

On the Role of Correlation and Abstraction in Cross-Modal Multimedia Retrieval

Jose Costa Pereira, *Student Member, IEEE*, Emanuele Coviello, Gabriel Doyle, Nikhil Rasiwasia, *Member, IEEE*, Gert R.G. Lanckriet, *Senior Member, IEEE*, Roger Levy, and Nuno Vasconcelos, *Senior Member, IEEE*

Abstract—The problem of cross-modal retrieval from multimedia repositories is considered. This problem addresses the design of retrieval systems that support queries *across* content modalities, for example, using an image to search for texts. A mathematical formulation is proposed, equating the design of cross-modal retrieval systems to that of isomorphic feature spaces for different content modalities. Two hypotheses are then investigated regarding the fundamental attributes of these spaces. The first is that low-level cross-modal correlations should be accounted for. The second is that the space should enable semantic abstraction. Three new solutions to the cross-modal retrieval problem are then derived from these hypotheses: correlation matching (CM), an unsupervised method which models cross-modal correlations, semantic matching (SM), a supervised technique that relies on semantic representation, and semantic correlation matching (SCM), which combines both. An extensive evaluation of retrieval performance is conducted to test the validity of the hypotheses. All approaches are shown successful for text retrieval in response to image queries and vice versa. It is concluded that both hypotheses hold, in a complementary form, although evidence in favor of the abstraction hypothesis is stronger than that for correlation.

Index Terms—Multimedia, content-based retrieval, multimodal, cross-modal, image and text, retrieval model, semantic spaces, kernel correlation, logistic regression



1 INTRODUCTION

CLASSICAL approaches to information retrieval are of a *unimodal* nature [1], [2], [3]. Text repositories are searched with text queries, image databases with image queries, and so forth. This paradigm is of limited use in the modern information landscape, where multimedia content is ubiquitous. Due to this, *multimodal* modeling, representation, and retrieval have been extensively studied in the multimedia literature [4], [5], [6], [7], [8], [9], [10], [11]. In multimodal retrieval systems, queries combining multiple content modalities (e.g., images and sound of a music video clip) are used to retrieve database entries with the same combination of modalities (e.g., other music video clips). These efforts have become increasingly widespread, due in part to large-scale research and evaluation efforts, such as TRECVID [12] and ImageCLEF [13], involving data sets that span multiple data modalities. However, much of this work

has focused on the straightforward extension of methods shown successful in the unimodal scenario. Typically, the different modalities are fused into a representation that does not allow individual access to any of them, for example, some form of dimensionality reduction of a large feature vector that concatenates measurements from images and text. Classical unimodal techniques are then applied to the low-dimensional representation.

In this work, we consider a richer interaction paradigm, which is denoted *cross-modal* retrieval. The goal is to build content models that enable interactivity with content *across* modalities. Such models can then be used to design cross-modal retrieval systems, where queries from one modality (e.g., video) can be matched to database entries from another (e.g., audio tracks). This form of retrieval can be seen as a generalization of current content labeling systems, where a primary modality is augmented with keywords, which can be subsequently searched. Examples include keyword-based image [14], [15], [16] and song [17], [18], [19] retrieval systems.

A defining property of cross-modal retrieval is the requirement that representations generalize across content modalities. This implies the ability to establish cross-modal links between the attributes (of different modalities) characteristic of each document or document class. Detecting these links requires deeper content understanding than what is obtained by classical matching of unimodal attributes. For example, while an image retrieval system can retrieve images of roses by matching red blobs, and a text retrieval system can retrieve texts about roses by matching the “rose” word, a cross-modal retrieval system must *understand* that the word “rose” matches the visual attribute “red blob.” This is much

- J.C. Pereira, E. Coviello, G.R.G. Lanckriet, and N. Vasconcelos are with the Department of Electrical and Computer Engineering, University of California, San Diego, EBU 1, Room 5101, Mail code 0409, 9500 Gilman Drive, La Jolla, CA 92093.
E-mail: {josecp, ecoviell}@ucsd.edu, {gert, nuno}@ece.ucsd.edu.
- G. Doyle and R. Levy are with the Department of Linguistics and the Department of Electrical and Computer Engineering, University of California, San Diego, 9500 Gilman Drive, La Jolla, CA 92093.
E-mail: {gdoyle, rlevy}@ucsd.edu.
- N. Rasiwasia is with the Yahoo!Labs, Bangalore, Karnataka 560037, India.
E-mail: nikhil.rasiwasia@gmail.com.

Manuscript received 16 Apr. 2013; accepted 2 July 2013; published online 11 Aug. 2013.

Recommended for acceptance by F. Fleuret.

For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number TPAMI-2013-04-0257.

Digital Object Identifier no. 10.1109/TPAMI.2013.142.

closer to what humans do than simple color or word matching. Hence, cross-modal retrieval is a better context than unimodal retrieval for the study of the fundamental hypotheses on multimedia modeling.

We exploit representations that generalize across content modalities to study two hypotheses on the joint modeling of images and text. The first, denoted the *correlation hypothesis*, is that explicit modeling of low-level correlations between the different modalities is important for the success of the joint models. The second, denoted the *abstraction hypothesis*, is that model benefits from semantic abstraction, i.e., the representation of images and text in terms of semantic (rather than low level) descriptors. These hypotheses are partly motivated by previous evidence that correlation, for example, correlation analysis on fMRI [20], and abstraction, for example, hierarchical topic models for text clustering [21] or hierarchical semantic representations for image retrieval [22], improve performance on unimodal retrieval tasks. Three joint image-text models that exploit low-level correlation, denoted *correlation matching* (CM), semantic abstraction, denoted *semantic matching* (SM), and both, denoted *semantic correlation matching* (SCM), are introduced.

The correlation and abstraction hypotheses are then tested by measuring the retrieval performance of these models on two reciprocal cross-modal retrieval tasks: 1) the retrieval of text documents in response to a query image, and 2) the retrieval of images in response to a query text. These are basic cross-modal retrieval problems, central to many applications of practical interest, such as finding pictures that effectively illustrate a given text (e.g., illustrate a page of a story book), finding the texts that best match a given picture (e.g., a set of vacation accounts about a given landmark), or searching using a combination of text and images. Model performance on these tasks is evaluated with two data sets: TVGraz [23] and a novel data set based on Wikipedia's featured articles. These experiments show that correlation modeling and abstraction yield independent benefits. In particular, the best results are obtained by a model that accounts for both low-level correlations—by performing a kernel canonical correlation analysis (KCCA) [24], [25]—and semantic abstraction—by projecting images and texts into a common semantic space [22] designed with logistic regression. This suggests that the hypotheses of abstraction and correlation are complementary, each improving the modeling in a different manner.

The paper is organized as follows. Section 2 discusses previous work in multimodal and cross-modal multimedia modeling. Section 3 presents a mathematical formulation for cross-modal modeling and discusses the two fundamental hypotheses analyzed in this work. Section 4 introduces the models underlying correlation, semantic, and semantic correlation matching. Section 5 summarizes an extensive experimental evaluation designed to test the hypotheses. Conclusions are presented in Section 6. A preliminary version of this work appeared in [26].

2 PREVIOUS WORK

The problems of image and text retrieval have been the subject of extensive research in the fields of information

retrieval, computer vision, and multimedia [2], [10], [12], [27], [28].

Unimodal Retrieval. In all these areas, the emphasis has been on unimodal approaches, where query and retrieved documents share a single modality [1], [2], [10], [29], [30]. For example, in [29] a query text, and in [30] a query image is used to retrieve similar text documents and images, based on low-level text (e.g., words) and image (e.g., DCTs) representations, respectively. However, this is not effective for all problems. The existence of a well-known *semantic gap*, between current image representations and those adopted by humans, severely hampers the performance of unimodal image retrieval systems [2].

Annotations. In general, successful retrieval from large-scale image collections requires that the latter be augmented with text metadata provided by human annotators. These manual annotations are typically in the form of a few keywords, a small caption, or a brief image description [12], [13], [27]. When this metadata is available, the retrieval operation tends to be unimodal and ignore the images—the text metadata of the query image is simply matched to the text metadata available for images in the database. Because manual image labeling is labor-intensive, recent research has addressed the problem of automatic image labeling.¹

Labeling. A common assumption is that images can be segmented into regions, which can be described by a small word vocabulary. The focus is then on learning a probability model that relates image regions and words. This can be done by learning a joint probability distribution for words and visual features, for example, using latent Dirichlet allocation (LDA) models [14], probabilistic latent semantic analysis (LSA) [31], histogramming methods [32], or a combination of Bernoulli distributions for text and kernel-based models for visual features [33], [34]. Alternatively, it is possible to use categorized images to train a dictionary of concept models, for example, Gaussian mixtures [16] or two-dimensional hidden Markov models [35], in a weakly supervised manner. The extent of association between images and concepts or words is measured by the likelihood of each image under these models. All these methods assume that each image or image region is associated with a single word.

Semantic Space. An alternative representation, where images are modeled as weighted combinations of concepts in a predefined vocabulary, is proposed in [22]. Statistical models of the distribution of low-level image features are first learned for each concept. The posterior probability of the features extracted from each image, under each of the concept models, is then computed. The image is finally represented by the vector of these posterior concept probabilities. This can be interpreted as a vector of semantic features, establishing a semantic feature space where each dimension is associated with a vocabulary concept. Fig. 1 illustrates how this descriptor, denoted a *semantic multinomial* (SMN), maps the image into the *semantic space*. All standard image analysis/classification tasks can then be conducted in the latter space, at a higher level of abstraction

1. Although not commonly perceived as being *cross-modal*, these systems support cross-modal retrieval, for example, by returning images in response to explicit text queries.

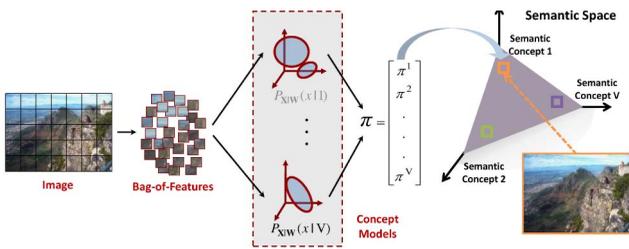


Fig. 1. Semantic space representation. An image is decomposed into a bag-of-features and represented by the vector of its posterior probabilities with respect to the concepts in a semantic vocabulary \mathcal{V} .

than that supported by low-level feature spaces. For example, image retrieval is formulated as retrieval by *semantic similarity*, by combining the semantic space with a suitable similarity function [22]. This allows assessments of image similarity in terms of weighted combinations of vocabulary words and substantially extends the range of concepts that can effectively be retrieved. It also increases the subjective quality of the retrieval results, even when the retrieval system makes mistakes, since images are retrieved by similarity of their content semantics rather than plain visual similarity [36].

