

Multimodal Preference Aggregation for Multimedia Information Retrieval

Eric Bruno, Stéphane Marchand-Maillet

Viper group - Computer Vision and Multimedia Lab - University of Geneva, Switzerland

Email: {name.surname}@unige.ch

Abstract—Representing and fusing multimedia information is a key issue to discover semantics in multimedia. In this paper we address more specifically the problem of multimedia content retrieval through the joint design of an original multimodal information representation and of a machine learning-based fusion algorithm. We first define a novel preference-based representation particularly adapted to the retrieval problem, and then, we investigate the RankBoost algorithm to combine those preferences to fulfill a user's query. Interestingly, it ends up being a flexible retrieval model that only manipulates ranking information and is blind to the intrinsic properties of the multimodal information input. The approach is tested on annotated images and on the complete TRECVID 2005 corpus and compared with SVM-based fusion strategies. The results show that our approach equals SVM performance but, contrary to SVM, is parameter free and faster.

I. INTRODUCTION

Determining semantic concepts by allowing users to iteratively and interactively refine their queries is a key issue in multimedia content-based retrieval. The Relevance Feedback loop allows us to build complex queries made out of documents marked as positive and negative examples. From this training set, a learning process has to create a model of the sought concept from a set of data features to finally provide relevant documents to the user. The success of this search strategy relies mainly on the representation spaces where data is embedded as well as on the learning algorithm operating in those spaces. These two issues are also intrinsically related to the problem of adequately fusing information arising from different sources. Various aspects of these problems have been studied with success for the last few years. This includes works on machine learning strategies such as active learning [1], imbalance classification algorithms [2], automatic kernel setting [3] or automatic labelling of training data [4]. Theoretical and experimental investigations have been achieved to determine optimal strategies for multimodal fusion: Kittler *et al* and R. Duin studied different rules for classifier combination [5], [6]; Wu *et al* propose the super-kernel fusion to determine

optimal combination of features for video retrieval [7]. In [8], Maximum Entropy, Boosting and SVM algorithms are compared to fuse audio-visual features. Multi-graph learning approaches [9] and latent semantic fusion [10] have been proposed recently for image and video retrieval and annotation. A number of further relevant references may be found into the Lecture Notes series on Multiple Classifier Systems [11].

The diversity of the features involved is a difficulty when dealing with fusion and learning. The multimedia descriptors may indeed be extracted from visual, audio or transcript streams using various operators providing outputs such as histograms, filter responses, statistical measures or symbolic labels. This heterogeneity imposes building complex learning setup that need to take into account all the variety of the features' mathematical and semantic properties [12][13].

We advocate for the definition of an *homogeneous representation* to store multimodal signals regardless their intrinsic dimensionality and scale. The fusion complexity would be then dramatically alleviated since a unique learning model can be indistinctly applied on multimodal information to determine document's semantic/relevance. It would allow to setup *fast* and *flexible* multimedia information retrieval systems. In our context, *fast* means that on-line learning is possible and *flexible* means that the system could handle any modalities blindly as far as they can be embedded into the homogeneous representation.

A first attempt for designing an homogeneous representation is to index documents according to their similarities (related to one or several features) to the other documents rather than to a feature vector. Considering a collection of documents, the similarity-based representation, stored in (dis)similarity matrices or some distance-based indexing structures [14], characterizes the content of an element of the collection relatively to a part of or the whole collection. Studies have been published for document retrieval and collection browsing by using pre-computed similarities. In [15], Boldareva *et al* proposed to index elements relatively to their closest neighbors, *i.e.* those who have the best probabilities to belong to the same class. This provides them with a sparse association graph structuring the multimedia collection and allowing fast retrieval of data. In [16], the idea of nearest neighbor networks is extended by creating edges for every combination of features. The resulting graph, called NN^k , allows to browse the data collection from various viewpoints

This work is partially supported by the Swiss NCCR (IM)2 (Interactive Multimodal Information Management) and the European Community under the Information Society Technologies (IST) programme of the 6th FP for RTD - project MultiMATCH contract IST-2005-2.5.10. The author is solely responsible for the content of this paper. It does not represent the opinion of the European Community, and the European Community is not responsible for any use that might be made of data appearing therein the IST European project Multimatch.

corresponding to the multiple features. In [17], we design a dissimilarity space [18] where elements are no longer represented by their multimodal features but by their relative dissimilarities with respect to a set of positive prototypes provided by users. As pointed out by authors, the similarity approach provides a convenient way for multimodal data fusion, since adding new features simply consists in adding new distances to the same representation framework. However, it still leaves open the problem of how to properly scale the similarity values to make them really comparable.

Following the similarity-based representation idea, we propose to simplify the similarity-based representation by retaining only ordering information from the distance measurements. The result is a *preference space* where every item is indexed through its relative ranking positions to a set of prototypes. The scaling issue is thus completely alleviated and we effectively obtain a unified representation of multimodal content, but at the price of losing an important amount of the initial information. In the following we address more specifically the problem of multimedia information retrieval using *query by example* and *relevance feedback* search paradigms. Problem position and terminologies are defined in section II). In section III we consider three multimodal information representations, namely *feature space*, *dissimilarity space* and the proposed *preference space*. Retrieving items from the preference space is then a ranking problem (section IV) that can be addressed using the RankBoost algorithm (section V). We end up with a multimedia search engine that a) builds its retrieval model upon multimodal information, b) is parameter free and c) is fast. Experiments on artificial and real data (annotated images and videos, see section VI) show that the preference space associated to RankBoost competes with SVM-based approaches in term of accuracy but speed up the retrieval by a factor greater than 10. Moreover, contrary to the SVM, our approach does not require to set query-sensitive parameters *a priori*. It is therefore a valid approach for on-line retrieval of multimedia information.

II. PROBLEM DEFINITION

A *multimedia document* is composed of multimodal contents (for instance visual, audio and textual content) and *multimedia information retrieval* will consist to determine the relevance of each document relatively to a given query. This relevance will reflect the adequation of the multimodal content to the query.

