

UC San Diego

UC San Diego Previously Published Works

Title

Bridging the gap: Query by semantic example

Permalink

<https://escholarship.org/uc/item/4bj3k85j>

Journal

IEEE Transactions on Multimedia, 9(5)

ISSN

1520-9210

Authors

Rasiwasia, Nikhil
Moreno, Pedro J.
Vasconcelos, Nuno

Publication Date

2007-08-01

Peer reviewed

Bridging the Gap: Query by Semantic Example

Nikhil Rasiwasia, *Student Member, IEEE*, Pedro J. Moreno, *Member, IEEE*, and Nuno Vasconcelos, *Member, IEEE*

Abstract—A combination of *query-by-visual-example (QBVE)* and *semantic retrieval (SR)*, denoted as *query-by-semantic-example (QBSE)*, is proposed. Images are labeled with respect to a vocabulary of visual concepts, as is usual in SR. Each image is then represented by a vector, referred to as a *semantic multinomial*, of posterior concept probabilities. Retrieval is based on the *query-by-example* paradigm: the user provides a query image, for which 1) a semantic multinomial is computed and 2) matched to those in the database. QBSE is shown to have two main properties of interest, one mostly practical and the other philosophical. From a practical standpoint, because it inherits the generalization ability of SR inside the space of known visual concepts (referred to as the *semantic space*) but performs much better outside of it, QBSE produces retrieval systems that are more accurate than what was previously possible. Philosophically, because it allows a direct comparison of visual and semantic representations under a common query paradigm, QBSE enables the design of experiments that explicitly test the value of semantic representations for image retrieval. An implementation of QBSE under the minimum probability of error (MPE) retrieval framework, previously applied with success to both QBVE and SR, is proposed, and used to demonstrate the two properties. In particular, an extensive objective comparison of QBSE with QBVE is presented, showing that the former significantly outperforms the latter both inside and outside the semantic space. By carefully controlling the structure of the semantic space, it is also shown that this improvement can only be attributed to the semantic nature of the representation on which QBSE is based.

Index Terms—Content-based image retrieval, Gaussian mixtures, image similarity, multiple instance learning, query by example, semantic retrieval, semantic space.

I. INTRODUCTION

CONTENT-BASED image retrieval, the problem of searching for digital images in large image repositories according to their content, has been the subject of significant research in the recent past [1]–[4]. Two main retrieval paradigms have evolved over the years: one based on visual queries, here referred to as *query-by-visual-example (QBVE)*, and the other based on text, here denoted as *semantic retrieval (SR)*. Early retrieval architectures were almost exclusively based on QBVE [5]–[7], [2], [3]. Under this paradigm, each

image is decomposed into a number of low-level visual features (e.g., a color histogram) and image retrieval is formulated as the search for the best database match to the feature vector extracted from a query image. It was, however, quickly realized that strict visual similarity is, in most cases, weakly correlated with the measures of similarity adopted by humans for image comparison.

This motivated the more ambitious goal of designing retrieval systems with support for semantic queries [8]. The basic idea is to annotate images with semantic keywords, enabling users to specify their queries through a natural language description of the visual concepts of interest. Because manual image labeling is a labor intensive process, SR research turned to the problem of the automatic extraction of semantic descriptors from images, so as to build models of visual appearance of the semantic concepts of interest. This is usually done by the application of machine learning algorithms. Early efforts targeted the extraction of specific semantics [9]–[12] under the framework of binary classification. More recently there has been an effort to solve the problem in greater generality, through the design of techniques capable of learning relatively large semantic vocabularies from informally annotated training image collections. This can be done with resort to both unsupervised [13]–[17] and weakly supervised learning [18], [19]. Advances in the joint use of QBVE and SR have also been demonstrated within TRECVID [20], a benchmark to promote progress in content-based retrieval from large video repositories, where recent research efforts have concentrated on the fusion of the retrieval results obtained with the two paradigms [21], [22].

In spite of these advances, the fundamental question of whether there is an intrinsic value to building models at a semantic level, remains poorly understood. On one hand, SR has the advantage of evaluating image similarity at a higher level of abstraction, and therefore better generalization¹ than what is possible with QBVE. On the other hand, the performance of SR systems tends to degrade for semantic classes that they were not trained to recognize. Since it is still difficult to learn appearance models for massive concept vocabularies, this could compromise the generalization gains due to abstraction. This problem is seldom considered in the literature, where most evaluations are performed with query concepts that are known to the retrieval system [13]–[17], [19].

In fact, it is not even straightforward to compare the two retrieval paradigms, because they assume different levels of query specification. While a semantic query is usually precise (e.g., “the White House”) a visual example (a picture of the “White House”) will depict various concepts that are irrelevant to the

Manuscript received October 9, 2006; revised April 5, 2007. This work was supported by NSF CAREER Award IIS-0448609, NSF Grant IIS-0534985, and a gift from Google. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Marcel Worring.

N. Rasiwasia is with the Statistical Visual Computing Laboratory, Department of Electrical and Computer Engineering, University of California, San Diego, CA 92093 USA (e-mail: nikux@ucsd.edu)

P. J. Moreno is with Google, Inc., New York, NY 10011 USA (e-mail: pedro@google.com)

N. Vasconcelos is with the Statistical Visual Computing Laboratory, Department of Electrical and Computer Engineering, University of California, San Diego, CA 92093 USA (e-mail: nuno@ece.ucsd.edu).

Digital Object Identifier 10.1109/TMM.2007.900138

¹Here, and throughout this work, we refer to the definition of “generalization” common in machine learning and content-based retrieval: the ability of the retrieval system to achieve low error rates outside of the set of images on which it was trained.

query (e.g., the street that surrounds the building, cars, people, etc.). It is, therefore, possible that better SR results could be due to a better interface (natural language) rather than an intrinsic advantage of representing images semantically. This may be of little importance when the goal is to build the next generation of (more accurate) retrieval systems. However, given the complexity of the problem, it is unlikely that significant further advances can be achieved without some understanding of the intrinsic value of semantic representations. If, for example, abstraction is indeed valuable, further research on appearance models that account for image taxonomies could lead to exponential gains in retrieval accuracy. Else, if the advantages are simply a reflection of more precise queries, such research is likely to be ineffective.

In this work, we introduce a framework for the objective comparison of the two formulations, by extending the query-by-example paradigm to the semantic domain. This consists of defining a semantic feature space, where each image is represented by the vector of posterior concept probabilities assigned to it by a SR system, and performing query-by-example in this space. We refer to the combination of the two paradigms as query-by-semantic-example (QBSE), and present an extensive comparison of its performance with that of QBVE. It is shown that QBSE has significantly better performance for both concepts known and unknown to the retrieval system, i.e., it can generalize beyond the vocabulary used for training. It is also shown that the performance gain is intrinsic to the semantic nature of image representation.

The paper is organized as follows. Section II briefly reviews previous retrieval work related to QBSE. Section III then reviews, in greater details, the minimum probability of error (MPE) formulation of retrieval, which has been successfully applied to both QBVE [23] and SR [24], and is adopted in this work. Section IV discusses the limitations of the QBVE and SR paradigms, motivating the adoption of QBSE. Section V proposes an implementation of QBSE, compatible with the MPE formulation. It is then argued, in Section VI, that the generalization ability of QBSE can significantly benefit from the combination of multiple queries, and various strategies are proposed to accomplish this goal. A thorough experimental evaluation of the performance of QBSE is presented in Section VII, where the intrinsic gains of semantic image representations (over strict visual matching) are quantified. Finally, Section VIII presents a conclusion and suggests some possibilities for future research.

II. RELATED WORK

The idea of representing documents as weighted combinations of the words in a pre-defined vocabulary is commonly used in information retrieval. In fact, the classic model for information retrieval is the vector space model of Salton [25], [26]. Under this model, documents are represented as collections of keywords, weighted by importance, and can be interpreted as points in the semantic space spanned by the vocabulary entries. In image retrieval, the major topic of recent interest has been that of learning semantic image representations, but few proposals have so far been presented on how to best exploit the semantic space for the design of retrieval systems. Nevertheless,

there have been some proposals to represent images as points in a semantic vector space. The earliest among these efforts [27], [28] were based on semantic information extracted from metadata—viz. origin, filename, image url, keywords from surrounding webpage text, manual annotations, etc.

A somewhat popular technique to construct content-based semantic spaces, is to resort to active learning based on user's relevance feedback [29]–[31]. The idea is to pool the images relevant to a query, after several rounds of relevance feedback, to build a model for the semantic concept of interest. Assuming that 1) these images do belong to a common semantic class and 2) the results of various relevance feedback sessions can be aggregated, this is a feasible way to incrementally build a semantic space. An example is given in [32], where the authors propose a retrieval system based on image embeddings. Using relevance feedback, the system gradually clusters images and learns a non-linear embedding which maps these clusters into a hidden space of semantic attributes. Cox *et al.* [33] also focus on the task of learning a predictive model for user selections, by learning a mapping between 1) the image selection patterns made by users instructed to consider visual similarity and 2) those of users instructed to consider semantic similarity.

These works have focused more on the issue of learning the semantic space than that of its application to retrieval. In fact, it is not always clear how the learned semantic information could be combined with the visual search at the core of the retrieval operation. Furthermore, the use of relevance feedback to train a semantic retrieval system has various limitations. First, it can be quite time consuming, since a sizable number of examples is usually required to learn each semantic model. Second, the assumption that all queries performed in a relevance feedback session are relative to the same semantic concept is usually not realistic, even when users are instructed to do so. For example, a user searching for pictures of “cafes in Paris” is likely to oscillate between searching for pictures of “cafes” and pictures of “Paris.”

The closest works, in the literature, to the QBSE paradigm adopted here, are the systems proposed by Smith *et al.* in [34], [35] and Lu *et al.* in [36]. To the best of our knowledge, [34] pioneered the idea of 1) learning a semantic space by learning a separate statistical model for each concept and 2) performing query by example in the space the resulting semantic concepts. The vector of semantic weights, denoted as the “model vector,” is learned from the image content and not from metadata or relevance feedback information. Each image receives a confidence score per semantic concept, based on the proximity of the image to the decision boundary of a support vector machine (SVM) trained to recognize the concept. Retrieval is then based on the L^2 similarity of the concept-score vectors, which play a role similar to that of the posterior concept probability vectors used in this work.²

While laying the foundations for QBSE, [34], [35] did not investigate any of the fundamental questions that we now consider. First, because there was no attempt to perform retrieval on databases not used for training, it did not address the problem

²In fact, in the machine learning literature, SVM scores are frequently converted into posterior probabilities by application of a simple sigmoidal transformation [37].

of generalization to concepts unknown to the retrieval system. As we will see, this is one of the fundamental reasons to adopt QBSE instead of the standard SR query paradigm. Second, although showing that QBSE outperformed a QBVE system, this work did not rely on the same image representation for the two query paradigms. While QBVE was based on either color or edge histogram matching, QBSE relied on a feature space composed of a multitude of visual features, including color and edge histograms, wavelet-based texture features, color correlograms and measures of texture co-occurrence. Because the representations are different, it is impossible to conclude that the improved performance of the QBSE system derives from an *intrinsic* advantage of semantic-level representations. In what follows, we preempt this caveat by adopting the same image representation and retrieval framework for the design of all systems.

III. MINIMUM PROBABILITY OF ERROR RETRIEVAL

The retrieval framework underlying all query paradigms discussed in this work is that of MPE retrieval, as introduced in [23]. In this section, we briefly review this framework, and how it can be used to implement various types of retrieval systems.

