

## PDF hosted at the Radboud Repository of the Radboud University Nijmegen

The following full text is a publisher's version.

For additional information about this publication click this link.

<https://hdl.handle.net/2066/228295>

Please be advised that this information was generated on 2021-11-02 and may be subject to change.

# Combining Information Sources for Video Retrieval

The Lowlands Team at TRECVID 2003

Thijs Westerveld<sup>◊</sup> Tzvetanka Ianeva<sup>†,\*</sup>, Lioudmila Boldareva<sup>‡</sup>, Arjen P. de Vries<sup>◊</sup>  
and Djoerd Hiemstra<sup>‡</sup>

<sup>◊</sup>Centrum voor Wiskunde en  
Informatica (CWI)  
Amsterdam, The Netherlands  
{thijs,arjen}@cwi.nl

<sup>†</sup>Departament  
d'Informàtica  
Universitat de València  
València, Spain  
tzveta.ianeve@uv.es

<sup>‡</sup>Department of Computer  
Science  
University of Twente  
Enschede, The Netherlands  
{boldarli,hiemstra}@cs.utwente.nl

## Abstract

The previous video track results demonstrated that it is far from trivial to take advantage of multiple modalities for the video retrieval search task. For almost any query, results based on ASR transcripts have been better than any other run. This year's main success in our runs is that a combination of ASR and visual performs better than either alone! In addition we experimented with dynamic shot models, combining topic examples, feature extraction and interactive search.

## 1 Introduction

Often it is stated that a successful video retrieval system should take advantage of information from all available sources and modalities. Merging knowledge from for example speech, vision, and audio would yield better results than using only one of them. But previous video track results demonstrated that it is far from trivial to take advantage of multiple modalities for a video retrieval search task.

For this year's TRECVID workshop, we experimented with combining different types of informa-

tion:

- combining different models/modalities
- combining multiple example images
- combining model similarity and human-judged similarity

The basic retrieval models we use to investigate the merging of different sources are the same models we used last year [12], they are described briefly in section 2 and more extensively in [13]. Section 2.1 describes a dynamic variant of the model that allows for describing spatio-temporal information as opposed to the spatial information captured in the basic models. These dynamic models can describe shots instead of still keyframes. For modelling the ASR transcripts, we use a basic language modelling approach (see Section 2.2). The interactive retrieval setup (Section 2.3) builds upon the visual models and language models. After the introduction of the separate models, we present experiments for combining textual and visual information (Section 4), combining different examples (Section 5) and the feature task (Section 6). Combining automatic similarity and interactive similarity is discussed in Section 7.

---

\*The work has been carried out while at CWI, supported by grants GV FPIB01.362 and CTESRR/2003/43

## 2 Generative Probabilistic Retrieval Model

The retrieval model we use to rank video shots is a generative model inspired by the language modelling approach to information retrieval [9, 6] and a similar probabilistic approach to image retrieval [10]. We present – concisely – the visual part of the model, referring the interested reader to [13] for more details. The visual model ranks images by their probability of generating the samples (pixel blocks) in one or more query example(s) [13]. The model is smoothed using background probabilities, calculated by marginalisation over the collection. So, a collection image  $\omega_i$  is compared to an example image  $\mathbf{x}$  consisting of  $N$  samples ( $\mathbf{x} = (x_1, x_2, \dots, x_N)$ ) by computing its ability to explain the samples  $\mathbf{x}$ . The retrieval status value (RSV) of an image  $\omega_i$  is defined as:

$$\text{RSV}(\omega_i) = \frac{1}{N} \sum_{j=1}^N \log [\kappa P(x_j|\omega_i) + (1 - \kappa)P(x_j)], \quad (1)$$

where  $\kappa$  is a mixing parameter.

