

Interactive Retrieval of Video Using Pre-computed Shot-Shot Similarities

Liudmila Boldareva Djoerd Hiemstra

University of Twente, EEMCS faculty,
P.O.Box 217, 7500 AE Enschede, The Netherlands
{L.Boldareva,D.Hiemstra}@utwente.nl

Abstract

A probabilistic framework for content-based interactive video retrieval is described. The developed indexing of video fragments is originated from the probability of the user's positive judgment about key-frames of video shots. Initial estimates of the probabilities are obtained from low-level feature representation. Only statistically significant estimate are picked out, the rest is replaced by an appropriate constant allowing efficient access at search time without loss of the search quality and leading to improvement in most experiments. With the time, these probability estimates are updated from the relevance judgment of users performing searches, resulting in further substantial increase in mean average precision.

1 Introduction

With rapid development of digital media, content-based multimedia retrieval has become an active research area. Having started its history in text documents, information retrieval quickly became dearly needed for other media such as still images and video.

The pioneering image retrieval systems used experience from the text retrieval domain, successfully adopting the vector space model [9, 15, 25]. Probabilistic approaches suggested for text retrieval [14, 22] gained less popularity with some notable exceptions [3, 32, 36]. One of the reasons for this lies in the difficulty of translating lower-level features that index the visual content into probability values. Often, content based retrieval systems rely on active participation of the searcher in the retrieval process, known as relevance feedback [23]. ‘Relevance feedback’ is a broad term sheltering various models of learning the user’s information need from the two-way communication where the searcher plays an active role. The early implementations of interactive visual retrieval systems are QBIC [9], MARS [25], MindReader [15], Viper [19], PicHunter [3]. More recently, machine learning methods have been successfully applied to visual information retrieval, e.g. self-organising maps [16] and support vector machines [5, 31].

One fundamental problem that needs to be solved in visual information retrieval, is the *semantic gap*—a mismatch between the human perception of visually rich documents and their representation in the storage. This has been an active research topic for decades [30, 29, 21, 11]. We believe that simple tools based on low-level feature representations, and matching functions on them, are unlikely to bridge the semantic gap in the near future, without additional functionality like relevance feedback and long-term learning of image representations.

Another difficult problem in multimedia retrieval is the *efficiency* of search algorithms. High dimensionality of the search space hinders effective indexing of it, introducing the problem of ‘dimensionality curse’ [7]: with many dimensions, that feature vectors usually are, it is hard to implement a similarity matching algorithm that is substantially faster than a linear scan over the data [35]. This puts a possible interaction away. The system has to rely on dimensionality-reduction indexing techniques (e.g. [8]) or other smart approaches to access objects most similar to the given examples [4]. Still, the performance of such

systems comes nowhere near to the performance of text retrieval systems.¹ The only viable solution may be processing as much of available information in advance as possible!

In this article we propose a framework for content-based indexing and retrieval, that

- can use any available technique for feature extraction and similarity matching, and allows easy combination of different sources of information;
- allows efficient interaction with the user, and is capable of learning from that interaction;

Therefore the proposed framework is spared of the two problems stated above.

We give a statistical interpretation to the data-driven similarity between elements of the video key-frames collection that fits into a probabilistic framework for efficient interactive retrieval. This framework accommodates both short-term learning within one retrieval session and long-term learning from relevance feedback gathered in multiple retrieval sessions.

In the next section a general Bayesian framework for interactive retrieval is introduced, and subsequently, two important components of it are discussed—the data representation and the user feedback. Section 5 talks about long-term learning from previous searches based on user’s relevance judgements. Experiments have been carried out to verify the benefits of the proposed methods and compare to other techniques. The data of TREC Video retrieval workshop [28] serves as a testbed. The experimental setup and results are reported in Section 7, after an overview of related work given in Section 6.

2 Interactive retrieval in Bayesian terms

Let \mathcal{I} be a collection of information objects x , e.g. key-frames for video shots, among which there is what the user is looking for, the search target denoted

¹At the time of writing this paper, <http://images.google.com> provides access to almost 1200 million images using text search techniques, orders of magnitudes more than any content-based approach could possibly manage.

here by T . During the search process, the system presents the user with intermediate retrieval results. The user can indicate which of the objects are relevant to his/her information need—those are positive examples. If an object is not relevant to the query, the user may indicate so, thus providing the system with negative examples. Given the feedback information, the retrieval system produces a new set of objects to be assessed by the user. There may be several loops of relevance feedback during one search session.

The probabilistic framework is introduced as follows. Consider disjoint random events of the user feedback regarding the relevance of an object x . Let δ_x be the corresponding indicator function taking the value one if the user marks the candidate object x as relevant and zero otherwise. Note that the present framework can be generalised for multiple choices of feedback, e.g. by introducing a third event of explicit negative judgement.

We want to use the concepts ‘relevant’ and ‘non-relevant’ without having to refer to lower-level features. Instead, the objects in the collection are related to each other according to the most likely user opinion about their relationship. For two objects x and y , the following conditional probability reflects their ‘measure of closeness’: $P(\delta_x = 1|T = y)$, the probability of object x being marked by the user as relevant given that y is referred to as the target for the search. When unambiguous, the shorthand notation $P(\delta_x|T)$ denotes the probability of a certain user action concerning the object x .

