

## The Lowlands team at TRECVID 2007

Robin Aly<sup>1</sup>, Claudia Hauff<sup>1</sup>, Willemijn Heeren<sup>1</sup>,  
Djoerd Hiemstra<sup>1</sup>, Franciska de Jong<sup>1</sup>, Roeland Ordelman<sup>1</sup>,  
Thijs Verschoor<sup>1</sup>, and Arjen de Vries<sup>2</sup>

<sup>1</sup>University of Twente, The Netherlands    <sup>2</sup> CWI, The Netherlands

February 26, 2008

### Abstract

Type	Run	Description	MAP
	<b>Official</b>		
A	UTen	English ASR	0.0031
A	UTt <sub>hs</sub> -t2-nm	Top-2 concepts from t <sub>hs</sub> graph method with neighbor multiply	0.0137
A	UTwiki-t2-nm	Top-2 Wikipedia concepts with neighbor multiply	0.0131
A	UTwiki-t2-en-nm	Top-2 Wikipedia concepts and English ASR with neighbor multiply	0.0107
A	UTwiki-t2-nl-nm	Top-2 Wikipedia concepts and Dutch ASR with neighbor multiply	0.0096
A	UTwordnet-t2-mult	Top-2 Wordnet concepts with neighbor multiply	0.0083
	<b>Additional</b>		
A	UTnl	Dutch ASR	0.0031
A	UTwikiS-t2-nT	Top-2 Wikipedia concepts on stemmed queries with neighbor using the concept detector scores from the B <sub>tsinghua</sub> -icrc <sub>5</sub> run	0.0410
A	UTt <sub>hs</sub> -t2-n	Top-2 concepts from t <sub>hs</sub> graph method of stemmed queries with neighbor the concept detector scores from the B <sub>tsinghua</sub> -icrc <sub>5</sub> run	0.0346
I	UTinter-wiki-nm	Interactive Search Task using Wikipedia concepts with neighbor multiply	0.0405
I	UTinter-en	Interactive Search Task using ASR based search	0.0338

In this report we summarize our methods and results for the search tasks in TRECVID 2007. We employ two different kinds of search: purely ASR based and purely concept based search. However, there is not significant difference of the performance of the two systems. Using neighboring shots for the combination of two concepts seems to be beneficial. General preprocessing of queries increased the performance and choosing detector sources helped. However, for all automatic search components we need to perform further investigations.

We also present a thorough analysis of the results of the TRECVID 2007 Interactive Search task that the Lowlands team participated in unofficially. This allowed for a comparison of our baseline system's functionality with that of other participants. Moreover, the analysis provides more insight into how users interact with the baseline search system. Finally, recommendations are made on how to improve both the baseline system's performance and the current user interface.

## 1 Introduction

Bridging the semantic gap is a key problem for multimedia information retrieval tasks such as video search [11]. It requires coupling of the well understood extraction methods for low level features from media files (e.g. color histograms or audio energy) and the semantically rich descriptions or concepts<sup>1</sup> in which users express their information needs (e.g. *Find me pictures of a sunrise*). In this paper we investigate how our concept combination methods [1] [4] perform against an ASR-only method<sup>2</sup>, and whether combining the two helps.

Concept detectors are commonly trained through positive and negative examples on a certain training dataset. For a particular domain appropriate sets of concepts and training data have to be selected. A less straightforward solvable problem is how to handle queries that do not correspond to exactly one concept from the selected set of concepts. Due to the lack of knowledge about the structure of the *semantic space*, it is not an option to simply increase the number of detectors up to the point where all requested concepts are covered. The hypothesis is that in order to support searching for *Condoleezza Rice* with a search system that only has the concepts *Face* and *Women* available, the uncovered concept has to be expressed as a combination of concepts for which detectors exist.

In this paper we describe three novel techniques to combine concept scores. The main innovations are in the score modification via the scores of preceding and following shots, and in combining the output for one detector with the output of other detectors. We also ran our IR system PF/Tijah [6] on the ASR output and investigated ways to integrate the results with the results from concept combination.

Furthermore, we developed a baseline video search interface and addressed its effectiveness and acceptance in unofficial interactive runs at TRECVID 2007<sup>3</sup>. Parts of this system will be developed further to study search in collections where the spoken content can be exploited as time-stamped metadata generated through e.g., ASR. This holds for audio collections and for video collections whose visual content mainly consists of talking heads; e.g., lecture recordings, meeting recordings, and interview collections. For speech-driven metadata the TRECVID tasks may be considered difficult as they target visual features in the video documents. On the other hand, this platform allows us to compare our baseline system's performance to that of other systems.

