vision problems, such as person re-identification, object detection, and image segmentation.

## References

## References

[1] D. G. Lowe, Distinctive image features from scale-invariant keypoints, International Journal of Computer Vision 60 (2) (2004)
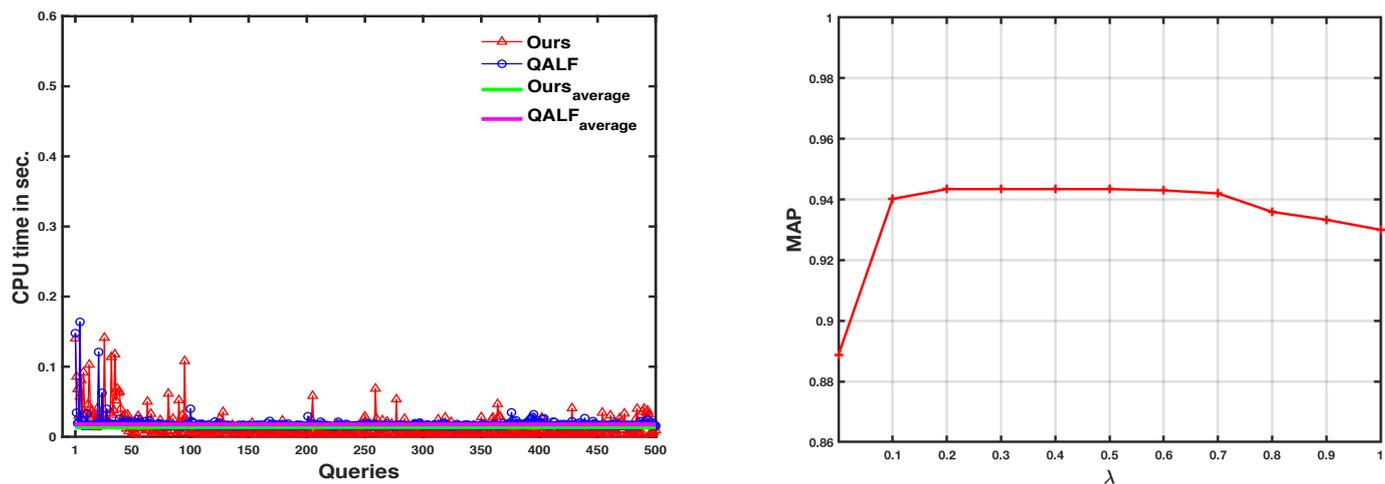
Figure 7: Left: Time complexity of our algorithm (red) and QALF[7] (blue) on Holiday dataset. Right: The impact of $\lambda$ on the retrieval performance, on Holiday dataset.

91–110.

[2] J. Sivic, A. Zisserman, Video google: A text retrieval approach to object matching in videos, in: 9th IEEE International Conference on Computer Vision (ICCV 2003), 14-17 October 2003, Nice, France, 2003, pp. 1470–1477.

[3] M. Jain, R. Benmokhtar, H. Jégou, P. Gros, embedding similarity-based image classification, in: International Conference on Multimedia Retrieval, ICMR '12, Hong Kong, China, June 5-8, 2012, 2012, p. 19.

[4] Y. Yang, F. Nie, D. Xu, J. Luo, Y. Zhuang, Y. Pan, A multimedia retrieval framework based on semi-supervised ranking and relevance feedback, IEEE Trans. Pattern Anal. Mach. Intell. 34 (4) (2012) 723–742.

[5] F. Yang, B. Matei, L. S. Davis, Re-ranking by multi-feature fusion with diffusion for image retrieval, in: 2015 IEEE Winter Conference on Applications of Computer Vision, WACV 2015, Waikoloa, HI, USA, January 5-9, 2015, 2015, pp. 572–579.

[6] S. Zhang, M. Yang, T. Cour, K. Yu, D. N. Metaxas, Query specific rank fusion for image retrieval, IEEE, Trans. Pattern Anal. Mach. Intell. 37 (4) (2015) 803–815.

[7] L. Zheng, S. Wang, L. Tian, F. He, Z. Liu, Q. Tian, Query-adaptive late fusion for image search and person re-identification, in: IEEE, CVPR, 2015, pp. 1741–1750.