Multimodal Retrieval. In parallel with these developments, advances have been reported in multimodal retrieval systems [8], [9], [10], [11], [12], [13], [27]. These are extensions of the classic unimodal systems, where a common retrieval system integrates information from various modalities. This can be done by fusing features from different modalities into a single vector [37], [38], [39], or by learning different models for different modalities and fusing their predictions [40], [41]. One popular approach is to concatenate features from different modalities and rely on unsupervised structure discovery algorithms, such as latent semantic analysis, to find multimodal statistical regularities. A good overview of these methods is given in [39], which also discusses the combination of unimodal and multimodal retrieval systems. Multimodal integration has also been applied to retrieval tasks including audio-visual content [42], [43]. In general, the inability to access each data modality individually (after the fusion of modalities) prevents the use of these systems for cross-modal retrieval.

Cross-Modal Retrieval. Recently, there has been progress toward cross-modal systems. This includes retrieval methods for corpora of images and text [8], [44], images and audio [45], [46], text and audio [47], images, text, and audio [46], [48], [49], [50], [51], or even other sources of data like EEG and fMRI [52]. One popular approach is to rely on manifold learning techniques [46], [48], [49], [50], [51], [52]. These methods learn a manifold from a matrix of distances between multimodal objects. The multimodal distances are formulated as a function of the distances between individual modalities, which allows us to single out particular modalities or ignore missing ones. Retrieval then consists of finding the nearest document, on the manifold, to a multimedia query (which can be composed of any subset of modalities). The main limitation of these methods is the lack of out-of-sample generalization. Since there is no computationally efficient way to project the query into the

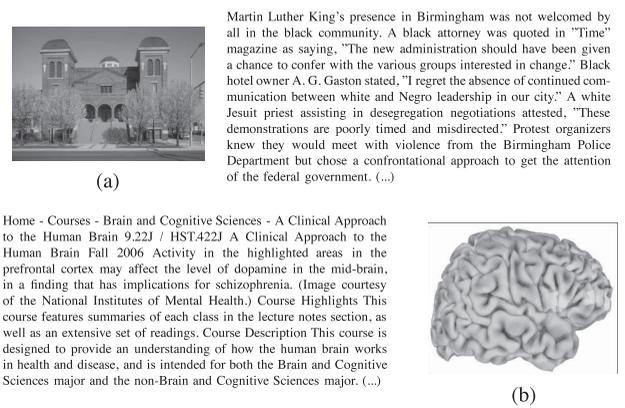


Fig. 2. Two examples of image-text pairs: (a) section from the Wikipedia article on the Birmingham campaign ("History" category), (b) part of a Cognitive Science class syllabus from the TVGraz data set ("Brain" category).

manifold, queries are restricted to the training set used to learn the latter. Hence, all unseen queries must be mapped to their nearest neighbors in this training set, defeating the purpose of manifold learning.

An alternative is to learn correlations between modalities [45], [53]. For example, Li et al. [45] compare canonical correlation analysis (CCA) and cross-modal factor analysis (CFA) in the context of audio-image retrieval. Both CCA and CFA perform a joint dimensionality reduction that extracts highly correlated features in the two data modalities. A kernelized version of CCA was also proposed in [53] to extract translation invariant semantics of text documents written in multiple languages. It was later used to model correlations between web images and corresponding captions in [20]. Another approach is *reranking*: unimodal retrieval is first performed using the query modality, and a second modality is used to rerank the results [54], [55].

Rich Annotation. Despite all these advances, current approaches tend to rely on a limited textual representation, in the form of keywords, captions, or small text snippets. We refer to these as forms of lighter annotation. This is at odds with the ongoing explosion of multimedia content on the web, where it is now possible to collect large sets of extensively annotated data. Examples include news archives, blog posts, or Wikipedia pages, where pictures are related to *complete* text articles, not just a few keywords. We refer to these data sets as *richly annotated*. While potentially more informative, rich annotation establishes a much more nuanced connection between images and text than light annotation. While keywords tend to be explicit image labels, many of the words in a rich text can be unrelated to the image used to illustrate it. For example, Fig. 2 shows a section of the Wikipedia article on the "Birmingham campaign," along with the associated image. Notice that, although related to the text, the image is clearly not representative of all the words in the article. The same is true for the webpage in Fig. 2b, from the TVGraz data set [23]. This is a course syllabus that, beyond the pictured brain, includes course information and other unrelated matters. A major long-term goal of modeling richly annotated data is to recover this *latent* relationship between the text and image components of a document, and exploit it in benefit of practical applications.

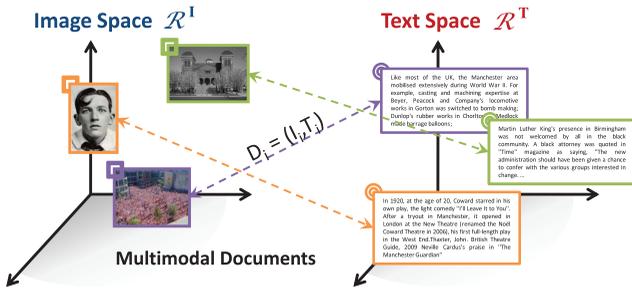


Fig. 3. A document (D_i) is a pair of an image (I_i) and a text (T_i) represented as vectors in feature spaces \mathcal{R}^I and \mathcal{R}^T , respectively. Documents establish a one-to-one mapping between points in \mathcal{R}^I and \mathcal{R}^T .

3 FUNDAMENTAL HYPOTHESES

In this section, we present a novel multimodal content modeling framework, which is flexible and applicable to rich content modalities. Although the fundamental ideas are applicable to any combination of modalities, we restrict the discussion to documents containing images and text.

3.1 The Problem

We consider the problem of information retrieval from a database $\mathcal{D} = \{D_1, \dots, D_{|D|}\}$ of documents comprising *image* and *text* components. Such documents can be quite diverse from a single text complemented by one or more images (e.g., a newspaper article) to documents containing multiple pictures and text sections (e.g., a Wikipedia page). For simplicity, we consider the case where each document consists of a single image and its accompanying text, i.e., $D_i = (I_i, T_i)$. Images and text are represented as vectors in feature spaces \mathcal{R}^I and \mathcal{R}^T , respectively, as illustrated in Fig. 3. In this way, documents establish a one-to-one mapping between points in \mathcal{R}^I and \mathcal{R}^T . Given a text (image) query $T_q \in \mathcal{R}^T$ ($I_q \in \mathcal{R}^I$), the goal of *cross-modal retrieval* is to return the closest match in the image (text) space \mathcal{R}^I (\mathcal{R}^T).

3.2 Multimodal Modeling

Whenever the image and text spaces have a natural correspondence, cross-modal retrieval reduces to a classical retrieval problem. Let

$$\mathcal{M} : \mathcal{R}^T \rightarrow \mathcal{R}^I$$

be an invertible mapping between the two spaces. Given a query T_q in \mathcal{R}^T , it suffices to find the nearest neighbor to $\mathcal{M}(T_q)$ in \mathcal{R}^I . Similarly, given a query I_q in \mathcal{R}^I , it suffices to find the nearest neighbor to $\mathcal{M}^{-1}(I_q)$ in \mathcal{R}^T . In this case, the design of a cross-modal retrieval system reduces to the design of an effective similarity function for determining the nearest neighbors.

In general, however, different representations are adopted for images and text, and there is no natural correspondence between \mathcal{R}^I and \mathcal{R}^T . In this case, the mapping \mathcal{M} has to be learned from examples. In this work, we map the two representations into intermediate spaces, \mathcal{V}^I and \mathcal{V}^T , that have a natural correspondence. This consists of learning two mappings

$$\mathcal{M}_I : \mathcal{R}^I \rightarrow \mathcal{V}^I \quad \mathcal{M}_T : \mathcal{R}^T \rightarrow \mathcal{V}^T$$

from each of the image and text spaces to two *isomorphic* spaces \mathcal{V}^I and \mathcal{V}^T , connected by an invertible mapping

$$\mathcal{M} : \mathcal{V}^T \rightarrow \mathcal{V}^I.$$

Given a text query T_q in \mathcal{R}^T , cross-modal retrieval reduces to finding the image I_r such that $\mathcal{M}_I(I_r)$ is the nearest neighbor of

$$\mathcal{M} \circ \mathcal{M}_T(T_q)$$

in \mathcal{V}^I . Similarly, given an image query I_q in \mathcal{R}^I , the goal is to find text T_r such that $\mathcal{M}_T(T_r)$ is the nearest neighbor of

$$\mathcal{M}^{-1} \circ \mathcal{M}_I(I_q)$$

in \mathcal{V}^T . Under this formulation, the main problem in the design of a cross-modal retrieval system is the design of the intermediate spaces \mathcal{V}^I and \mathcal{V}^T (and the corresponding mappings \mathcal{M}_I and \mathcal{M}_T).

3.3 The Fundamental Hypotheses

Since the goal is to design representations that generalize across content modalities, the solution of this problem requires some ability to derive a more *abstract* representation than the sum of the parts (low-level features) extracted from each content modality. Given that such abstraction is the hallmark of true image or text understanding, this problem enables the exploration of some central questions in multimedia modeling. Consider, for example, a query for a “swan.” While 1) a unimodal image retrieval system can successfully retrieve images of “swans” in that they are the only white objects in a database, 2) a text retrieval system can successfully retrieve documents about “swans” because they are the only documents containing the word “swan,” and 3) a multimodal retrieval system can simply match “white” to “white” and “swan” to “swan,” a cross-modal retrieval system cannot solve the task without understanding that “white is a visual attribute of swan.” Hence, cross-modal retrieval is a more effective paradigm for testing fundamental hypotheses in multimedia representation than unimodal or multimodal retrieval.

In this work, we exploit the cross-modal retrieval problem to test two such hypotheses regarding the joint modeling of images and text:

- \mathcal{H}_1 (*correlation hypothesis*). Low-level cross-modal correlations are important for joint image-text modeling.
- \mathcal{H}_2 (*abstraction hypothesis*). Semantic abstraction is important for joint image-text modeling.