This paper is structured as follows. In Section 2 we introduce the system we used for our experiments. In Section 3 we elaborate on our concept combination methods.

---

<sup>1</sup>In TRECVID terminology high level features

<sup>2</sup>ASR: automatic speech recognition

<sup>3</sup><http://www-nlpir.nist.gov/projects/t01v/>

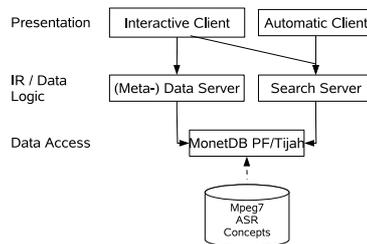


Figure 1: Architecture

Section 4 briefly outlines the PF/Tijah system. Section 5 shows the setup for the interactive search task. Section 6 describes the experiments we undertook to verify our methods. Section 7 concludes the paper.

## 2 Minos System Overview

We named the IR System which we used to carry out the runs *Minos*<sup>4</sup>. It is designed to allow several search strategies as well as to combine them. The system architecture is shown in Figure 1.

In the data access layer we use the XML database, MonetDB-XQuery, with PF/Tijah as a Text IR extension. The data is stored in MPEG7 documents which contain time interval, English and Dutch ASR output and the scores of the concept detectors from the University of Amsterdam. MonetDB-XQuery provides a method to execute queries using an XML remote procedure call (XRPC).

In the IR logic layer, the search server is concerned with encapsulating the information retrieval logic and to hide the system's complexity from the presentation layer. It has the ability to use different search modules. The two search modules implemented at the moment provide concept based and text (ASR) based search. The data server provides a unified interface to deliver (meta-) data to the user. Bulk binary data, such as key frames are provided through a URL. The protocol from the IR logic layer to the presentation layer is using Web Services defined by the web service description language (WSDL) to ensure interoperability.

In the presentation layer we implemented two clients. One client is designed to carry out automatic search tasks. It gets passed the TRECVID topic file and automatically executes one topic after each other for all system configurations using the web service. The interactive client allows a human user to interact with the search system. At start up, the client program is told which search module it should use. This setting, together with the text query, gets passed to the search server. The server returns a list of shot identifiers together with a rank and a score. For all the shot identifiers the needed metadata is retrieved from data server. The key frame pictures get loaded from a potentially independent web server.

<sup>4</sup>Minos is a mythologic King of Crete who created a un-escapable labyrinth

### 3 Concept Combination

As was mentioned earlier concept combination is carried out because one concept is unlikely to be enough to answer a user’s query. Our notion of combination [1] focuses mainly on the co-occurrence of concepts. Unlike techniques mentioned in [15] we do not take relationships between concepts into account. Therefore the two concepts *Animal* and *Dog* would be treated the same for a query “Find me dogs”. This allows the *Animal* concept to introduce noise (e.g. *Cats*) into the result. A big advantage is that there is no need for an ontology to represent those relationships.

#### 3.1 Query To Concepts

Users cannot be expected to know the concepts that are available to the system. User queries usually either consist of a few keywords (e.g. *Beach*) or more elaborate natural language requests (e.g. *Find me pictures of a beach with people.*). In the best case, the query contains one or more concept names and syntactic matching is sufficient. However, often this will not be the case. For instance, the set of concepts included in TRECVID include *Outdoor*, *Waterscape* and *People* but not *Beach*. Hence, the first task is the extraction of TRECVID concepts underlying the queries. The natural language query and the concepts available for the collection are matched and a ranking of relevant concepts is derived that shall resemble the information need expressed in the query as close as possible. We implemented two query to concept approaches: one is based on WordNet [3] glosses and Wikipedia pages, the second is based on WordNet’s graph structure.

In the gloss (Wikipedia) approach, we consider WordNet glosses (Wikipedia pages) describing a concept as substitutes of the concepts. The relevant concepts to a query can then be found by using Text IR methods on the collection of the documents describing the concepts.

In the second approach, WordNet’s graph structure is exploited. TRECVID concepts are mapped to synsets in WordNet. The distances between query terms and concepts on the graph are used to rank the concepts.