[8] C. Deng, R. Ji, W. Liu, D. Tao, X. Gao, Visual reranking through weakly supervised multi-graph learning, in: IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013, 2013, pp. 2600–2607.

[9] S. Zhang, M. Yang, X. Wang, Y. Lin, Q. Tian, Semantic-aware co-indexing for image retrieval, in: IEEE International Conference on Computer Vision, ICCV 2013, Sydney, Australia, December 1-8, 2013, 2013, pp. 1673–1680.

[10] S. Bai, Z. Zhou, J. Wang, X. Bai, L. J. Latecki, Q. Tian, Ensemble diffusion for retrieval, in: IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017, 2017, pp. 774–783.

[11] E. Zemene, M. Pelillo, Interactive image segmentation using constrained dominant sets, in: Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII, 2016, pp. 278–294.

[12] D. Qin, S. Gammeter, L. Bossard, T. Quack, L. J. V. Gool, Hello neighbor: Accurate object retrieval with k-reciprocal nearest neighbors, in: 2011 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2011), 2011, pp. 777–784.

[13] O. Chum, A. Mikulík, M. Perdoch, J. Matas, Total recall II: query expansion revisited, in: The 24th IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Colorado Springs, CO, USA, 20-25 June 2011, 2011, pp. 889–896.

[14] G. Tolias, H. Jégou, Visual query expansion with or without geometry: Refining local descriptors by feature aggregation, Pattern Recognition 47 (10) (2014) 3466–3476.

[15] J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Object retrieval with large vocabularies and fast spatial matching, in: 2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition CVPR, 18-23 June 2007, Minneapolis, Minnesota, USA, 2007.

[16] Y. S. Avrithis, Y. Kalantidis, Approximate gaussian mixtures for large scale vocabularies, in: Computer Vision - ECCV 2012 - 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part III, 2012, pp. 15–28.

[17] E. Zemene, L. T. Alemu, M. Pelillo, Constrained dominant sets for retrieval, in: 23rd International Conference on Pattern Recognition, ICPR 2016, Cancún, Mexico, December 4-8, 2016, 2016, pp. 2568–2573.

[18] H. Jégou, F. Perronnin, M. Douze, J. Sánchez, P. Pérez, C. Schmid, Aggregating local image descriptors into compact codes, IEEE Trans. Pattern Anal. Mach. Intell. 34 (9) (2012) 1704–1716.

[19] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: NIPSAdvances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States., 2012, pp. 1106–1114.

[20] F. Perronnin, C. R. Dance, Fisher kernels on visual vocabularies for image categorization, in: 2007 IEEE Computer Society Conference on Computer Vision and Pattern Recognition CVPR, 18-23 June 2007, Minneapolis, Minnesota, USA, 2007.

[21] H. Jégou, A. Zisserman, Triangulation embedding and democratic aggregation for image search, in: 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Columbus, OH, USA, June 23-28, 2014, pp. 3310–3317.

[22] T. Do, Q. D. Tran, N. Cheung, Faemb: A function approximation-based embedding method for image retrieval, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR , Boston, MA, USA, June 7-12, 2015, 2015, pp. 3556–3564.

[23] Y. Kalantidis, C. Mellina, S. Osindero, Cross-dimensional weighting for aggregated deep convolutional features, in: Computer Vision - ECCV 2016 Workshops - Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part I, 2016, pp. 685–701.

[24] A. Iscen, G. Tolias, Y. S. Avrithis, T. Furon, O. Chum, Efficient diffusion on region manifolds: Recovering small objects with compact CNN representations, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Honolulu, HI, USA, July 21-26, 2017, 2017, pp. 926–935.

[25] E. Zemene, L. T. Alemu, M. Pelillo, Dominant sets for "constrained" image segmentation, CoRR abs/1707.05309.

[26] R. Sznitman, C. J. Becker, F. Fleuret, P. Fua, Fast object detection with entropy-driven evaluation, in: 2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland CVPR, OR, USA, June 23-28, 2013, 2013, pp. 3270–3277.