In the following, we consider a collection \mathcal{X} containing l multimedia documents x . The terms item, element or object are also used to refer to x . The *query by example* search paradigm consists in gathering user's judgements indicating, for some objects, whether they are relevant or irrelevant to the user request. This set, denoted \mathcal{Q} , is called the *query* and is composed of positive and negative subsets, respectively

$$\mathcal{P} = \{x_i^+\}_{i=1}^p \text{ and } \mathcal{N} = \{x_i^-\}_{i=1}^n.$$

The query \mathcal{Q} is then used to train a machine that will produce a decision function ranking documents according to their relevance to the query.

This paradigm might be embedded in the *Relevance Feedback* (RF) strategy, where these two steps (user judgement and ranking estimation) are iterated until the search converges to a satisfactory result.

III. MULTIMODAL CONTENT REPRESENTATIONS

Expressing multimodal content involves first to extract various descriptors from the multimedia objects. Ideally, each descriptor depicts an appropriate aspect of the multimodal features of the documents. Assuming such descriptors are available, we discuss in the following how efficient representations may be derived to store descriptors and to facilitate their fusion.

A. Feature-based representation

Assuming m distinct descriptors are designed (and extraction procedures implemented), the multimodal representation of an object x is the set of m feature vectors $\{\mathbf{x}^k\}_{k=1}^m$ living respectively in feature spaces $\{\mathcal{F}^k\}_{k=1}^m$. The dimension of each feature space intrinsically depends of the descriptor they express. The feature-based representation is rather straightforward, but not really convenient since it mixes heterogeneous vectors of various dimensions and scales. Fusion and ranking algorithms need to manage the diversity of the representation, thus making them more dependent on complex parameter setting procedures and less flexible to handle new descriptors.

To avoid this situation, modality-independent representations are desirable. For that purpose, (dis)similarity-based representations have been recently proposed [15], [17], [19], [20]. As pointed out by these authors, similarities are convenient to manipulate multimodal information since they form a homogeneous representation of the content. Moreover, similarity representations are generally made such as their dimensionality remain much lower than their feature counterparts.

B. Dissimilarity-based representation

In [17], we proposed a *Query-based Dissimilarity Space* (QDS), derived from the dissimilarity spaces introduced by Pekalska *et al* [21]. For a given feature space \mathcal{F}^k , the corresponding QDS, denoted $\mathcal{D}_{\mathcal{P}}^k$, is defined relatively to the positive set \mathcal{P} by the mapping $\mathbf{d}^k(x, \mathcal{P}) \in \mathbb{R}^p$

$$\mathbf{d}^k(x, \mathcal{P}) = [d^k(x, x_1^+), d^k(x, x_2^+), \dots, d^k(x, x_p^+)]^T, \quad (1)$$

where $d^k(x, x_i^+) \in \mathbb{R}^+$ is the dissimilarity from any object $x \in \mathcal{X}$ to the prototype x_i^+ when the measure is done in \mathcal{F}^k . Using QDS, an object x is thus represented with a set of m dissimilarity vectors $\{\mathbf{d}^k\}_{k=1}^m$ living in p -dimensional dissimilarity spaces $\{\mathcal{D}_{\mathcal{P}}^k\}_{k=1}^m$,

$$\mathcal{D}_{\mathcal{P}}^k = \begin{pmatrix} d^k(x_1, x_1^+) & d^k(x_2, x_1^+) & \dots & d^k(x_l, x_1^+) \\ d^k(x_1, x_2^+) & d^k(x_2, x_2^+) & \dots & d^k(x_l, x_2^+) \\ \vdots & \vdots & \ddots & \vdots \\ d^k(x_1, x_p^+) & d^k(x_2, x_p^+) & \dots & d^k(x_l, x_p^+) \end{pmatrix}. \quad (2)$$

The QDS presents two decisive advantages relatively to feature spaces: 1) It provides a unified representation of multimodal information channels, and 2) is particularly adapted to the class asymmetry typically exhibited by the positive and negative classes. This asymmetry corresponds to a $(1+x)$ class setup where the one class, presumably well-clustered in the feature space, encompasses the sought documents (positive class), while an unknown number x of classes, partially represented by negative examples, is supposed to model all irrelevant documents. Classical learning approaches, by applying a symmetric treatment to all classes are not really efficient for such a setup. Learning the negative classes, while being feasible using traditional non-linear learning machines, becomes challenging when only few samples are available. Nevertheless, we show in [17] how a built-in property of $\mathcal{D}_{\mathcal{P}}$ is to transform the asymmetric classification setup such that it becomes linearly separable.

However, the issue of how properly scaling dissimilarity spaces so that modalities become easily comparable still remains. This problem might be left out to the fusion and ranking algorithms [17], but a more elegant solution would be to end up with a fully homogeneous multimodal representation.

C. Preference-based representation

We propose to simplify the QDS representation by replacing the dissimilarity components $d^k(x, x_i^+)$ with the ranking position $\pi^k(x, x_i^+) \in \mathbb{N}$ of an object x with respect to the prototype x_i^+ according to the dissimilarity measure d^k and the collection \mathcal{X} ,

$$\pi^k(x, x_i^+) = \sum_{x_j \in \mathcal{X}} \llbracket d^k(x_j, x_i^+) \leq d^k(x, x_i^+) \rrbracket. \quad (3)$$

The notation $\llbracket \kappa \rrbracket$ is defined to be 1 if predicate κ holds and 0 otherwise. Considering the p positive prototypes and the m dissimilarity measures, the multimodal representation of an object x may be represented as a *unique* $(p * m)$ -dimensional vector of *preferences*

$$\begin{aligned} \pi(x) = & [\pi^1(x, x_1^+), \pi_1^2(x, x_1^+) \dots, \pi^m(x, x_1^+), \\ & \pi^1(x, x_2^+), \pi_1^2(x, x_2^+) \dots, \pi^m(x, x_2^+), \\ & \vdots \\ & \pi^1(x, x_p^+), \pi_1^2(x, x_p^+), \dots, \pi^m(x, x_p^+)]^T. \end{aligned}$$

For the sake of readability, the notation $\pi^k(\cdot, x_i^+)$ is simplified to $\pi_j(\cdot)$, $j = k + m * (i - 1)$, $j \in [1, p * m]$, with i iterating over all objects $x_i^+ \in \mathcal{P}$ and k over the m modalities. The multimodal *preference space* embedding all objects $x \in \mathcal{X}$ is therefore

$$\Pi_{\mathcal{P}} = \begin{pmatrix} \pi_1(x_1) & \pi_1(x_2) & \dots & \pi_1(x_l) \\ \pi_2(x_1) & \pi_2(x_2) & \dots & \pi_2(x_l) \\ \vdots & \vdots & \ddots & \vdots \\ \pi_{pm}(x_1) & \pi_{pm}(x_2) & \dots & \pi_{pm}(x_l) \end{pmatrix}. \quad (4)$$