The hypotheses are tested by comparing three possibilities for the design of the intermediate spaces \mathcal{V}^I and \mathcal{V}^T of cross-modal retrieval. In the first case, two feature transformations map \mathcal{R}^I and \mathcal{R}^T onto *correlated* d -dimensional *subspaces* denoted as \mathcal{U}^I and \mathcal{U}^T , respectively, which act as \mathcal{V}^I and \mathcal{V}^T . This maintains the level of semantic abstraction of the representation while maximizing the correlation between the two spaces. We refer to this matching technique as *correlation matching*. In the second case, a pair of transformations is used to map the image and text spaces into a pair of semantic spaces \mathcal{S}^I and \mathcal{S}^T , which then act as \mathcal{V}^I and \mathcal{V}^T . This increases the semantic

TABLE 1
Taxonomy of Proposed Approaches to Cross-Modal Retrieval

	correlation hypothesis	abstraction hypothesis
CM	✓	
SM		✓
SCM	✓	✓

abstraction of the representation without directly seeking correlation maximization. The spaces \mathcal{S}^I and \mathcal{S}^T are made isomorphic by using the same set of semantic concepts for both modalities. We refer to this as *semantic matching*. Finally, a third approach combines the previous two techniques: project onto maximally correlated subspaces \mathcal{U}^I and \mathcal{U}^T , and then project again onto a pair of semantic spaces \mathcal{S}^I and \mathcal{S}^T , which act as \mathcal{V}^I and \mathcal{V}^T . We refer to this as *semantic correlation matching*.

Table 1 summarizes which hypotheses hold for each of the three approaches. The comparative evaluation of the performance of these approaches on cross-modal retrieval experiments provides indirect evidence for the importance of the above hypotheses to the joint modeling of images and text. The intuition is that a better cross-modal retrieval performance results from a more effective joint modeling.

4 CROSS-MODAL RETRIEVAL

In this section, we present the three approaches in detail.

4.1 Correlation Matching

The design of a mapping from \mathbb{R}^T and \mathbb{R}^I to the correlated spaces \mathcal{U}^T and \mathcal{U}^I requires a combination of dimensionality reduction and some measure of correlation between the text and image modalities. In both text and vision literature, dimensionality reduction is frequently accomplished with methods such as latent semantic indexing (LSI) [56] and principal component analysis (PCA) [57]. These are members of a broader class of learning algorithms, denoted subspace learning, which are computationally efficient and produce linear transformations that are easy to conceptualize, implement, and deploy. Furthermore, because subspace learning is usually based on second-order statistics, such as correlation, it can be easily extended to the multimodal setting and kernelized. This has motivated a number of multimodal subspace methods. In this work, we consider *cross-modal factor analysis*, *canonical correlation analysis*, and *kernel canonical correlation analysis*. All these

methods include a training stage, where the subspaces \mathcal{U}^I and \mathcal{U}^T are learned, followed by a projection stage, where images and text are projected into these spaces. Fig. 4 illustrates this process. Cross-modal retrieval is performed in the low-dimensional subspaces.

4.1.1 Linear Subspace Learning

CFA seeks transformations that best represent coupled patterns between different subsets of features (e.g., different modalities) describing the same objects [45]. It finds the orthonormal transformations Ω_I and Ω_T that project the two modalities onto a shared space, $\mathcal{U}^I = \mathcal{U}^T = \mathcal{U}$, where the projections have minimum distance

$$\|X_I\Omega_I - X_T\Omega_T\|_F^2. \quad (1)$$

X_I and X_T are matrices containing corresponding features from the image and text domains, and $\|\cdot\|_F^2$ is the Frobenius norm. It can be shown that this is equivalent to maximizing

$$\text{trace}(X_I\Omega_I\Omega_T'X_T'), \quad (2)$$

and the optimal matrices Ω_I, Ω_T can be obtained by a singular value decomposition of the matrix $X_I'X_T$, i.e.,

$$X_I'X_T = \Omega_I\Lambda\Omega_T', \quad (3)$$

where Λ is the matrix of singular values of $X_I'X_T$ [45].

CCA [58] learns the d -dimensional subspaces $\mathcal{U}^I \subset \mathbb{R}^I$ (image) and $\mathcal{U}^T \subset \mathbb{R}^T$ (text), where the correlation between the two data modalities is maximal. It is similar to principal components analysis, in the sense that it learns a basis of canonical components, directions $w_i \in \mathbb{R}^I$ and $w_t \in \mathbb{R}^T$, but seeks directions along which the data are maximally correlated

$$\max_{w_i \neq 0, w_t \neq 0} \frac{w_i' \Sigma_{IT} w_t}{\sqrt{w_i' \Sigma_I w_i} \sqrt{w_t' \Sigma_T w_t}}, \quad (4)$$

where Σ_I and Σ_T are the empirical covariance matrices for images $\{I_1, \dots, I_{|D|}\}$ and text $\{T_1, \dots, T_{|D|}\}$, respectively, and $\Sigma_{IT} = \Sigma_{TI}'$ the cross covariance between them. Repeatedly solving (4) for directions that are orthogonal to all previously obtained solutions provides a series of canonical components. It can be shown that the canonical components in the image space can be found as the eigenvectors of $\Sigma_I^{-1/2} \Sigma_{IT} \Sigma_T^{-1} \Sigma_{TI} \Sigma_I^{-1/2}$, and in the text space as the eigenvectors of $\Sigma_T^{-1/2} \Sigma_{TI} \Sigma_I^{-1} \Sigma_{IT} \Sigma_T^{-1/2}$. The first d

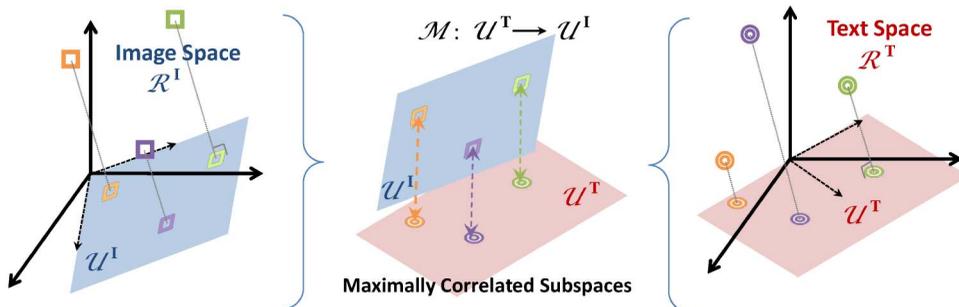


Fig. 4. Correlation matching. Text (\mathbb{R}^T) and images (\mathbb{R}^I) are projected onto two maximally correlated isomorphic subspaces \mathcal{U}_T and \mathcal{U}_I , respectively.

eigenvectors $\{w_{i,k}\}_{k=1}^d$ and $\{w_{t,k}\}_{k=1}^d$ define a basis of the subspaces \mathcal{U}^I and \mathcal{U}^T .

4.1.2 Nonlinear Subspace Learning

CCA and CFA can only model linear dependences between image and text features. This limitation can be avoided by mapping these features into high-dimensional spaces, with a pair of nonlinear transformations $\phi_T: \mathbb{R}^T \rightarrow \mathcal{F}^T$ and $\phi_I: \mathbb{R}^I \rightarrow \mathcal{F}^I$. Application of CFA or CCA in these spaces can then recover complex patterns of dependence in the original feature space. As is common in machine learning, the transformations $\phi_T(\cdot)$ and $\phi_I(\cdot)$ are computed only implicitly by the introduction of two kernel functions $\mathcal{K}_T(\cdot, \cdot)$ and $\mathcal{K}_I(\cdot, \cdot)$, specifying the inner products in \mathcal{F}^T and \mathcal{F}^I , i.e., $\mathcal{K}_T(T_m, T_n) = \langle \phi_T(T_m), \phi_T(T_n) \rangle$ and $\mathcal{K}_I(I_m, I_n) = \langle \phi_I(I_m), \phi_I(I_n) \rangle$, respectively.

KCCA [24], [25] implements this type of extension for CCA, seeking directions $w_i \in \mathcal{F}^I$ and $w_t \in \mathcal{F}^T$, along which the two modalities are maximally correlated in the transformed spaces. The canonical components can be found by solving

$$\max_{\alpha_i \neq 0, \alpha_t \neq 0} \frac{\alpha_i' K_I K_T \alpha_t}{V(\alpha_i, K_I) V(\alpha_t, K_T)}, \quad (5)$$

where $V(\alpha, K) = \sqrt{(1 - \kappa)\alpha' K^2 \alpha + \kappa \alpha' K \alpha}$, $\kappa \in [0, 1]$, is a regularization parameter, and K_I and K_T are the kernel matrices of the image and text representations, for example, $(K_I)_{mn} = \mathcal{K}_I(I_m, I_n)$. Given optimal α_i and α_t for (5), w_i and w_t are obtained as linear combinations of the training examples $\{\phi_I(I_k)\}_{k=1}^{|\mathcal{D}|}$ and $\{\phi_T(T_k)\}_{k=1}^{|\mathcal{D}|}$, with α_i and α_t as weight vectors, i.e., $w_i = \Phi_I(X_I)^T \alpha_i$ and $w_t = \Phi_T(X_T)^T \alpha_t$, where $\Phi_I(X_I)(\Phi_T(X_T))$ is the matrix whose rows contain the high-dimensional representation of the image (text) features. To optimize (5), we solve a generalized eigenvalue problem using the software package of [25]. The first d generalized eigenvectors, where $1 \leq d \leq |\mathcal{D}|$, are the d weight vectors $\{\alpha_{i,k}\}_{k=1}^d$ and $\{\alpha_{t,k}\}_{k=1}^d$ that define the bases $\{w_{i,k}\}_{k=1}^d$ and $\{w_{t,k}\}_{k=1}^d$ of the two maximally correlated d -dimensional subspaces $\mathcal{U}^I \subset \mathcal{F}^I$ and $\mathcal{U}^T \subset \mathcal{F}^T$.