[27] T. Deselaers, T. Weyand, H. Ney, Image retrieval and annotation using maximum entropy, in: Evaluation of Multilingual and Multi-modal Information Retrieval, 7th Workshop of the Cross-Language Evaluation Forum, CLEF 2006, Alicante, Spain, September 20-22, 2006, Revised Selected Papers, 2006, pp. 725–734.

[28] L. Ma, J. Lu, J. Feng, J. Zhou, Multiple feature fusion via weighted entropy for visual tracking, in: 2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015, 2015, pp. 3128–3136.

[29] Shannon, A mathematical theory of communication., Bell Syst. Tech. J. (1948) 27, 379423.

[30] D. Nistér, H. Stewénius, Scalable recognition with a vocabulary tree, in: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition CVPR, 17-22 June 2006, New York, NY, USA, 2006, pp. 2161–2168.

[31] H. Jegou, M. Douze, C. Schmid, Hamming embedding and weak geometric consistency for large scale image search, in: Computer Vision - ECCV 2008, 10th European Conference on Computer Vision, Marseille, France, October 12-18, 2008, Proceedings, Part I, 2008, pp. 304–317.

[32] J. Philbin, O. Chum, M. Isard, J. Sivic, A. Zisserman, Lost in quantization: Improving particular object retrieval in large scale image databases, in: 2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition CVPR, 24-26 June 2008, Anchorage, Alaska, USA, 2008.

[33] K. R. Mopuri, R. V. Babu, Object level deep feature pooling for compact image representation, in: 2015 IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops, Boston, MA, USA, June 7-12, 2015, 2015, pp. 62–70.

[34] L. Zheng, S. Wang, Z. Liu, Q. Tian, Packing and padding: Coupled multi-index for accurate image retrieval, in: 2014 IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Columbus, OH, USA, June 23-28, 2014, 2014, pp. 1947–1954.

[35] R. Arandjelovic, A. Zisserman, Three things everyone should know to improve object retrieval, in: 2012 IEEE Conference on Computer Vision and Pattern Recognition, Providence, RI, USA, June 16-21, 2012, 2012, pp. 2911–2918.

[36] L. Zheng, S. Wang, Z. Liu, Q. Tian, Lp-norm IDF for large scale image search, in: 2013 IEEE Conference on Computer Vision and Pattern Recognition, Portland CVPR, OR, USA, June 23-28, 2013, 2013, pp. 1626–1633.

[37] H. Jegou, M. Douze, C. Schmid, On the burstiness of visual elements, in: 2009 IEEE Computer Society Conference on Computer Vision and Pattern Recognition CVPR, 20-25 June 2009, Miami, Florida, USA, 2009, pp. 1169–1176.

[38] R. Arandjelovic, P. Gronát, A. Torii, T. Pajdla, J. Sivic, Netvlad: CNN architecture for weakly supervised place recognition, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR, Las Vegas, NV, USA, June 27-30, 2016, 2016, pp. 5297–5307.

[39] S. Bai, Z. Zhou, J. Wang, X. Bai, L. J. Latecki, Q. Tian, Ensemble diffusion for retrieval, in: IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017, 2017, pp. 774–783.

[40] A. Gordo, J. Almazán, J. Revaud, D. Larlus, Deep image retrieval: Learning global representations for image search, in: Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VI, 2016, pp. 241–257.

[41] F. Radenovic, G. Tolias, O. Chum, CNN image retrieval learns from bow: Unsupervised fine-tuning with hard examples, in: Computer Vision - ECCV 2016 - 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part I, 2016, pp. 3–20.

[42] J. Xu, C. Shi, C. Qi, C. Wang, B. Xiao, Part-based weighting aggregation of deep convolutional features for image retrieval, CoRR abs/1705.01247.

[43] A. Babenko, A. Slesarev, A. Chigorin, V. S. Lempitsky, Neural codes for image retrieval, in: Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I, 2014, pp. 584–599.

[44] S. Rota Bulò, M. Pelillo, I. M. Bomze, Graph-based quadratic optimization: A fast evolutionary approach, Computer Vision and Image Understanding 115 (7) (2011) 984–995.