#### 3.2 Concept Preprocessing

Given the ranked list of concepts that are returned for a text query the system still has to select some concepts from this list for their combination. Using the whole list is not advisable as the query to concept step might return all concepts available to the system, although the irrelevant ones only with very small score. In [4] we performed studies on various strategies. Taking the top-2 concepts from the list showed the best performance. We used this setting in all experiments throughout this paper.

We used the concept detector scores from the *A\_uva.Coeus\_4* run of the high level feature detection task. We chose this run because we used the detector results from the University of Amsterdam [13]. Because we used these detectors in earlier experiments [1, 4], we expect better comparability. As our methods need scores within the interval  $[0..1[$  we linearly scaled the scores to the desired interval. We had to take this decision as probabilistic scores were not available.

Functions on single concept:

$$r(c, s_j) = \frac{\text{rank}(s_j) - \text{minRank}(c)}{\text{maxRank}(c) - \text{minRank}(c)} \quad (1)$$

$$\text{smooth}(c, s_j) = \frac{\sum_{i=j-nh}^{j+nh} r(c, s_i)}{2nh + 1} \quad (2)$$

Functions on multiple concepts:

$$\text{mult}(C, s_j) = \exp\left(\sum_{c \in C} \log(r_c(s_j))\right) \quad (3)$$

$$n(C, s_j) = \frac{\sum_{c \in C} r_c(s_j) \frac{\sum_{c' \in C \setminus c} \text{smooth}(c, s_j)}{|C|-1}}{|C|} \quad (4)$$

$$\text{nm}(C, s_j) = \frac{\sum_{c \in C} r_c(s_j) \exp\left(\sum_{c' \in C \setminus c} \log(\text{smooth}(c, s_j))\right)}{|C|} \quad (5)$$

Figure 2: Combination Functions

### 3.3 Combination of Concept Scores

In the following we describe the combination methods we used to calculate a joint score from the output of multiple detectors.

Figure 2 shows the definition of all used combination functions. The function  $r$  (1) returns the previous described derived score of the shot  $s_j$  as calculated from the rank. The function  $\text{smooth}$  (2) assumes that it is more likely that a concept  $c$  appears in the shot  $s_j$  if it also appears in previous or following shots. Similar approaches have been investigated using the text from automatic speech recognition associated with shots [5]. We define a surrounding neighborhood as a fixed number  $nh$  of shots before and after the actual shot  $s_j$  that contribute to the score of  $s_j$ .

The function  $\text{mult}$  (3) multiplies adds the logarithm of the scores of all concept detectors. At the end it applies the  $\exp()$  function to bring the resulting score back into the interval  $[0..1]$ .

The Neighbor function  $n$  (4) considers all base scores multiplied with the average of the smoothed scores of the other concepts to apply.  $\text{nm}$  (5) is an extension of the  $\text{mult}$  function which weighs the individual scores by the  $\log()$  of averaged smoothed scores of other concepts.

## 4 PF/Tijah TextIR

We kept all information in an MPEG7 conform documents. To store the scores of the feature donations we extended the `mpeg7:VideoSegmentType` to include Concepts

subelement which in turn contains all concept scores of each subject.

Because the unit of retrieval was a shot, we used all ASR and automatic speech translation [7] from speaker segments overlapping with the shot segment to retrieve a shot. In this way the text associated with the shot could be a little more than what was actually spoken during the shot. Neighboring shots are considered to have a similar relevance; therefore this is not problematic.

In order to keep the data format to MPEG7 we extended the available vocabulary to also contain concept scores. This was done through creating a new schema on top of the existing MPEG7 schemas extending the existing type *VideoSegmentType* to allow definition of concepts.

We used the protocol the MonetDB's XML Remote Procedure Calls (XRPC) to communicate with the data layer. We implemented three such XRPC functions: (i) one which gets passed the query text and the language returning a ranking of shots, (ii) one which gets passed the query text and returns a list of concepts and (iii) a function which retrieved all metadata for a list of shot identifiers.

To see if a joint result of ASR output and concept combination could be beneficial we use the score from the shots found from ASR as "artificial" concept that could get combined like the others.

## 5 Interactive Search

We developed a baseline video search interface and addressed its effectiveness and acceptance in unofficial interactive runs. The system will be developed further to study search in collections where the spoken content can be exploited as time-stamped metadata generated through e.g., ASR. This holds for audio and video collections whose visual content mainly consists of talking heads; e.g., lecture recordings, meeting recordings, and interview collections. For speech-driven metadata the TRECVID tasks may be considered difficult as they target visual features in the video documents. However, this platform allows us to compare our baseline system's performance to that of other systems.