Collection images  $\omega_i$  are modelled as mixtures of Gaussians with a fixed number of components  $C$  [10, 13]:

$$P(\mathbf{x}|\omega_i) = \sum_{c=1}^{N_C} P(C_{i,c}) \mathcal{G}(\mathbf{x}, \boldsymbol{\mu}_{i,c}, \boldsymbol{\Sigma}_{i,c}), \quad (2)$$

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

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

where  $n$  is the dimensionality of the feature space and  $(\mathbf{x} - \boldsymbol{\mu})^T$  is the matrix transpose of  $(\mathbf{x} - \boldsymbol{\mu})$ .

The samples are 8 by 8 pixel blocks, described by their DCT coefficients and their position in the image plane; the models are trained on these using standard EM [3], assuming a diagonal co-variance matrix.

### 2.1 Dynamic Model

The Dynamic model is a Gaussian Mixture Model in DCT-space-time domain [7]. Instead of modelling just a single image (keyframe) we model a one-second video sequence around the keyframe as a single entity. The dynamic model is an extension of the static model, the sampling process is very similar. We take 29 frames around the keyframe and cut them in distinct blocks of 8 by 8 pixels. Each block is then described by its DCT coefficients, its x and y position in the image plane and its position in time (normalised between 0 and 1). Given this setup, the static model can be seen as a special case of the dynamic model where the temporal feature takes a fixed value of 0.5. The intuitive explanation in this case is: we do not know what is happening before and after the central time moment matching the keyframe.

The number of samples for training the dynamic model is much larger than for the static model (i.e. 29 times as large). The training process remains the same: feature vectors are fed to the EM algorithm to find the parameters  $\mu_c$ ,  $\Sigma_c$ , and  $P(C_c)$ . Because we use diagonal covariance matrices ( $\Sigma_c$ ), the components are aligned to the axes. The resulting models capture the appearance and disappearance of objects, but not all temporal information. Figure 1 shows an example video sequence and a visualisation of the resulting model. It is easy to see how well the model fits the data, the *tree* in the top left corner is only visible at the beginning of the sequence. The corresponding component in the model also disappears at about  $t = 0.5$ . In other words, the GMM captures the disappearing of the tree. This effect is impossible to be captured from the static model where time is not taken into account.

### 2.2 ASR

The ASR approach used the same hierarchical language model as last year [12]. We model video as a sequence of scenes, each consisting of a sequence of shots. The generative model mixes the different levels of the hierarchy, with models for shots and scenes. Given a query with  $N_t$  terms  $\mathbf{q} = (q_1, q_2, \dots, q_{N_t})$ , the RSV of a shot  $\omega_i$  is defined as:

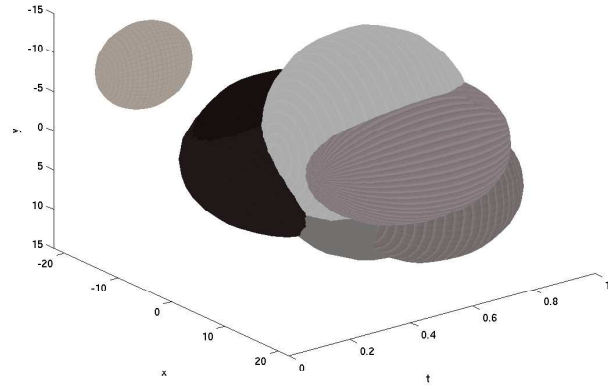


Figure 1: A shot represented by 29 frames around the keyframe (top) and a 3D visualisation of dynamic GMM computed from it (bottom). The bottom image shows mean colour and mean texture of the components where the standard deviation from the mean position in the  $x-y-t$  is below 2. Prior probabilities and variance in colour and texture are not visualised.

$$\text{RSV}(\omega_i) = \frac{1}{N_t} \sum_{j=1}^{N_t} \log[\lambda_{\text{Shot}} P(q_j | \text{Shot}_i) + \lambda_{\text{Scene}} P(q_j | \text{Scene}_i) + \lambda_{\text{Coll}} P(q_j)]$$

with  $\lambda_{\text{Coll}} = 1 - \lambda_{\text{Shot}} - \lambda_{\text{Scene}}$  (3)

$\text{Shot}_i$  and  $\text{Scene}_i$  are the shot and scene to which  $\omega_i$  belongs. 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 more 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. Because scene boundaries are not known, we assume (pragmatically) that each sequence of 5 consecutive shots forms a scene.