A. Visual and Semantic-Level Retrieval Systems

The starting point for any retrieval system is an image database $\mathcal{D} = \{\mathcal{I}_1, \dots, \mathcal{I}_D\}$. Images are observations from a random variable \mathbf{X} , defined on some feature space \mathcal{X} . In the absence of labels, each image is considered an observation from a different class, determined by a random variable Y defined on $\{1, \dots, D\}$. In this case, the retrieval system is said to operate at the *visual-level*. Given a query image \mathcal{I}_q , the MPE retrieval decision is to assign it to the class of largest posterior probability, i.e.,

$$y^* = \arg \max_y P_{Y|\mathbf{X}}(y|\mathcal{I}_q). \quad (1)$$

A semantic-level retrieval system augments the database \mathcal{D} with a vocabulary $\mathcal{L} = \{w_1, \dots, w_L\}$ of semantic concepts or keywords w_i , and each image \mathcal{I}_i with a pre-specified caption \mathbf{c}_i , making $\mathcal{D} = \{(\mathcal{I}_1, \mathbf{c}_1), \dots, (\mathcal{I}_D, \mathbf{c}_D)\}$. Note that \mathbf{c}_i is a binary L -dimensional vector such that $\mathbf{c}_{i,j} = 1$ if the i th image was annotated with the j th keyword in \mathcal{L} . The database is said to be weakly labeled if the absence of a keyword from caption \mathbf{c}_i does not necessarily mean that the associated concept is not present in \mathcal{I}_i . For example, an image containing “sky” may not be explicitly labeled with that keyword. This is usually the case in practical scenarios, since each image is likely to be annotated with a small caption that only identifies the semantics deemed as most relevant to the labeler. We assume weak labeling throughout this work. Many possibilities exist for image representation. At the *visual level*, the image can be represented as a set of n feature vectors $\mathcal{I} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, where $\mathbf{x}_i \in \mathcal{X}$. Although any type of visual features are acceptable, we only consider *localized features*, i.e., features of limited spatial support. At the *semantic level*, concepts are drawn from a random variable W , which takes values in $\{1, \dots, L\}$, so that $W = i$ if

and only if \mathbf{x} is a sample from the concept w_i . Each concept induces a probability density $\{P_{\mathbf{X}|W}(\mathbf{x}|i)\}_{i=1}^L$ on \mathcal{X} , from which feature vectors are drawn. The generative model for a feature vector \mathbf{x} consists of two steps: a concept w_i is first selected, with probability $P_W(i) = \pi_i$, and the vector then drawn from $P_{\mathbf{X}|W}(\mathbf{x}|i)$. Both concepts and feature vectors are drawn independently, with replacement.

Given a new image \mathcal{I} , the MPE annotation rule is to label it with the concept of largest posterior probability

$$w^* = \arg \max_w P_{W|\mathbf{X}}(w|\mathcal{I}). \quad (2)$$

Once all images in an unlabeled database are annotated in this way, it becomes possible to support retrieval from that database with natural language queries. Given a query concept w_q , the optimal retrieval decision (in the MPE sense) is to select the image for which w_q has the largest posterior annotation probability, i.e., the image of index

$$i^* = \arg \max_i P_{\mathbf{X}|W}(\mathcal{I}_i|w_q). \quad (3)$$

B. Query by Visual Example

A QBVE system operates at the visual level and assumes that the feature vectors which compose any image \mathcal{I} are sampled independently

$$P_{\mathbf{X}|Y}(\mathcal{I}|y) = \prod_j P_{\mathbf{X}|Y}(\mathbf{x}_j|y). \quad (4)$$

Some density estimation [38] procedure is used to estimate the distributions $P_{\mathbf{X}|Y}(\mathbf{x}|y)$. This produces a vector of parameters Γ_y per image, e.g., $\Gamma_y = \{\mu_y^j, \Sigma_y^j, \alpha_y^j\}$, $y = 1, \dots, D$, when

$$P_{\mathbf{X}|Y}(\mathbf{x}|y; \Gamma_y) = \sum_j \alpha_y^j \mathcal{G}(\mathbf{x}, \mu_y^j, \Sigma_y^j) \quad (5)$$

is a mixture of Gaussians. Here, α_y is a probability mass function such that $\sum_j \alpha_y^j = 1$, $\mathcal{G}(\mathbf{x}, \mu, \Sigma)$ a Gaussian density of mean μ and covariance Σ , and j an index over the mixture components. In this work, we adopt the Gaussian mixture representation, and all parameter vectors Γ_y are learned by maximum likelihood, using the well known expectation-maximization (EM) algorithm [39]. Image retrieval is based on the mapping $g : \mathcal{X} \rightarrow \{1, \dots, D\}$ of (1), implemented by combining (4), (5) and Bayes rule. Although any prior class distribution $P_Y(i)$ can be supported, we assume a uniform distribution in what follows.

C. Semantic Retrieval

Under the MPE framework, SR is, in many ways, similar to QBVE. Images are assumed to be independently sampled from concept distributions

$$P_{\mathbf{X}|W}(\mathcal{I}|w) = \prod_j P_{\mathbf{X}|W}(\mathbf{x}_j|w) \quad (6)$$

and some density estimation procedure used to estimate the distributions $P_{\mathbf{X}|W}(\mathbf{x}|w)$. This produces a vector of parameters Ω_w per concept, e.g., $\Omega_w = \{\nu_w^j, \Phi_w^j, \beta_y^j\}$ when

$$P_{\mathbf{X}|W}(\mathbf{x}|w; \Omega_w) = \sum_j \beta_y^j \mathcal{G}(\mathbf{x}, \nu_w^j, \Phi_w^j) \quad (7)$$

is a mixture of Gaussians. As in the QBVE case, the parameter vectors Ω_w are learned by maximum likelihood. Image labeling is based on the mapping $g : \mathcal{X} \rightarrow \{1, \dots, L\}$ of (2), implemented by combining (6), (7), Bayes rule, and a uniform prior concept distribution $P_W(w)$. Image retrieval is based on the mapping $g : \{1, \dots, L\} \rightarrow \{1, \dots, D\}$ of (3), implemented by combining (6) with (7).

IV. QUERY BY SEMANTIC EXAMPLE

Both the QBVE and SR implementations of MPE retrieval have been extensively evaluated in [23] and [19], [24]. Although these evaluations have shown that the two implementations are among the best known techniques for visual and semantic-level retrieval, the comparison of the two retrieval paradigms is difficult. We next discuss this issue in greater detail, and motivate the adoption of an alternative retrieval paradigm, QBSE, that combines the best properties of the two approaches.

A. Query by Visual Example versus Semantic Retrieval

Both QBVE and SR have advantages and limitations. Because concepts are learned from collections of images, SR can *generalize* significantly better than QBVE. For example, by using a large training set of images labeled with the concept “sky,” containing both images of sky at daytime (when the sky is mostly blue) and sunsets (when the sky is mostly orange), a SR system can learn that “sky” is sometimes blue and others orange. This is a simple consequence of the fact that a large set of “sky” images populate, with high probability, the blue and orange regions of the feature space. It is, however, not easy to accomplish with QBVE, which only has access to two images (the query and that in the database) and can only perform direct matching of visual features. We refer to this type of abstraction, as *generalization inside the semantic space*, i.e., inside the space of concepts that the system has been trained to recognize.

While better generalization is a strong advantage for SR, there are some limitations associated with this paradigm. An obvious difficulty is that most images have multiple semantic interpretations. Fig. 1 presents an example, identifying various semantic concepts as sensible annotations for the image shown. Note that this list, of relatively salient concepts, is a small portion of the keywords that could be attached to the image. Other examples include colors (e.g., “yellow” train), or objects that are not salient in an abstract sense but could become very relevant in some contexts (e.g., the “paint” of the markings on the street, the “letters” in the sign, etc.). In general, it is impossible to predict all annotations that may be relevant for a given image. This is likely to compromise the performance of a SR system. Furthermore, because queries are specified as text, a SR system is



Fig. 1. Image containing various concepts: “train,” “smoke,” “road,” “sky,” “railroad,” “sign,” “trees,” “mountain,” and “shadows,” with variable degrees of presence.

usually limited by the size of its vocabulary³. In summary, SR can generalize poorly *outside the semantic space*.

Since visual retrieval has no notion of semantics, it is not constrained by either vocabulary or semantic interpretations. When compared to SR, QBVE systems can generalize better outside the semantic space. In the example of Fig. 1, a QBVE would likely return the image shown as a match to a query depicting an industrial chimney engulfed in dark smoke (a more or less obvious query prototype for images of “pollution”) despite the fact that the retrieval system knows nothing about “smoke,” “pollution,” or “chimneys.” Obviously, there are numerous examples where QBVE correlates much worse with perceptual similarity than SR. We have already seen that when the latter is feasible, i.e., inside the semantic space, it has better generalization. Overall, it is sensible to expect that SR will perform better inside the semantic space, while QBVE should fare better outside of it. In practice, however, it is not easy to compare the two retrieval paradigms. This is mostly due to the different forms of query specification. While a natural language query is usually precise (e.g., “train” and “smoke”), a query image like that of Fig. 1 always contains a number of concepts that are not necessarily relevant to the query (e.g., “mountain,” or even “yellow” for the train color). Hence, the better performance of SR (inside the semantic space) could be simply due to higher query precision. A fair comparison would, therefore, require the optimization of the precision of visual queries (e.g., by allowing the QBVE system to rely on image regions as queries) but this is difficult to formalize.

Overall, both the engineering question of how to design better retrieval systems (with good generalization inside and outside of the semantic space) and the scientific question of whether there is a real benefit to semantic representations, are difficult to answer under the existing query paradigms. To address this problem we propose an alternative paradigm, which is denoted as *query by semantic example* (QBSE).

³It is, of course, always possible to rely on text processing ideas based on thesauri and ontologies like WordNet [40] to mitigate this problem. For example, query expansion can be used to replace a query for “pollution” by a query for “smoke,” if the latter is in the vocabulary and the former is not. While such techniques are undeniably useful for practical implementation of retrieval systems, they do not reflect an improved ability, by the retrieval system, to model the relationships between visual features and words. They are simply an attempt to fix these limitations *a posteriori* (i.e., at the language level) and are, therefore, beyond the scope of this work. In practice, it is not always easy to perform text-based query expansion when the vocabulary is small, as is the case for most SR systems, or when the queries report to specific instances (e.g., a person’s name).

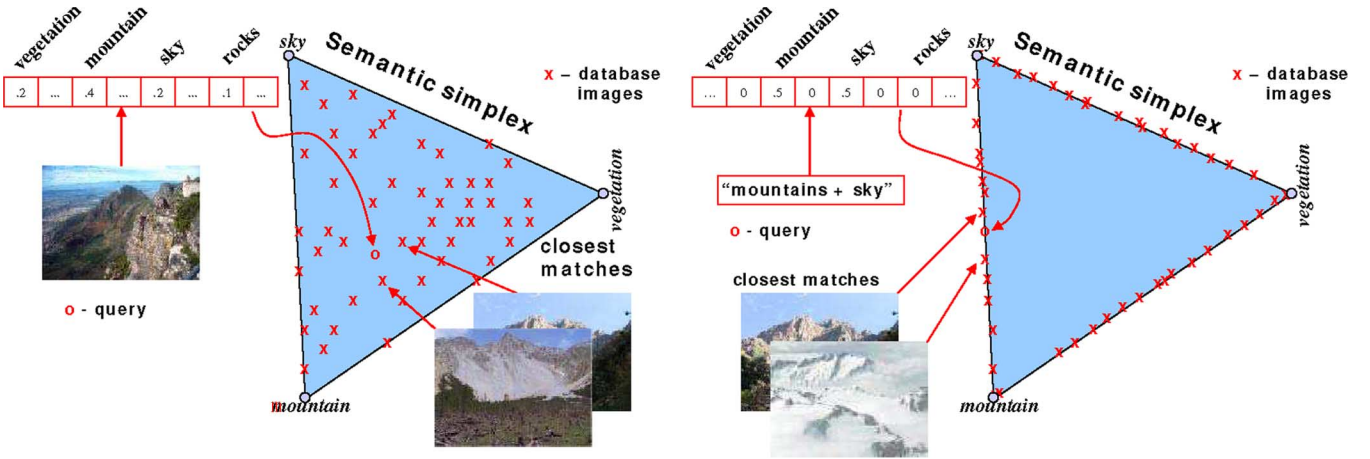


Fig. 2. Semantic image retrieval. Left: Under QBSE the user provides a query image, probabilities are computed for all concepts, and the image represented by the concept probability distribution. Right: Under the traditional SR paradigm, the user specifies a short natural language description, and only a small number of concepts are assigned a nonzero posterior probability.