4.1.3 Image and Text Projections

Images and text are represented by their projections p_I and p_T onto the subspaces \mathcal{U}^I and \mathcal{U}^T , respectively. p_I (p_T) is obtained by computing the dot products between the vector representing the image (text) $I \in \mathbb{R}^I$ ($T \in \mathbb{R}^T$) and the image (text) basis vectors spanning \mathcal{U}^I (\mathcal{U}^T). For CFA, the basis vectors are the columns of Ω_I and Ω_T , respectively. For CCA, they are $\{w_{i,k}\}_{k=1}^d$ and $\{w_{t,k}\}_{k=1}^d$. In the case of KCCA, an image $I \in \mathbb{R}^I$ is first mapped into \mathcal{F}^I and subsequently projected onto $\{w_{i,k}\}_{k=1}^d$, i.e., $p_I = \mathcal{P}_I(\phi_I(I))$ with

$$\begin{aligned} p_{I,k} &= \langle \phi_I(I), w_{i,k} \rangle \\ &= \langle \phi_I(I), [\phi_I(I_1), \dots, \phi_I(I_{|\mathcal{D}|})] \alpha_{i,k} \rangle \\ &= [\mathcal{K}_I(I, I_1), \dots, \mathcal{K}_I(I, I_{|\mathcal{D}|})] \alpha_{i,k}, \end{aligned} \quad (6)$$

where $k = 1, \dots, d$. Analogously, a text $T \in \mathbb{R}^T$ is mapped into \mathcal{F}^T and then projected onto $\{w_{t,k}\}_{k=1}^d$, i.e., $p_T = \mathcal{P}_T(\phi_T(T))$, using $\mathcal{K}_T(\cdot, \cdot)$.

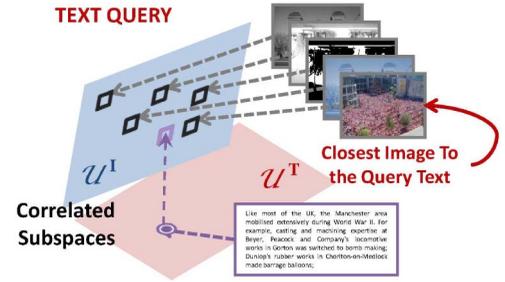


Fig. 5. Example of cross-modal retrieval using CM. Here, CM is used to find the images that best match a query text.

4.1.4 Correlation Matching

For all methods, a natural invertible mapping between the projections onto \mathcal{U}^I and \mathcal{U}^T follows from the correspondence between the d -dimensional bases of the subspaces, as $w_{i,1} \leftrightarrow w_{t,1}, \dots, w_{i,d} \leftrightarrow w_{t,d}$. This results in a compact, efficient representation of both modalities, where vectors p_I and p_T are coordinates in two isomorphic d -dimensional subspaces, as shown in Fig. 4. Given an image query I with projection p_I , the text $T \in \mathbb{R}^T$ that most closely matches it is that for which p_T minimizes

$$D(I, T) = d(p_I, p_T), \quad (7)$$

for some suitable distance measure $d(\cdot, \cdot)$ in a d -dimensional vector space. Similarly, given a query text T with projection p_T , the closest image match $I \in \mathbb{R}^I$ is that for which p_I minimizes $d(p_I, p_T)$. An illustration of cross-modal retrieval using CM is given in Fig. 5.

4.2 Semantic Matching

An alternative to subspace learning is to map images and text to representations at a higher level of abstraction, where a natural correspondence can be established. This is obtained by augmenting the database \mathcal{D} with a vocabulary $\mathcal{V} = \{v_1, \dots, v_K\}$ of semantic concepts. These can be generic or application dependent, ranging from generic document attributes, such as “Long” or “Short,” to specific topics such as “History” or “Biology,” or any other categories that are deemed relevant. Individual documents are grouped into these semantic concepts. Two mappings \mathcal{L}_T and \mathcal{L}_I are then implemented using classifiers of text and images, respectively. \mathcal{L}_T maps a text $T \in \mathbb{R}^T$ into a vector π_T of posterior probabilities $P_{V|T}(v_j|T)$, $j \in \{1, \dots, K\}$ with respect to each of the concepts in \mathcal{V} . The space S^T of these vectors is referred to as the *semantic space for text*, and the probabilities in π_T as the *semantic text features*. Similarly, \mathcal{L}_I maps an image I into a vector π_I of *semantic image features* in a *semantic space for images* S^I .

Semantic representations have two advantages for cross-modal retrieval. First, they provide a higher level of abstraction. While features in \mathbb{R}^T and \mathbb{R}^I frequently have no obvious interpretation (e.g., image features tend to be edges, edge orientations or frequency bases), the features in S^T and S^I are (semantic) concept probabilities (e.g., the probability that the image belongs to the “History” or “Biology” document classes). Previous work has shown that increased feature abstraction can lead to substantially better generalization for tasks such as image retrieval [22]. Second,

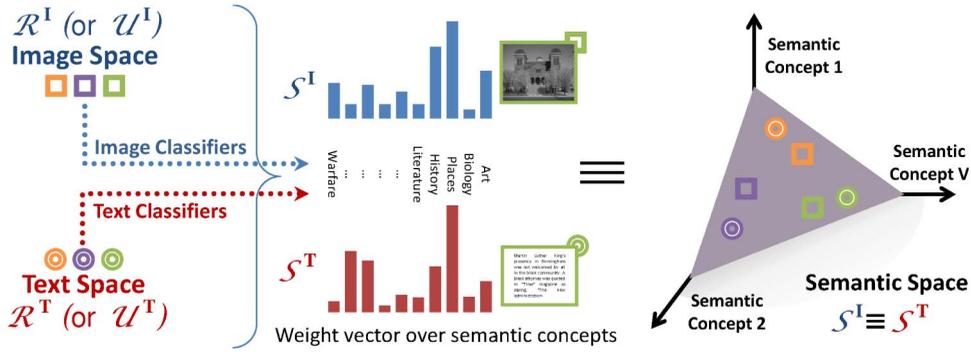


Fig. 6. Semantic matching. Text and images are mapped into a common semantic space, using the posterior class probabilities produced by a multiclass text or image classifier.

the semantic spaces \mathcal{S}^T and \mathcal{S}^I are isomorphic, since both images and text are represented as vectors of posterior probabilities with respect to the *same* set of semantic concepts. Hence, the spaces can be treated as being the same, i.e., $\mathcal{S}^T = \mathcal{S}^I$, leading to the representation of Fig. 6.

4.2.1 Learning

Many classification techniques can be used to learn the mappings \mathcal{L}_T and \mathcal{L}_I . In this work, we consider three popular methods. Logistic regression computes the posterior probability of a particular class by fitting image (text) features to a logistic function. Parameters are chosen to minimize the loss function,

$$\min_w \frac{1}{2} w'w + C \sum_i \log(1 + \exp(-y_i w'x_i)), \quad (8)$$

where y_i is the class label, x_i the feature vector in the input space, and w a vector of parameters. A multiclass logistic regression can be learned for the image and text modalities, by making x_i the image and text representation, $I \in \mathbb{R}^I$ and $T \in \mathbb{R}^T$, respectively. In our implementation, this is done with the Liblinear software package [59].

Support vector machines (SVMs) learn the separating hyperplane of largest margin between two classes, using

$$\begin{aligned} \min_{w,b,\xi} \quad & \frac{1}{2} w'w + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i(w'x_i + b) \geq 1 - \xi_i, \quad \forall i \\ & \xi_i \geq 0, \end{aligned} \quad (9)$$

where w and b are the hyperplane parameters, y_i the class label, x_i input feature vectors, ξ_i slack variables that allow outliers, and $C > 0$ a penalty on the number of outliers. Although the SVM output does not have a probabilistic interpretation, a sigmoidal transformation of the SVM scores $y_i w'x_i$ is often taken as a proxy for the posterior class probabilities. This is, for example, supported by the LibSVM [60] package, which we use in our implementation.

Boosting methods combine weak learners into a strong decision rule. Many boosting algorithms have been proposed in the literature. In this work, we adopt the multiclass boosting method of [61]. This is based on multidimensional codewords (y^k) and predictors (f). Each class k is mapped to a distinct class label y^k , and the strong classifier, $F(x)$, is a mapping from examples $x_i \in \mathcal{X}$ into class labels y^k

$$F(x) = \arg \max_k y^k f^*(x), \quad (10)$$

where $f^*(x) : \mathcal{X} \rightarrow \mathbb{R}$ is the continuous valued predictor that maximizes the classification margin. Posterior class probabilities can then be recovered by applying a nonlinear transformation to the classifier output. In our implementation, this is done with recourse to the multiclass boosting software package of [61].

4.2.2 Retrieval

Given a query image I (text T), represented by $\pi_I \in \mathcal{S}^I$ ($\pi_T \in \mathcal{S}^T$), SM-based cross-modal retrieval returns the text T (image I), represented by $\pi_T \in \mathcal{S}^T$ ($\pi_I \in \mathcal{S}^I$), that minimizes

$$D(I, T) = d(\pi_I, \pi_T), \quad (11)$$

for some suitable distance measure d between probability distributions. An illustration of cross-modal retrieval using SM is given in Fig. 7.

4.3 Semantic Correlation Matching

CM and SM are not mutually exclusive. In fact, a corollary to the two hypotheses discussed above is that there may be a benefit in combining CM and SM. CM extracts maximally correlated features from \mathbb{R}^T and \mathbb{R}^I . SM builds semantic spaces using original features to gain semantic abstraction. When the two are combined by building semantic spaces using the feature representation produced by correlation maximization, it may be possible to improve on the individual performances of both CM and SM. To combine the two approaches, the maximally correlated subspaces \mathcal{U}^I and \mathcal{U}^T are first learned and the projections (p_I, p_T) of each image-text pair (I, T) computed, as discussed in Section 4.1. The transformations \mathcal{L}_I and \mathcal{L}_T are then learned in each of these subspaces to produce the semantic spaces \mathcal{S}^I and \mathcal{S}^T , respectively. Retrieval is finally based on the image-text distance $D(I, T)$ of (11), based on the semantic mappings $\pi_I = \mathcal{L}_I(p_I)$ and $\pi_T = \mathcal{L}_T(p_T)$.

5 EXPERIMENTS

In this section, we describe an extensive experimental evaluation of the proposed cross-modal retrieval framework.