The features in the ASR model are simply the word tokens from the LMSI transcripts [5]. We estimate the foreground ( $P(q_j | \text{Shot}_i)$ ) and background ( $P(q_j)$ ) probabilities in Equation 3 by taking the term frequency and document frequency respectively [6]. 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  $\lambda_{\text{Coll}} = 0.700$ .

### 2.3 Interactive retrieval

In an interactive setting, the models above might be used as follows: First the user enters a text query that uses the speech recognition transcripts. As a result, he/she gets a screen full of video key frames. The user might now: *a)* rephrase his/her query, *b)* walk down the ranked list, i.e. look at the next screen of key frames, or *c)* use relevance feedback, i.e. retrieve key frames that are similar to one or more of the key frames currently on the screen. The text queries (see Section 2.2) can be evaluated almost instantly by the system. Whenever the system has to find similar images, however, the Gaussian mixture models of Equation 1 are used. It is well-known that similarity search using low level features does not scale very well. Given one example image, ranking the 32,000 key frames of the collection would take about 15 minutes on a fairly fast personal computer. Ranking 3.2 million key frames would take approximately 100 times as long.

To get reasonable on-line performance, we used a brute-force solution by pre-computing for each key frame  $\mathbf{x}$  its nearest neighbours  $\omega_i$  (see e.g. [4]), according to the Gaussian mixture models. In practice, this means we have to calculate all image-image similarity pairs, and then take those image pairs that are nearest neighbours. The image pairs  $(\mathbf{x}, \omega_i)$  are stored together with the probability  $P(\mathbf{x}|\omega_i)$  defined by Equation 1 in an access structure called *associative matrix*. Using the associative matrix, we avoid expensive run time distance computations [1]. To search for images that are similar to two or more example images  $\mathbf{x}_j$ , the model’s probabilities  $P(\mathbf{x}|\omega_i)$  are combined by assuming conditional independence between selected images  $\mathbf{x}_1, \mathbf{x}_2, \dots$  given the hypothesised target image  $\omega_i$ :

$$P(\mathbf{x}_1, \mathbf{x}_2, \dots | \omega_i) = \prod_{j=1}^n P(\mathbf{x}_j | \omega_i) \quad (4)$$

For images that are not among the nearest neighbours (that is, for which  $P(\mathbf{x}_j|\omega_i)$  is not stored in the associative matrix), we use a back-off mechanism, effectively replacing  $P(\mathbf{x}_j|\omega_i)$  with a smoothing constant  $\bar{p}$ .

Besides efficient real time access to similar images, the associative matrix serves another purpose: It can be trained from user interaction. By collecting user feedback, we hope to adapt the matrix in such a way that it represents user’s associations between images. A possible interpretation of an image pair  $(\mathbf{x}_j, \omega_i)$  in the matrix is the following: “Users that were looking for image  $\omega_i$  selected image  $\mathbf{x}_j$ ”<sup>1</sup>, or more precisely, 80 % of the users that sought image  $\omega_i$  selected image  $\mathbf{x}_j$  as relevant if  $P(\mathbf{x}_j|\omega_i) = 0.8$ .

## 3 Experimental Setup

### 3.1 Building Models and Queries

The search collection is indexed using the procedures described above. For each shot, we build a static model a dynamic model and a Language Model.

<sup>1</sup>As in collaborative filtering used by e.g. Amazon.com: “Customers who bought this book also bought:...”

Building queries from the topic descriptions is mostly automatic. The only manual action in constructing visual queries was selecting one or more image or video examples to be used for ranking. A textual query was constructed manually for each topic by taking only the content words from the topic description. From there on, the whole retrieval process was automatic, except for the experiments described in section 5.1. All image examples are rescaled to at most 272x352 pixels and then jpg compressed with a quality level of 20% to match size and quality of the collections videos. Further experiments are needed to test if this scaling and degrading of queries influences retrieval results. The interactive experiments only used textual queries; not the image examples from the topics.