It consists in a unique pm -dimensional natural number space providing a fully homogeneous representation of

multimodal information. Reading $\Pi_{\mathcal{P}}$ column-wise gives the preference vectors $\pi(x)$ of every object $x \in \mathcal{X}$, while reading row-wise yields the complete ordering of \mathcal{X} relatively to a given positive example x_i^+ and a given modality k (i and j are given by the relation $j = k + m * (i - 1)$). Similarly to the QDS approach, $\Pi_{\mathcal{P}}$ represents the two classes \mathcal{P} and \mathcal{N} asymmetrically since every element is evaluated relatively to the positive instances only.

It is worth noting however that we obtain this modality-independent representation at the price of losing most information about the initial feature distributions; only ordering information is actually preserved. Our objective now is to define a machine learning effectively able to learn from preferences as efficiently as learning directly in feature spaces or in dissimilarity spaces.

IV. THE RANKING PROBLEM

The ranking problem could be formulated as follows: For each item $x \in \mathcal{X}$, it exists ranking features π_1, \dots, π_{pm} , where each π_j defines a linear ordering of the instances $x \in \mathcal{X}$. In our formulation, $\pi_j \in \mathbb{N}$ and $\pi_j(x_1) < \pi_j(x_0)$ means x_1 preferred to x_0 .

Additionally to the ranking features, there exists a *feedback* function $\Phi : \mathcal{X} \times \mathcal{X}$ which provides to the learner the desired form of the final ranking. Formally $\Phi(x_1, x_0) > 0$ means that x_1 should be ranked above x_0 while $\Phi(x_1, x_0) < 0$ means the opposite. $\Phi(x_1, x_0) = 0$ means no preferences between x_0 and x_1 and the magnitude of $|\Phi(x_1, x_0)|$ indicates how important is to rank x_1 above or below x_0 . The *bipartite* feedback function is special but common case in document retrieval: the function is said bipartite if there exists two disjoint set \mathcal{X}_1 and \mathcal{X}_0 such that Φ ranks all instances x_1 of \mathcal{X}_1 above instances x_0 of \mathcal{X}_0 . These subsets are respectively the positive and negative subsets \mathcal{P} and \mathcal{N} we defined in section II.

Learning such a feedback function implies estimating a ranking $H : \mathcal{X} \rightarrow \mathbb{R}$ through the optimization of a ranking loss function penalizing every miss-ordered pair of items. We consider the loss proposed in [22]

$$\sum_{\substack{x^+ \in \mathcal{N} \\ x^- \in \mathcal{P}}} \Phi(x^+, x^-) [H(x^+) - H(x^-)]. \quad (5)$$

The function $H(x)$ is a ranking of items x stating that x^+ is ranked higher than x^- whenever $H(x^+) > H(x^-)$. Interestingly, in case of bipartite feedback, the problem becomes separable and the ranking loss simplifies to [22]

$$\sum_{x \in \mathcal{Q}} w(x) s(x) H(x), \quad (6)$$

where the user feedback is carried by both

$$s(x) = \begin{cases} +1 & \text{if } x \in \mathcal{P} \\ -1 & \text{if } x \in \mathcal{N} \end{cases}, \quad (7)$$

and $w(x)$ a weight giving the importance of the rank of the item x .

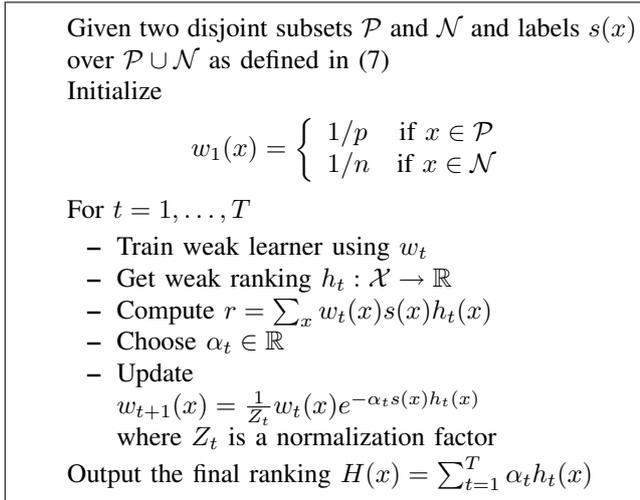


Fig. 1. The RankBoost algorithm for bipartite feedback

V. RANKBOOST

Following the boosting principle, the final ranking H results from a weighted sum of weak rankings $h_t : \mathcal{X} \rightarrow \mathbb{R}$

$$H(x) = \sum_{t=1}^T \alpha_t h_t(x), \quad (8)$$

which is estimated through an Adaboost-like algorithm, namely RankBoost [22] (see Figure 1). This greedy coordinate-wise search algorithm aims at iteratively minimizing the normalization factor Z_t by choosing at each round an appropriate pair $\{\alpha_t, h_t\}$. For a given weak hypothesis $h_t \in [-1, 1]$, it has been shown [23] that Z_t is minimized for

$$\alpha_t = \frac{1}{2} \ln \frac{1 + r_t}{1 - r_t}, \quad (9)$$

where r is the weighted classification rate

$$r_t = \sum_{x \in \mathcal{Q}} w_t(x) s(x) h_t(x). \quad (10)$$

The algorithm is run over a number T of iterations which is predefined or may depend on the training error. In our implementation, the loop is stopped whenever the training error is equal to 0, with a maximum of $2pm$ iterations.