B. Query by Semantic Example

A QBSE system operates at the semantic level, representing images by vectors of concept counts $\mathcal{I} = (c_1, \dots, c_L)^T$. Each feature vector of the image is assumed to be sampled from the probability distribution of a semantic class (concept) and c_i is the number of feature vectors drawn from the i th concept. The count vector for the y th image is drawn from a multinomial variable \mathbf{T} of parameters $\boldsymbol{\pi}_y = (\pi_y^1, \dots, \pi_y^L)^T$

$$P_{\mathbf{T}|Y}(\mathcal{I}|y; \boldsymbol{\pi}_y) = \frac{n!}{\prod_{k=1}^L c_k!} \prod_{j=1}^L (\pi_y^j)^{c_j} \quad (8)$$

where π_y^i is the probability that a feature vector is drawn from the i th concept. The random variable \mathbf{T} can be seen as the result of a feature transformation from the space of visual features \mathcal{X} to the L -dimensional probability simplex \mathcal{S}_L . This mapping, $\boldsymbol{\Pi} : \mathcal{X} \rightarrow \mathcal{S}_L$ such that $\boldsymbol{\Pi}(\mathbf{X}) = \mathbf{T}$, maps the Gaussian mixtures $P_{\mathbf{X}|Y}(\mathcal{I}|y)$ into the multinomials $P_{\mathbf{T}|Y}(\mathcal{I}|y)$, and establishes a correspondence between images and points $\boldsymbol{\pi}_y \in \mathcal{S}_L$, as illustrated by Fig. 2 (left). Since the entries of $\boldsymbol{\pi}_y$ are the posterior probabilities of the semantic concepts ω_i , $i = 1, \dots, L$ given the y th image, we refer to the probability simplex \mathcal{S}_L as the *semantic simplex*, and to the probability vector $\boldsymbol{\pi}_y$ itself as the *semantic multinomial* (SMN) that characterizes the image. A QBSE system operates on the simplex \mathcal{S}_L , according to a similarity mapping $f : \mathcal{S}_L \rightarrow \{1, \dots, D\}$ such that

$$f(\boldsymbol{\pi}) = \arg \max_y s(\boldsymbol{\pi}, \boldsymbol{\pi}_y) \quad (9)$$

where $\boldsymbol{\pi}$ is the query SMN, $\boldsymbol{\pi}_y$ the SMN that characterizes the y th database image, and $s(\cdot, \cdot)$ an appropriate similarity function. The user provides a query image, for which a SMN $\boldsymbol{\pi}$ is computed, and compared to all the SMNs $\boldsymbol{\pi}_y$ previously stored for the images in the database.

This query paradigm has a number of interesting properties. First, the mapping of the visual features into the probability simplex \mathcal{S}_L can be seen as an abstract mapping of the image to a *semantic space* where each concept probability π_y^i , $i = 1, \dots, L$ is

a *semantic feature*. Semantic features, or concepts, outside the vocabulary simply define directions orthogonal to the learned semantic space. This implies that, by projecting these dimensions onto the simplex, the QBSE system can generalize beyond the known semantic concepts. In the example of Fig. 1, the mapping of the image onto the semantic simplex assigns high probability to (known) concepts such as “train,” “smoke,” “railroad,” etc. This makes the image a good match for other images containing large amounts of “smoke,” such as those depicting industrial chimneys or “pollution” in general. The system can therefore establish a link between the image of Fig. 1 and “pollution,” despite the fact that it has no *explicit* knowledge of the “pollution” concept.⁴ Second, when compared to QBVE, QBSE complements all the advantages of query by example with the advantages of a semantic representation. Moreover, since in both cases queries are specified by the same examples, any differences in their performance can be directly attributed to the semantic versus visual nature of the associated image representations.⁵ This enables the objective comparison of QBVE and QBSE.

V. PROPOSED QUERY BY SEMANTIC EXAMPLE SYSTEM

QBSE is a generic retrieval paradigm and, as such, can be implemented in many different ways. Any implementation must specify a method to estimate the SMN that describes each image, and a similarity function between SMNs. We next describe an implementation which is compatible with MPE retrieval.

A. Semantic Labeling System

SMN parameter vectors $\boldsymbol{\pi}_i$ are learned with a semantic labeling system, which implements the mapping $\boldsymbol{\Pi}$, by computing

⁴Note that this is different from text-based query expansion, where the link between “smoke” and “pollution” must be *explicitly* defined. In QBSE, the relationship is instead inferred automatically, from the fact that both concepts have commonalities of visual appearance.

⁵This assumes, of course, that a common framework, such as MPE, is used to implement both the QBSE and QBVE systems.

an estimate of posterior concept probabilities given the observed feature vectors

$$\pi_w = P_{W|X}(w|\mathcal{I}). \quad (10)$$

While this can be done with any system that produces posterior concept probabilities, we adopt the weakly supervised method of [19]. This method formulates semantic image labeling as an L -ary classification problem. A semantic class density $P_{X|W}(x|w)$ is learned for each concept w from the set \mathcal{D}_w of all training images labeled with the w^{th} in \mathcal{L} , using a *hierarchical estimation* procedure first proposed, in [41], for image indexing. This procedure is itself composed of two steps.

First, a Gaussian mixture is learned for each image in \mathcal{D}_w , producing a sequence of mixture densities

$$P_{X|S,W}(x|s,w) = \sum_k \alpha_{w,s}^k \mathcal{G}(x, \mu_{w,s}^k, \Sigma_{w,s}^k) \quad (11)$$

where S is a hidden variable that indicates the index of the image in \mathcal{D}_w . Note that, if a QBVE system has already been implemented, these densities are just replicas of the ones of (5). In particular, if the mapping $\mathcal{T} : \{1, \dots, L\} \times \{1, \dots, D\} \rightarrow \{1, \dots, D\}$ translates the index (w, s) of the s th image in \mathcal{D}_w into the image's index y on \mathcal{D} , i.e., $y = \mathcal{T}(w, s)$, then

$$P_{X|S,W}(x|s,w) = P_{X|Y}(x|\mathcal{T}(w, s)).$$

Omitting, for brevity, the dependence of the mixture parameters on the semantic class w , assuming that each mixture has K components, and that the cardinality of \mathcal{D}_w is D_w , this produces $D_w K$ mixture components of parameters $\{\alpha_s^k, \mu_s^k, \Sigma_s^k\}$, $s = 1, \dots, D_w, k = 1, \dots, K$. The second step is an extension of the EM algorithm, which clusters the Gaussian components into the mixture distribution of (7), using a hierarchical estimation technique (see [24], [41] for details). Because the number of parameters in each image mixture is orders of magnitude smaller than the number of feature vectors extracted from the image, the complexity of estimating concept mixtures is negligible when compared to that of estimating the individual image mixtures. It follows that the overall training complexity is equivalent to that required to train a QBVE system based on (1). In [24] it is shown that this labeling method achieves better performance than a number of other state-of-the-art methods available in the literature [16], [17].

B. Semantic Multinomial

Given an image $\mathcal{I} = \{x_1, \dots, x_n\}$ the posterior concept probabilities of (10) are computed by combining (6), (7), and Bayes rule, assuming a uniform prior concept distribution $P_W(w)$. As is usual in probability estimation, these posterior probabilities can be inaccurate for concepts with a small number of training images. Of particular concern are cases where some of the π_i are very close to zero, and can become ill-conditioned during retrieval, where noisy estimates are amplified by ratios or logs of probabilities. A common solution is to introduce a prior distribution to regularize these parameters.

For this, it is worth considering an alternative procedure for the estimation of the π_i . Instead of (10), this consists of com-

puting the posterior concept probabilities $P_{W|X}(w|x_k)$, $w \in \{1, \dots, L\}$ of *each* feature vector x_k , assign x_k to the concept of largest probability, and count the number c_w of vectors assigned to each concept. The maximum likelihood estimate of the probabilities is then given by [38]

$$\pi_w^{\text{ML}} = \arg \max_{\pi_w} \prod_{j=1}^L \pi_j^{c_j} = \frac{c_w}{\sum_j c_j} = \frac{c_w}{n}. \quad (12)$$

Regularization can then be enforced by adopting a Bayesian parameter estimation viewpoint, where the parameter π is considered a random variable, and a prior distribution $P_{\Pi}(\pi)$ introduced to favor parameter configurations that are, *a priori*, more likely. Conjugate priors are frequently used, in Bayesian statistics [42], to estimate parameters of distributions in the exponential family, as is the case of the multinomial. They lead to a closed-form posterior (which is in the family of the prior), and *maximum a posteriori probability* parameter estimates which are intuitive. The conjugate prior of the multinomial is the Dirichlet distribution

$$\pi \sim \text{Dir}(\alpha) = \frac{\Gamma(\sum_{j=1}^L \alpha_j)}{\prod_{j=1}^L \Gamma(\alpha_j)} \prod_{j=1}^L \pi_j^{\alpha_j - 1} \quad (13)$$

of *hyper-parameters* α_i , and where $\Gamma(\cdot)$ is the Gamma function. Setting⁶ $\alpha_i = \alpha$, the maximum a posteriori probability estimates are

$$\begin{aligned} \pi_w^{\text{posterior}} &= \arg \max_{\pi_w} P_{\mathcal{T}|\Pi}(c_1, \dots, c_L | \pi) P_{\Pi}(\pi) \\ &= \arg \max_{\pi_w} \prod_{j=1}^L \pi_j^{c_j} \prod_{j=1}^L \pi_j^{\alpha - 1} \\ &= \frac{c_w + \alpha - 1}{\sum_{j=1}^L (c_j + \alpha - 1)}. \end{aligned} \quad (14)$$

This is identical to the maximum likelihood estimates obtained from a sample where each count is augmented by $\alpha - 1$, i.e., where each image contains $\alpha - 1$ more feature vectors from each concept. The addition of these vectors prevents zero counts, regularizing π . As α increases, the multinomial distribution tends to uniform.

Thresholding the individual feature vector posteriors and counting is likely to produce worse probability estimates than those obtained, with (10), directly from the entire collection of feature vectors. Nevertheless, the discussion above suggests a strategy to regularize the probabilities of (10). Noting, from (12), that $c_w = n\pi_w^{\text{ML}}$, the regularized estimates of (14) can be written as

$$\pi_w^{\text{posterior}} = \frac{\pi_w^{\text{ML}} + \pi_0}{\sum_j (\pi_j^{\text{ML}} + \pi_0)}.$$

with $\pi_0 = (\alpha - 1)/(n)$. Hence, regularizing the estimates of (10) with

$$\pi_w^{\text{reg}} = \frac{\pi_w + \pi_0}{1 + L\pi_0} \quad (15)$$

⁶Different hyper-parameters could also be used for the different concepts.

is equivalent to using maximum a posteriori probability estimates, in the thresholding plus counting paradigm, with the Dirichlet prior of (13).

C. Similarity Function

There are many known methods to measure the distance between two probability distributions, all of which can be used to measure the similarity of two SMNs. Furthermore, because the latter can also be interpreted as normalized vectors of counts, this set can be augmented with all measures of similarity between histograms. We have compared various similarity functions for the purpose of QBSE.

