

CWI at the TREC-2002 video track

Thijs Westerveld, Arjen de Vries, Alex van Ballegooij
CWI
PO Box 94079, 1090 GB Amsterdam
The Netherlands
{thijs, arjen, alexb}@cwi.nl

1 Introduction

We present a probabilistic model for the retrieval of multimodal documents. The model is based on Bayesian decision theory and combines models for text based search with models for visual search. The textual model, applied to the LIMSI transcripts, is based on the language modelling approach to text retrieval. The visual model, a mixture of Gaussian densities, describes keyframes selected from shots. Both models have proved successful on media specific retrieval tasks. Our contribution is the combination of both techniques in a unified model, ranking shots on ASR-data and visual features simultaneously.

Using this model, we tried to answer the following questions.

- Is it useful to identify important parts in query images?
- Can using (additional) query images from outside the search collection¹ help improve retrieval results?
- Does it help to have *multiple* image examples for a query, or are we better of using only *one* good example?
- Can a combination combined textual and visual query perform better than queries in a single modality?

¹Throughout this document, we refer to the search collection used in the TREC-2002 video track as *the search collection*.

Because of problems with the similarity measure we used in the submitted runs, we mainly report on post-hoc experiments on the TREC-2002 data in which we used a different measure. Both measures are discussed in Section 2, where we also present our retrieval model. Section 3 reports on the post-hoc experiments and Section 4 summarises our main findings. The official results can be found in appendix A.

2 Probabilistic Multimedia Retrieval

In a probabilistic retrieval setting, the goal is to find the document D^* with highest probability given a query Q :

$$D^* = \arg \max_i P(D_i|Q) = \arg \max_i \frac{P(Q|D_i)P(D_i)}{P(Q)} \quad (1)$$

Usually, (1) is used as a scoring function and a ranked list is returned rather than the one most probable document.

If we assume that all documents have equal prior probability, (1) reduces to the maximum likelihood (ML) criterion, which is approximated by the minimum KL-divergence between query model and document model: $D^* = \arg \min_i \text{KL}[P_q(\mathbf{x})||P_i(\mathbf{x})]$.

$$\begin{aligned} \text{KL}[P_q(\mathbf{x})||P_i(\mathbf{x})] &= \int P(\mathbf{x}|D_q) \log \frac{P(\mathbf{x}|D_q)}{P(\mathbf{x}|D_i)} d\mathbf{x} \\ &= \int P(\mathbf{x}|D_q) \log P(\mathbf{x}|D_q) d\mathbf{x} - \int P(\mathbf{x}|D_q) \log P(\mathbf{x}|D_i) d\mathbf{x}, \end{aligned}$$

where \mathbf{x} are feature vectors describing the documents.

The first integral is independent of D_i and can be ignored, thus

$$D^* = \arg \max_i \int P(\mathbf{x}|D_q) \log P(\mathbf{x}|D_i) d\mathbf{x} \quad (2)$$

Now suppose query and document models generate a mixture of textual features \mathbf{x}_t and visual features \mathbf{x}_v :²

$$P(\mathbf{x}|D_i) = P(\mathbf{x}_t|D_i)P(t) + P(\mathbf{x}_v|D_i)P(v).$$

We can then integrate over these different feature sets separately and arrive at the following ranking formula for multimodal retrieval [5].

$$D^* = \arg \max_i [P(t) \int_{\mathbf{x}_t} P(\mathbf{x}_t|D_q) \log P(\mathbf{x}_t|D_i) d\mathbf{x}_t + P(v) \int_{\mathbf{x}_v} P(\mathbf{x}_v|D_q) \log P(\mathbf{x}_v|D_i) d\mathbf{x}_v] \quad (3)$$

2.1 Text Model

To describe the probability distributions of the textual terms, we take a language modelling approach to information retrieval [2]. Such a model operates on discrete signals (i.e. words), thus we can replace the integral from (3) by a sum. Moreover, the query model D_q is usually nothing more than the empirical distribution of the query, therefore we only need to sum over the words in the query. The document model is usually taken to be a mixture of foreground ($P(x_{t,j}|D_i)$) and background ($P(x_{t,j})$) probabilities for the query terms $x_{t,j}$, interpolated using mixing parameter λ (cf. Section 2.1.1). If our textual query consists of N_t terms $\mathbf{x}_t = (x_{t,1}, x_{t,2}, \dots, x_{t,N_t})$ then the textual part of our ranking formula is the following.

$$D_i^* = \arg \max_i \frac{1}{N_t} \sum_{j=1}^{N_t} \log [\lambda P(x_{t,j}|D_i) + (1 - \lambda)P(x_{t,j})] \quad (4)$$

² $P(t)$ and $P(v)$ are the prior probabilities of drawing respectively textual or visual features from a document; assumed uniform across documents.