2.1 Estimating probability of relevance

The goal is to predict, or identify, the set of objects relevant to the user’s information need, based on his/her request accompanied by feedback and the existing data representation. In a Bayesian framework [22, 3] the problem is restated as estimation of the probability of relevance $P(T)$ given user’s relevance judgements $\{\delta_{x_1}, \dots, \delta_{x_n}\}$ on the set of candidate objects $\{x_1 \dots x_n\}$ and the data indexing. We write it down in the following iterative form, with the assumption that the user actions $\{\delta_{x_1}, \dots, \delta_{x_n}\}$ are conditionally independent

given the target T .

$$P^{\text{new}}(T) = P(T|\delta_{x_1}, \dots, \delta_{x_n}) = \frac{P^{\text{old}}(T) \prod_{s=1}^n P(\delta_{x_s}|T)}{P(\delta_{x_1}, \dots, \delta_{x_n})}. \quad (1)$$

The conditional independence assumption used here states that the user’s judgement about the relevance of a certain item is not affected by the relevance of other displayed items.

$P^{\text{new}}(T)$ becomes $P^{\text{old}}(T)$ for the next iteration; $\{(\delta_{x_1}, \dots, \delta_{x_n})\}$ are provided by the user. $P(\delta_{x_s}|T)$ represents conditional dependency between elements x and T . According to Equation (1), in order to retrieve the most likely relevant answers T , one needs to know all δ_{x_s} and $P(\delta_{x_s}|T)$. They are the subjects of the following two sections.

3 Association-based data representation

If we were to use a graphic model, the collection of objects and the corresponding conditional probabilities $P(\delta_x|T)$ could be visualised as a directed graph with nodes $x \in \mathcal{I}$, and weighted arcs connecting them. Each object x is described by its associations with a number of other objects. The strength of the association is the weight of the arc, $P(\delta_x|T)$. The whole structure is called here ‘association matrix’ denoted further by \mathbf{M} .

Preferably for each item there needs to be need only a few associations, which refer to high-level semantics and agree with the observed users’ acts of relevance feedback. We arrive at these associations as follows. Starting at the point when we do not have knowledge about human perception of similarities between objects, the initial associations are derived from a similarity measure on lower-level features, such as colours, textures, shapes. Typically such similarity measures take values in the range of (non-negative) real numbers and thus cannot be directly used as an initial estimate for $P(\delta_x|T)$.

As a first step, the pair-wise similarities are transformed by fitting a probability distribution on it. Since a priori we cannot prefer some objects from

the collection to others in the sense of the distribution of estimates of $P(\delta_x|T)$, the underlying pair-wise similarities of the whole collection are assumed to be conform the Normal distribution. This gives equal emphasis of the alike similarities and spreads the observations evenly on the interval $[0, 1]$ according to their probability of occurrence and not to the magnitude of the similarity measure. As a result it reduces the influence of outliers and preserves the scale of the similarities between objects, i.e. ‘improves the discrimination capabilities of the similarity measure’ [1]. The result of this step is a square table with all possible pair-wise estimates of conditional probabilities, also containing errors induced by the mismatch between data-driven similarity measures and human perception.

The value of $P(\delta_x = 0|T) \equiv 1 - P(\delta_x = 1|T)$ computed in this way can be interpreted as a *P-value*, the probability that a variable assumes a value greater than or equal to the observed one strictly by chance. That is, the *P-value* is the probability that the computed similarity between two objects purely by chance does not exceed the ‘true’ one, therefore x will not be found relevant to T by the user. As a second step, we specify some α , the upper bound for the *P-value*, so that only statistically significant pair-wise similarities and their corresponding $P(\delta_x|T)$ are taken into account. Probability estimates for not significant similarities are replaced by an appropriate constant further denoted by \bar{p} . That means, while updating $P(T)$ for each object in Equation (1), the conditional probabilities $P(\delta_x|T)$ is substituted by \bar{p} if its estimate is below $1-\alpha$. Here $1-\alpha$ serves as a cut-off threshold for the right tail of the distribution of pair-wise similarities. Similarly, a threshold for the left tail of the distribution will differentiate between significant and non-significant estimates of dis-similarity. The corresponding threshold for the dis-similarity-based estimates for the rest of the article is set to zero, because judging non-relevance from low-level features is less practical. An object x that has a significant $P(\delta_x = 1|T = y)$ is called a *neighbour* of y . Our idea is that the left out estimates contain more noise than useful information and removing them will not harm retrieval quality, while

improving efficiency.

The data representation in the form of association matrix as described here has the following advantages:

- Different sources of information on the (estimated) relevance of objects are on the same scale and can be efficiently combined;
- The relevance information obtained from the searchers can be used to improve conditional probability estimates in the association matrix.

4 Input provided by the user

Another factor that plays a role in an interactive retrieval session is the input from the user. One question is how to interpret user actions. Another related question is how to select candidate objects to be presented to the user.

4.1 Interpretation of user feedback

During a search session, the current probability of an object to satisfy the user's information need $P(T)$ is updated according to Equation (1). Every object can be marked relevant/non-relevant to the user's information need (or not marked at all), and these events are disjoint. If the user is not supposed to ignore the objects presented for relevance assessment, the objects that are not marked by the user as relevant, take part in the probability update as if they are explicitly rejected by the user, and for their neighbours, $P(\delta_x = 0|T) \equiv 1 - P(\delta_x = 1|T)$ is used in Equation (1) to update $P(T)$.