### 5.1 User Interface

Since most users, i.e. non-expert users, usually formulate text queries when using search engines, we only included query-by-keyword search (as opposed to query-by-example or query-by-concept search) in our baseline search system. As opposed to more advanced video search systems such as MediaMill [13], Informedia [2], and Físchlár [12], we have not (yet) included ways to use relations between shots for result presentation, such as temporal relations or stories (e.g., [12, 13]), or semantic relations ([2, 13]). This was due to lack of development time.

We tested two manners of query processing for retrieval: (i) ASR-based search (UTinter\_en), and (ii) concept-based search (UTinter\_wiki\_nm). These differences currently do not affect the type or manner of information presentation in the UI.

A screen shot of the UI is given in Figure 3. After processing a query entered at the top of the screen, the total number of results found is reported and retrieved shots

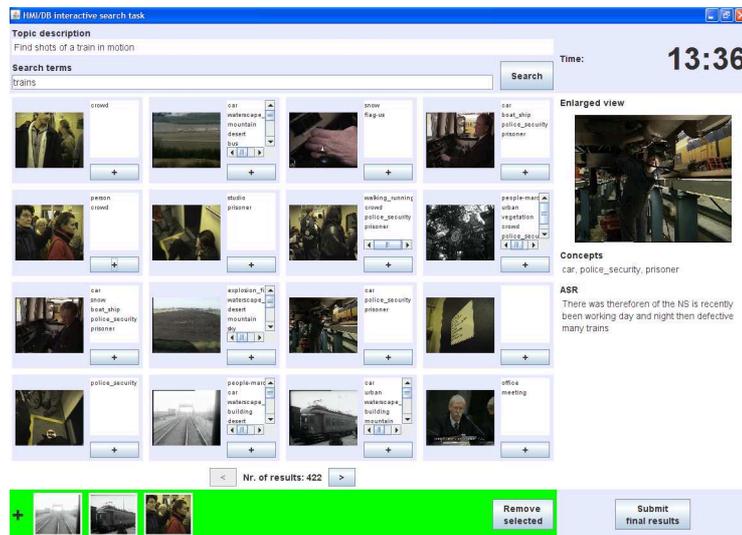


Figure 3: Screen shot of the search interface.

are presented. Results are shown in sets of 16 key-frames per page (in a 4 x 4 matrix). For each key-frame the concepts most strongly associated with it are given as well as the option to move that particular shot to the list of results that users definitely want to keep. This is done by clicking the plus-button next to a shot. The definitive selection is shown in the green bar at the bottom of the screen. Clicking on a key-frame gives more precise information on that frame on the right hand side of the screen: an enlarged view of the shot, the list of concepts associated with it, and the machine-translated English version of the Dutch ASR text associated with the shot.

## 5.2 Method

Each search run in TRECVID consists of 24 topics (or: search tasks). In the interactive task these do not have to be completed by a single person. Since we tested two system variants, 48 topics had to be searched for. We had six Dutch participants (age range=21-27; 1 female, 5 males) each complete eight topics, four on each system variant, see Table 1 for the design. All participants used search engines on a daily basis and three out of six indicated to also search for videos. They furthermore regularly searched online library catalogs. All were novice users of the system.

Topics, queries and results were in English, the second language of our users. Tests were run on PCs with 19" monitors in a quiet room. Before the actual test, users filled out a demographic questionnaire, which was followed by written instructions and practice with the search system on a topic that was not part of that searcher's test set. This lasted about 20 minutes. During testing, system variant and topic order was counterbalanced across participants. Between performing the search tasks on the two system variants, participants got a short break, and after each topic they filled in a post-

	$S_{ASR}$	$S_{Concept}$
$T_{0197} - T_{0200}$	1	2
$T_{0201} - T_{0204}$	3	4
$T_{0205} - T_{0208}$	5	6
$T_{0209} - T_{0212}$	2	1
$T_{0213} - T_{0216}$	4	3
$T_{0217} - T_{0220}$	6	5

Table 1: Design scheme showing which participant completed which topics on which system variant  $S$ . Subjects with uneven numbers first got the system variant with ASR-based search, for subjects with even numbers it was the other way around.