### 3.2 Guarding Covariance

In estimating GMMs from images using EM, it sometimes happens that one of the components collapses onto a single point in features space. This means, the covariance of such a component tends to zero and the component explains nothing except for this particular point in feature space. For example in the TRECVID dataset, it happens that a component collapses onto perfectly black or perfectly white blocks. In retrieval such singular solutions can cause underflow division-by-zero problems and queries that happen to have samples that fall into the exact center of a singular component (e.g. queries with perfectly black blocks), will get infinity scores. To work around this, in our official submissions, we set the prior probability of components with a covariance smaller than some threshold to zero. Effectively, this means we ignored these components during retrieval. An alternative, better, solution is to guard the covariances during training of the models, and to make sure they do not get too small. This is what we did for the static model in post hoc experiments. In the following, both the official scores and the post hoc scores are listed, they are labeled *official* and *COVguard* respectively.

Description	MAP official	MAP COVguard
ASR only	.130	
static only	.022	.025
static + ASR	.105	.134
dynamic only	.022	
dynamic + ASR	.132	
ARS + noAnchor	.133	
static + noAnchor	.022	
static + ASR + noAnchor	.106	
dynamic + noAnchor	.022	
dynamic + ASR + noAnchor	.135	

Table 1: Mean Average Precision (MAP) for ASR only, visual only and combined runs (static and dynamic models).

## 4 Combining Modalities

If we (unrealistically) assume textual and visual information are independent, we can easily compute a joint probability of observing query text and visual example from a document. Under this assumption, we can simply multiply the individual probabilities (or sum the log probabilities from Equations 1 and 3). In previous TREC experiments, we found that such a combination strategy works well provided that the individual runs each have useful results [12]. This year, this seems to be the case: a combination of the dynamic models described in Section 2.1 and ASR (Section 2.2) outperforms the individual runs (see Table 1). In the official results, the static run performs worse than ASR only. But, after retraining the static models and guarding the covariances in the process (see Section 3.2), the static+ASR combination run also outperforms the monomodal variants. However, rebuilding the dynamic models in the same way is expected to give a similar increase in performance.<sup>2</sup> Therefore, the comparison in the remainder of this section is based on the official results. The dynamic combination outperforms the static one, even though the MAP score for static only is the same as that for

<sup>2</sup>At the moment, we are recomputing the dynamic models to verify this.

dynamic only. MAP scores hide a lot of information though. The dynamic run has a higher initial precision (See Figure 2), and thus in some ways, the dynamic run is better than the static run and therefore more useful in a combination. While the dynamic models represent some of the temporal aspects of a shot, like objects (dis-)appearing, the main advantage over the static models seems to result from two (related) aspects: more training data describing the visual content, and less dependency on choosing an appropriate keyframe.

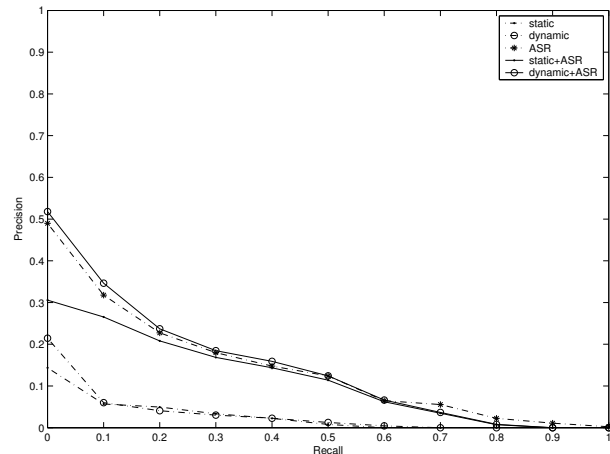


Figure 2: Recall-Precision graph for static, dynamic and ASR runs as well as multimodal combinations.