Deploying explicit negative relevance judgements is less obvious, however. Differ from giving positive examples, explaining non-relevance is harder for the user [26, Section 3]. Excessive amount of negative feedback may have negative effect on retrieval [19]. Therefore associations to the neighbours of the negative examples are not considered and \bar{p} is used in the update of their $P(T)$ in Equation (1). This scenario roughly corresponds to nearest neighbour search,

effectively eliminating seen non-marked examples from further consideration. Here \bar{p} plays the role of a smoothing constant. When it equals zero, the search space is limited to the neighbours of all positive examples.

4.2 New display for the next iteration

After updating $P(T)$, a new set of objects should be presented to the user for relevance judgement. Selection of candidate objects (display update) is an important part of the search process, since it determines what the system will learn from the interaction. Each iteration should bring the user closer to his/her target object. ‘Closer to the target’ may have various interpretations, such as: the posterior probability $P(T)$ of the desired information object tends to 1; or the target object(s) approach the top of the ranked list. In this article we describe experiments performed with the following three display update strategies.

Best-Target. Following the probability ranking principle [22], $P(T)$ is considered as a score that the element receives during a retrieval session. The next display set consists of (new) objects that have largest values of $P(T)$. This ‘Best-Target’ strategy is plausible for a user unfamiliar with content-based retrieval (thus, the majority of potential users). The screen often contains objects that are the neighbours of good examples provided by the user. The user is able to observe the immediate result of his/her action. This display update strategy does not intent to explore *new* objects that, in the selected similarity measure, differ from the relevant ones already seen.

Non-deterministic strategies. In order to diversify the set of displayed elements, we introduce non-deterministic strategies as extensions to the ‘Best-Target’ one. First, we consider weighted selection of display candidates, or *Proportional sampling*: The chance to be displayed for an object is proportional to its probability of relevance. When the distribution of $P(T)$ is peaked, proportional sampling converges to the ‘Best-Target’ method.

Second, instead of selecting the candidates with the highest score, the *Sample-of-Best* strategy makes the selection among those objects of which the probability of relevance increased since several previous iterations. This includes mainly neighbours of the relevant examples. Occasionally, elements with low but consistently growing $P(T)$ may be selected for display, giving the user a chance to see potentially relevant objects that may be very different from what he/she has already seen. Ideally, the number of elements of which $P(T)$ increases should shrink on to the group of objects that satisfy the user’s information need.

5 Learning from past retrieval sessions

As stated before, there may be more than one association matrix to be used in the retrieval process. As we later show in experiments, combination of information sources may result in a better retrieval quality. Still, it is more promising to improve the existing feature estimates stored in the index.

At the end of a successful retrieval session the system is in possession of the list of objects displayed to the user $\{x_1, x_2, \dots, x_m\}$ and the corresponding relevance judgements $\{\delta_{x_1}, \delta_{x_2}, \dots, \delta_{x_m}\}$. This information can be used for improving the corresponding estimates of the probabilities $P(\delta_{x_1}|T)$, $P(\delta_{x_2}|T)$, \dots , $P(\delta_{x_m}|T)$.

The event of selecting x by the user as relevant/non-relevant should result in an improved estimate of the corresponding probability $P(\delta_x|T)$. To update the involved estimates, we use the maximum likelihood (ML) principle, which boils down to counting events. Let $P(\delta_x|T)$ should be updated after observing one feedback action on x while seeking T . The following equation corresponds to frequency-based update for the case when δ represents binary choice:

$$P^{\text{new}}(\delta_x|T) = \frac{\kappa \cdot P^{\text{old}}(\delta_x|T) + i}{\kappa + 1}, \quad i = \begin{cases} 1 & \text{for positive feedback} \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

Here κ is the number of observations prior to the current retrieval session. In the beginning, when there are no observations, it denotes the weight of the

feature-based estimate of $P(\delta_x|T)$. Every new observation adds a unit to the denominator, and in the case of positive judgement, a unit is also added to the numerator.

Such a method of updating estimates of $P(\delta|T)$ enhances subsequent retrieval sessions. It requires an extensive interaction history, which can be achieved in, for instance, the World Wide Web environment where the number of potential users is large.

6 Related work

Maron and Kuhns in [18] talked about indexing of library documents with respect to the actions taken by the searcher. The set of possible user actions includes among others (a) providing index terms for an information request, including specification of fields of interest, and (b) marking a document relevant, given the indexing terms. They suggest that this index can be refined based on the judgements of searchers. We used this work as a source of inspiration adapting it to the case when objects in the collection are atomic elements that can serve as descriptions for other members of the collection.

For modelling the interaction between the user and the system, we adapted a Bayesian framework similar to that of the PicHunter retrieval system [3]. The image relevance is determined based on the judgements of the user and (a model for) conditional dependencies between the elements. There is a large body of work dedicated to learning from relevance feedback in content-based retrieval, by far not restricted by probabilistic models. The main learning objective is to separate the answers to the query from the rest of the collection. This can be achieved in many ways: e.g., through (variations of) feature space transformation [25, 15, 27, 24], or through partitioning the data set [17, 5, 31]. The Bayesian framework turns out very suitable for modelling learning from the interaction, as it accommodates both short-term and long-term learning capabilities.