1) *Kullback-Leibler (KL) Divergence*: The KL divergence between two distributions π and π' is

$$s_{\text{KL}}(\pi, \pi') = \text{KL}(\pi \| \pi') = \sum_{i=1}^L \pi_i \log \frac{\pi_i}{\pi'_i}. \quad (16)$$

It is nonnegative, and equal to zero when $\pi = \pi'$. For retrieval, it also has an intuitive interpretation as the asymptotic limit of (1) when Y is uniformly distributed [43]. However, it is not symmetric, i.e., $\text{KL}(\pi \| \pi') \neq \text{KL}(\pi' \| \pi)$. A symmetric version can be defined as

$$s_{\text{symmKL}}(\pi, \pi') = \text{KL}(\pi \| \pi') + \text{KL}(\pi' \| \pi) \quad (17)$$

$$= \sum_{i=1}^L \pi_i \log \frac{\pi_i}{\pi'_i} + \sum_{i=1}^L \pi'_i \log \frac{\pi'_i}{\pi_i}. \quad (18)$$

2) *Jensen-Shannon Divergence*: The Jensen-Shannon divergence (JS) is a measure of whether two samples, as defined by their empirical distributions, are drawn from the same source distribution [44]. It is defined as

$$s_{\text{JS}}(\pi, \pi') = \text{KL}(\pi \| \hat{\pi}) + \text{KL}(\pi' \| \hat{\pi}) \quad (19)$$

where $\hat{\pi} = (1/2)\pi + (1/2)\pi'$. This divergence can be interpreted as the average distance (in the KL sense) between each distribution and the average of all distributions.

3) *Correlation*: The correlation between two SMNs is defined as

$$s_{\text{CO}}(\pi, \pi') = \pi^T \pi' = \sum_i \pi_i \times \pi'_i. \quad (20)$$

Unlike the KL or JS divergence, which attain their minimum value (zero) for equal distributions, correlation is maximum in this case. The maximum value is, however, a function of the distributions under consideration. This limitation can be avoided by the adoption of the *normalized correlation*

$$s_{\text{NC}}(\pi, \pi') = \frac{\pi^T \pi'}{\|\pi\| \|\pi'\|} = \frac{\sum_i \pi_i \times \pi'_i}{\sqrt{\sum_j \pi_j^2} \sqrt{\sum_j \pi_j'^2}}. \quad (21)$$

4) *Other Similarity Measures*: A popular set of image similarity metrics is that of L^p distances

$$s_{L^p}(\pi, \pi') = \left(\sum_{i=1}^L |\pi_i - \pi'_i|^p \right)^{\frac{1}{p}}. \quad (22)$$

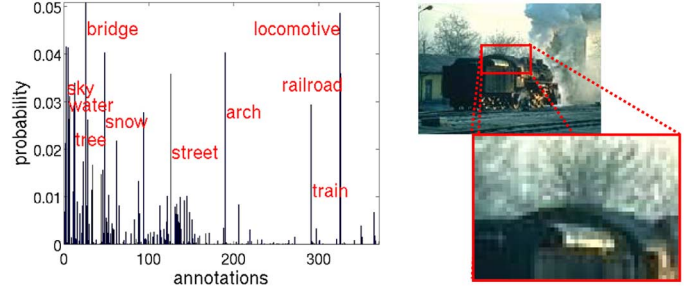


Fig. 3. Image and its associated SMN. Note that, while most of the concepts of largest probability are present in the image, the SMN assigns significant probability to “bridge” and “arch.” These are due to the geometric structure shown on the image close-up.

These distances are particularly common in color-based retrieval, where they are used as metrics of similarity between color histograms. Another popular metric is the histogram intersection (HI) [45],

$$s_{\text{HI}}(\pi, \pi') = \sum_{i=1}^L \min(\pi_i, \pi'_i) \quad (23)$$

the maximization of which is equivalent to minimizing the L^1 norm.

VI. MULTIPLE IMAGE QUERIES

A QBSE system can theoretically benefit from the specification of queries through multiple examples. We next give some reasons for this and discuss various alternatives for query combination.

A. The Benefits of Query Fusion

Semantic image labeling is, almost by definition, a noisy endeavor. This is a consequence of the fact that various interpretations are usually possible for a given arrangement of image intensities. An example is given in Fig. 3 where we show an image and the associated SMN. While most of the probability mass is assigned to concepts that are present in the image (“railroad,” “locomotive,” “train,” “street,” or “sky”), two of the concepts of largest probability do not seem related to it: “bridge” and “arch.” Close inspection of the image (see close-up presented in the figure), provides an explanation for these labels: when analyzed locally, the locomotive’s roof actually resembles the arch of a bridge. This visual feature seems to be highly discriminant, since when used as a query in a QBVE system, most of the top matches are images with arch-like structures, not trains (see Fig. 7). While these types of errors are difficult to avoid, they are *accidental*. In particular, the arch-like structure of Fig. 3 is the result of viewing a particular type of train, at a particular viewing angle, and a particular distance. It is unlikely that similar structures will emerge consistently over a set of train images. There are obviously other sources of error, such as classification mistakes for which it is not possible to encounter a plausible explanation. But these are usually even less consistent, across a set of images, than those due to accidental visual resemblances. A pressing question is then whether it is possible to exploit the lack

of consistency of these errors to obtain a better characterization of the query image set?

We approach this question from a *multiple instance* learning perspective [46], [47], formulating the problem as one of learning from *bags of examples*. In QBSE, each image is modeled as a bag of feature vectors, which are drawn from the different concepts according to the probabilities π_i . When the query consists of multiple images, or bags, the negative examples that appear across those bags are inconsistent (e.g., the feature vectors associated with the arch-like structure which is prominent in Fig. 3 but does not appear consistently in all train images), and tend to be spread over the feature space (because they also depict background concepts, such as roads, trees, mountains, etc., which vary from image to image). On the other hand, feature vectors corresponding to positive examples are likely to be concentrated within a small region of the space. It follows that, although the distribution of positive examples may not be dominant in any individual bag, the consistent appearance in all bags makes it dominant over the entire query ensemble. This suggests that a better estimate of the query SMN should be possible by considering a set of multiple query images.

In addition to higher accuracy, a set of multiple queries is also likely to have better generalization, since a single image does not usually exhibit all possible visual manifestations of a given semantic class. For example, images depicting “bikes on roads” and “cars in garage” can be combined to retrieve images from the more general class of “vehicles.” A combination of the two query image sets enables the retrieval system to have a more complete representation of the vehicle class, by simultaneously assigning higher weights to the concepts “bike,” “cars,” “road,” and “garage.” This enables the retrieval of images of “bikes in garage” and “cars on roads,” matches that would not be possible if the queries were used individually.

B. Query Combination

Under MPE retrieval, query combination is relatively straightforward to implement by QBVE systems. Given two query images $\mathcal{I}_q^1 = \{\mathbf{x}_1^1, \mathbf{x}_2^1, \dots, \mathbf{x}_n^1\}$ and $\mathcal{I}_q^2 = \{\mathbf{x}_1^2, \mathbf{x}_2^2, \dots, \mathbf{x}_n^2\}$, the probability of the composite query $\mathcal{I}_q^C = \{\mathbf{x}_1^1, \mathbf{x}_2^1, \dots, \mathbf{x}_n^1, \mathbf{x}_1^2, \mathbf{x}_2^2, \dots, \mathbf{x}_n^2\}$ given class $Y = y$ is

$$\begin{aligned} P_{\mathbf{X}|Y}(\mathcal{I}_q^C|y) &= \prod_{k=1}^n P_{\mathbf{X}|Y}(\mathbf{x}_k^1|y) \prod_{l=1}^n P_{\mathbf{X}|Y}(\mathbf{x}_l^2|y) \\ &= P_{\mathbf{X}|Y}(\mathcal{I}_q^1|y) P_{\mathbf{X}|Y}(\mathcal{I}_q^2|y). \end{aligned} \quad (24)$$

The MPE decision of (1) for the composite query is obtained by combining (24) with (5) and Bayes rule. In the context of QBSE, there are at least three possibilities for query combination. The first is equivalent to (24), but based on the probability of the composite query \mathcal{I}_q^C given semantic class $W = w$

$$\begin{aligned} P_{\mathbf{X}|W}(\mathcal{I}_q^C|w) &= \prod_{k=1}^n P_{\mathbf{X}|W}(\mathbf{x}_k^1|w) \prod_{l=1}^n P_{\mathbf{X}|W}(\mathbf{x}_l^2|w) \\ &= P_{\mathbf{X}|W}(\mathcal{I}_q^1|w) P_{\mathbf{X}|W}(\mathcal{I}_q^2|w) \end{aligned} \quad (25)$$

TABLE I
RETRIEVAL AND QUERY DATABASE

Database	Semantic Space	Source	# Retrieval Images	# Query Images	# Classes
Corel50	Inside	Corel Stock Photo CDs	4500	500	50
Corel15	Outside	Corel Stock Photo CDs	1200	300	15
Flickr18	Outside	www.flickr.com	1440	360	18

which is combined with (7) and Bayes rule to compute the posterior concept probabilities of (10). We refer to (25) as the “LKLD combination” strategy for query combination. It is equivalent to taking a geometric mean of the probabilities of the individual images given the class.

A second possibility is to represent the query as a mixture of SMNs. This relies on a different generative model than that of (25): the i th query is first selected with probability λ_i and a count vector is then sampled from the associated multinomial distribution. It can be formalized as

$$P_{\mathbf{T}}(\mathcal{I}_q^C; \boldsymbol{\pi}_q) = \frac{n!}{\prod_{k=1}^L c_k!} \prod_{j=1}^L (\lambda_1 \pi_1^j + \lambda_2 \pi_2^j)^{c_j} \quad (26)$$

where $P_{\mathbf{T}}(\mathcal{I}_q^C; \boldsymbol{\pi}_q)$ is the multinomial distribution for the query combination, of parameter $\boldsymbol{\pi}_q = \lambda_1 \boldsymbol{\pi}_1 + \lambda_2 \boldsymbol{\pi}_2$. $\boldsymbol{\pi}_1$ and $\boldsymbol{\pi}_2$ are the parameters of the individual multinomial distribution, and $\boldsymbol{\lambda} = (\lambda_1, \lambda_2)^T$ the vector of query selection probabilities. If $\lambda_1 = \lambda_2$, the two SMNs are simply averaged. We adopt the uniform query selection prior, and refer to this strategy as “SMN combination.” Geometrically, it sets the combined SMN to the centroid of the simplex that has the SMNs of the query images as vertices. This ranks highest the database SMN which is closest to this centroid.

The third possibility, henceforth referred to as “KL combination,” is to execute the multiple queries separately, and combine the resulting image rankings. For example, when similarity is measured with the KL divergence, the divergence between the combined image SMN, $\boldsymbol{\pi}_q$, and database SMNs $\boldsymbol{\pi}_y$ is

$$s_{\text{KL}}(\boldsymbol{\pi}_q, \boldsymbol{\pi}_y) = \frac{1}{2} \text{KL}(\boldsymbol{\pi}_1 \| \boldsymbol{\pi}_y) + \frac{1}{2} \text{KL}(\boldsymbol{\pi}_2 \| \boldsymbol{\pi}_y). \quad (27)$$

It is worth noting that this combination strategy is closely related to that used in QBVE. Note that the use of (24) is equivalent to using the arithmetic average (mean) of log-probabilities which, in turn, is identical to combining image rankings, as in (27). For QBVE the two combination approaches are identical.

VII. EXPERIMENTAL EVALUATION

In this section we report on an extensive evaluation of QBSE. We start by describing the evaluation procedure and the various databases used. This is followed by some preliminary tuning of the parameters of the QBSE system and the analysis of a number of retrieval experiments, that can be broadly divided into two classes. Both compare the performance of QBSE and QBVE, but while the first is performed inside the semantic space, the second studies retrieval performance outside of the latter.

A. Evaluation Procedure

In all cases, performance is measured with *precision* and *recall*, a classical measure of information retrieval performance

[25], which is also widely used by the image retrieval community [48], and one of the metrics adopted by the TRECVID evaluation benchmark. Given a query and the top “ N ” database matches, also called as *scope*, if “ r ” of the retrieved objects are relevant (where relevant means belonging to the class of the query), and the total number of relevant objects in the database is “ R ,” then precision is defined as “ r/N ,” i.e., the percentage of N which are relevant and recall as “ r/R ,” which is the percentage of all relevant images contained in the retrieved set. Precision-recall is commonly summarized by the mean average precision (MAP)[16]. This consists of averaging the precision at the ranks where recall changes, and taking the mean over a set of queries. Because some authors [49] consider the characterization of retrieval performance by curves of *precision-scope* more expressive for image retrieval, we also present results with this measure.