Results of the dynamic+ASR are further improved by filtering out shots with anchor persons. These shots are identified by computing the likelihood of generating the set of frames annotated with the labels *news\_person\_face* and *studio\_setting* from the shot models.<sup>3</sup> The 1077 shots<sup>4</sup> with the highest likelihood are assumed to contain anchor persons and are filtered out. The resulting MAPs are shown in Table 1. A detailed analysis of the combination of ASR and visual results can be found in [2].

<sup>3</sup>For computational reasons, we computed the likelihood of a subset of the samples from the frames, instead of the likelihood of the full set.

<sup>4</sup>After 1077 documents, we noticed a drop in Anchor person likelihood scores.

## 5 Combining Topic Examples

### 5.1 Merging Run Results

Last year, we found combining multiple examples in one query is far from trivial [12]. Conflicting examples (like different views under different weather conditions) can easily cause bad results. Still, if we use just one 'best' topic example as a query there's a risk of missing a lot of relevant material. Suppose, for example, we selected a close up shot of a point being scored in basketball, and most of the relevant shots in the collection happen to be overview shots of the playing field. This year, instead of combining multiple examples in a single query, we experiment with running separate queries for each example and merging run results afterwards. For each topic, we manually select a set of 'good' examples, run separate queries for each example, and then merge the results using a simple round-robin approach. Duplicates are filtered out afterwards.

In addition to this within topic merging of examples, we also experiment with merging results from different models. The goal here is to find out if it is possible to decide in advance what would be a good model to use for a given topic. When there are only still image examples, one could decide to use the static models, whilst the dynamic model might be used for video examples. For each topic, we manually selected a visual model to be used. The main strategy is the one described above: dynamic models for video examples and static ones for still image examples (more precisely, 'Topic Models', see Section 5.2). In addition, for topics where this seems appropriate, we filter out unwanted shots from the results set. The approach is similar to the one described in [8]. We use *Anchor Person*, *people*, *studio setting* and *weather maps* filters and remove those shots that have a high likelihood<sup>5</sup> of explaining annotated samples of the given feature (See also section 6).

The MAP of this merging run is .039, the highest among the runs that use only visual information, even though, for some topics, it contains results from the disappointing Topic Model run (See Section 5.2).

<sup>5</sup>High means within some top  $K$ , where  $K$  is set using separate training data.

It may be interesting to see if topic specific model selection can be useful.

### 5.2 Building Topic Models

Instead of explaining the query samples from the document models (and combining the results for different examples afterwards), we could go the other way and explain document samples from query models (called 'Topic Models' in this paper). In this setting, documents that have a high likelihood of being generated from the query model are assumed to be relevant. For a full comparison of this topic model approach and the original document model approach from section 2, see [11].

In the topic model approach, all available examples for a given topic are used to build a GMM. The assumption here is that different components in the GMM can capture different aspects of the information need, thus contrasting examples (e.g. day and night shots of city views) are no longer a problem. During training of the model we allowed topic samples to be explained by a background model; this way very common, non distinguishing query samples do not contribute as much to the model parameters. After building the models we can use our normal ranking function (Equation 1), only now the roles of query and document are reversed; we have a set of document samples  $\mathbf{x}$  that need to be explained from query model  $\omega_i$ .

Smoothing with background probabilities as in Equation 1 is not wise in this case. When we use smoothing, common samples get a high score. This is useful in the standard query likelihood approach (Section 2), since these high background scores are independent of the document model under consideration; the contribution of the individual document models to scores of common samples, is therefore less important and the influence of these samples on the ranking is relatively small. However with the reversed likelihood that we are using here, smoothing means that common *document* samples get high scores. Obviously, this is *not* document independent and thus documents with a lot of common samples will get high scores and end up at the top of the ranked list. For this reason, we do not smooth the scores for the

Topic Models and rank the documents using:

$$\text{RSV}(\mathbf{x}^D) = \frac{1}{N} \sum_{j=1}^N \log P(x_j^D | \omega_Q), \quad (5)$$

where  $\mathbf{x}^D$  is the set of  $N$  samples from document  $D$  and  $\omega_Q$  is the model built from the samples from topic  $Q$ . Using this approach we obtain a MAP of only .005 and for many topics we retrieve a lot of black frames. Apparently, even without the smoothing, we retrieve mainly shots with common samples. This is not too surprising, since shots with common samples which are easily explained by *any* model, are often also well explained by a specific query model.