TABLE 2
MAP Scores (Validation Set) of Different Distance Measures

	measure	$d(p, q)$	TVGraz			Wikipedia		
			img query	txt query	avg	img query	txt query	avg
CM	ℓ_1	$\sum_i p_i - q_i $	0.376	0.418	0.397	0.193	0.234	0.214
	ℓ_2	$\sum_i (p_i - q_i)^2$	0.391	0.444	0.417	0.199	0.243	0.221
	NC	$\frac{p^T q}{\ p\ \ q\ }$	0.498	0.476	0.487	0.288	0.239	0.263
	NC_c	$\frac{(p-\mu_p)^T (q-\mu_q)}{\ p-\mu_p\ \ q-\mu_q\ }$	0.486	0.462	0.474	0.287	0.239	0.263
SM	KL	$\sum_i p_i \log \frac{p_i}{q_i}$	0.362	0.564	0.463	0.206	0.274	0.240
	ℓ_1	$\sum_i p_i - q_i $	0.525	0.573	0.549	0.220	0.274	0.247
	ℓ_2	$\sum_i (p_i - q_i)^2$	0.492	0.570	0.531	0.205	0.276	0.241
	NC	$\frac{p^T q}{\ p\ \ q\ }$	0.582	0.581	0.582	0.301	0.276	0.289
	NC_c	$\frac{(p-\mu_p)^T (q-\mu_q)}{\ p-\mu_p\ \ q-\mu_q\ }$	0.598	0.578	0.588	0.352	0.272	0.312
SCM	KL	$\sum_i p_i \log \frac{p_i}{q_i}$	0.560	0.623	0.592	0.311	0.270	0.291
	ℓ_1	$\sum_i p_i - q_i $	0.623	0.633	0.628	0.334	0.273	0.304
	ℓ_2	$\sum_i (p_i - q_i)^2$	0.605	0.615	0.610	0.315	0.267	0.291
	NC	$\frac{p^T q}{\ p\ \ q\ }$	0.665	0.632	0.649	0.371	0.279	0.325
	NC_c	$\frac{(p-\mu_p)^T (q-\mu_q)}{\ p-\mu_p\ \ q-\mu_q\ }$	0.669	0.633	0.651	0.382	0.281	0.332

μ_p and μ_q are the sample averages for p and q , respectively.

5.2.2 Correlation Matching

A set of experiments was performed to compare the performance of CFA, CCA, and KCCA. In all cases, the number of canonical components was validated in each retrieval experiment. As shown in Table 3, KCCA had the top performance. Best results were achieved with a chi-square radial basis function kernel⁴ for images, a histogram intersection kernel for text [65], [66], and regularization constants $\kappa = 10\%$ on TVGraz and $\kappa = 50\%$ on Wikipedia. To verify the importance of modeling correlations, we considered two alternative representations. The first implemented dimensionality reduction but no correlation modeling. The two modalities were independently projected into subspaces of the same dimension, learned with PCA. The second investigated the benefits of complementing correlation with discriminant modeling, by introducing a linear discriminant analysis on the correlated subspaces discovered by KCCA. It is denoted *linear discriminant kernel canonical correlation analysis* (LD-KCCA). As shown in Table 3, neither alternative improved on the average MAP scores of KCCA. This shows that there are benefits to correlation matching beyond dimensionality reduction and that further gains are not trivial to achieve, supporting the hypothesis that correlation modeling is important for cross-modal retrieval. Given its good performance, KCCA was used in all remaining CM experiments.

5.2.3 Semantic Matching

A set of experiments was performed to evaluate the impact of the classification architecture used to design the semantic space on retrieval accuracy. Three architectures were compared: logistic regression, boosting, and SVMs. As shown in Table 4, the semantic space obtained with logistic regression performed best for both cross-modal retrieval tasks. It was, thus, chosen to implement SM in all remaining experiments.

4. $\mathcal{K}(x, y) = \exp\left(\frac{d_{\chi^2}(x, y)}{\gamma}\right)$, where $d_{\chi^2}(x, y)$ is the chi-square distance between x and y , and γ is the average chi-square distance among training points.

5.2.4 Optimization

The experiments above lead to a retrieval architecture that combines KCCA for learning correlated subspaces, logistic regression to learn the semantic space, and the centered normalized correlation NC_c distance measure to evaluate

TABLE 3
MAP Scores (Validation Set) under the CM Hypothesis

	img q.	txt q.	avg.
TVGraz			
LD-KCCA	0.428	0.471	0.450
KCCA	0.486	0.462	0.474
CCA	0.284	0.254	0.269
CFA	0.195	0.179	0.187
PCA	0.162	0.144	0.153
Wikipedia			
LD-KCCA	0.242	0.241	0.242
KCCA	0.287	0.239	0.263
CCA	0.210	0.174	0.192
CFA	0.195	0.156	0.176
PCA	0.208	0.132	0.170

TABLE 4
MAP Scores (Validation Set) under the SM Hypothesis

	img q.	txt q.	avg.
TVGraz			
Log. Reg.	0.598	0.578	0.588
SVM	0.556	0.548	0.552
Boosting	0.567	0.476	0.522
Wikipedia			
Log. Reg.	0.352	0.272	0.312
SVM	0.318	0.237	0.278
Boosting	0.322	0.207	0.265

TABLE 5
Optimal Parameters (Validation Set) for
Best Retrieval Architecture

	Cbk size	LDA topics	KCCA comps
TVGraz			
SCM		400	700
SM	4096	100	n/a
CM		200	8
Wikipedia			
SCM		200	100
SM	4096	600	n/a
CM		20	10

(7) and (11). Using this architecture, a final round of experiments was used to determine the best combination of 1) BOW codebook size for image representation, 2) number of LDA topics for text representation, and 3) number of KCCA components, for each of the CM, SM, and SCM retrieval regimes and data set. Table 5 summarizes the optimal parameter configuration, which was used in the remaining experiments.

5.3 Testing the Fundamental Hypotheses

This architecture was used on a set of experiments aimed to test the fundamental hypotheses of Section 3. In these experiments, MAP scores were measured on the test set.

5.3.1 Overall Performance

Table 6 compares the scores of cross-modal retrieval with CM, SM, SCM, and the baseline TTI method. The table provides evidence in support of the two hypotheses of Section 3.3, both joint dimensionality reduction and semantic abstraction are beneficial for multimodal modeling, leading to a nontrivial improvement over TTI. For example, in TVGraz, the average MAP score of CM is more than double that of TTI. For SM, the improvement is more than threefold. Overall, the best performance is achieved by SCM. Similar conclusions can be drawn for Wikipedia, although the average gains of SCM are slightly lower than in TVGraz. This is not surprising, since the retrieval scores are generally lower on Wikipedia than on TVGraz. As discussed in Section 5.1, this is explained by the broader scope of the Wikipedia categories.

Fig. 8 presents a more detailed analysis of the retrieval performance, in the form of PR curves. CM, SM, and SCM again achieve large improvements over TTI. These improvements tend to occur at all levels of recall, indicating better generalization, and often involve substantial increases in precision, indicating higher accuracy. Overall, these results suggest that the contributions of cross-modal correlation and semantic abstraction are *complementary*: not only is there an independent benefit to both correlation modeling and abstraction, but the *best performance is achieved when the two are combined*.

TABLE 6
MAP Scores (Test Set) of CM, SM, SCM, and TTI,
on TVGraz and Wikipedia

	img. query	txt. query	avg.	gain
TVGraz				
SCM	0.664	0.649	0.657	—
SM	0.619	0.585	0.602	9%
CM	0.460	0.450	0.455	44%
TTI [44]	0.216	0.153	0.185	255%
Wikipedia				
SCM	0.362	0.273	0.318	—
SM	0.350	0.249	0.300	6%
CM	0.267	0.219	0.243	31%
TTI [44]	0.237	0.137	0.187	70%

5.3.2 Per-Class Performance

Fig. 8 shows the per-class MAP scores of all methods. SCM has higher MAP than CM and SM on all classes of TVGraz and is either comparable to or better than CM and SM on the majority of Wikipedia classes. TTI does very poorly in general and seems biased toward one class. This is evident from Figs. 8c and 8f, where it achieves a very high score on one class—“Frog” on TVGraz and “Warfare” on Wikipedia—and very low scores in the remaining. In both cases, the favored class has a large number of training examples.

Two examples of text queries and corresponding retrieval results, using SCM, are shown in Fig. 10. The text query is presented along with its probability vector π_T and the ground-truth image. The top five image matches are shown below the text, along with their probability vectors π_I . Finally, Fig. 11 shows some examples of image-to-text retrieval. Since displaying the retrieved texts would require too much space, we present the associated ground-truth images instead. The query images are framed in the left column, and the images associated with the four best text matches are shown on the right.

5.4 Robustness

The previous experiments indicate that semantic spaces are beneficial for cross-modal retrieval. However, in each experiment, the semantic space was designed with a vocabulary \mathcal{V} identical to the ground-truth semantics. This could be argued to give an unfair advantage to SM and SCM. To evaluate this possibility, we performed a number of additional experiments that evaluated the robustness of SM to mismatches between semantic vocabulary and ground-truth semantics. Two classes of experiments were performed.

5.4.1 Extended Semantics

This set of experiments tested the impact of the size of the vocabulary \mathcal{V} on SM performance. It was based on an *extended vocabulary* \mathcal{V}' , which was *shared by the two data sets*. This contained the 10 classes from TVGraz, the 10 classes of Wikipedia, the 20 classes of Wikipedia featured articles that were not used in the Wikipedia data set, and 20 categories from the *Pascal-Sentences* [67] data set (50 image/text pairs per class). Overall, \mathcal{V}' contained 60 classes. The ground-truth semantics were as before, i.e., the classes in the second column of Tables 7 and 8.