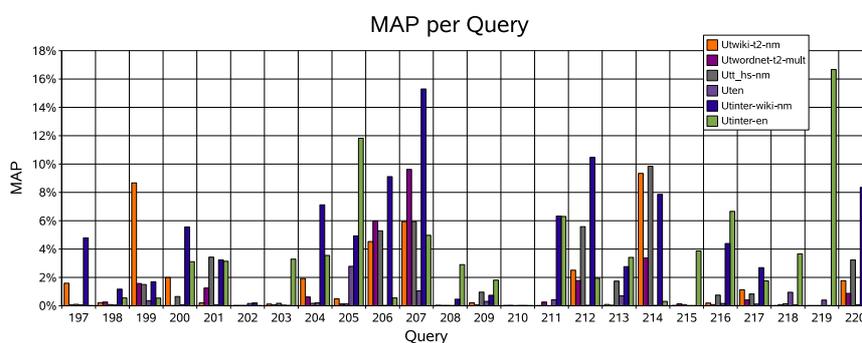


Figure 4: Per Query Average Precision

topic questionnaire (translated to Dutch from the CMU2006 example<sup>5</sup>). They received monetary compensation for their efforts.

Participants mostly used the full 15 minutes per topic; they stopped early if no more matching shots were found. In these cases 10 to 14 minutes were used, except for topic 0219 that was abandoned after 6 minutes (see section 6.2.1 for more information on this topic). During testing we measured the interaction with the system by logging user actions. After the interactive task a post-test questionnaire was administered on the system's general usability.

For score computation, result sets were filled to 1000 results. If the user's result set was not large enough it was completed with the results from his/her last query, and if necessary the set was further completed with the results from the automatic run for that topic<sup>6</sup>.

## 6 Experiments

In this section we describe the experiments we did to verify our methods. First in Section 6.1 Runs according to the automatic search task description of TRECVID are described. The following Section 6.2 describes the outcome of our interactive user studies with the search system.

### 6.1 Automatic Runs

All our official runs are automatic runs. For the six runs we used the text IR based method with the Wikipedia and WordNet corpus and the graph based query to concept method hierarchical shortest path. We left out other graph based methods as they did not help increasing the performance in [4]. The given topics were then fully automated executed by the system.

Overall one of our runs reached the median of all submitted runs. Later we found out that there were some simple changes of our methods which improved the results significantly.

To compare the different Query to Concept mechanisms we compare the two official runs UTwiki-t2-nm and UTt\_hs-nm together with the unofficial run UTwordnet-t2-nm (MAP 0.0139) it is not possible to conclude whether graph or text based methods are to be preferred.

A comparison between the combination methods based on the official run is problematic. There is an indication that the neighbor multiply method is better to the multiply method. To what extent this is true would have to be verified by runs using the same Query to Text method but varying combination methods.

We also compared the performance of our system when using Dutch and English language. For Dutch we used the direct ASR output and human translated topics. The result of this unofficial run UTnl was 0.0031 and therefore exactly the same as the one from English, which was machine translated.

Furthermore, we investigated whether using text scores, as another concept, helps. From the listed runs we have to conclude that using ASR - at least in this manner - is decreasing performance.

Additional checks on the returned concepts from the Query to Concept phase revealed that very often the same concepts were chosen. Investigations showed that this was due the nearly constant beginning of the textual topic "Find shots of". Introducing a stop word mechanism which removed this bit yielded significant improvements. Hereafter all reported results were achieved using this stop wording.

To see whether the chosen source of concept detector scores matters we ran the combination UTwiki-t2-nm on all available sources, see Figure 5. It can be seen that the achieved MAP is significantly different depending on the source. The source we chose for the official runs (A\_uva.Coeus\_4) performed within the upper third of the sources. The run B\_tsinghua-icrc\_5 yielded the best results. We used this detector source for another intensive investigation of the performance of all query to concept and

---

<sup>5</sup>Last visited Oct. 22 2007: [http://www-nlpir.nist.gov/projects/tvpubs/tv6.papers/cmu\\_talk\\_search.slides.pdf](http://www-nlpir.nist.gov/projects/tvpubs/tv6.papers/cmu_talk_search.slides.pdf)

<sup>6</sup>Double entries were removed.

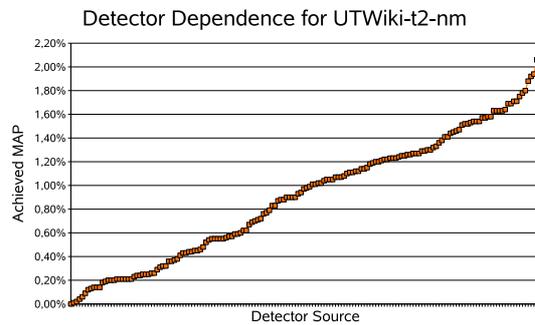


Figure 5: Dependence on Base Detectors

combination methods. As reference we report the run which resulted in 0.0410 MAP, which was using the Wiki explanation of the concept, and the graph based method t.hs.