When using a large collection, to achieve a near-real time system response, it is not unusual that as much data as possible is processed beforehand. For example, using pre-computed sets of nearest neighbours in (several) feature spaces has proved to speed up database performance during information search and browsing [6]. Strategies to combine relevance information from different sources are also studied in that work. We envision the proposed data organisation as a directed graph, but this is not unique. Stemming from a different paradigm, in [12] the collection of images is represented as a graph where images represent vertices connected to a small number of other vertices's. Those are determined as nearest neighbours in numerous weighted combinations of available feature spaces. The optimal weights in a given query context is determined by the positive examples selected by the user when browsing. By in advance quantifying the feature vectors into tree-structured self-organising maps [16], relevance for images is determined by their distance on the map to the user-judged examples. Later in time, the relevance judgements make up for a separate self-organising map that is used along the feature-based maps. In [13] it is proposed to represent images with a vector of relevance judgements collected from past retrieval sessions. By applying latent semantic indexing technique, a set of closely associated images is determined and used during retrieval. This very similar in spirit to our approach described here.

7 Experiments and evaluation

7.1 Video collection, data pre-processing and experiment setup

The experimental evaluation is performed in the framework of 2003 TREC Video retrieval evaluation workshop² [28]. The video materials are CNN, ABC and C-SPAN news programs recorded between 1998 and 2001. The videos are seg-

²Some of the experiments performed with the collection of TREC Video-02 are presented in [2].

mented into shots, and from each shot a representative key-frame is extracted. The key-frames and shot boundaries are part of the data set, as well as speech transcripts from a large-vocabulary automatic speech recognition system [10]. In total the test collection contains about 60 hours of recordings. Video shots represented by key-frames are the objects the system deals with. There are 24 search tasks, or topics, based on real user requests. Each search task consists of a short text description and, in most cases, few image and/or video examples of the user’s information need. Experiments reported here are carried out with the following association matrices:

1. \mathbf{M}^t . Similarity between shots is computed using language models built on the text from the speech transcripts. The language model used in the experiments is described in [36].
2. Similarity between shots is based on their visual properties, namely:
 - (a) \mathbf{M}^c . Weighted L_1 distance on 3 colour moments (mean, variance and skewness) in Hue, Saturation, Value colour space of the whole image. Implemented according to [29].
 - (b) \mathbf{M}^b ‘Bag of blocks’ likelihood as described in [37]. Similarity between two key-frames (g, f) is computed as likelihood that samples taken from g could serve as a model to explain f . A draw of 100 blocks of 8x8 pixels represents the key-frame.
3. \mathbf{M}^{t+b} . Run-time combination of the association matrices based on textual and visual features. The combined relevance score for an object is computed as sum of the scores it achieves when using each of the matrices. The combination uses the assumption that the distribution of text-based pair-wise similarities is independent of the visual-based one. By counting neighbours in each of the two matrices to combine, we find that this assumption cannot be rejected for the text-based association matrix \mathbf{M}^t with either of \mathbf{M}^c and \mathbf{M}^b .

The threshold $(1-\alpha)$ is set such that on average 1.2–1.4% of possible values need to be stored. The value of \bar{p} is set to 0.10 for all experiments which appears to be close to the optimum.

In addition to comparison within TREC evaluation framework, we implemented Rui’s relevance feedback algorithm as described in [25]. The feature space is the one used for computation of \mathbf{M}^c , that is, three colour moments for each of Hue, Saturation, Value channels and L_1 norm as the distance measure. The vector components are normalised, as described in the paper, so that they have equal emphasis on the resulting similarity. The initial weights for the vector components are set uniform, to be consequently updated by the intra-weight update algorithm. A random draw of six key-frames³ from the collection served as six image queries for each topic. The results in the graphs are averaged over these six queries.

We perform an empirical study on performance differences caused by the prior distribution of the probability of relevance. In order to provide for a number of experiments a better than arbitrary estimate of the prior probability, the text from search topics serves as a text query to match against the speech transcripts. Such a scenario is quite specific to video collections, where the speech is aligned with the image and can be seen as a surrogate annotation. In an un-annotated still image collection, a text query is of little use, as well as in the case when the text query has no match. For that situations the prior probability of relevance is determined by the image queries that we used in the Rui’s setup.

A retrieval session starts with browsing a display set of 12 key-frames generated by the text or an image query. The key-frames are ranked by their probability of relevance. A standard TREC evaluation metric, mean average precision (MAP), is used as a measure of user’s satisfaction (see [34, Appendix]). Wilcoxon signed rank test [20] determines whether the performance figures between two methods differ significantly.

³The number of answers to the 24 topics varied from 6 to 665

7.2 Automated experiments

In the series of experiments referred to as ‘automated’, the user input is replaced with relevance judgements of the TREC assessors who play the role of a ‘generic user’. The experiments are carried out on a subset of the collection selected such that half of the key-frames are relevant to at least one of the 24 topics. This yields a set of 4096 shots of the collection of 2003. In this way we could test the proposed probabilistic framework, and find the best setup to be used in the experiments with real users. A selection of automated experiments has been repeated on the whole 2003 TREC Video test data set containing about 32 000 key-frames. The results are consistent; note that mean average precision number is lower when all key-frames are used, due to the lower proportion of relevant shots.

7.2.1 Effect of truncation the association matrix

Mean average precision curves for \mathbf{M}^t are shown in Fig. 1. Two situations, when using one of the six images as a query, and when querying with words from the TREC topics descriptions, are shown. Keeping only significant $P(\delta_x|T)$ leads to the increase of mean average precision compared to using all pairs of conditional probabilities, both with visual- and text-based matrices. This evidence confirms the choice of the α -value to determine significant pair-wise similarities (Section 3).