Using the statistical language modelling approach for video retrieval, we would like to exploit the hierarchical data model of video, in which a video is subdivided in scenes, which are subdivided in shots, which are in turn subdivided in frames. Statistical language models are particularly well-suited for modelling such complex representations of the data. We can simply extend the mixture to include the different levels of the hierarchy, with models for shots and scenes.³

$$\text{Shot}^* = \arg \max_i \frac{1}{N_t} \sum_{j=1}^{N_t} \log [\lambda_{\text{Shot}} P(x_{t,j}|\text{Shot}_i) + \lambda_{\text{Scene}} P(x_{t,j}|\text{Scene}_i) + \lambda_{\text{Coll}} P(x_{t,j})] \quad (5)$$

with $\lambda_{\text{Coll}} = 1 - \lambda_{\text{Shot}} - \lambda_{\text{Scene}}$

The main idea behind this approach is that a good shot contains the query terms and is part of a scene having more occurrences of the query terms. Also, by including scenes in the ranking function, we hope to retrieve the shot of interest, even if the video's speech describes it just before it begins or just after it is finished. Depending on the information need of the user, we might use a similar strategy to rank scenes or complete videos instead of shots, that is, the best scene might be a scene that contains a shot in which the query terms (co-)occur.

2.1.1 Estimating Parameters

The features in the textual part of our model are simply the words themselves. For the textual part of our retrieval function (5), we only need to estimate foreground ($P(x_{t,j}|D_i)$) and background ($P(x_{t,j})$) probabilities. Both measures are estimated in the standard way, by taking the term frequency and document frequency respectively [2]. We used the TREC-2002 video search collection to find the optimal values for the mixing parameters: $\lambda_{\text{Shot}} = 0.090$, $\lambda_{\text{Scene}} = 0.210$, and $\lambda_{\text{Coll}} = 0.700$. Since we trained these parameters on the test collection, we cannot say anything about how well these numbers generalise across

³We assume each shot is a separate class and replace ω_i with Shot_i .

collections.⁴ Yet, for each of the mixing parameters, there is quite a large range of values for which the scores are close to optimal. In this work we do not look into the stability of these parameters across collections, we are only interested in finding the optimal settings for this collection and evaluating the retrieval model with these optimal settings.

2.2 Image Model

We use a Gaussian Mixture Model for describing document densities [4].

$$P(\mathbf{x}_v|D_i) = \sum_{c=1}^C P(\theta_{i,c}) \mathcal{G}(\mathbf{x}_v, \boldsymbol{\mu}_{i,c}, \boldsymbol{\Sigma}_{i,c}),$$

where C is the number of components in the mixture model, $\theta_{i,c}$ is component c of document model D_i and $\mathcal{G}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the Gaussian density with mean vector $\boldsymbol{\mu}$ and co-variance matrix $\boldsymbol{\Sigma}$:

$$\mathcal{G}(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^n |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2} \|\mathbf{x} - \boldsymbol{\mu}\|_{\boldsymbol{\Sigma}}}, \quad (6)$$

$$\text{where } \|\mathbf{x} - \boldsymbol{\mu}\|_{\boldsymbol{\Sigma}} = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$$

and n is the length of the feature vector \mathbf{x} .