A. Weak ranking

The weak ranking h_t is produced through a *weak learner*. It has to provide a new ranking from ranking features π_i conforming the best to the bipartite feedback. For example, the weak learner proposed in [22] selects at each iteration the ranking feature π_i minimizing the training error. The output preserves only relative-ordering information so as to be independent of specific preference values,

$$h(x) = \begin{cases} +1 & \text{if } \pi_i(x) < \theta \\ -1 & \text{if } \pi_i(x) \geq \theta \end{cases}. \quad (11)$$

As illustrated in Figure 2, this weak learner consists in fitting a step function to the user feedbacks $\{s(x_j)\}_{j=1}^q$

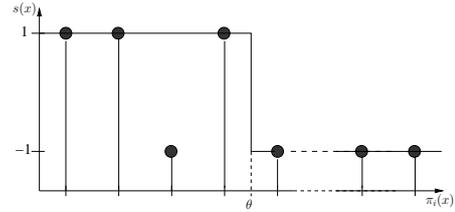


Fig. 2. Binary weak ranking. The $\pi_i(x)$'s are ordered in increasing order.

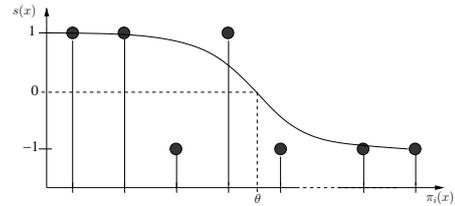


Fig. 3. Soft weak ranking. The $\pi_i(x)$'s are ordered increasing order.

sorted by increasing order of $\pi_i(x_j)$. The best weak ranking is the one maximizing equation (10) over the q candidate weak rankings for the pm preferences π_i . The evaluation of all candidates is done in $O(qpm)$.

As defined in (11), the function $h(x)$ provides at each iteration a binary ranking. The final ranking H (eq. (8)) is thus an injection on \mathcal{X} whose image has at most cardinality 2^T , ie $H(x) : \mathcal{X} \rightarrow \{v^1, \dots, v^{2^T}\}$. Typically when the training set is small or when the ranking problem is simple, RankBoost converges in a few iterations (T small) and consequently provides a coarse ranking partitioning the collection \mathcal{X} in few blocks. To get a finer ranking, we propose to use the a soft ranking function,

$$h(x) = 2e^{-\gamma \pi_i^2(x)} - 1. \quad (12)$$

Learning this weak ranking consists of choosing the pair (π_i, γ) that maximize the classification rate r_t (10). Given a ranking feature π_i , a grid search on γ is achieved rather than a time-consuming non-linear regression. The grid vertices are positioned at the middle of the q ranking intervals (see Figure 3). With this approximation, the weak learner complexity remains $O(qpm)$.

VI. EXPERIMENTS

The behavior and performance of ranking data in the three representation spaces (feature, dissimilarity and preference) are studied here. As stated before, RankBoost (soft and binary weak ranking) will be used to learn preferences. As far as feature and dissimilarity spaces are concerned, ranking are produced with the SVM algorithm as it is considered as an effective and standart technique for multimedia retrieval [24], [25], [26]. Depending of experiment, linear or non-linear (eg using RBF kernel) SVM is used.

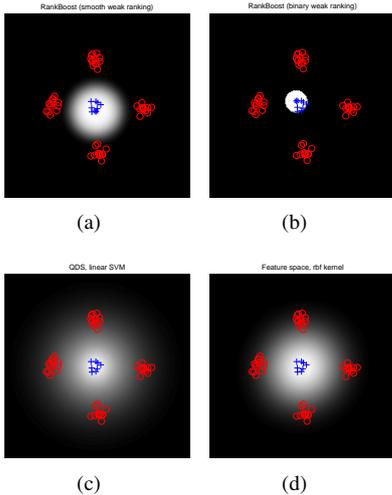


Fig. 4. Cross toy example

A. Toy examples

Artificial data allows us to concretely illustrate how rankings are learned in the various representation spaces. The following toy examples are made so as to be representative of the class asymmetry we generally meet in real applications. For every learning technique, the learned ranking is superimposed to the items; white areas correspond to top ranks and black areas to the last rank. Moreover, prototypes selected by Rankboost are indicated with the * marker.

The first example (Figure 4) corresponds to an ideal separable case, where all the positive instances (cross marker) belong to the same cluster, while the negative samples are distributed around (circle marker). The corresponding dissimilarity space is built using pairwise Euclidean distances while the preference space is derived by ordering dissimilarities. Linear SVM is used to learn in QDS while a RBF-SVM with an appropriate scale parameter operates in feature space.

In each case (preferences, dissimilarities and features, respectively in Figure 4.a, b, c and d), a perfect ranking has been estimated. As the class setup is simple, only one weak ranking (indicated by the selected prototype) is necessary for RankBoost (Figure 4.a and b). It implies that the final ranking is binary when using the binary weak ranking function, while learning with the soft weak ranking provides us with a more convenient continuous ranking. As mention in section III-B, a linear function is also enough to catch the positive class within dissimilarity space, while a non-linear RBF-based ranking function is needed in feature space.

The second example (Figure 5) depicts a less obvious problem, the XOR configuration. The classes are no longer linearly separable neither in preference space nor in dissimilarity space. In that case, the linear SVM used in QDS failed in estimating the ranking. On contrary, RankBoost succeeds in finding the two positive clusters and selects one prototype per cluster.

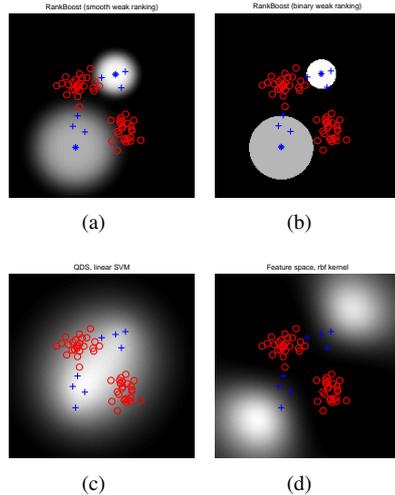


Fig. 5. XOR toy example

B. Real data

1) *Corel image collection*: The studied image collection is a subset of the Corel collection. It contains 1159 images annotated with 1 to 10 keywords per image (including some non-sense descriptions). The images are categorized into 49 classes. Textual and visual features are considered for fusing experiments: The vector space model $\mathcal{F}^{\text{text}}$ containing tf-idf weights is built from keywords (2035 terms). The color space $\mathcal{F}^{\text{color}}$ contains 166 bins HSV histograms and the texture space $\mathcal{F}^{\text{texture}}$ is made of Gabor filter bank outputs (120 dimensions). Cosine distance is considered for textual features, while Euclidean is used in visual feature spaces.