Fig. 2 shows search progress with the colour moments-based visual feature. Clearly, learning the optimal weights for the vector components does not perform better than the learning within the Bayesian framework. It is interesting to note, that on average, the set of optimal weights found by the algorithm, emphasizes the same vector components that have larger weights in [29], which we used in the implementation. It is worth mentioning that the vector-based model cannot straightforward take advantage of a text query. Using a text feature vector is certainly possible, but the retrieval efficiency would be impaired by extremely high dimensionality of the text feature space. However, the main

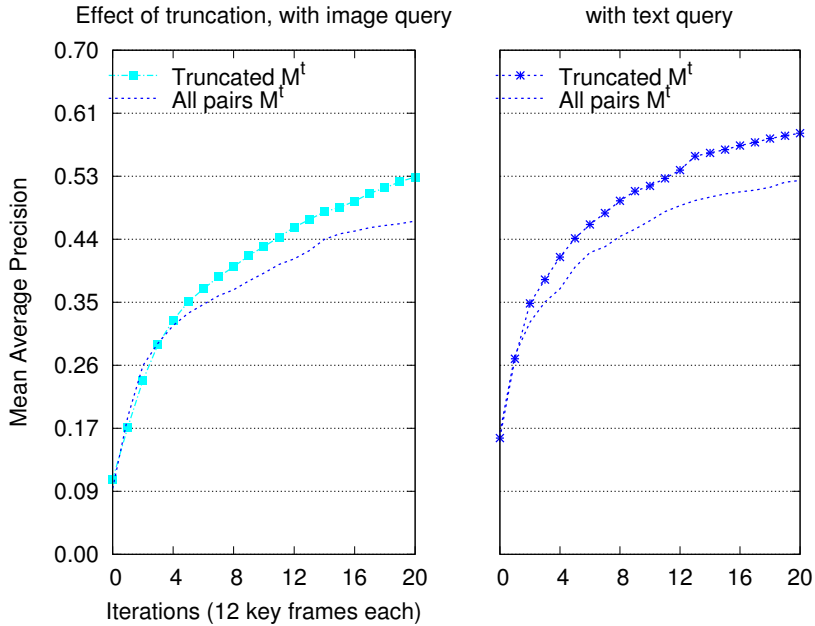


Figure 1: All pairs and truncated conditional probabilities for M^t , using image queries (left) and the text query (right).

advantage of the proposed framework over vector space models is in its potential to use complex visual content representations, that do not necessarily fit into the vector space models.

7.2.2 Display update strategies.

The ‘Best-Target’ display update with an ad-hoc tuned value of \bar{p} offers great improvement over iterations, both with the text and visual queries. By making sure that the user does not see the same object twice, the danger of getting stuck in an isolated island of nearest neighbours.

Proportional sampling does not differ substantially from the ‘Best-Target’ strategy. When the probability of relevance distribution becomes peaked, so that few elements share most of the probability mass, the system tends to select the next display set among those few elements. Because they also occupy the top of the ranked list, the display turns out to be quite similar to that of the ‘Best-

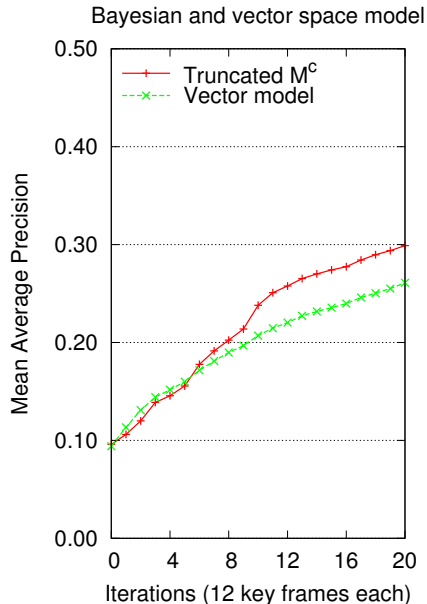


Figure 2: All pairs and truncated conditional probabilities using \mathbf{M}^b , and Rui’s relevance feedback model. Image queries used.

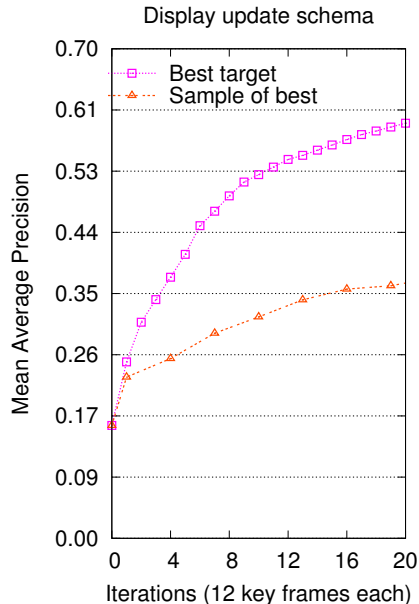


Figure 3: Display update strategies for \mathbf{M}^b matrix, with the text query. Proportional sampling is not plotted.

Target’ strategy (on average they share 68% of displayed objects when using text query, and 43% with image query). Consequently, there is no significant difference in mean average precision.

The ‘Sample-of-Best’ strategy in practice does not perform better than the deterministic ‘Best-Target’ method (see Fig. 3). This does not support the consideration that sampling among the promising candidates with increasing $P(T)$ can potentially give a better collection representation.

7.2.3 Combination of different modalities

Figure 4 shows mean average precision curves as a result of combining the scores from different matrices. From the graphs one can conclude that combining different sources of information improves the retrieval quality, in terms of mean average precision, by several percentage points. The combination of two ma-

trices both based on visual appearance of the key-frames, namely M^b and M^c , does not have such effect and at best results in mean average precision that not exceeding the one that performs better. This is an expected result, since the score based on the colour statistics does not bring in new information compared to the Bag-of-Blocks likelihood score, that already has the color information in it.