2.2.1 Bags of Blocks

Just like in our textual approach, for the query model, we can simply take the empirical distribution of the query samples. If a query-image \mathbf{x}_v consists of N_v samples: $\mathbf{x}_v = (x_{v,1}, x_{v,2}, \dots, x_{v,N_v})$ then $P(x_{v,i}|D_q) = \frac{1}{N_v}$. For the document model, we take a mixture of foreground and background probabilities, i.e. the (foreground) probability of drawing a query sample from the document's Gaussian mixture model, and the (background) probability of drawing it from any Gaussian mixture in the collection. In other words, the query image is viewed as a bag of blocks (BoB), and its probability is estimated as the joint probability of all its blocks. The BoB measure

⁴Obviously, the official runs have used different mixing parameter values. see Appendix A.

for query images then becomes:

$$D_v^* = \arg \max_i \frac{1}{N_v} \sum_{j=1}^{N_v} \log [\kappa P(x_{v,j}|i) + (1 - \kappa) P(x_{v,j})], \quad (7)$$

where κ is a mixing parameter and the background probability $P(x_{v,j})$ can be found by marginalising over all M documents in the collection:

$$P(x_{v,j}) = \sum_{i=1}^M P(x_{v,j}|D_i) P(D_i).$$

Again we assume uniform document priors ($P(D_i) = \frac{1}{M}$ for all i). In text retrieval, one of the reasons for mixing the document model with a collection model is to assign non-zero probabilities to words that are not observed in a document. Smoothing is not necessary in the visual case, since the documents are modelled as mixtures of Gaussians, having infinite support. Another motivation for mixing is to weight term importance: a common sample \mathbf{x} (i.e., a sample that occurs frequently in the collection) has a relatively high probability $P(\mathbf{x})$ (equal for all documents), and therefore $P(\mathbf{x}|D)$ has only little influence on the probability estimate. In other words, relatively common terms and common blocks influence the final ranking only marginally.

2.2.2 Asymptotic Likelihood Approximation

A disadvantage of using the BoB measure is its computational complexity. In order to rank the collection given a query, we need to compute the posterior probability $P(\mathbf{x}_v|\omega_i)$ of each image block \mathbf{x}_v in the query for each document ω_i in the collection. For evaluating a retrieval method this is fine, but for an interactive retrieval system, optimisation is necessary.

An alternative is to represent the query image, like the document image, as a Gaussian model (instead of by its empirical distribution as a bag of blocks), and then compare these two models using the KL-divergence. Yet, if we use Gaussians to model the class conditional densities of the mixture components, there is no closed-form solution for the visual

part of the resulting ranking formula (3). As a solution, Vasconcelos assumes that the Gaussians are well separated and derives an approximation, ignoring the overlap between the mixture components: the asymptotic likelihood approximation (ALA) [4]. The ALA is the measure we used in our official TREC-2002 runs, (see Appendix A). However, in post hoc analysis, we found that one of the assumptions underlying the ALA is not plausible for the collection at hand and, moreover, using it decreases performance compared to the BoB measure (for details see [5]). In the remainder of this work we will concentrate on the BoB measure.

2.2.3 Estimating Parameters

For estimating the parameters of the Gaussian mixture model, we used the EM algorithm [1]. We described a document as a set of samples, where each sample is described by a number of DCT coefficients in the YCbCr colour space⁵. Then we used EM to fit a mixture of 8 Gaussian (for details see [5]). Finally, we described the position in the image plane of each component as a 2D-Gaussian with mean and covariance computed from the positions of the samples assigned to this component. We evaluated different values for mixing parameter κ on the TREC-2002 video search collection and found the optimal value: $\kappa = 0.9$.

3 Experiments

Fully automatic creation of queries from topic descriptions was not required in this year’s video track. However, there was a distinction between manual and interactive runs. In an interactive run a user can interact with a system to locate relevant shot. In a manual run a user has one go at creating a query from a topic descriptions and then submits this query to the system to retrieve relevant shots. All our runs are manual runs in which we experimented with different ways of creating queries from topic statements. In

⁵We use the first 10 coefficients from the Y channel and the two DC coefficients of the Cb and the Cr channels.

the following subsections, we investigate the following questions:

- Is it useful to identify important parts in query images?
- Can using (additional) query images from outside the search collection help improve retrieval results?
- Does it help to have *multiple* image examples for a query, or are we better off using only *one* good example?
- Can a combination combined textual and visual query perform better than queries in a single modality?