B. Databases

The complete evaluation of a QBSE system requires three different databases. The first is a *training database*, used by the semantic labeling system to learn concept probabilities. The second is a *retrieval database* from which images are to be retrieved. The third is a database of *query images*, which do not belong to either the training or retrieval databases. In the first set of experiments, the training and retrieval databases are identical, and the query images are inside the semantic space. This is the usual evaluation scenario for semantic image retrieval [14], [17], [16]. In the second, designed to evaluate generalization, both query and retrieval databases are outside the semantic space.

1) *Training Database:* We relied on the Corel dataset, used in [14], [16], [17] as the training database for all experiments. This dataset, henceforth referred to as *Corel50*, consists of 5000 images from 50 Corel Stock Photo CDs, divided into a training set of 4500 images (which was used to learn the semantic space), and a test set of 500 images (which did not play a role in the learning stage). Each CD includes 100 images of a common topic, and each image is labeled with 1-5 semantic concepts. Overall there are 371 keywords in the data set, leading to a 371-dimensional semantic simplex. With respect to image representation, all images were normalized to size 181×117 or 117×181 and converted from RGB to the YBR color space. Image observations were derived from 8×8 patches obtained with a sliding window, moved in a raster-scan fashion. A feature transformation was applied to this space by computing the 8×8 discrete cosine transform (DCT) of the three color components of each patch. The parameters of the semantic class mixture hierarchies were learned in the subspace of the resulting 192-dimension feature space composed of the first 21 DCT coefficients from each channel. In all experiments, the SMN associated with each image was computed with these semantic class-conditional distributions.

2) *Retrieval and Query Database:* To evaluate retrieval performance we carried out tests on three databases *Corel50*, *Flickr18*, and *Corel15*.⁷

a) *Inside the Semantic Space:* Retrieval performance inside the semantic space was evaluated by using *Corel50* as both

⁷The dataset used for experiments is available from <http://www.svcl.ucsd.edu/projects/qbse/>

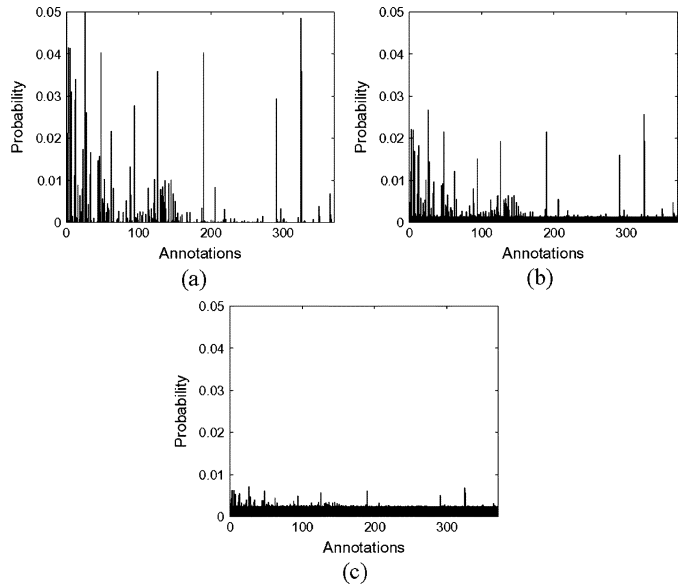


Fig. 4. SMN of the *train* query of Fig. 7 as a function of the ratio $L(\alpha - 1)/(n)$ adopted for its regularization.

TABLE II
EFFECT OF SMN REGULARIZATION ON THE MAP SCORE OF QBSE

Ratio	100	10	1	0.1	0.01	0.001	0.0001	0.00001
<i>Corel50</i>	0.1544	0.1744	0.1833	0.1768	0.1709	0.1683	0.1672	0.1667
<i>Corel15</i>	0.1878	0.2030	0.2156	0.2175	0.2160	0.2150	0.2144	0.2141
<i>Flickr18</i>	0.1447	0.1557	0.1625	0.1615	0.1594	0.1578	0.1569	0.1564

retrieval and query database. More precisely, the 4500 training images served as the *retrieval database* and the remaining 500 as the *query database*. This experiment relied on clear ground truth regarding the relevance of the retrieved images, based on the theme of the CD to which the query belonged.

b) *Outside the Semantic Space:* To test performance outside the semantic space, we relied on two additional databases. The first, *Corel15*, consisted of 1500 images from 15⁸ previously unused Corel CDs. Once again, the CD themes (nonoverlapping with those of *Corel50*) served as the ground truth. To address some criticism that “Corel is easy” [50], [51], we collected a second database from the online photo sharing website <http://www.flickr.com>. The images in this database were extracted by placing queries on the flickr search engine, and manually pruning images that appeared irrelevant to the specified queries. Note that the judgments of relevance did not take into account how well a content-based retrieval system would perform on the images, simply whether they appeared to be search errors (by flickr) or not. The images are shot by flickr users, and hence differ from the Corel Stock photos, which have been shot professionally. This database, *Flickr18*, contains 1800 images divided into 18 classes according to the manual annotations provided by the online users. For both databases, 20% of randomly selected images served as *query images* and the remaining 80% as the *retrieval database*. Table I summarizes the composition of the databases used.

⁸“Adventure Sailing,” “Autumn,” “Barnyard Animals,” “Caves,” “Cities of Italy,” “Commercial Construction,” “Food,” “Greece,” “Helicopters,” “Military Vehicles,” “New Zealand,” “People of World,” “Residential Interiors,” “Sacred Places,” and “Soldier.”

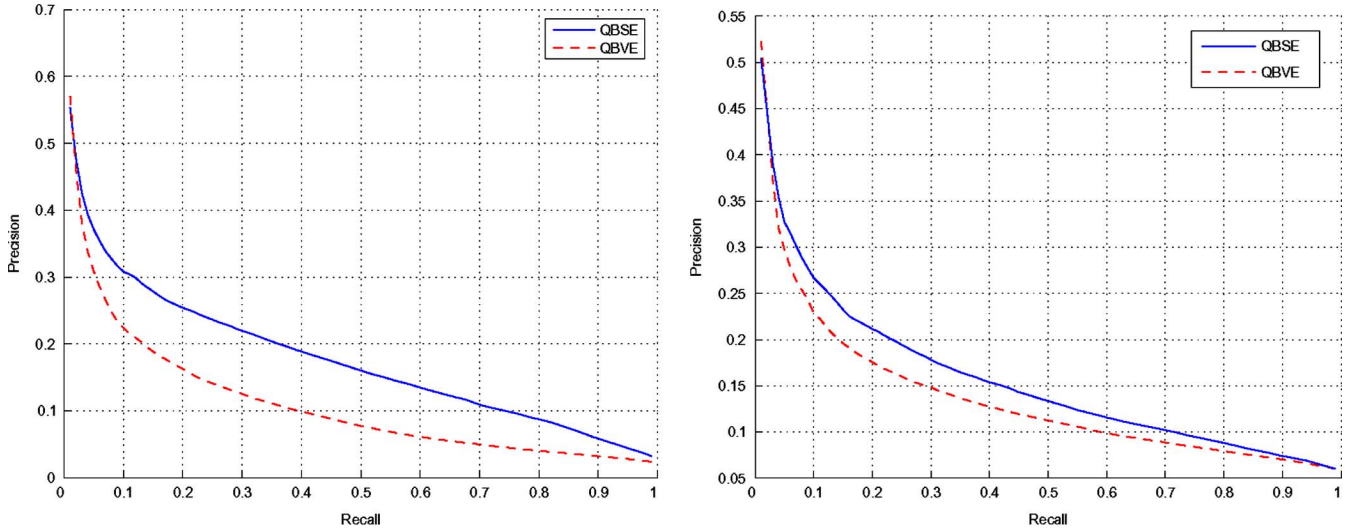


Fig. 5. Average precision-recall of single-query QBSE and QBVE. Left: Inside the semantic space (*Core150*). Right: Outside the semantic space (*Flickr18*).

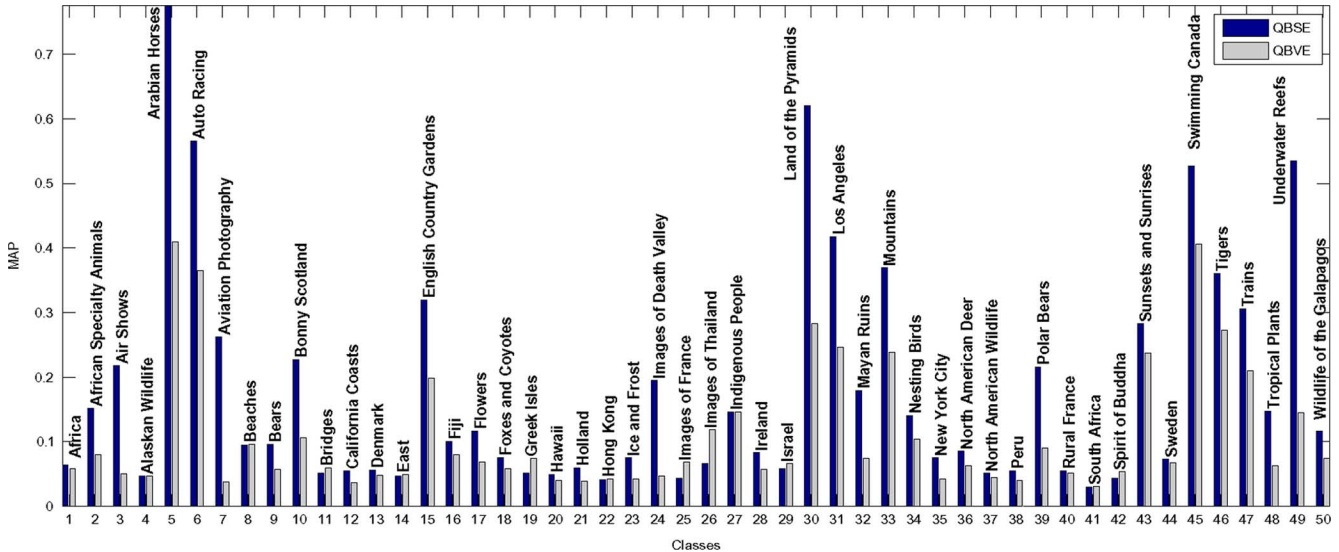


Fig. 6. MAP scores of QBSE and QBVE across the 50 classes of *Core150*.

A QBVE system only requires a query and a retrieval database. In all experiments, these were made identical to the query and retrieval databases used by the QBSE system. Since the performance of QBVE does not depend on whether queries are inside or outside the semantic space, this establishes a benchmark for evaluating the generalization of QBSE.

C. Model Tuning

All parameters of our QBVE system have been previously optimized, as reported in [23]. Here, we concentrate on the QBSE system, reporting on the impact of 1) SMN regularization and 2) choice of similarity function on the retrieval performance. The parameters resulting from this optimization were used in all subsequent experiments.

1) *Effect of Regularization on QBSE*: Table II presents the MAP obtained with values of $L\pi_0$ (15), ranging from 10^{-5} to 100. Fig. 4 presents the SMN of the *train* query of Fig. 7, for some of the values of $L\pi_0$. It can be seen that very large values of

α force the SMN towards a uniform distribution, e.g., Fig. 4(c), and almost all semantic information is lost. Fig. 4(b) shows the SMN regularized with the optimal value of $\pi_0 = 1/L$, where exceedingly low concept probabilities are lower-bounded by the value of 0.001. This regularization is instrumental in avoiding very noisy distance estimates during retrieval.