In an additional (not submitted) experiment, we calculate the odds rather than the likelihood of the document samples:

$$\begin{aligned} \text{RSV}(\mathbf{x}^D) &= \frac{\sum_{j=1}^N \log P(x_j^D | \omega_Q)}{\sum_{j=1}^N \log P(x_j^D | \neg \omega_Q)} \\ &= \frac{\sum_{j=1}^N \log P(x_j^D | \omega_Q)}{\sum_{j=1}^N \log P(x_j^D)}. \end{aligned} \quad (6)$$

Using this approach, we get a MAP of .013, significantly higher than for the document likelihood ranking discussed above, but lower than the original query likelihood ranking. Perhaps, we still have a background influence and are now confronted with a preference for uncommon documents. In post workshop experiments, we found that indeed it is important to balance the influence of common and uncommon shots [11]. In those experiments, we use the likelihood ratio of Equation 6, but we smooth the numerator to balance the influence of common and uncommon shots. Using this smoothed version of the likelihood ratio, we get a MAP of .033. We have to admit that these new experiments are based on the COVguard version of the training, but even then the score without smoothing is much lower (.015).

## 6 Feature Extraction

In a way, feature extraction is the same as retrieval. There is a more or less coherent set of examples (annotated images) for a specific information need (feature to be detected). Thus, it is possible to use the

techniques developed for search on the feature extraction task. The results could be regarded as a baseline for systems that put more effort into building detectors for specific features.

Starting from the results of the collaborative annotation effort, we construct sets of examples for each specific feature, by grouping labels. For example everything annotated as either *building* or *house* is used in the example set for the *building* feature. The set of examples is then converted into feature vectors using the usual procedure (computing DCTs from 8x8 pixel blocks). From each set of samples, we take a random subset<sup>6</sup> and use that as the set of samples to use our search techniques on. We used the sample likelihood approach (See Section 2) and the reversed likelihood approach based on feature models both with and without using the background probabilities during training (Section 5.2). For some features some of the approaches got results above median, but overall the results are disappointing.<sup>7</sup> Possible reasons for this are the high diversity in the set of annotated examples for a given feature and again the influence of the background probabilities (cf. Section 5.2). Further experiments with odds rather than likelihoods and more carefully selected examples are needed to test whether feature extraction as search is in principle a reasonable option.

## 7 Interactive Experiments

The interactive experiments are built on results from the automatic models as described in Section 2.3. Since computing the associative matrix is computationally very intensive, we computed the likelihood using a random subset of 100 image samples<sup>8</sup>. Based on experiments on the TREC 2002 video collection on which we simulated user feedback using the assessments [1], we identified the following four exper-

<sup>6</sup>We take at most 10,000 samples per feature, to keep everything tractable. For sets with fewer than 10,000 annotated samples, all available samples are used.

<sup>7</sup>Even a random system scores above median on some topics.

<sup>8</sup>A small experiment indicated that using a subset of 100 samples only, already gives a ranking similar to the one obtained by using the full set of samples.



iments:

1. **Randomised** The user gets random key frames on the screen;
2. **Text-Only** The next screen of key frames is only based on the textual query, not on the user’s feedback;
3. **Feedback** The next screen of key frames is based on the user’s feedback using the associative matrix based on the mixture models;
4. **Feedback-trained** The next screen of key frames is based on the user’s feedback using an associative matrix that was trained on the user feedback gathered by experiments 1 to 3.

Only systems 3 and 4 actually used the feedback; the other two systems recorded the feedback for the resulting ranked list to make use of the users’ notion of relevance, but ignored it when updating the screen. Systems 2–4 started out with textual query, which was automatically taken from the topic descriptions, thus emulating the realistic situation when the (inexperienced) user copies the topic description into the query field. For the same topic both systems started with the same set of images.