3.1 Selecting Query Images

In general, it is hard to guess what would be a good example image for a specific query. If we look for shots of the *Golden Gate bridge*, we might not care from what angle the bridge was filmed, or if the clip was filmed on a sunny or a cloudy day; visually however, such examples may be very different (Figure 1). If a user has presented three examples and no additional information, the best we can do is try to find documents that describe all example images well. Unfortunately, a document may be ranked low even though it models the samples from one example image well, as it may not explain the samples from the other images.

For each topic, we computed which of the example images would have given the best results if it had been used as the only example for that topic. We compared these *best example* results to the *full topic* results in which we used all available visual examples. In the *full topic* case, the set of available topics was regarded as one large bag of blocks. We ranked documents by their probability of generating all blocks in all query images. For the single image queries in the *best example*, we used all samples from the single visual example to rank documents.

Since it is problematic to use multiple examples in a query, we wanted to see if it is possible to guess in advance what would be a good example for a specific

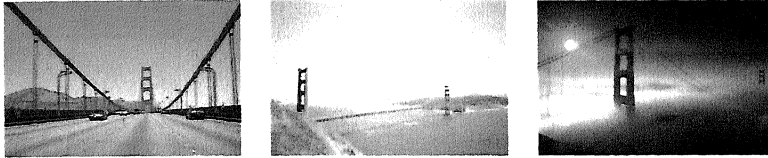


Figure 1: Visual examples of the Golden Gate bridge.

topic. Therefore, we hand-picked for each topic a single representative from the available examples and compared these *manual example* results to the other two result sets.

The results for the different settings are listed in Table 1. A first thing to notice is that all scores are rather low. When we take a closer look at the topics with higher average precision scores, we see that these mainly contain examples from the search collection. In other words, we can find similar shots from within the same video, but generalisation is a problem.

The fact that using the best image example outperforms the use of all examples shows that indeed combining results from different visual examples can degrade results. Looking at the results, manually selecting good examples seems a non-trivial task, but the drop in performance is partly due to the generalisation problem. If one of the image examples happens to come from the collection it scores high. If we fail to select that particular example, the score for the manual example run drops. Simply counting how often the manually selected example was the same as the best performing example, we see that this was the case for 8 out of 13 topics.⁶

3.2 Selecting Important Regions

In last year’s video track, we saw that query articulation, i.e. the manual identification of important parts in a query image, can help improve retrieval results [3]. We also noticed that this requires an enormous effort from a user. In our probabilistic setting, selecting important (and coherent) regions is much

	full topic	best example	manual example
vt075	.0038	.2438	.2438
vt076	.4854	.4323	.1760
vt077	.0000	.0000	.0000
vt078	.0000	.0000	.0000
vt079	.0000	.0040	.0000
vt080	.0048	.0977	.0977
vt081	.0000	.0000	.0000
vt082	.0330	.0234	.0234
vt083	.0000	.0000	.0000
vt084	.0046	.0046	.0046
vt085	.0000	.0000	.0000
vt086	.0053	.0704	.0704
vt087	.0000	.0000	.0000
vt088	.0046	.0069	.0069
vt089	.0000	.0000	.0000
vt090	.0000	.0305	.0305
vt091	.0095	.0095	.0095
vt092	.0003	.0106	.0000
vt093	.0006	.0006	.0000
vt094	.0021	.0021	.0021
vt095	.0000	.0000	.0000
vt096	.0323	.0323	.0323
vt097	.1312	.1408	.0000
vt098	.0000	.0003	.0003
vt099	.0000	.0000	.0000
MAP	.0287	.0444	.0279

Table 1: MAP for Full Topics, Best Examples and Manual Examples

⁶We ignored the topics for which there is only one example and the ones for which none of the examples retrieved relevant documents.

easier. After building a query-model (like we build document models) a user can simply select one or more meaningful components from the query-model. In retrieval, we can then use only the Bag of Blocks corresponding to the selected component(s). For example in Figure 2a, we selected the components that together form the US flag. Similarly, we can indicate we want multiple parts to be present in the target shots, e.g. *boat* and *water* and *sky* (Figure 2b). Note that even though the union of the sets of samples is in this case the full image, this differs from simply taking the using all samples as a query. If the full image were used, we would have looked for shots with relatively few water samples; the selection of components compensates for that and looks for documents that explain all 3 concepts equally well.