2) *Effect of the Similarity Function on QBSE*: Table III presents a comparison of the seven similarity functions discussed in the text. It is clear that L^2 distance and histogram intersection do not perform well. All information theoretic measures, KL divergence, symmetric KL divergence and Jensen-Shanon divergence, have superior performance, with an average improvement of 15%. Among these the KL divergence performs the best. The closest competitors to KL divergence are the correlation and normalized correlation metrics. Although, they outperform KL divergence inside the semantic space (*Core150*), their performance is inferior for databases outside the semantic space (*Flickr18*, *Core115*). This indicates

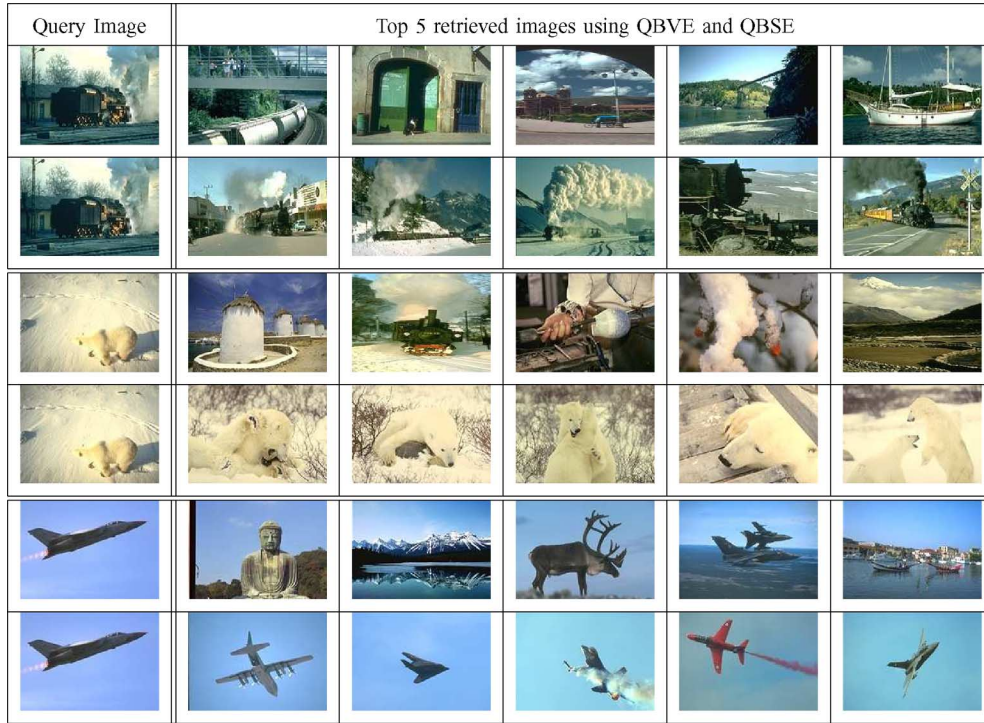


Fig. 7. Some examples where QBSE performs better than QBVE. The second row of every query shows the images retrieved by QBSE.

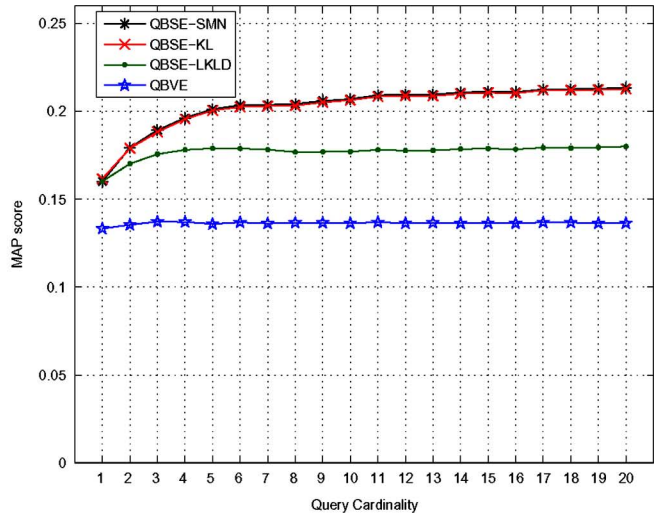
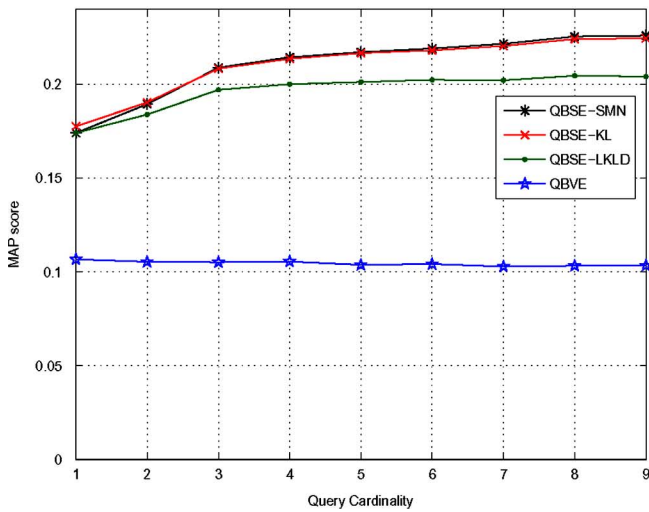


Fig. 8. MAP as a function of query cardinality for multiple image queries. Comparison of QBSE, with various combination strategies, and QBVE. Left: Inside the semantic space (*Corel50*). Right: Outside the semantic space (*Flickr18*).

TABLE III
EFFECT OF THE SIMILARITY FUNCTION ON THE MAP SCORE OF QBSE

Similarity Function	KL	symmKL	JS	CO	NC	L2	HI
<i>Corel50</i>	0.1768	0.1733	0.1740	0.2108	0.1938	0.1461	0.1692
<i>Corel15</i>	0.2175	0.2164	0.2158	0.1727	0.2041	0.1830	0.2119
<i>Flickr18</i>	0.1615	0.1602	0.1611	0.1392	0.1595	0.1408	0.1600

that the KL divergence is likely to have better generalization. While further experiments will be required to reach definitive conclusions, this has led us to adopt the KL divergence in the remaining experiments.

D. Performance Within the Semantic Space

Fig. 5 (left) presents the precision-recall curves obtained on *Corel50* with QBVE and QBSE. It can be seen that the precision of QBSE is significantly higher than that of QBVE, at most levels of recall. The competitive performance of QBVE at low recall can be explained by the fact that there are always some database images which are visually similar to the query. However, performance decreases much more dramatically than that of QBSE, as recall increases, confirming the better generalization ability of the latter. The MAP scores for QBSE and QBVE

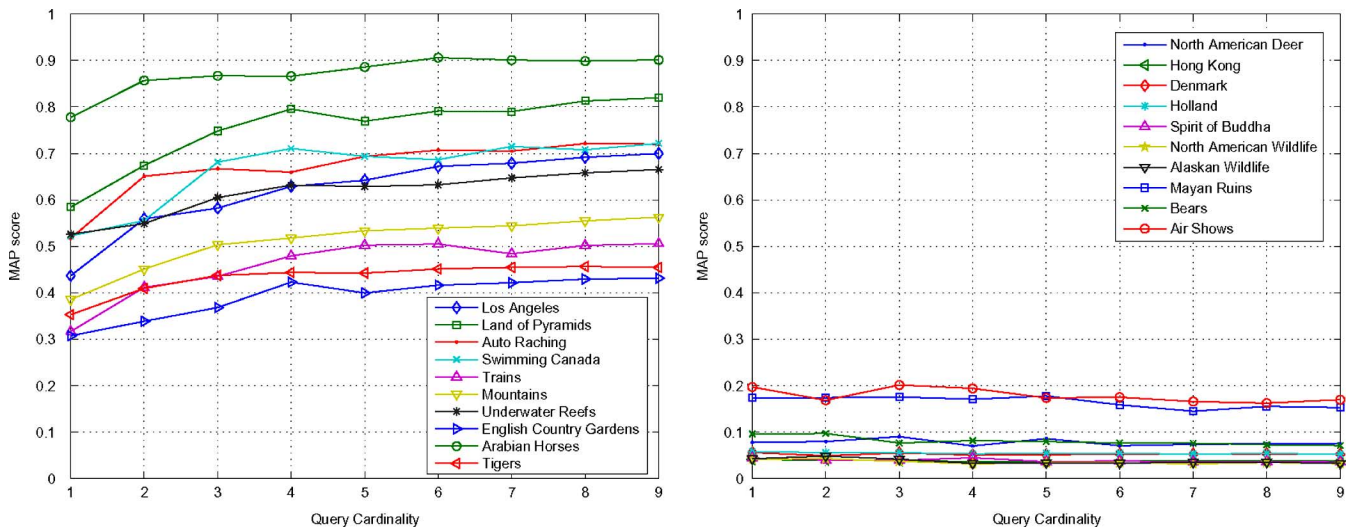


Fig. 9. Effect of multiple image queries on the MAP score of various classes from *Core150*. Left: Classes with highest MAP gains. Right: Classes with lowest MAP gains.

are 0.1665 and 0.1094, respectively, and the chance MAP performance is 0.0200. Fig. 6 presents a comparison of the performance on individual classes, showing that QBSE outperforms QBVE in almost all cases.

The advantages of QBSE are also illustrated by Fig. 7, where we present the results of some queries, under both QBVE and QBSE. Note, for example, that for the query containing *white smoke* and a large area of *dark train*, QBVE tends to retrieve images with *whitish* components, mixed with *dark* components, that have little connection to the *train* theme. Furthermore, the arch-like structure highlighted in Fig. 3 seems to play a prominent role in visual similarity, since three of the five top matches contain arches. Due to its higher level of abstraction, QBSE is successfully able to generalize the main semantic concepts of *train*, *smoke* and *sky*, realizing that the white color is an irrelevant attribute to this query (as can be seen in the last column, where an image of *train with black smoke* is successfully retrieved).

E. Multiple Image Queries

Since the test set of *Core150* contains 9 to 11 images from each class, it is possible to use anywhere from 1 to 9 images per query. When the number of combinations was prohibitively large (for example, there are close to 13 000 combinations of 5 queries), we randomly sampled a suitable number of queries from the set. Fig. 8 (left) shows the MAP values for multiple image queries, as a function of query cardinality, under both QBVE and QBSE for *Core150*. In the case of QBSE, we also compare the three possible query combination strategies: “*LKLD*,” “*SMN*,” and “*KL Combination*.” It is clear that, inside the semantic space, the gains achieved with multiple QBSE queries are unparalleled on the visual domain. In [52], the authors have experimented with multiple query images on a

QVBE system. They show that, using two examples, precision increases by around 15% at 10% recall (over single example queries) but no further improvements are observed for three or more images. We have found that, while the MAP of QBSE increases with the number of images, no gain is observed under QBVE. For QBSE, among the various combination methods, combining SMNs yields best results, with a gain of 29.8% over single image queries. “*LKLD*” and “*KL Combination*” exhibit a gain of 17.3% and 26.4%, respectively. This increase of precision with query cardinality is experienced at all levels of recall.

Fig. 9 shows the performance of 1–9 image queries for the best and the worst ten classes, sorted according to the gain in MAP score. It is interesting to note that in all of the best 10 classes, single image query performs well above chance, while the opposite holds for the worst 10. This means that moderate performance of a QBSE system can be considerably enhanced by using multiple query images, but this is not a cure for fundamental failures. Overall, the MAP score increases with the number of queries for 76% of the classes. For the classes with unsatisfactory MAP score, poor performance can be explained by 1) significant inter-concept overlap (e.g., “*Air Shows*” versus “*Aviation Photography*”), 2) incongruous concepts that would be difficult even for a human labeler (e.g., “*Holland*” and “*Denmark*”), or 3) failure to learn semantic homogeneity among the images, e.g., “*Spirit of Buddha*.” Nevertheless, for 86% of the classes QBSE outperforms QBVE by an average MAP score of 0.136. On the remaining QBVE is only marginally better than QBSE, by an average MAP score of 0.016. Fig. 10 (Left) presents the average precision-recall curves, obtained with the number of image queries that performed best, for QBSE and QBVE on *Core150*. It is clear that QBSE significantly outperforms QBVE at all levels of recall, the average MAP gain being of 111.73%.