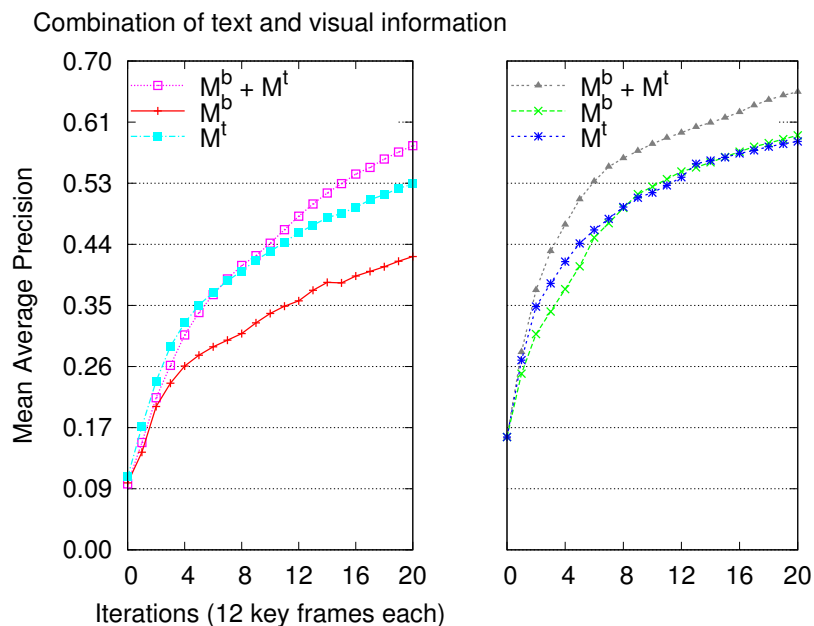


Figure 4: Combinations of different information sources, using image queries (left) and the text query (right).

7.3 Live experiments

In the *live* experiments, the search tasks have been performed on 2003 TREC Video collection by real users. All of them are students of University of Twente aged between 19 and 26. Each search task took at most 15 minutes. None of the users was familiar with the search system or was related to the development of it. A large proportion of the users' positive feedback turns out to be relevant according to the ground truth (see 'Agreement' in Table 1), thus the described automated experiments can indeed serve as an approximation to real life (see

[33] for an analysis of agreement between the TREC assessors).

The set-up for live search sessions is similar to the automated experiments, using ‘Best-Target’ display update strategy, \mathbf{M}^b as the access index. Words from the descriptions of search tasks served as the text query. The users retrieved key-frames (images), and not the corresponding videos. The resulting mean average precision at the end of the ‘Best-Target’ experiment is 0.245 (see Table 1), which is seventh best mean average precision for that year. For this run, 78% of the shots selected by the user were relevant according to TREC. At the same time, 48% of all relevant shots that have been displayed, were not marked as such. The users tend to miss some relevant key-frames from the display sets that contained many of those. Partially, the relevant items are missed due to the fact that the user observed still frames, and not the video shots themselves. Therefore the key-frames of the relevant shots that do not show the required object or scene have not been marked as relevant. This issue can be resolved by enabling video display in the interface.

In the other experiment that showed the user screens sampled uniformly (mean average precision 0.026), the proportion of missed shots is much lower (31%), as well as the agreement with TREC committee (55%). This is an indication that the users are inclined to mark a larger proportion of the displayed key-frames as relevant when little of those are on the screen, i.e. the independence assumption used in Equation (1) apparently does not hold.

Table 1: TRECVID experiments with real users

System type	MAP	Agreement with NIST	Missed Relev.
Best-Target	0.245	78.74%	48.98%
Uniform Sampling	0.026	55.00%	31.25%

7.4 Learning from live experiments

We use the feedback data collected from six live retrieval sessions to conduct a preliminary experiment with training the association matrix from user feedback.

To be able to use a controlled environment, these experiments are automated. However, the training data itself comes from real user sessions and contains user erroneous responses regarding the ground truth provided by TREC, mentioned in Section 7.3. The estimates of conditional probabilities for the key-frames that have been displayed to the users in the search sessions are updated using a ML method described in Section 5, Equation (2). The key-frames marked as relevant serve as ‘targets’ in the training. After the training, the automated experiment is repeated using the same set-up, but with the new, updated association matrix, in this case M^b . As shown in Fig. 5, mean average precision is substantially increased, and the improvement is statistically significant. The most increase in mean average precision is observed at the beginning of the search, so that the user sees the desired objects earlier. The difference in mean average precision is not only due to more favourable re-arrangement of the relevant shots.⁴ In the experiments with the trained association matrix more relevant objects have been ‘displayed’ to the user (shown in the figure).

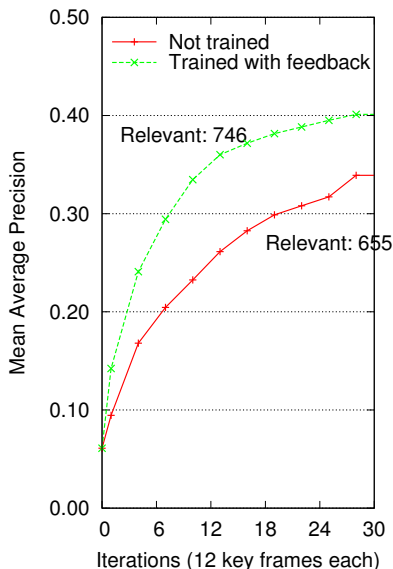


Figure 5: Automated experiments before and after training on real user feedback, using complete 2003 TREC Video test collection.