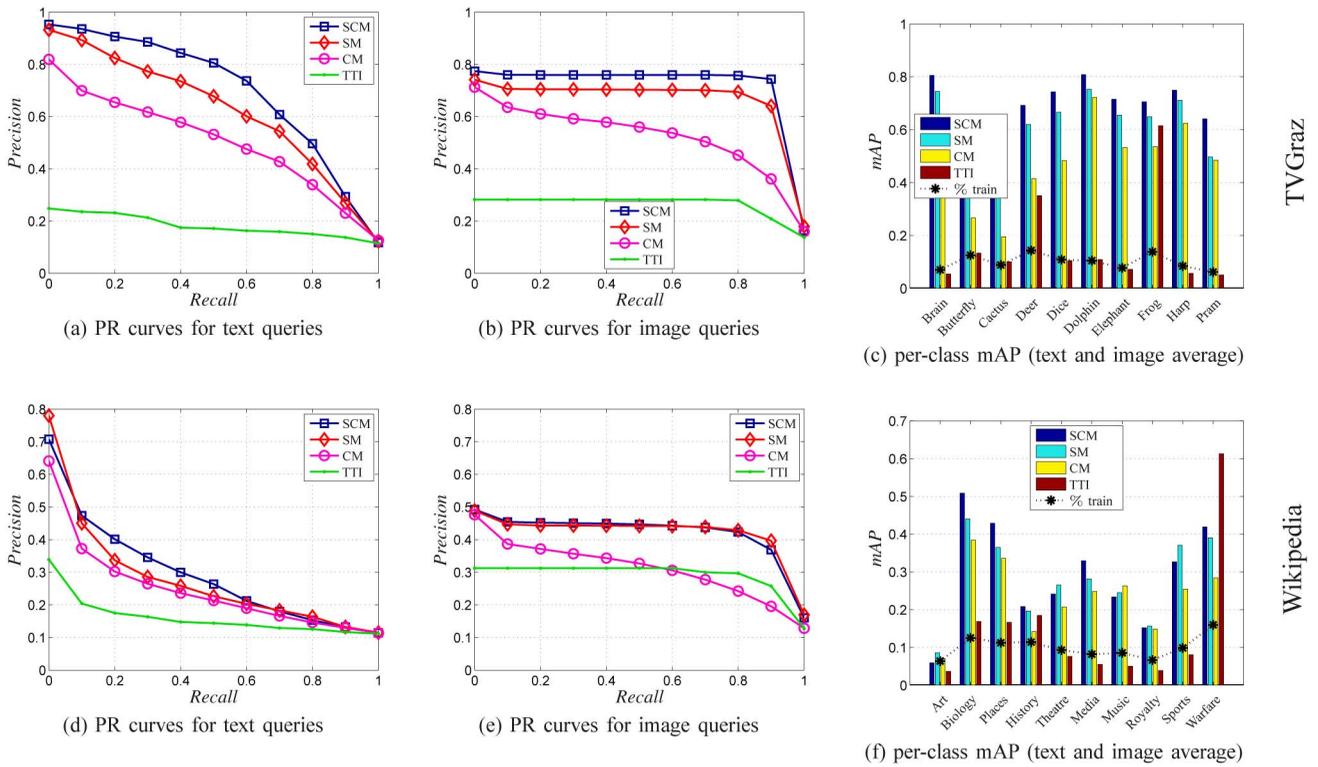


Fig. 8. PR curves of cross-modal retrieval using both text (a), (d) and image (b), (e) queries on TVGraz (top) and Wikipedia (bottom). Average (across image and text queries) per-class mAP scores also shown in (c) and (f).

To evaluate the impact of the composition of the semantic space on retrieval scores, we repeated the retrieval experiment using multiple subsets of \mathcal{V} as vocabulary \mathcal{V} . Starting with \mathcal{V} containing the 10 ground-truth classes, we sequentially added one of the remaining classes in \mathcal{V} to \mathcal{V} . This produced a sequence of semantic spaces with between 11 and 60 dimensions. To introduce randomness, the whole experiment was repeated five times, using a sequence of randomly selected classes to add at each step. Fig. 9 presents the MAP scores as a function of the vocabulary size, for image and text queries on the two data sets. The straight horizontal lines are the scores obtained when \mathcal{V} contained the 10 original classes. The image query task appears to be slightly more affected than its text counterpart; this is a natural consequence of the noisier semantic

descriptor of images when compared to that of texts [68]. While there is some degradation of performance as the vocabulary grows, the effect is small. This indicates that the performance of SM is fairly insensitive to the size of the vocabulary \mathcal{V} .

5.4.2 Alternative Semantics

In the previous experiments, the vocabulary \mathcal{V} always included the ground-truth semantics. To further test the robustness of SM to the make-up of the semantic space, a final set of experiments was performed with ground-truth semantics that are only loosely related to the vocabulary \mathcal{V} . For this, we defined a new set of ground-truth semantics for each data set, according to Tables 7 and 8. In all

TABLE 7
TVGraz Semantics

Alternative	Vocabulary
Anatomy	1. Brain
Pollination	2. Butterfly
	3. Cactus
Land Animals	4. Deer
	7. Elephant
Marine Animals	6. Dolphin
	8. Frog
Objects	5. Dice
	9. Harp
	10. Pram

TABLE 8
Wikipedia Semantics

Alternative	Vocabulary
Humanities	1. Art & architecture
	3. Geography & places
	4. History
	5. Literature & theatre
	2. Biology
Nature	6. Media
	7. Music
Entertainment	9. Sport & recreation
	8. Royalty & nobility
Honor	10. Warfare

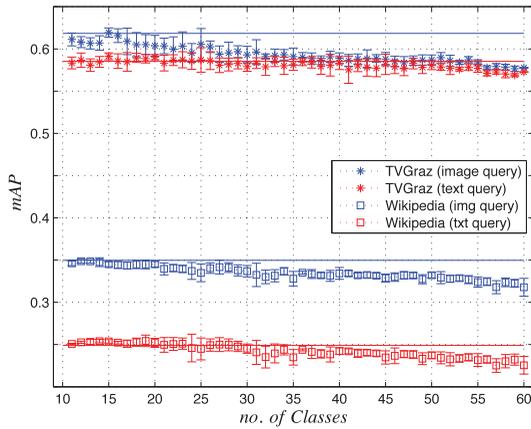


Fig. 9. MAP scores under SM. The solid horizontal line is the score obtained with the 10 original data set categories.

experiments, the vocabulary \mathcal{V} consisted of the original data set classes, also shown in the tables.

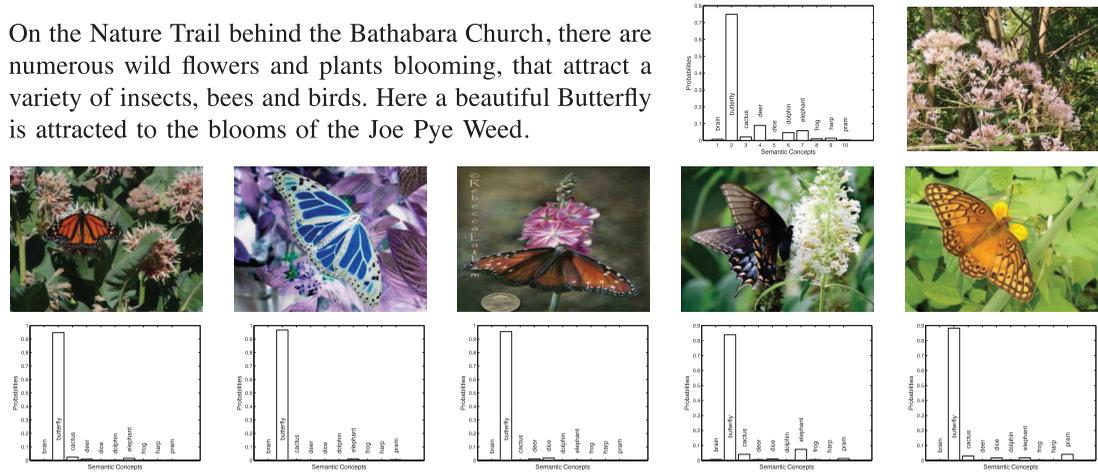
Table 9 presents a comparison of the average MAP scores achieved with the alternative ground-truth semantics of Tables 7 and 8 (denoted “alt. semantics”) and with the original data set classes (denoted “vocabulary”). Since there are fewer classes in the alternative semantics, the retrieval

TABLE 9
Average MAP Scores (Test Set) under the Original (“Vocabulary”) and Alternative Semantics

	Retrieval Evaluation	
	Alternative	Vocabulary
TVGraz		
SCM	0.584	0.657
SM	0.568	0.602
CM	0.492	0.455
TTI [44]	0.292	0.185
Wikipedia		
SCM	0.448	0.318
SM	0.436	0.300
CM	0.413	0.243
TTI [44]	0.347	0.187

performance is expected to improve. However, the fact that these classes are more abstract could also lead to a degradation. The two behaviors are visible in the table. On Wikipedia, where the original classes are already quite abstract, all methods have improved performance under the alternative semantics. On TVGraz, where the alternative semantics are much more abstract than the vocabulary classes, performance decreases for SM and SCM. Note,

On the Nature Trail behind the Bathabara Church, there are numerous wild flowers and plants blooming, that attract a variety of insects, bees and birds. Here a beautiful Butterfly is attracted to the blooms of the Joe Pye Weed.



Between October 1 and October 17, the Japanese delivered 15,000 troops to Guadalcanal, giving Hyakutake 20,000 total troops to employ for his planned offensive. Because of the loss of their positions on the east side of the Matanikau, the Japanese decided that an attack on the U.S. defenses along the coast would be prohibitively difficult. Therefore, Hyakutake decided that the main thrust of his planned attack would be from south of Henderson Field. His 2nd Division (augmented by troops from the 38th Infantry Division), under Lieutenant General Masao Maruyama and comprising 7,000 soldiers in three infantry regiments of three battalions each was ordered to march through the jungle and attack the American defenses from the south near the east bank of the Lunga River. Shaw, “First Offensive”, p. 34, and Rottman, “Japanese Army”, p. 63. (...)

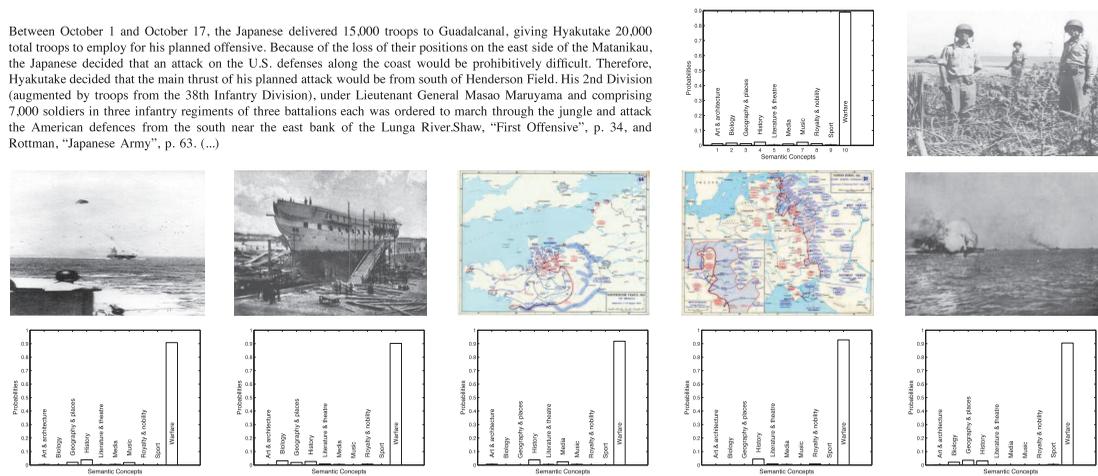


Fig. 10. Two examples of text-based cross-modal retrieval using SCM. The first example is from TVGraz and the second example from Wikipedia. The query text, associated probability vector, and ground-truth image are shown on the top; retrieved images are presented at the bottom.