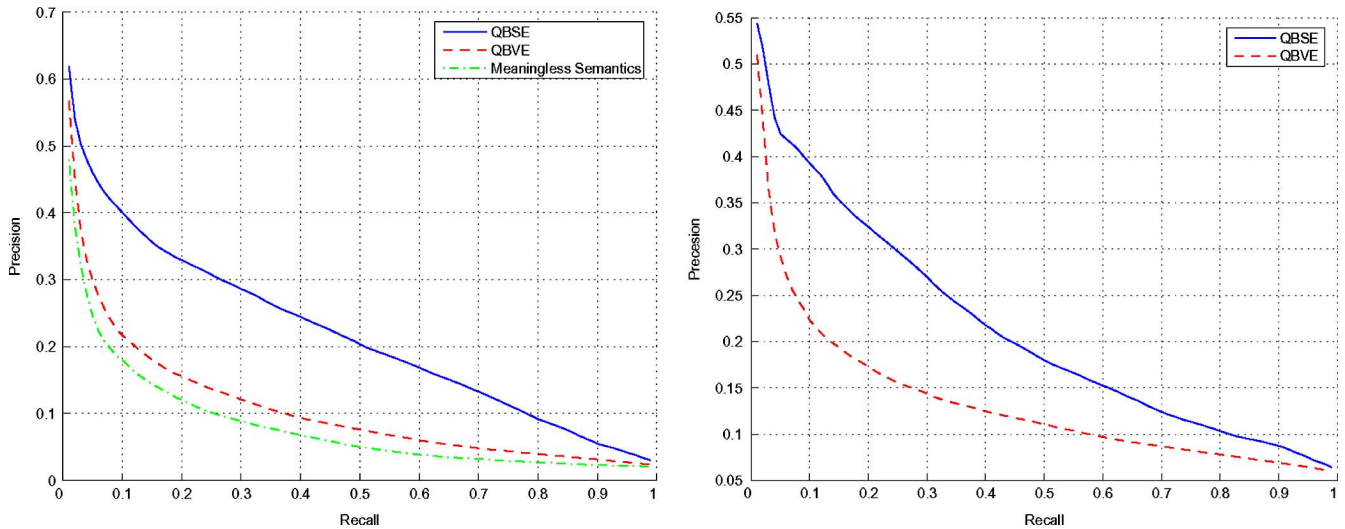


Fig. 10. Best precision-recall curves achieved with QBSE and QBVE on *Corel50*. Left: Inside the semantic space (*Corel50*); also shown is the performance with meaningless semantic space. Right: Outside the semantic space (*Flickr18*).

Query Image	Multiple Image Query					
township						
Helicopter						

Fig. 11. Examples of multiple-image QBSE queries. Two queries (for “Township” and “Helicopter”) are shown, each combining two examples. In each case, two top rows presents the single-image QBSE results, while the third presents the combined query.

F. Performance Outside the Semantic Space

Fig. 5 (Right) presents precision-recall curves obtained on *Flickr18*⁹, showing that outside the semantic space single-query QBSE is marginally better than QBVE. When combined with Fig. 5 (Left), it confirms that, overall, single-query QBSE has better generalization than visual similarity: it is substantially better inside the semantic space, and has slightly better performance outside of it. For multiple image queries we performed experiments with up to 20 images per query (both databases

⁹For brevity, we only document the results obtained with *Flickr18*, those of *Corel15* were qualitatively similar

contain 20 test images per class). As was the case for *Corel50*, multiple image queries benefit QBSE substantially but have no advantage for QBVE. This is shown in Fig. 8 (Right), where we present the MAP score as a function of query cardinality. With respect to the combination strategy, “SMN” once again outperforms “KL” (slightly) and “LKLD Combination” (significantly).

An illustration of the benefits of multiple image queries is given in Fig. 11. The two top rows present query images from the class “Township” (*Flickr18*) and single-query QBSE retrieval results. The third row presents the result of combining the two queries by “SMN combination.” It illustrates the wide

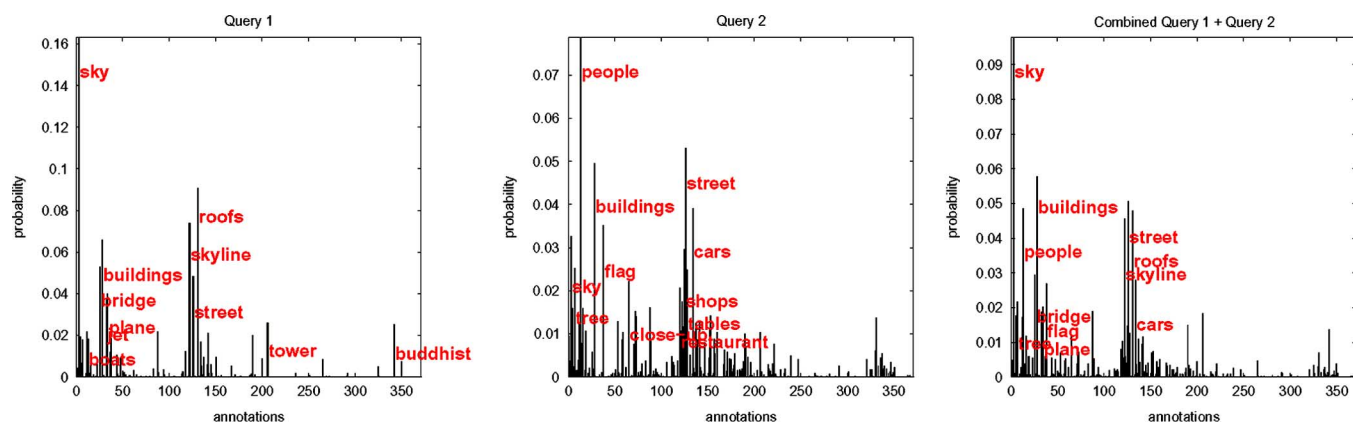


Fig. 12. SMN of individual and combined queries from class “Township” of Fig. 11. Left column shows the first query SMN, center the second and, right the combined query SMN.

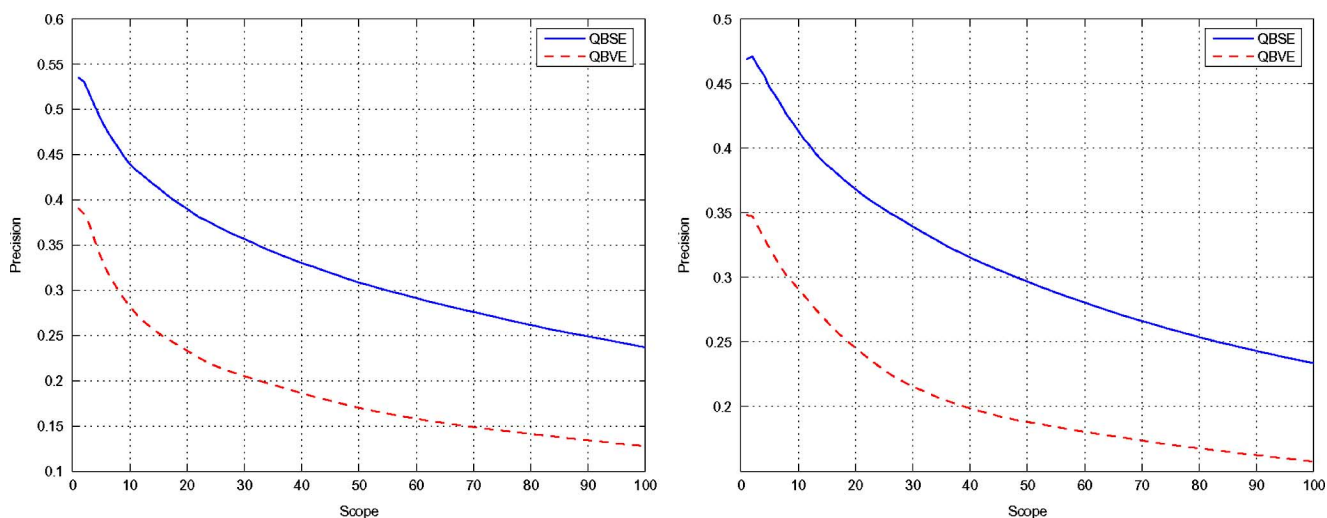


Fig. 13. Performance of QBSE compared to QBVE, based on precision-scope curve for $N = 1$ to 100. Left: Inside the semantic space (*Corel50*). Right: Outside the semantic space (*Flickr18*).

variability of visual appearance of the images in the “Township” class. While single-image queries fail to express the semantic richness of the class, the combination of the two images allows the QBSE system to expand “indoor market scene” and “buildings in open air” to an “open market street” or even a “railway platform.” This is revealed, by the SMN of the combined query, presented in Fig. 12 (right), which is a semantically richer description of the visual concept “Township,” containing concepts (like “sky,” “people,” “street,” “skyline”) from both individual query SMNs. The remaining three rows of Fig. 11 present a similar result for the class “Helicopter” (*Corel15*).

Finally, Fig. 10 presents the best results obtained with multiple queries under both the QBSE and QBVE paradigms. A similar comparison, using the precision-scope curve is shown in Fig. 13. It is clear that, when multiple image queries are adopted, QBSE significantly outperforms QBVE, even outside the semantic space. Table IV summarizes the MAP gains of QBSE, over QBVE, for all datasets considered. In the case of *Flickr18* the gain is of 55.47%. Overall, the table emphatically points out that QBSE significantly outperforms QBVE, both inside and outside the semantic space. Since the basic visual representation (DCT features and Gaussian mixtures) is shared by the two

approaches, this is strong indication that *there is a benefit* to the use of semantic representations in image retrieval. To further investigate this hypothesis we performed a final experiment, based on QBSE with a semantically meaningless space. Building on the fact that all semantic models are learned by grouping images with a common semantic concept, this was achieved by replicating the QBSE experiments with random image groupings. That is, instead of a semantic space composed of concepts like “sky” (learned from images containing sky), we created a “semantic space” of nameless concepts learned from random collections of images. Fig. 10 (left) compares (on *Corel50*) the precision-recall obtained with QBSE on this “meaningless semantic space,” with the previous results of QBVE and QBSE. It is clear that, in the absence of semantic structure, QBSE has *very poor* performance, and is *clearly inferior* to QBVE.

VIII. CONCLUSIONS AND FUTURE WORK

The results above provide *strong* support to the conclusion that *semantic representations have an intrinsic benefit for image retrieval*. While this could be dismissed as a trivial conclusion, we believe that doing so would be unwise, for two main reasons. First, it had not been previously shown that SR systems can

TABLE IV
MAP OF QBVE AND QBSE ON ALL DATASETS CONSIDERED

Database	Chance	QBVE	QBSE	% increase
Corel50	0.0200	0.1067	0.2259	111.73
Corel15	0.0667	0.2176	0.2980	36.95
Flickr18	0.0556	0.1373	0.2134	55.47

generalize beyond the restricted vocabulary on which they are trained. This is certainly not the case for the current standard query paradigm in the SR literature. Second, the results above suggest some interesting hypotheses for future research, which could lead to long-term gains that are more significant than simple out-of-vocabulary generalization. For example, given that the higher abstraction of the semantic representation enables better performance than visual matching, it appears likely that larger gains could be achieved by considering additional semantic layers.

It is currently unclear how significant these gains could be, or where they would stop. On one hand, it does not appear that simply building hierarchical spaces will be sufficient to overcome all the difficulties. On the other, humans seem to be able to perform extremely fast decisions for classification problems involving concepts similar to those considered in this work (e.g., “food” versus “not food,” or “animal” versus “not animal”) [53]. The puzzling properties of these decisions are that 1) the amount of time they require, approximately 150 ms, only allows for strictly feed forward processing with a small number of neural layers, and 2) by human standards the classification results are not great (close to 95% accuracy, on average). It could be that the visual stimulus is classified by a very large number of simple semantic detectors, which identify the concepts that are most likely to be present in the scene. Exact classification would only be done in a second-stage, guided by *attentional* mechanisms, which must decide between a small number of hypothesis and can rely on more complex processing.