⁴Mean average precision is sensitive to having relevant shots on top of the ranked list.

8 Discussion and open questions

We found that the proposed image feature normalisation and smoothing by replacing similarities below a certain threshold with a constant, in the investigated visual-based and text-based feature spaces, results in higher mean average precision compared to the method that uses all pairwise similarity values. To perform the normalisation and truncation of the association matrix, we used statistics of a particular collection we wanted to search in. On a collection of video frames, our probabilistic model that uses colour moments for indexing, performs slightly better than a vector-space model using the same feature set, although we are aware of a limited nature of this comparison. The advantage of the proposed framework is that it can deploy any available technique of image understanding, to create an initial association matrix.

Truncation of the matrix enables efficient combination of different similarity measures, such as visual information from key-frames and transcripts of the speech occurring in video shots. Combination of independent sources of information has positive effect on the retrieval. We used equal weights when adding up the scores, but there should possibly be better weighting schemes—this needs to be investigated.

So far we did not observe any improvement when attempting to efficiently diversify the set of displayed objects as compared to common strategy of showing the most relevant candidates. Although, from the point of view of optimal learning, the ‘Best-Target’ is not the best choice, it is hard to compete with when real life data is used in the collection.

Keeping only the significant values allows an interactive retrieval system to ensure fast response time, which is a necessary condition for an interactive retrieval system. This in turn will provide vast amount of training data. Learning from the history of relevance judgements in retrieval sessions with the real users substantially improves the successive searches. The improvement is observed not only due to higher ranking of the previous positive examples, but also due

to a larger number of relevant key-frames that are displayed to the user. More advanced learning techniques are needed to update the conditional probability estimates between the unseen objects.

Acknowledgement

The authors would like to thank Thijs Westerveld of CWI Amsterdam for providing BoB likelihood scores for this collection.

References

- [1] S. Aksoy and R. M. Haralick. Feature normalization and likelihood-based similarity measures for image retrieval. *Pattern Recognition Letters*, 22(5): 563–582, 2001.
- [2] L. Boldareva and D. Hiemstra. Interactive content-based retrieval using pre-computed object-object similarities. In *Proceedings International Conference on Image and Video Retrieval, CIVR*, pages 308–316, Dublin, Ireland, July 2004.
- [3] I. J. Cox, M. L. Miller, T. P. Minka, T. V. Papatomas, and P. N. Yianilos. The bayesian image retrieval system, PicHunter: Theory, implementation, and psychophysical experiments. *IEEE Tran. On Image Processing*, 9(1): 20–37, 2000.
- [4] A. P. de Vries, N. Mamoulis, N. Nes, and M. L. Kersten. Efficient k-NN search on vertically decomposed data. In *ACM SIGMOD International Conference on Management of Data*, Madison, WI, USA, June 2002.
- [5] Harris Drucker, Behzad Shahrory, and David C. Gibbon. Relevance feedback using support vector machines. In *Proc. 18th International Conference on Machine Learning*, pages 122–129. Morgan Kaufmann, San Francisco, CA, 2001.

- [6] R. Fagin. Combining fuzzy information from multiple systems. *Journal of Computer and System Sciences*, 58:83–99, 1999.
- [7] C. Faloutsos. *Searching Multimedia Databases By Content*. Kluwer Academic Publishers, Boston, USA, 1996.
- [8] C. Faloutsos and K.-I. Lin. FastMap: A fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. In Michael J. Carey and Donovan A. Schneider, editors, *Proc. of the 1995 ACM SIGMOD International Conference on Management of Data*, pages 163–174. ACM Press, 1995.
- [9] M. Flickner, H. Sawhney, W. Niblack, and J. Ashley. Query by image and video content: the QBIC system. *IEEE Computer*, 28(9):310–15, 1995.
- [10] J. L. Gauvain, L. Lamel, and G. Adda. The LIMSI broadcast news transcription system. *Speech Communication*, 37(1-2):89–108, 2002.
- [11] V. N. Gudivada and V. V. Raghavan. Design and evaluation of algorithms for image retrieval by spatial similarity. *ACM Transactions on Information Systems*, 13(2):115–144, 1995.
- [12] Daniel Heesch and Stefan M. Ruger. NN^k networks for content-based image retrieval. In Sharon McDonald and John Tait, editors, *26th European Conference on IR Research, ECIR 2004*, volume 2997 of *Lecture Notes in Computer Science*, pages 253–266, Sunderland, UK, April 2004. Springer. ISBN 3-540-21382-1.
- [13] D. R. Heisterkamp. Building a latent semantic index of an image database from patterns of relevance feedback. In *Proceedings of 16th International Conference on Pattern Recognition (ICPR 2002)*, volume 4, pages 134–145, Quebec, Canada, 2002.
- [14] D. Hiemstra. *Using language models for information retrieval*. PhD thesis,

Centre for Telematics and Information Technology, University of Twente, 2001.