Fusion is operated in feature space, dissimilarity space and preference space. In feature and dissimilarity space we have considered a state-of-the-art hierarchical fusion scheme [17], [7]. At the first level, base classifiers are trained in each monomodal space. At the second level, a super classifier is used to fuse soft-outputs of all base classifiers. Base classifiers and super classifier are RBF SVM. Optimal classifier parameters have been determined through a leave-one-out cross validation.

Retrieval performance is given in terms of Mean Average Precision (MAP). Average Precision (AP) is the sum of the precision at each relevant hit in the retrieved list, divided by the minimum between the number of relevant documents in the collection and the length of the list. The MAP is simply the AP averaged over several classes. Additionally to the algorithm performance, a baseline consisting in retrieving randomly documents is always provided. All results are displayed in Figure 6.

Multimodal retrieval (Figure 6.b) and text-only search (Figure 6.a) are studied. In both cases we observe that for RankBoost, soft ranking outperforms largely binary ranking. Moreover, the soft ranking performs similarly to the SVM approaches whereas it uses only a degraded version of the original features. The second observation we can make is that the multimodal retrieval outperforms only very slightly the keyword-only search, whatever

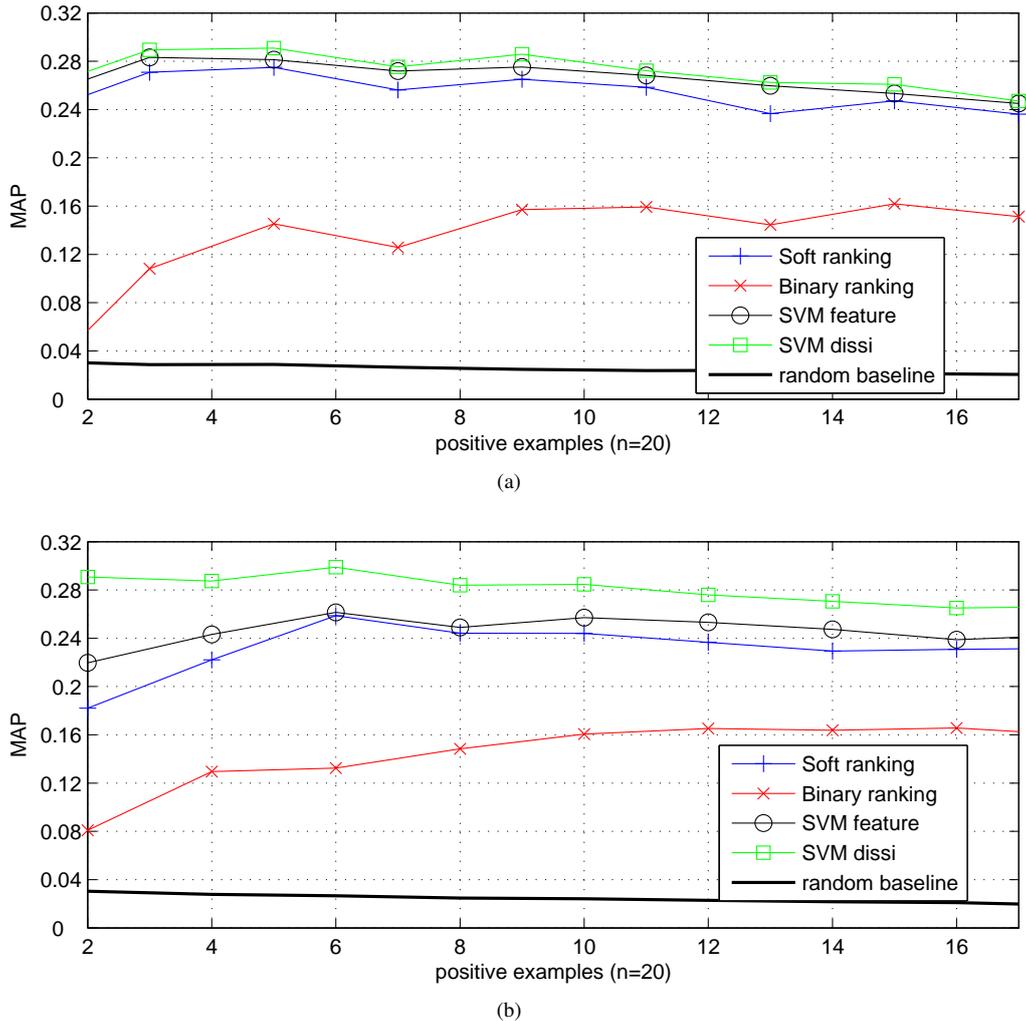


Fig. 6. Image retrieval results with a) multimodal fusion and b) keywords-only search

the approach considered. This result seems to indicate that keywords bring much of the category information, and that color or texture low-level information are of little help in that case. This observation is confirmed by analyzing the ranking features selected by RankBoost to build retrieval models: among all retrieval instances, text information is used for 93% of them, while color and texture features are only used for respectively 36% and 17% of the cases.

Concerning the computational time on this particular example (Table I), we observe that soft ranking is slightly faster than binary ranking. It is also interesting to note that RankBoost is around 20 times faster than the hierarchical SVM approaches.

2) *TRECVID video corpus*: We now consider the TRECVID 2005 benchmark. In our setup, videos are segmented into around 89'500 segments using the common shot reference [27]. These shots are considered as individual and independent documents. This means that no contextual information is taken into account and that shot description is restricted to its audiovisual content (*eg*

TABLE I
COMPUTATIONAL TIME
(IN SECOND, INTEL XEON 2.80GHZ)

p+n	SVM in \mathcal{F}	SVM in $\mathcal{D}_{\mathcal{P}}$	binary rkg	soft rkg
20	0.46	0.36	0.009	0.01
30	0.80	0.68	0.028	0.024
60	2.95	2.75	0.17	0.12
100	9.43	9.23	0.56	0.52

visual, audio and speech¹ information).