From each of the query images, we selected meaningful components and we used the corresponding samples as queries. If we take a look at the individual components and their results, we see that the components are often homogeneous in colour and/or texture and that results are often meaningful (Figure 3) or, if there is little semantics in the component, at least visually similar (Figure 4). It is not clear yet how this can be used for highly specific queries like the video track queries.

3.3 Using Query Examples from Outside the Collection

In Section 3.1, we argued that selecting the right query image is important. On the one hand therefore, one would like to expand a query to have as many different query images as possible. On the other hand, we saw that it is difficult to combine multiple examples in one query (Section 3.1). We investigate whether using (additional) examples from outside the collection can improve retrieval effectiveness. We expect that this is not the case; in previous experiments [5, 3] we saw that we can only find relevant shots if the query images are highly similar to the relevant shots, i.e. if they are from the same collection and preferably from the same video.

First of all, we had a look at the original examples provided by NIST. Most, if not all, of the video examples in this set come from either the search col-

lection itself, or the highly comparable⁷ feature train or feature test set. We found that the topics that contributed most to our MAP score were the ones with examples from the search collection. If we remove videos from which shots are used as examples from the relevance judgements, our MAP score for a purely visual run (using full examples for all queries) drops from .0287 to .0029; purely visual runs from other groups show a similar drop in performance. This indicates that visual retrieval systems are able to locate the query examples in the collection, but generalisation seems problematic. Furthermore, the best examples as reported in table 1 are mainly video examples from either the search collection or the comparable training data. Only for three topics, the best scoring example was an image example from outside these collections. Yet, for these three topics no video examples were available.

We experimented with query expansion by adding additional example images found using Google image search⁸. We manually created short queries from the topic descriptions and submitted these to Google image search. From the result list we selected images based that we thought were good examples for the topic. This way we expanded topics with up to 7 additional image examples. We ran these new examples as queries against the collection and recomputed the best scoring examples for each topic. For 5 out of 25 topics none of the examples retrieved any relevant documents. The best scoring examples for the remaining 20 topics were video examples in 12 cases and image examples from Google in 8 cases. Clearly, if we try more examples we have a better chance of having a good example among them, yet the problem remains how to combine multiple examples or how to identify a good example without knowledge of the relevant documents in the collection.

3.3.1 Combining Textual and Visual runs

We combined textual and visual runs using our combined ranking formula (3). Since we had no data to estimate the parameters for mixing textual and visual information we used $P(t) = P(v) = 0.5$. For the

⁷In fact, these are distinct subsets of one larger collection.

⁸<http://images.google.com>

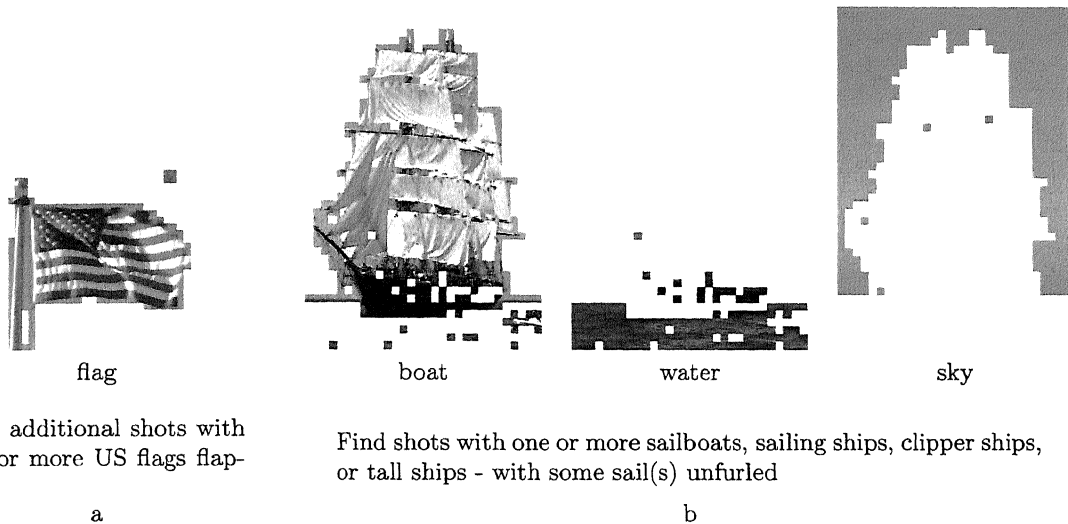


Figure 2: Selecting components from images

textual part we tried both short and long queries, for the visual part we used full queries and best-example queries. Table 2 shows the results for combinations with the BoB measure. We also experimented with combinations with the ALA measure, but we found that in the ALA case it is difficult to combine textual and visual scores, because they are on different scales (see also Appendix A). The BoB measure is closer to the KL-divergence and, on top of that, more similar to our textual approach, and thus easier to combine with the textual scores.