The same user interface was used for all systems. The screen contained twelve scaled-down images (3 rows by 4), which could be magnified to their real size by hovering mouse pointer over an icon. Above each image, the user could check a box to indicate its relevance to the topic. The topic description was available to the user at any time, and could be removed from the screen to save space for key frames. To support the users in their action, recent positive examples were displayed at the side of the screen, even when the feedback was not used by the system for the display update.

For Systems 1–3 there were two groups of three subjects performing 15 minute search sessions for each topic, with  $3 \times 3$  Latin square experiment design. The supplemental experiment with System 4 was performed by one person. The ranked lists for the submission to TREC were produced by taking the images that were selected by the user and concatenating them with the ranking produced by the model at the last iteration step. All systems used the

collected relevance feedback to fill up the result list up to the allowed 1000 shots.

Table 2 contains in column 2 MAP achieved with each method after at most 15 minutes of interaction<sup>9</sup>. It does not come as a surprise that the prior information from speech transcripts and text queries from the topic descriptions, in combination with user’s notion of similarity, has great impact on the retrieval quality (System 2), further improved by the relevance feedback based on visual contents of the key frames (in Systems 3 and 4). Remarkably low value for the Randomised system indicates, that for a large collection such display update strategy is not practical. The users got to see very little relevant key frames on the screen, and in almost half of the cases the system was unable to produce any one.

Column “Relevant found by users” of Table 2 shows the proportion of relevant shots identified by the users during the search, from all relevant shots that occurred in the submitted results. By its definition, the MAP measure is sensitive to the number relevant shots put on top of the ranked list. The difference in MAP between Systems 3,4 and System 2 is substantial, but still it hides the fact that the users of the systems with feedback identified 49% and 60% of all detected relevant shots, whereas for the text-only system this number is about 38%. Preliminary experiments with training the associative matrix (System 4) indicate that the relevance information obtained from the user is valuable and should be used to complement existing techniques for feature extraction.

Table 2: Statistics of interactive runs

System type	MAP	Rel.found by users	Agreement with NIST	Missed Relev.
Randomised (1)	0.026	5.79%	55.00%	31.25%
Text only (2)	0.214	37.83%	85.02%	51.91%
Feedback (3)	0.245	49.18%	78.74%	48.98%
Trained (4)	0.240	60.65%	64.03%	32.07%

<sup>9</sup>For some topics the system returned no relevant shots which lead to a higher MAP reported in the TREC summary table.

Next to “Relevant found by users” are shown the agreement between the users and the ground truth, and the percentage of missed relevant shots. From the numbers we suspect that the user tends to give his/her relevance judgements relative to other images present on the screen: In the Randomised system the user clicked larger proportion of irrelevant key frames, than in other cases. Conversely, when more relevant frames are displayed, the user missed a larger proportion of them than when the screen contained arbitrary key frames (see in Table 2 column “Missed relevant”, which shows proportion of key frames of relevant shots, that have not been marked by the users as relevant). Partially the relevant documents are not recognised as such due to the fact that the user saw single frames representing a shot which may not have contained the relevant objects, and not the video shots themselves. This is an indication that there is a need for a better presentation of a shot to the user.

Figure 3 shows the recall-precision graphs of systems 1–4. On most recall levels, the feedback run is better than the text run. The feedback run on the trained matrix showed slightly worse performance than the feedback run on the visual feature matrix, but at higher recall levels there is a noticeable advantage.

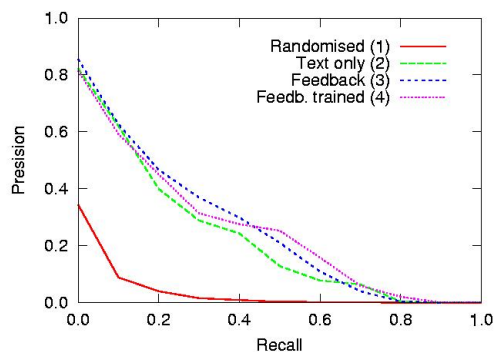


Figure 3: Precision at various recall levels for interactive runs.