The Search Task, as defined in TRECVID-05, consists in retrieving shots that are relevant to some predefined queries (called topics). There are 24 topics concerning people (person-X queries), objects (specific or generic), locations, sports and combinations of the former. For each topic, keywords, pictures and several video shots (4-10) are provided as positive examples. Further details about the Search Task may be found in [28]. During the experiments, we only considered video shots as positive examples. The positive examples are completed with ten

¹the speech transcripts extracted by Automatic Speech Recognition (ASR) are also available.

negative examples randomly selected within the test set. Starting with this initial query, a *relevance feedback* loop is initiated by adding to the query up to 10 new positive and negative examples returned in the 1000-entries hit-list. The process is repeated ten times. Following the TRECVID evaluation protocol, the performance was measured at each iteration by MAP at 1000. Additionally to the algorithm performance, a baseline consisting of retrieving randomly documents is always provided.

The multimodal features are derived from the six following text and audiovisual descriptors:

- Color histogram, $4 \times 4 \times 4$ bins in YCbCr space
- Motion vector histogram, 66 bins quantization of the MPEG block motion vectors [29]
- Local features, SIFT descriptors extracted around the Lowe salient points [30],
- Face detection [31],
- Word occurrence histogram (vector space model) computed from ASR,
- Dominant audio features [32] extracted from the audio stream.

The distance measures used are Euclidean for color and motion histograms. An approximation of the minimal matching distance is applied on local features to determine partial similarities [33]. Euclidean distance in the 30-dimensional eigenface space gives the similarity between the detected faces. Cosine distance is used for the vector space model and finally the audio similarity measure proposed in [32] is used for audio features.

The fusion strategies remain the hierarchical RBF SVM approach in feature spaces and dissimilarity spaces. For feature space however, we adapt the RBF-kernel to the distances used, $k_d(x, y) = e^{-\frac{d(x, y)}{2\sigma^2}}$ (it is worth noting that k_d is strictly a RBF-kernel when d is an Euclidean distance). Optimal classifier parameters have been cross-validated using the TRECVID development set.

MAP results are given in Figure 7.a. We compare multimodal retrieval techniques with the best monomodal search (hierarchical SVM in \mathcal{D}_P^{ASR}). The overall RankBoost performance remains very close to the best retrieval result provided by the hierarchical SVM in dissimilarity space. Soft ranking and binary ranking have now similar performance and the latter is even slightly better when the training set becomes large. However in that case, the soft ranking is around three times faster than the binary ranking and ten times faster the SVM learning² (Figure 7.b). This rapidity is explained by the fact that soft ranking systematically selects less features than binary ranking to produce the final ranking (Figure 7.c) and thus converges faster and provides simpler retrieval models. In all cases, the retrieval accuracy benefits from multimodal fusion and largely outperforms the ASR-only search. On the contrary to the Corel experiment, we observe now that the retrieval models produced by RankBoost (soft ranking) are fully multimodal. As shown in Figure 8, the modality usage,

²The computational complexity to learn the SVM in \mathcal{F} and in \mathcal{D}_P is equivalent. Only the computational time for SVM in \mathcal{D}_P is thus reported in Figure 7.b

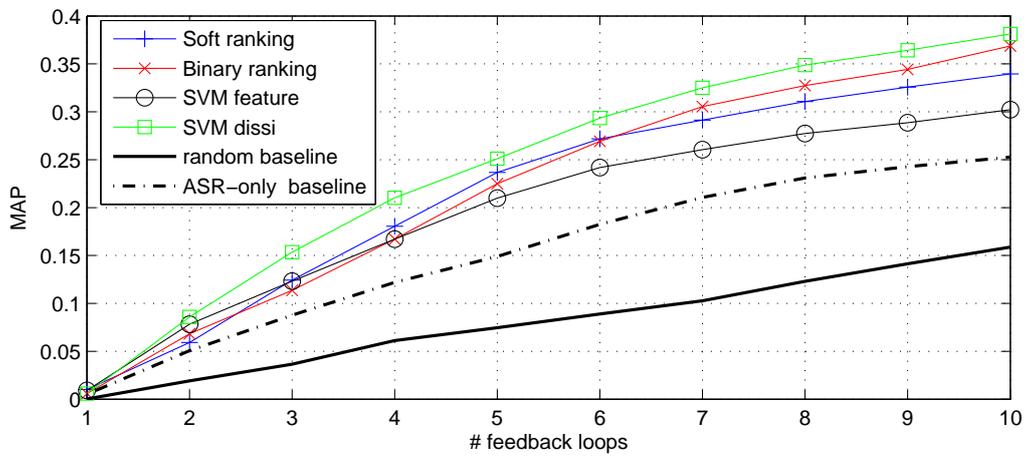
ie the frequency of the selection of each descriptor to build the final ranking over the 24 queries, is almost 100% for every modality. This indicate that all multimodal information sources are needed to fulfil the semantic level required by the TRECVID queries.

VII. CONCLUSION

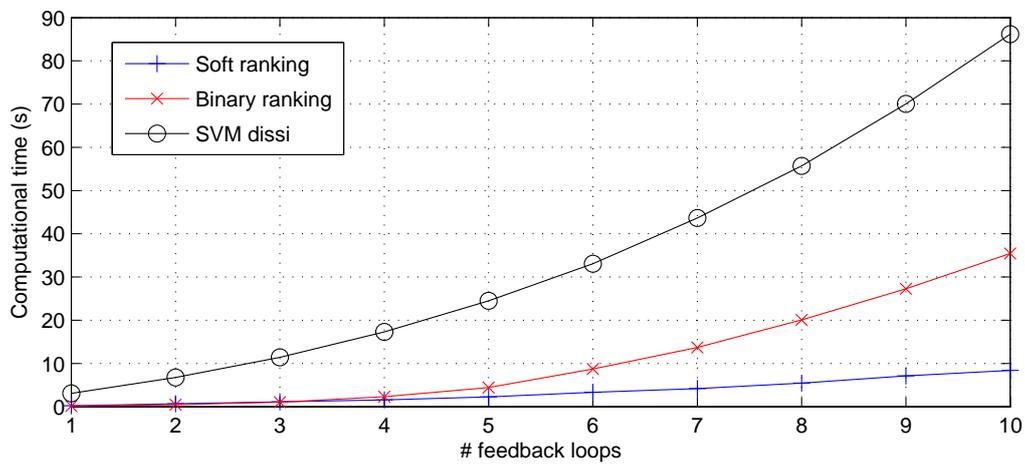
The preference space we introduced in this paper is a degraded but lightweight representation of the original feature space where all information relative to multimedia content is stored. The preferences have the strong advantage to completely abstract multimodal content from dimensionality and scaling issues, and thus to facilitate fusion of heterogeneous descriptors. The challenge is then how to implement retrieval algorithms in preference space that are as effective as techniques based on more traditional representations (eg feature space). The RankBoost algorithm offers us a very convenient solution, especially when considering the soft ranking function as a weak ranking. The performance is very close to state of the art SVM-based fusion algorithm operating in feature or dissimilarity spaces. The algorithm is parameter free and thus avoid any lengthy and hazardous parameters estimation. Finally, RankBoost is really fast compared to SVM-based approaches which is a crucial argument for online retrieval systems.