Both our results and these observations suggest that there may be a benefit in designing retrieval systems with large concept taxonomies. These could be learned automatically, or exploit the structure of natural language. The QBSE paradigm now proposed could be easily extended to the multi-resolution semantic spaces that are likely to result from a hierarchical concept representation. Furthermore, it would allow an objective characterization of the gains achievable at the different levels of the taxonomy. We intend to explore these questions in future work.

ACKNOWLEDGMENT

The authors would like to thank K. Barnard for providing the Corel dataset used in [14]. Finally, they would like to thank the reviewers for insightful comments that helped to improve the paper.

REFERENCES

[1] A. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, “Content-based image retrieval: The end of the early years,” *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 22, no. 12, pp. 1349–1380, 2000.

[2] W. Niblack *et al.*, “The QBIC Project: Querying Images by Content Using Color, Texture, and Shape,” in *Storage and Retrieval for Image and Video Databases*. San Jose, CA: SPIE, Feb. 1993, pp. 173–181.

[3] A. Pentland, R. Picard, and S. Sclaroff, “Photobook: Content-based manipulation of image databases,” *Int. J. Comput. Vis.*, vol. 18, no. 3, pp. 233–254, Jun. 1996.

[4] N. Vasconcelos and M. Kunt, “Content-based retrieval from image databases: Current solutions and future directions,” in *Proc. Int. Conf. Image Processing*, Thessaloniki, Greece, 2001.

[5] A. Jain and A. Vailaya, “Image retrieval using color and shape,” *Pattern Recognit. J.*, vol. 29, Aug. 1996.

[6] J. Smith and S. Chang, *VisualSEEK: A Fully Automated Content-Based Image Query System*. Boston, MA: ACM, 1996, pp. 87–98.

[7] R. Manmatha and S. Ravela, “A syntactic characterization of appearance and its application to image retrieval,” *Proc. SPIE*, vol. 3016, 1997.

[8] R. Picard, “Digital libraries: Meeting place for high-level and low-level vision,” in *Proc. Asian Conf. Computer Vision*, Singapore, USA, Dec. 1995.

[9] M. Szummer and R. Picard, “Indoor-outdoor image classification,” in *Proc. Workshop in Content-based Access to Image and Video Databases*, Bombay, India, 1998.

[10] A. Vailaya, A. Jain, and H. Zhang, “On image classification: City versus landscape,” *Pattern Recognit.*, vol. 31, pp. 1921–1936, Dec. 1998.

[11] N. Haering, Z. Myles, and N. Lobo, “Locating deducuous trees,” in *Proc. Workshop in Content-based Access to Image and Video Libraries*, San Juan, Puerto Rico, 1997, pp. 18–25.

[12] D. Forsyth and M. Fleck, “Body plans,” in *Proc. IEEE CVPR*, San Juan, Puerto Rico, 1997, pp. 678–683.

[13] K. Barnard and D. Forsyth, “Learning the semantics of words and pictures,” in *Proc. Int. Conf. Computer Vision*, Vancouver, 2001, vol. 2, pp. 408–415.

[14] P. Duygulu, K. Barnard, N. Freitas, and D. Forsyth, “Object recognition as machine translation: Learning a lexicon for a fixed image vocabulary,” in *Proc. ECCV*, Copenhagen, Denmark, 2002.

[15] D. Blei and M. Jordan, “Modeling annotated data,” in *Proc. ACM SIGIR Conf. Research and Development in IR*, 2003.

[16] S. Feng, R. Manmatha, and V. Lavrenko, “Multiple Bernoulli relevance models for image and video annotation,” in *Proc. IEEE CVPR*, Washington, DC, 2004.

[17] V. Lavrenko, R. Manmatha, and J. Jeon, *A model for learning the semantics of pictures*. Vancouver, BC, Canada: NIPS, 2003.

[18] P. C. H. Kueck and N. Freitas, “A constrained semi-supervised learning approach to data association,” in *Proc. ECCV*, Prague, Czech Republic, 2004.

[19] G. Carneiro and N. Vasconcelos, “Formulating semantic image annotation as a supervised learning problem,” in *Proc. IEEE CVPR*, San Diego, 2005.

[20] W. Kraaij, P. Over, and A. Smeaton, “TRECVID 2006: An introduction,” in *Proc. TRECVID*, 2006.

[21] C. Snoek, M. Worring, D. Koelma, and A. Smeulders, “Learned lexicon-driven interactive video retrieval,” in *Proc. CIVR*, 2006, pp. 11–20.

[22] A. Amir *et al.*, “IBM research TRECVID-2005 video retrieval system,” in *Proc. NIST TRECVID Workshop*, Gaithersburg, MD, Nov. 2005.

[23] N. Vasconcelos, “Minimum probability of error image retrieval,” *IEEE Trans. Signal Processing*, vol. 52, pp. 2322–2336, Aug. 2004.

[24] G. Carneiro, A. Chan, P. Moreno, and N. Vasconcelos, “Supervised learning of semantic classes for image annotation and retrieval,” *IEEE Pattern Anal. Machine Intell.*, vol. 29, no. 3, pp. 394–410, Mar. 2007.

[25] G. Salton and J. McGill, *Introduction to Modern Information Retrieval*. New York: McGraw-Hill, 1983.

[26] G. Salton, A. Wong, and C. S. Yang, “A vector space model for automatic indexing,” *Commun. ACM*, vol. 18, no. 11, pp. 613–620, 1975.

[27] Y. Kiyoki, T. Kitagawa, and T. Hayama, “A metadatabase system for semantic image search by a mathematical model of meaning,” *SIGMOD Rec.*, vol. 23, no. 4, pp. 34–41, 1994.

[28] J. Han and L. Guo, “A new image retrieval system supporting query by semantics and example,” in *Proc. Int. Conf. Image Processing*, 2002.

[29] N. Vasconcelos and A. Lippman, “Learning over multiple temporal scales in image databases,” in *Proc. ECCV*, Dublin, Ireland, 2000.

[30] L. W. Y. Lu, H. J. Zhang, and C. Hu, “Joint semantics and feature based image retrieval using relevance feedback,” *IEEE Trans. Multimedia*, vol. 5, no. 3, pp. 339–347, 2003.

[31] X. He, O. King, W. Ma, M. Li, and H. Zhang, “Learning a semantic space from user’s relevance feedback for image retrieval,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 13, no. 1, pp. 39–48, Jan. 2003.

- [32] C. S. Lee, W.-Y. Ma, and H. Zhang, "Information embedding based on user's relevance feedback for image retrieval," *Proc. SPIE*, vol. 3846, pp. 294–304, 1999.
- [33] I. J. Cox, J. Ghosn, T. V. Pappas, and P. N. Yianilos, "Hidden annotation in content based image retrieval," in *Proc. IEEE CVPR*, 1997.
- [34] J. Smith, M. Naphade, and A. Natsev, "Multimedia semantic indexing using model vectors," *ICME*, pp. 445–448, 2003.
- [35] J. R. Smith, C.-Y. Lin, M. R. Naphade, A. Natsev, and B. L. Tseng, "Validity-weighted model vector-based retrieval of video," *Proc. SPIE*, vol. 5307, pp. 271–279, 2003.
- [36] M. Z. J. Lu and S.-P. Ma, "Automatic image annotation based-on model space," in *Proc. IEEE NLP-KE*, 2005.
- [37] J. Platt, "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods," *Adv. Large Margin Classifiers*, vol. 6174, 1999.
- [38] R. Duda, P. Hart, and D. Stork, *Pattern Classification*. New York: Wiley, 2001.
- [39] A. Dempster, N. Laird, and D. Rubin, "Maximum-likelihood from incomplete data via the EM algorithm," *J. R. Statist. Soc.*, vol. B-39, 1977.
- [40] C. Fellbaum, *Wordnet: An Electronic Lexical Database*. Cambridge, MA: MIT Press, 1998.
- [41] N. Vasconcelos, "Image indexing with mixture hierarchies," in *Proc. IEEE CVPR*, Kauai, HI, 2001.
- [42] A. Gelman, J. Carlin, H. Stern, and D. Rubin, *Bayesian Data Analysis*. London, U.K.: Chapman & Hall, 1995.
- [43] N. Vasconcelos, "A unified view of image similarity," in *Proc. Int. Conf. Pattern Recognition*, Barcelona, Spain, 2000.
- [44] T. Cover and J. Thomas, *Elements of Information Theory*. New York: Wiley, 1991.
- [45] M. Swain and D. Ballard, "Color Indexing," *Int. J. Comput. Vis.*, vol. 7, no. 1, pp. 11–32, 1991.
- [46] O. Maron and T. Lozano-Perez, *A Framework For Multiple Instance Learning*, 1998.
- [47] P. Auer, "On learning from multi-instance examples: Empirical evaluation of a theoretical approach," in *Proc. 14th Int. Conf. Machine Learning*, 1997, pp. 21–29.
- [48] J. Smith, "Image retrieval evaluation," in *Proc. IEEE Workshop on Content-Based Access of Image and Video Libraries*, 1998.
- [49] Y. Rui and T. Huang, "Optimizing learning in image retrieval," *Proc. IEEE Conf. Computer Vision Pattern Recognition*, 2000.
- [50] H. Muller, S. Marchand-Maillet, and T. Pun, "The truth about core-evaluation in image retrieval," in *Proc. Int. Conf. Image and Video Retrieval*, 2002, pp. 38–49.
- [51] T. Westerveld and A. P. de Vries, "Experimental evaluation of a generative probabilistic image retrieval model on 'easy' data," in *Proc. Multimedia Information Retrieval Workshop*, Toronto, ON, Canada, Aug. 2003.
- [52] S. M. M. Tahaghoghi, J. A. Thom, and H. E. Williams, "Are two pictures better than one?," in *Proc. 12th Australasian Database Conf.*, Washington, DC, 2001, pp. 138–144.

- [53] S. Thorpe, D. Fize, and C. Marlot, "Speed of processing in the human visual system," *Nature*, pp. 520–522, 1996.



Nikhil Rasiwasia (S'03) received the B.Tech. degree in electrical engineering from the Indian Institute of Technology, Kanpur, in 2005. He is currently pursuing the Ph.D. degree in the Statistical Visual Computing Laboratory, Electrical and Computer Engineering Department, University of California, San Diego.

His main interests are in computer vision and machine learning. He is also interested in photography and computer graphics.

Mr. Rasiwasia was awarded the Certificate of Merit for Academic Excellence in Undergraduate Programs in 2002 and 2004.



Pedro J. Moreno (M'96) received the Ph.D. degree in electrical and computer engineering from Carnegie Mellon University, Pittsburgh, PA.

He is a Senior Research Scientist at the Google Research and Development Center, New York, NY. He joined Google in May 2004. Previously, he was with HP Labs for more than six years. His main interests are in the practical applications of machine learning techniques in several fields, such as audio indexing, image retrieval, text classification, and noise robustness.



Nuno Vasconcelos (M'91) received the Licenciatura degree in electrical engineering and computer science from the Universidade do Porto, Portugal, in 1988, and the M.S. and Ph.D. degrees from the Massachusetts Institute of Technology, Cambridge, in 1993 and 2000, respectively.

From 2000 to 2002, he was a Member of Research Staff at the Compaq Cambridge Research Laboratory, which in 2002 became the HP Cambridge Research Laboratory. In 2003, he joined the Electrical and Computer Engineering Department, University of California, San Diego, where he heads the Statistical Visual Computing Laboratory. He has authored more than 50 peer-reviewed publications. His work spans various areas, including computer vision, machine learning, signal processing and compression, and multimedia systems.

Dr. Vasconcelos is the recipient of a U.S. National Science Foundation CAREER Award and a Hellman Fellowship.