For most of the topics, textual runs give the best results, however for some topics using the visual examples is useful. This is mainly the case when either the topics come from the search collection or when the relevant documents are outliers in the collection. This illustrates how difficult it is to search a generic video collection using visual information only. We only succeed if the relevant documents are either highly similar to the examples provided or very dissimilar from the other documents in the collection (and therefore relatively similar to the query examples). When both textual and visual runs have reasonable scores, combining the runs can improve on the individual runs, however, when one of them has

inferior performance, a combination only adds noise and lowers the scores.

4 Conclusions

We presented a probabilistic framework for multi-modal retrieval in which textual and visual retrieval models are integrated seamlessly and evaluated the framework using the search task from the TREC-2002 video track. We found that even though the topics were specifically designed for content-based retrieval, and relevance was defined visually, a textual search outperforms visual search for most topics. The main conclusion in this work is that visual retrieval using the presented model for specific queries does not generalise very well. The model could retrieve shots that are highly similar to the query examples (i.e. shots from the same video), but other similar shots were found mostly by coincidence, because they happened to have for example the same colour sky or grass. For more general queries, the model seems useful. When we select a single component from an example, results are intuitive, i.e. visually similar. It is unclear how this helps in retrieving relevant documents for highly

Topic	Tshort	Tlong	BoBfull	BoBbest	BoBfull +Tshort	BoBfull +Tlong	BoBbest +Tshort	BoBbest +Tlong
vt075	.0000	.0082	.0038	.2438	.0189	.0569	.2405	.3537
vt076	.4075	.6242	.4854	.4323	.5931	.7039	.5757	.6820
vt077	.1225	.5556	.0000	.0000	.0000	.0000	.0000	.0000
vt078	.1083	.2778	.0000	.0000	.0000	.0000	.0000	.0000
vt079	.0003	.0006	.0000	.0040	.0003	.0000	.0063	.0050
vt080	.0000	.0000	.0048	.0977	.0066	.0059	.0845	.0931
vt081	.0154	.0333	.0000	.0000	.0037	.0000	.0000	.0000
vt082	.0080	.0262	.0330	.0234	.0181	.0335	.0145	.0210
vt083	.1669	.1669	.0000	.0000	.0962	.0962	.0078	.0078
vt084	.7500	.7500	.0046	.0046	.6875	.6875	.6875	.6875
vt085	.0000	.0000	.0000	.0000	.0000	.0000	.0000	.0000
vt086	.0554	.0676	.0053	.0704	.0536	.0215	.0791	.0600
vt087	.0591	.0295	.0000	.0000	.0052	.0003	.0052	.0003
vt088	.0148	.0005	.0046	.0069	.0052	.0046	.0069	.0069
vt089	.0764	.0764	.0000	.0000	.0503	.0503	.0045	.0045
vt090	.0229	.0473	.0000	.0305	.0006	.0075	.0356	.0477
vt091	.0000	.0000	.0095	.0095	.0000	.0086	.0000	.0086
vt092	.0627	.0687	.0003	.0106	.0191	.0010	.0078	.0106
vt093	.1977	.1147	.0006	.0006	.0099	.0021	.0071	.0012
vt094	.0232	.0252	.0021	.0021	.0122	.0036	.0122	.0036
vt095	.0034	.0021	.0000	.0000	.0008	.0012	.0011	.0010
vt096	.0000	.0000	.0323	.0323	.0161	.0161	.0323	.0323
vt097	.1002	.0853	.1312	.1408	.1228	.1752	.1521	.1474
vt098	.0225	.0086	.0000	.0003	.0068	.0000	.0004	.0003
vt099	.0726	.0606	.0000	.0000	.0000	.0000	.0000	.0000
MAP	.0916	.1212	.0287	.0444	.0691	.0750	.0784	.0870