REFERENCES

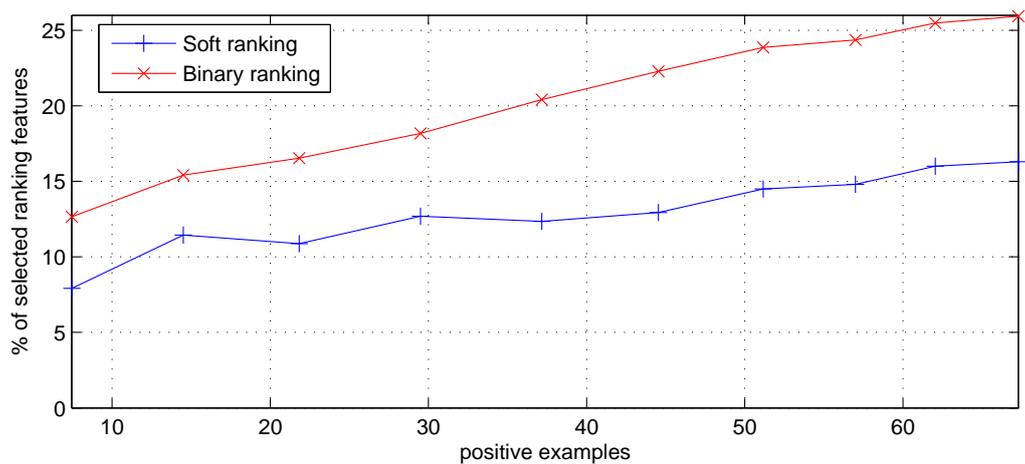
- [1] E. Y. Chang, B. Li, G. Wu, and K. Go, "Statistical learning for effective visual information retrieval," in *Proceedings of the IEEE International Conference on Image Processing*, 2003.
- [2] X. Zhou and T. Huang, "Small sample learning during multimedia retrieval using biasmap," in *Proceedings of the IEEE Conference on Pattern Recognition and Computer Vision, CVPR'01*, vol. 1, Hawaii, 2004, pp. 11–17.
- [3] X. Zhou, A. Garg, and T. Huang, "A discussion of nonlinear variants of biased discriminant for interactive image retrieval," in *Proc. of the 3rd Conference on Image and Video Retrieval, CIVR'04*, 2004, pp. 353–364.
- [4] R. Yan, A. Hauptmann, and R. Jin, "Negative pseudo-relevance feedback in content-based video retrieval," in *Proceedings of ACM Multimedia (MM2003)*, Berkeley, USA, 2003.
- [5] J. Kittler, M. Hatef, R. Duin, and J. Matas, "On combining classifiers," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, no. 3, pp. 226–239, 1998.
- [6] R. Duin, "The combining classifier: To train or not to train?" in *Proceedings of the 16th International Conference on Pattern Recognition, ICPR'02*, vol. II. Quebec City: IEEE Computer Society Press, 2004, pp. 765–770.
- [7] Y. Wu, E. Y. Chang, K.-C. Chang, and J. Smith, "Optimal multimodal fusion for multimedia data analysis," in *Proceedings of ACM Int. Conf. on Multimedia*, New York, 2004.
- [8] W. H. Hsu and S.-F. Chang, "Generative, discriminative, and ensemble learning on multi-modal perceptual fusion toward news video story segmentation," in *ICME*, Taipei, Taiwan, June 2004.
- [9] M. Wang, X.-S. Hua, X. Yuan, Y. Song, and L.-R. Dai, "Optimizing multi-graph learning: towards a unified video annotation scheme," in *MULTIMEDIA '07: Proceedings of the 15th international conference on Multimedia*. New York, NY, USA: ACM, 2007, pp. 862–871.
- [10] T.-T. Pham, N. E. Maillot, J.-H. Lim, and J.-P. Chevallet, "Latent semantic fusion model for image retrieval and annotation," in *CIKM '07: Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*. New York, NY, USA: ACM, 2007, pp. 439–444.
- [11] M. Haindl, J. Kittler, and F. Roli, "Multiple classifier systems," in *Series: Lecture Notes in Computer Science*, vol. 4472/2007. Springer, 2007.



(a)



(b)



(c)

Fig. 7. Multimodal video retrieval using a Relevance Feedback strategy with a) Mean Average Precision, b) Computational time for Hard ranking and soft ranking, c) percentage of ranking features selected by RankBoost, and d) modality usage for the 24 queries.

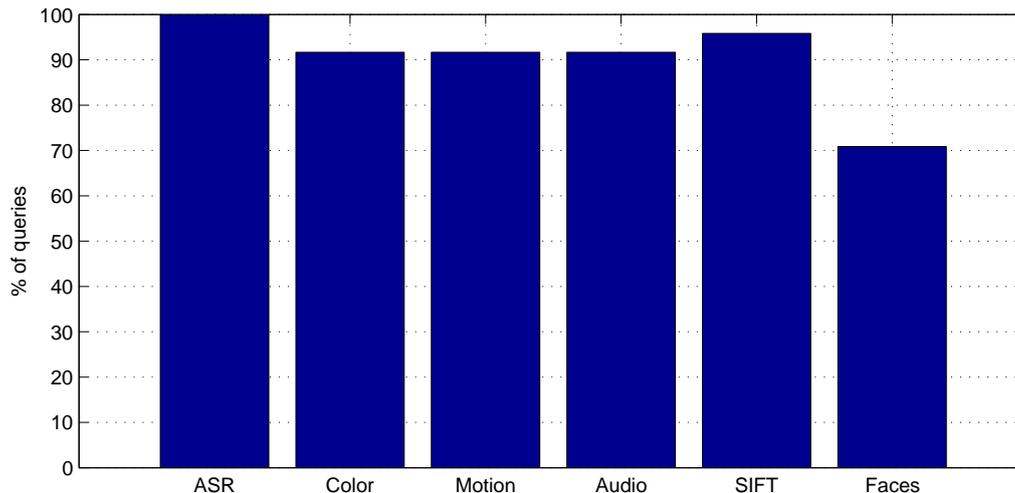


Fig. 8. Modality usage: Selected ranking feature frequencies over the 24 TRECVID queries (in %).