In the questionnaires that the users filled before and after the search sessions, the users merely preferred to make many fast iterations to waiting between the search steps, regardless the fact that a session never exceeded 15 minutes. The average number of iteration steps per search was 65, with the maximum of 150. Thus, an algorithm that quickly produces the next set of examples, is desired. At the same time, in our small user study we found that the users did not like the Randomised system because of the lack of good examples on the screen, and, conversely, liked when the next screen seemed to be the consequence of their feedback, giving the user the feeling of being in control of the process. Thus, when choosing an algorithm for display update, one should keep a trade off between speed and quality of that.

## 8 Conclusions

This paper presented the models used and experiments carried out for the TRECVID 2003 workshop. We focused on combining information on different levels. We combined information from different sources and modalities, information from different visual examples and human and model based information.

We extended our generative probabilistic models to include temporal information and found this improves the results. Combining these results with results from a ASR run gives another improvement. Whenever both text results (from ASR language models) and visual results (from GMM models) do something useful, a combination gives even better results.

Manually selecting good visual examples and useful models for a given topic gave the best results among the visual runs. The selection of (multiple) good examples and their combination are the main cause for this success.

Intuitively, building a model from all (or some) visual topic examples and ranking the documents by their probability of being generated from such a topic model, seems at least as good an approach as the reversed process (trying to explain query samples from document models). In practise it turns out that in

fact this topic model approach works better, provided that smoothing is used (See also [11]).

The interactive search results are obviously far better, than the manual ones. Also here we noticed that a combination of textual and visual information performed better than text alone. The user action remains an indispensable component of an information retrieval system.

## References

- [1] L. Boldareva, D. Hiemstra, and W. Jonker. Relevance feedback in probabilistic multimedia retrieval. In *DELOS Workshop on Multimedia Contents in Digital Libraries*, 2003. <http://www.music.tuc.gr/MDL/>.
- [2] A. P. de Vries, T. Westerveld, and T. Ianeva. Combining multiple representations on the TRECVID search task. In *International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2004)*, 2004. to appear.
- [3] 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.
- [4] R. Fagin. Fuzzy queries in multimedia database systems. In *ACM Symposium on Principles of Database Systems (PODS)*, pages 1–10, 1998.
- [5] J. Gauvain, L. Lamel, and G. Adda. The LIMSI broadcast news transcription system. *Speech Communication*, 37(1–2):89–108, 2002.
- [6] D. Hiemstra. *Using language models for information retrieval*. PhD thesis, Centre for Telematics and Information Technology, University of Twente, 2001.
- [7] T. Ianeva, A. P. de Vries, and T. Westerveld. A dynamic probabilistic retrieval model. In *IEEE International Conference on Multimedia and Expo (ICME)*, 2004. to appear.
- [8] T. I. Ianeva, A. P. de Vries, and H. Röhrig. Detecting cartoons: a case study in automatic video-genre classification. In *2003 IEEE International Conference on Multimedia & Expo*, 2003.
- [9] J. M. Ponte and W. B. Croft. A language modeling approach to information retrieval. In *Proceedings of the 21st ACM Conference on Research and Development in Information Retrieval (SIGIR'98)*, pages 275–281, 1998.
- [10] N. Vasconcelos. *Bayesian Models for Visual Information Retrieval*. PhD thesis, Massachusetts Institut of Technology, 2000.
- [11] T. Westerveld and A. P. de Vries. Multimedia retrieval using multiple examples. Technical Report INS-E0403, CWI, Amsterdam, The Netherlands, March 2004.
- [12] T. Westerveld, A. P. de Vries, and A. van Ballegooij. CWI at the TREC-2002 video track. In E. M. Voorhees and D. K. Harman, editors, *The Eleventh Text REtrieval Conference (TREC-2002)*, 2003.
- [13] 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*, 2003(2):186–198, 2003. special issue on Unstructured Information Management from Multimedia Data Sources.