Fig. 11. Image-to-text retrieval on TVGraz (top row) and Wikipedia (bottom). Query images are framed in the far-left column. The four most relevant texts, represented by their ground-truth images, are shown in the remaining columns.

however, that these variations do not affect the relative performance of the different methods. In both cases, CM and SM achieve significant improvements over TTI and the best overall performance is obtained when they are combined (SCM). In summary, this experiment confirms all the conclusions reached above.

6 CONCLUSION

The increasing availability of multimodal information demands novel representations for content-based retrieval. In this work, we proposed models applicable to cross-modal retrieval. This entails the retrieval of database entries from one content modality in response to queries from another. While the emphasis was on cross-modal retrieval of images and rich text, the proposed models support many other content modalities. By requiring representations that can generalize across modalities, cross-modal retrieval establishes a suitable context for the objective investigation of fundamental hypotheses in multimedia modeling.

We have considered two such hypotheses, regarding the importance of low-level cross-modal correlations and semantic abstraction in multimodal modeling. The hypotheses were objectively tested by comparing the performance of three methods: 1) CM, based on the correlation hypothesis, 2) SM, based on the abstraction hypothesis, and 3) SCM, based on the combination of the two. All of these map objects from different native spaces (e.g., rich text and images) to a pair of isomorphic spaces, where a natural correspondence can be established for cross-modal retrieval purposes. The retrieval performance of the three solutions was tested on two data sets, "Wikipedia" and "TVGraz," which combine images and rich text, and compared to a *state-of-the-art* cross-modal retrieval method (TTI).

While the two fundamental hypotheses were shown to hold for the two data sets, where both CM and SM achieved significant improvements over TTI, SM achieved overall better performance than CM. This implies stronger evidence for the abstraction than for the correlation hypothesis. However, the two hypotheses were also found to be complementary, with SCM achieving the best results of all methods considered.

ACKNOWLEDGMENTS

This work was funded by FCT graduate Fellowship SFRH/BD/40963/2007 and US National Science Foundation grant CCF-0830535. The authors would like to thank Malcolm Slaney for helpful discussions.

REFERENCES

- [1] G. Salton and M. McGill, *Introduction to Modern Information Retrieval*. McGraw-Hill, 1983.
- [2] A. Smeulders, M. Worring, S. Santini, A. Gupta, and R. Jain, "Content-Based Image Retrieval at the End of the Early Years," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 12, pp. 1349-1380, Dec. 2000.
- [3] B. Logan and A. Salomon, "A Music Similarity Function Based on Signal Analysis," *Proc. IEEE Int'l Conf. Multimedia and Expo*, pp. 745-748, 2001.
- [4] S. Sclaroff, M. Cascia, S. Sethi, and L. Taycher, "Unifying Textual and Visual Cues for Content-Based Image Retrieval on the World Wide Web," *J. Computer Vision and Image Understanding*, vol. 75, no. 1, pp. 86-98, 1999.
- [5] C. Frankel, M. Swain, and V. Athitsos, "Webseer: An Image Search Engine for the World Wide Web," technical report, Computer Science Dept., Univ. of Chicago, 1996.
- [6] W. Li, K. Candan, and K. Hirata, "SEMCOG: An Integration of SEMantics and COGNition-Based Approaches for Image Retrieval," *Proc. ACM Symp. Applied Computing*, pp. 136-143, 1997.
- [7] K. Barnard and D. Forsyth, "Learning the Semantics of Words and Pictures," *Proc. IEEE Int'l Conf. Computer Vision*, vol. 2, pp. 408-415, 2001.
- [8] L. Denoyer and P. Gallinari, "Bayesian Network Model for Semi-Structured Document Classification," *Information Processing and Management*, vol. 40, no. 5, pp. 807-827, 2004.
- [9] C. Snoek and M. Worring, "Multimodal Video Indexing: A Review of the State-of-the-Art," *J. Multimedia Tools and Applications*, vol. 25, no. 1, pp. 5-35, 2005.
- [10] R. Datta, D. Joshi, J. Li, and J. Wang, "Image Retrieval: Ideas, Influences, and Trends of the New Age," *ACM Computing Surveys*, vol. 40, no. 2, pp. 1-60, 2008.
- [11] J. Iria, F. Ciravegna, and J. Magalhães, "Web News Categorization Using a Cross-Media Document Graph," *Proc. ACM Int'l Conf. Image and Video Retrieval*, pp. 1-8, 2009.
- [12] A.F. Smeaton, P. Over, and W. Kraaij, "Evaluation Campaigns and TRECVID," *Proc. Eighth ACM Int'l Workshop Multimedia Information Retrieval*, pp. 321-330, 2006.
- [13] T. Tsirikia and J. Kludas, "Overview of the Wikipedia Multimedia Task at ImageCLEF 2008," *Evaluating Systems for Multilingual and Multimodal Information Access*, pp. 539-550, Springer, 2009.
- [14] K. Barnard, P. Duygulu, D. Forsyth, N. De Freitas, D. Blei, and M. Jordan, "Matching Words and Pictures," *J. Machine Learning Research*, vol. 3, pp. 1107-1135, 2003.
- [15] Y. Mori, H. Takahashi, and R. Oka, "Automatic Word Assignment to Images Based on Image Division and Vector Quantization," *Proc. Recherche d'Information Assistée par Ordinateur*, 2000.