- [15] Y. Ishikawa, R. Subramanya, and C. Faloutsos. MindReader: Querying databases through multiple examples. In *Proc. 24th Int. Conference Very Large Data Bases, VLDB*, pages 218–227, 1998.
- [16] J. Laaksonen, M. Koskela, S. Laakso, and E. Oja. Self-organizing maps as a relevance feedback technique in contentbased image retrieval. *Pattern Analysis and Applications*, 4(2-3):140–152, 2001.
- [17] S. MacArthur, C. Brodley, and C. Shyu. Relevance feedback decision trees in content-based image retrieval. In *Proc. of IEEE Workshop on Content-based Access of Image and Video Libraries (CBAIVL'00)*, South Carolina, 2000.
- [18] M. E. Maron and J. L. Kuhns. On relevance, probabilistic indexing and information retrieval. *Journal of the Association for Computing Machinery*, 7:216–244, 1960.
- [19] H. Müller, W. Müller, S. Marchand-Maillet, T. Pun, and et al. Strategies for positive and negative relevance feedback in image retrieval. In *Proc. of International Conference on Pattern Recognition*, Barselona, Spain, Sept 2000.
- [20] L. Ott. *An introduction to statistical methods and data analysis*. Boston : PWS-KENT, 3 edition, 1988.
- [21] A. Pentland, R. Picard, and S. Sclaroff. Photobook: Content-based manipulation of image databases. *International Journal of Computer Vision*, 18(3):233–254, Feb 1996.
- [22] S. E. Robertson. The probability ranking principle in IR. *Journal of Documentation*, 33(4):294–304, 1977.

- [23] J. J. Rocchio. Relevance feedback in information retrieval. In G. Salton, editor, *The Smart Retrieval System: Experiments in Automatic Document Processing*, pages 313–323. Prentice Hall, 1971.
- [24] Y. Rui and T. Huang. Optimizing learning in image retrieval. In *Proceeding of IEEE int. Conference On Computer Vision and Pattern Recognition*, June 2000.
- [25] Y. Rui, T. Huang, S. Mehrotra, and M. Ortega. A relevance feedback architecture in content-based multimedia information retrieval systems. In *Proc. of IEEE Workshop on Content-based Access of Image and Video Libraries, in conjunction with IEEE CVPR*, 1997.
- [26] Ian Ruthven and Mounia Lalmas. A survey on the use of relevance feedback for information access systems. *Knowledge Engineering Review*, 18(2):95–145, 2003.
- [27] Simone Santini and Ramesh Jain. Integrated browsing and querying for image databases. *IEEE MultiMedia*, 7(3):26–39, 2000.
- [28] A. Smeaton, W. Kraaij, and P. Over. TRECVID 2003 - an introduction. In *Text Retrieval Conference TRECVID Workshop*. National Institute of Standards and Technology, Nov 2003. URL <http://www-nlpir.nist.gov/projects/tv2003/tv2003.html>.
- [29] M. A. Stricker and M. Orengo. Similarity of color images. In *Storage and Retrieval for Image and Video Databases (SPIE)*, pages 381–392, San Diego/La Jolla, California, USA, February 1995.
- [30] H. Tamura, S. Mori, and T. Yamawaki. Texture features corresponding to visual perception. *IEEE Trans. on Systems, Man, and Cybernetics*, 6(8):460–473, 1978.
- [31] Simon Tong and Edward Chang. Support vector machine active learning for

- image retrieval. In *Proceedings of the ninth ACM international conference on Multimedia*, pages 107–118. ACM Press, 2001. ISBN 1-58113-394-4.
- [32] N. M. Vasconcelos. *Bayesian Models for Visual Information Retrieval*. PhD thesis, Massachusetts Institute of Technology, 2000.
- [33] Ellen M. Voorhees. Variations in relevance judgments and the measurement of retrieval effectiveness. *Information Processing Management*, 36(5):697–716, 2000.
- [34] E.M. Voorhees and D.K. Harman, editors. *The Tenth Text Retrieval Conference (TREC-2001)*, Gaithersburg, MD, USA, Nov 2002. Department of Commerce, National Institute of Standards and Technology.
- [35] Roger Weber, Hans-Jörg Schek, and Stephen Blott. A quantitative analysis and performance study for similarity-search methods in high-dimensional spaces. In *Proc. of the 24rd International Conference on Very Large Data Bases*, pages 194–205. Morgan Kaufmann Publishers Inc., 1998. ISBN 1-55860-566-5.
- [36] T. Westerveld, A.P. de Vries, A.R. van Ballegooij, F.M.G. de Jong, and D. Hiemstra. A probabilistic multimedia retrieval model and its evaluation. *EURASIP Journal on Applied Signal Processing, special issue on Unstructured Information Management from Multimedia Data Sources*, pages 186–198, February 2003. URL <http://www.cwi.nl/~arjen/pub/EJASP-preprint.pdf>.
- [37] T. Westerveld, T. Ianeva, L. Boldareva, A.P. de Vries, and D. Hiemstra. Combining information sources for video retrieval. In *Text Retrieval Conference TRECVID Workshop*. National Institute of Standards and Technology, Nov 2003. URL <http://www-nlpir.nist.gov/projects/tv2003/tv2003.html>.

List of Figures

1	All pairs and truncated conditional probabilities for \mathbf{M}^t , using image queries (left) and the text query (right).	15
2	All pairs and truncated conditional probabilities using \mathbf{M}^b , and Rui's relevance feedback model. Image queries used.	16
3	Display update strategies for \mathbf{M}^b matrix, with the text query. Proportional sampling is not plotted.	16
4	Combinations of different information sources, using image queries (left) and the text query (right).	17
5	Automated experiments before and after training on real user feedback, using complete 2003 TREC Video test collection. . . .	19