Table 2: Average precision per topic, for Textual runs, BoB runs and combined runs

Q:



Top 5:



Figure 3: Top 5 results for a homogeneous query with clear semantics ('Sky')

specific topics like the video track topics, but it helps in gaining insight in the models performance. In future work, we will further investigate the influence of individual components on retrieval results. In addition, we intend to look at how incorporating different sources of additional information (e.g. contextual frames, the movement in video or user interaction) can help improve results across collections. Combining multiple examples in one query is still problematic, but combining textual and visual runs seems possible using the presented framework. When one of the runs is poor, a combined run, including the noise, is less effective than the single best run. However, when the individual runs have reasonable scores, combining them improves retrieval effectiveness.

References

- [1] A. Dempster, N. Laird, and D. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society, series B*, 39(1):1–38, 1977.
- [2] D. Hiemstra. *Using language models for information retrieval*. PhD thesis, Centre for Telematics and Information Technology, University of Twente, 2001.

- [3] The Lowlands team. Lazy users and automatic video retrieval tools in (the) lowlands. In *The 10th Text Retrieval Conference (TREC-2001)*, 2002.
- [4] N. Vasconcelos. *Bayesian Models for Visual Information Retrieval*. PhD thesis, Massachusetts Institut of Technology, 2000.
- [5] T. Westerveld, A. P. de Vries, A. 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*, 2003.

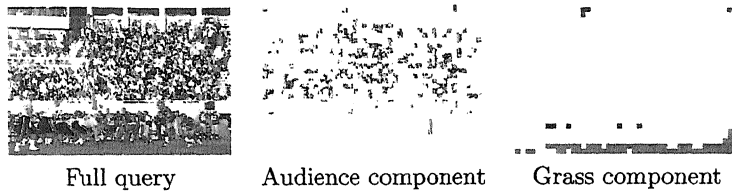
A Official Results

In the official runs, we used the Asymptotic Likelihood Approximation (see section 2.2.2). We distinguished between the NIST images (the visual examples from the official topics) and Google images (additional examples we found with manual query expansion using Google and submitted four runs:

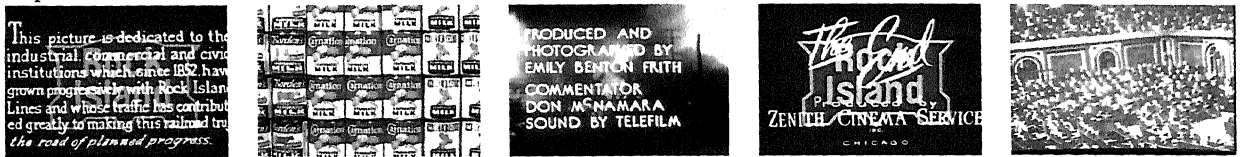
run1 Text only.

run2 Text + NIST images.

run3 Text + selected components from NIST images.



Top 5 audience:



Top 5 grass:



Figure 4: Top 5 results for homogeneous queries with unclear or false semantics

run4 Text + selected components from both NIST properly. and Google images.

The text model’s mixing parameters have been optimized using the TREC-2001 corpus, giving $\lambda_{\text{shot}} = 0.015$, $\lambda_{\text{scene}} = 0.135$, and $\lambda_{\text{coll}} = 0.850$. For run3 and run4, we manually selected *important* components from the query model (cf. Section 3.2). In all runs that involved visual examples, we computed a single new (8 component) Gaussian mixture model from all available visual blocks and we used that model in our ALA ranking formula. The results for the official runs and for the same runs after fixing some bugs⁹ are shown in Table 3. We see that also with the ALA measure text only results are by far the best (run 1). Combinations that also used visual information scored lower, not only on MAP, but also on average precision for each individual topic. In contrast to our findings with the BoB measure we were not able to combine textual and visual information

runName	MAP
run1	.0917
run2	.0016
run3	.0022
run4	.0038
run1 fixed	.1212
run2 fixed	.0082
run3 fixed	.0137
run4 fixed	.0069

Table 3: Official results and same runs after bug fix.

⁹A normalisation error in the training of the models and exchanging a few videos from the search and feature detection collections.