- [16] G. Carneiro, A. Chan, P. Moreno, and N. Vasconcelos, "Supervised Learning of Semantic Classes for Image Annotation and Retrieval," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 29, no. 3, pp. 394-410, Mar. 2007.
- [17] G. Tzanetakis and P. Cook, "Musical Genre Classification of Audio Signals," *IEEE Trans. Speech and Audio Processing*, vol. 10, no. 5, pp. 293-302, July 2002.
- [18] D. Turnbull, L. Barrington, D. Torres, and G. Lanckriet, "Semantic Annotation and Retrieval of Music and Sound Effects," *IEEE Trans. Audio, Speech, and Language Processing*, vol. 16, no. 2, pp. 467-476, Feb. 2008.
- [19] D. Eck, P. Lamere, T. Bertin-Mahieux, and S. Green, "Automatic Generation of Social Tags for Music Recommendation," *Proc. Advances in Neural Information Processing Systems*, vol. 20, pp. 385-392, 2008.
- [20] D. Hardoon, S. Szedmak, and J. Shawe-Taylor, "Canonical Correlation Analysis: An Overview with Application to Learning Methods," *J. Neural Computation*, vol. 16, no. 12, pp. 2639-2664, 2004.
- [21] D. Blei, A. Ng, and M. Jordan, "Latent Dirichlet Allocation," *J. Machine Learning Research*, vol. 3, pp. 993-1022, 2003.
- [22] N. Rasiwasia, P. Moreno, and N. Vasconcelos, "Bridging the Gap: Query by Semantic Example," *IEEE Trans. Multimedia*, vol. 9, no. 5, pp. 923-938, Aug. 2007.
- [23] I. Khan, A. Saffari, and H. Bischof, "TVGraz: MultiModal Learning of Object Categories by Combining Textual and Visual Features," *Proc. 33rd Workshop Austrian Assoc. for Pattern Recognition*, 2009.
- [24] J. Shawe-Taylor and N. Cristianini, *Kernel Methods for Pattern Analysis*. Cambridge Univ. Press, 2004.
- [25] A. Vinokourov, D. Hardoon, and J. Shawe-Taylor, "Learning the Semantics of Multimedia Content with Application to Web Image Retrieval and Classification," *Proc. Fourth Int'l Symp. Independent Component Analysis and Blind Source Separation*, 2003.
- [26] N. Rasiwasia, J.C. Pereira, E. Coviello, G. Doyle, G. Lanckriet, R. Levy, and N. Vasconcelos, "A New Approach to Cross-Modal Multimedia Retrieval," *Proc. ACM Int'l Conf. Multimedia*, pp. 251-260, 2010.
- [27] M. Paramita, M. Sanderson, and P. Clough, "Diversity in Photo Retrieval: Overview of the ImageCLEF 2009 Photo Task," *Multilingual Information Access Evaluation: Multimedia Experiments*, pp. 45-59, Springer, 2010.
- [28] C. Meadow, B. Boyce, D. Kraft, and C. Barry, *Text Information Retrieval Systems*. Emerald Group, 2007.
- [29] G. Salton, *The SMART Retrieval System—Experiments in Automatic Document Processing*. Prentice Hall, 1971.
- [30] N. Vasconcelos, "Minimum Probability of Error Image Retrieval," *IEEE Trans. Signal Processing*, vol. 52, no. 8, pp. 2322-2336, Aug. 2004.
- [31] F. Monay and D. Gatica-Perez, "Modeling Semantic Aspects for Cross-Media Image Indexing," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 29, no. 10, pp. 1802-1817, Oct. 2007.
- [32] J. Jeon, V. Lavrenko, and R. Manmatha, "Automatic Image Annotation and Retrieval Using Cross-Media Relevance Models," *Proc. 26th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval*, pp. 119-126, 2003.
- [33] V. Lavrenko, R. Manmatha, and J. Jeon, "A Model for Learning the Semantics of Pictures," *Proc. Advances in Neural Information Processing Systems*, vol. 16, 2004.
- [34] S. Feng, R. Manmatha, and V. Lavrenko, "Multiple Bernoulli Relevance Models for Image and Video Annotation," *Proc. IEEE Conf. Computer Vision on Pattern Recognition*, vol. 2, pp. 1002-1009, 2004.
- [35] J.Z. Wang and J. Li, "Learning-Based Linguistic Indexing of Pictures with 2-D MHMMs," *Proc. ACM Int'l Conf. Multimedia*, pp. 436-445, 2002.
- [36] N. Vasconcelos, "From Pixels to Semantic Spaces: Advances in Content-Based Image Retrieval," *IEEE Trans. Computers*, vol. 40, no. 7, pp. 20-26, July 2007.
- [37] T. Westerveld, "Image Retrieval: Content versus Context," *Proc. Content-Based Multimedia Information Access at Recherche d'Information Assistée par Ordinateur*, pp. 276-284, 2000.
- [38] T. Pham, N. Maillot, J. Lim, and J. Chevallet, "Latent Semantic Fusion Model for Image Retrieval and Annotation," *Proc. ACM Int'l Conf. Information and Knowledge Management*, pp. 439-444, 2007.
- [39] H. Escalante, C. Hernández, L. Sucar, and M. Montes, "Late Fusion of Heterogeneous Methods for Multimedia Image Retrieval," *Proc. ACM Int'l Conf. Multimedia Information Retrieval*, pp. 172-179, 2008.
- [40] G. Wang, D. Hoiem, and D. Forsyth, "Building Text Features for Object Image Classification," *Proc. IEEE Conf. Computer Vision on Pattern Recognition*, pp. 1367-1374, 2009.
- [41] T. Klieger, K. Chandramouli, J. Nemrava, V. Svatek, and E. Izquierdo, "Combining Image Captions and Visual Analysis for Image Concept Classification," *Proc. Workshop Neural Networks for Signal Processing at ACM SIG Int'l Conf. Knowledge Discovery and Data Mining*, pp. 8-17, 2008.
- [42] S. Nakamura, "Statistical Multimodal Integration for Audio-Visual Speech Processing," *IEEE Trans. Neural Networks*, vol. 13, no. 4, pp. 854-866, July 2002.
- [43] J. Fisher III, T. Darrell, W. Freeman, and P. Viola, "Learning Joint Statistical Models for Audio-Visual Fusion and Segregation," *Proc. Advances in Neural Information Processing Systems*, pp. 772-778, 2001.
- [44] G. Qi, C. Aggarwal, and T. Huang, "Towards Semantic Knowledge Propagation from Text Corpus to Web Images," *Proc. ACM Int'l Conf. World Wide Web*, pp. 297-306, 2011.
- [45] D. Li, N. Dimitrova, M. Li, and I. Sethi, "Multimedia Content Processing through Cross-Modal Association," *Proc. ACM Int'l Conf. Multimedia*, pp. 604-611, 2003.
- [46] H. Zhang, Y. Zhuang, and F. Wu, "Cross-Modal Correlation Learning for Clustering on Image-Audio Dataset," *Proc. ACM Int'l Conf. Multimedia*, pp. 273-276, 2007.
- [47] M. Slaney, "Semantic-Audio Retrieval," *Proc. IEEE Int'l Conf. Acoustics Speech, and Signal Processing*, vol. 4, pp. 4108-4111, 2002.
- [48] Y. Yang, D. Xu, F. Nie, J. Luo, and Y. Zhuang, "Ranking with Local Regression and Global Alignment for Cross Media Retrieval," *Proc. ACM Int'l Conf. Multimedia*, pp. 175-184, 2009.
- [49] Y. Zhuang, Y. Yang, and F. Wu, "Mining Semantic Correlation of Heterogeneous Multimedia Data for Cross-Media Retrieval," *IEEE Trans. Multimedia*, vol. 10, no. 2, pp. 221-229, Feb. 2008.
- [50] Y. Zhuang, Y. Yang, F. Wu, and Y. Pan, "Manifold Learning Based Cross-Media Retrieval: A Solution to Media Object Complementary Nature," *J. VLSI Signal Processing Systems*, vol. 46, no. 2, pp. 153-164, 2007.
- [51] Y. Yang, Y. Zhuang, F. Wu, and Y. Pan, "Harmonizing Hierarchical Manifolds for Multimedia Document Semantics Understanding and Cross-Media Retrieval," *IEEE Trans. Multimedia*, vol. 10, no. 3, pp. 437-446, Apr. 2008.
- [52] V. Mahadevan, C.W. Wong, J.C. Pereira, T.T. Liu, N. Vasconcelos, and L.K. Saul, "Maximum Covariance Unfolding: Manifold Learning for Bimodal Data," *Proc. Advances in Neural Information Processing Systems*, vol. 24, pp. 918-926, 2011.
- [53] A. Vinokourov, J. Shawe-Taylor, and N. Cristianini, "Inferring a Semantic Representation of Text Via Cross-Language Correlation Analysis," *Proc. Advances in Neural Information Processing Systems*, vol. 15, pp. 1473-1480, 2003.
- [54] W. Hsu, T. Mei, and R. Yan, "Knowledge Discovery over Community-Sharing Media: From Signal to Intelligence," *Proc. IEEE Int'l Conf. Multimedia and Expo*, pp. 1448-1451, 2009.
- [55] T. Mei, W. Hsu, and J. Luo, "Knowledge Discovery from Community-Contributed Multimedia," *IEEE Trans. Multimedia*, vol. 17, no. 4, pp. 16-17, Oct.-Dec. 2010.
- [56] S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and R. Harshman, "Indexing by Latent Semantic Analysis," *J. Am. Soc. for Information Science*, vol. 41, no. 6, pp. 391-407, 1990.
- [57] I. Jolliffe, *Principal Component Analysis*. John Wiley & Sons, 2005.
- [58] H. Hotelling, "Relations between Two Sets of Variates," *Biometrika*, vol. 28, pp. 321-377, 1936.
- [59] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin, "LIBLINEAR: A Library for Large Linear Classification," *J. Machine Learning Research*, vol. 9, pp. 1871-1874, 2008.
- [60] C.-C. Chang and C.-J. Lin, "LIBSVM: A Library for Support Vector Machines," *ACM Trans. Intelligent Systems and Technology*, vol. 2, pp. 27:1-27:27, 2011.
- [61] M.J. Saberian and N. Vasconcelos, "Multiclass Boosting: Theory and Algorithms," *Proc. Advances in Neural Information Processing Systems*, vol. 24, pp. 2124-2132, 2011.
- [62] G. Griffin, A. Holub, and P. Perona, "The Caltech-256," technical report, Caltech, 2006.
- [63] C. Manning, P. Raghavan, and H. Schütze, *An Introduction to Information Retrieval*. Cambridge Univ. Press, 2008.

- [64] G. Doyle and C. Elkan, "Accounting for Word Burstiness in Topic Models," *Proc. ACM Int'l Conf. Machine Learning*, pp. 281-288, 2009.
- [65] M. Swain and D. Ballard, "Color Indexing," *Int'l J. Computer Vision*, vol. 7, no. 1, pp. 11-32, 1991.
- [66] S. Boughorbel, J. Tarel, and N. Boujemaa, "Generalized Histogram Intersection Kernel for Image Recognition," *Proc. IEEE Int'l Conf. Image Processing*, vol. 3, pp. 161-164, 2005.
- [67] C. Rashtchian, P. Young, M. Hodosh, and J. Hockenmaier, "Collecting Image Annotations Using Amazon's Mechanical Turk," *Proc. NAACL HLT Workshop Creating Speech and Language Data with Amazon's Mechanical Turk*, pp. 139-147, 2010.
- [68] J.C. Pereira and N. Vasconcelos, "On the Regularization of Image Semantics by Modal Expansion," *Proc. IEEE Conf. Computer Vision on Pattern Recognition*, pp. 3093-3099, 2012.



Jose Costa Pereira received the licenciatura in computer science and engineering in 2000 and the MS degree in computational methods in 2003 from the Universidade do Porto, Portugal. He is currently working toward the PhD degree in the Statistical Visual Computing Laboratory, Department of Electrical and Computer Engineering, University of California, San Diego. He received the PhD fellowship from the Ministry of Sciences and Education, Portugal, for 2008-2012. His research interests include computer vision, multimedia, and machine learning. He is a student member of the IEEE.



Emanuele Coviello received the laurea triennale in information engineering in 2006, and the laurea specialistica in telecommunications engineering in 2008 from the University of Padova, Italy. He is working toward the PhD degree in electrical and computer engineering in the Computer Audition Laboratory, University of California, San Diego. His main research interests include machine learning applied to (music) information retrieval and multimedia data modeling.



Gabriel Doyle received the AB degree in mathematics from Princeton University in 2005 and the MA degree in linguistics from the University of California, San Diego, in 2011, where he is currently working toward the PhD degree in linguistics. His primary research area is computational psycholinguistics, especially models of language acquisition that integrate multiple sources of information.



Nikhil Rasiwasia received the BTech degree in electrical engineering from the Indian Institute of Technology Kanpur, India, in 2005, and the MS and PhD degrees from the University of California, San Diego, in 2007 and 2011, respectively. He is currently a scientist at Yahoo! Labs Bangalore. His research interests include areas of computer vision and machine learning. He was recognized as an "Emerging Leader in Multimedia" in 2008 by the IBM T.J. Watson Research Center. He also received the Best Student Paper Award at ACM Multimedia 2010. He is a member of the IEEE.



Gert R.G. Lanckriet received the MS degree in electrical engineering from the Katholieke Universiteit Leuven, Belgium, in 2000, and the MS and PhD degrees in electrical engineering and computer science from the University of California, Berkeley, in 2001 and 2005, respectively. In 2005, he joined the Department of Electrical and Computer Engineering, University of California, San Diego, where he is the head of the Computer Audition Laboratory. His research interest focuses on the interplay of convex optimization, machine learning, and signal processing, with applications in computer audition, music information retrieval, and personalized health. He was awarded the SIAM Optimization Prize in 2008 and has received a Hellman fellowship, an IBM Faculty Award, a US National Science Foundation CAREER Award, and an Alfred P. Sloan Foundation Research fellowship. In 2011, *MIT Technology Review* named him one of the 35 top young technology innovators in the world (TR35). He is a senior member of the IEEE.



Roger Levy received the BS degree in mathematics from the University of Arizona in 1996, the MS degree in anthropological sciences from Stanford University in 2002, and the PhD degree in linguistics from Stanford University in 2005. He is currently an associate professor in the Department of Linguistics, University of California, San Diego, where he is the head of the Computational Psycholinguistics Laboratory. In 2005-2006, he was a postdoctoral fellow in Informatics at the University of Edinburgh and joined the University of California, San Diego, in 2006. His research focuses on theoretical and applied questions in the processing of natural language. He received a Hellman fellowship, a US National Science Foundation Career Award, and an Alfred P. Sloan fellowship. He is currently an associate editor for the *Journal of Cognitive Science*. In 2013-2014, he will be a fellow at the Center for Advanced Study in the Behavioral Sciences at Stanford University.



Nuno Vasconcelos received the licenciatura in electrical engineering and computer science from the Universidade do Porto, Portugal, in 1988, and the MS and PhD degrees from the Massachusetts Institute of Technology in 1993 and 2000, respectively. From 2000 to 2002, he was a member of the research staff at the Compaq Cambridge Research Laboratory, which in 2002 became the HP Cambridge Research Laboratory. In 2003, he joined the Department of Electrical and Computer Engineering at the University of California, San Diego, where he is the head of the Statistical Visual Computing Laboratory. He is the recipient of a US National Science Foundation CAREER Award, a Hellman Fellowship, and has authored more than 150 peer-reviewed publications. His work spans various areas, including computer vision, machine learning, signal processing and compression, and multimedia systems. He is a senior member of the IEEE.

► For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.