- [12] J. R. Smith, A. Jaimes, C.-Y. Lin, M. Naphade, A. Natsev, and B. Tseng, "Interactive search fusion methods for video database retrieval," in *IEEE International Conference on Image Processing (ICIP)*, 2003.
- [13] J. Yang and A. Hauptmann, "Multi-modality analysis for person type classification in news video," in *Electronic Imaging'05 - Conference on Storage and Retrieval Methods and Applications for Multimedia*, San Jose, USA, Jan 2005.
- [14] E. Chávez, G. Navarro, R. Baeza-Yates, and J. Marroquin, "Searching in metric spaces," *ACM Computing Surveys*, vol. 33, no. 3, pp. 273–321, Sep. 2001.
- [15] L. Boldareva and D. Hiemstra, "Interactive content-based retrieval using pre-computed object-object similarities," in *Conference on Image and Video Retrieval, CIVR'04*, Dublin, Ireland, 2004, pp. 308–316.
- [16] D. Heesch, A. Yavlinsky, and S. Rüger, "Nnk networks and automated annotation for browsing large image collections from the World Wide Web," in *Proc ACM Int'l Conference Multimedia*, 2006, pp. 220–224.
- [17] E. Bruno, N. Moenne-Loccoz, and S. Marchand-Maillet, "Design of multimodal dissimilarity spaces for retrieval of multimedia documents," *To appear in IEEE Transaction on Pattern Analysis and Machine Intelligence*, 2008.
- [18] E. Pekalska, R. P. W. Duin, and P. Paclík, "Prototype selection for dissimilarity-based classifiers," *Pattern Recogn.*, vol. 39, no. 2, pp. 189–208, 2006.
- [19] D. Heesch and S. Rueger, "NNk networks for content-based image retrieval," in *26th European Conference on Information Retrieval*, Sunderland, UK, 2004.
- [20] G. P. Nguyen, M. Worring, and A. W. M. Smeulders, "Similarity learning via dissimilarity space in CBIR," in *MIR '06: Proceedings of the 8th ACM international workshop on Multimedia information retrieval*. New York, NY, USA: ACM Press, 2006, pp. 107–116.
- [21] E. Pekalska, P. Paclík, and R. Duin, "A generalized kernel approach to dissimilarity-based classification," *Journal of Machine Learning Research*, vol. 2, pp. 175–211, December 2001.
- [22] Y. Freund, R. Iyer, R. E. Schapire, and Y. Singer, "An efficient boosting algorithm for combining preferences," *Journal of Machine Learning Research*, vol. 4, pp. 933–969, November 2003.
- [23] R. E. Schapire and Y. Singer, "Improved boosting algorithms using confidence-rated predictions," *Journal of Machine Learning*, vol. 37, no. 3, pp. 297–336, 1999.
- [24] J. Cheng and K. Wang, "Active learning for image retrieval with co-svm," vol. 40, no. 1, pp. 330–334, January 2007.
- [25] S. Hoi, R. Jin, J. Zhu, and M. Lyu, "Semi-supervised svm batch mode active learning for image retrieval," 2008, pp. 1–7.
- [26] R. Liu, Y. Wang, T. Baba, D. Masumoto, and S. Nagata, "Svm-based active feedback in image retrieval using clustering and unlabeled data," vol. 41, no. 8, pp. 2645–2655, August 2008.
- [27] C. Petersohn, "Fraunhofer HHI at TRECVID 2004: Shot boundary detection system," in *TREC Video Retrieval Evaluation Online Proceedings*, 2004.
- [28] A. F. Smeaton, P. Over, and W. Kraaij, "Evaluation campaigns and trecvid," in *MIR '06: Proceedings of the 8th ACM international workshop on Multimedia information retrieval*. New York, NY, USA: ACM Press, 2006, pp. 321–330.
- [29] A. Jain, A. Vailaya, and X. Wei, "Query by video clip," *Multimedia Syst.*, vol. 7, no. 5, pp. 369–384, 1999.
- [30] D. Lowe, "Object recognition from local scale invariant features," in *Proceedings of the International Conference in Computer Vision, ICCV'99*, Corfu, 1999, pp. 1150–1157.
- [31] P. Viola and M. Jones, "Robust real-time face detection," *International Journal of Computer Vision (IJCV)*, vol. 57, no. 2, pp. 137–154, 2004.
- [32] J. Gu, L. Lu, H. Zhang, and J. Yang, "Dominant feature vectors based audio similarity measure," in *PCM*, no. 2, 2004, pp. 890–897.
- [33] N. Moënné-Loccoz, E. Bruno, and S. Marchand-Maillet, "Interactive partial matching of video sequences in large collections," in *IEEE International Conference on Image Processing*, Genova, Italy, 11-14 September 2005.

Eric Bruno received his M.S. degree from the Engineers School of Physics in Strasbourg, France in 1995, and his Ph.D in signal processing from the Joseph Fourier University, Grenoble, France in 2001. Since 2002, he is working at the Computer Vision and Multimedia Laboratory, University of Geneva, Switzerland, as a research associate. His research interests focus on multimedia information retrieval, machine learning and statistical approaches for fusion.

Stéphane Marchand-Maillet received his PhD on theoretical image processing from Imperial College, London in 1997. He then joined the Institut Eurecom at Sophia-Antipolis (France) where he worked on automatic video indexing techniques based on human face localization and recognition. Since 1999, he is Assistant Professor in the Computer Vision and Multimedia Lab at the University of Geneva, where he is working on content-based multimedia retrieval as head of the Viper research group. He has authored several publications on image analysis and information retrieval, including a book on low-level image analysis. His current research interests are in the study of semantic extraction from collaborative interaction.