## 6.2 Interactive Search

The results presented in this section were gathered from a small number of users, as is usually the case for this type of evaluation. This implies that effects must be interpreted with caution.

### 6.2.1 Evaluation results

The UTinter\_en interactive run on average got a MAP of 3.5% and the UTinter\_wiki\_nm run got 4.63%. This difference was not significant; UTinter\_en scored higher on some topics whereas UTinter\_wiki\_nm got better results on others. For topics 0197, 0207, 0212, 0214 and 0220 concept-based search scored higher than ASR-based search. For topics 0205, 0215, 0218 and 0219 it was the other way around.

Most noticeable are the results for topic 0219 (*Find shots that contain the cook character in the Klokhuis series*), where ASR-based interactive search outperforms all other conditions (both interactive and automatic). Given that the content as well as our searchers are Dutch, they could have used their knowledge of the TV show during search. Analysis of the users' logs showed that this was not necessary: the high score on topic 0219 was the result from one user finding a single shot for the query 'cook klokhuis' as opposed to another user who found no relevant shots at all using different queries.

In comparison with the official interactive runs of other groups, our baseline system ranks among the lowest scoring interactive systems. This may be considered unsurprising given the basic nature of our UI and the fact that we had novice users. In comparison with the corresponding automatic runs (0.31% and 1.37%, respectively) an improvement was found with users in the loop.

User	Text	Concept	Average
1	0.51	0.48	0.50
2	0.63	0.93	0.80
3	0.47	0.79	0.61
4	0.67	0.35	0.53
5	0.69	0.60	0.63
6	0.65	0.86	0.74
avg.	0.60	0.67	0.63

Table 2: Average precision at the number of saved shots, per user and per search type.

### 6.2.2 User Performance and Usability

Besides system performance we monitored how the system was used by logging user actions. In the `UTinter_wiki_nm` run (i.e. concept-based search) participants on average formulated almost 17 queries, looked at 25 previews and saved about 12 shots per topic. Average query length was 2.8 words. In the `UTinter_en` (i.e. ASR-based search) run participants on average formulated almost 27 queries, looked at 25 previews and saved about 12 shots per topic. Average query length was 1.7 words.

One might hypothesize that the shots saved by participants are relevant. To assess this expectation we examined the agreement of the participants' choices with the NIST judgments. Table 2 shows the results. The users' precision at the number of saved shots was 63% on average (shots not judged by NIST were excluded from the analysis). This number may seem rather low, but it is comparable to earlier results [14]. Moreover, given the average number of 12 saved shots per topic, the impact of one or two irrelevant shots on the precision of the saved shots is relatively high. In the following we call the disagreement between the users in our study and the NIST assessors error. However, we remind the reader that there is the possibility that the users in our study could identify the relevance of shot better than the official assessors, as they natively spoke the same language as in the video content.

We examined three cases to better understand why users in some cases may have selected irrelevant shots. These cases were chosen such that the total number of saved shots was over 20, and the number of irrelevant shots saved in comparison to the number of relevant shots saved was high: (i) 18 irrelevant v. 7 relevant shots for topic 0212 with text-based search, (ii) 10 irrelevant v. 21 relevant shots for topic 0213 with text-based search, and (iii) 17 irrelevant v. 26 relevant shots for topic 0220 with concept-based search. Topic 0212 asked to *Find shots in which a boat moves past*. Most errors were made when shots were selected that contained non-moving boats, lying along quays. Topic 0213 was *Find shots of a woman talking toward the camera in an interview - no other people visible*. Shots of women involved in a dialog, but currently not speaking, gave the most frequent errors. Furthermore, shots with other people visible were saved. Thirdly, topic 0220 asked searchers to *Find gray-scale shots of a street with one or more buildings and one or more people*. Errors contained e.g., indoor shots with house-like walls and paths rather than streets. In these three example cases, we think that users were unaware of the many mistakes they made, since they each indi-

	$S_{ASR}$	$S_{Concept}$
EaseToFind	3-7-5-2-7	11-2-2-5-4
Time	2-1-3-5-13	2-1-5-6-10
Satisfaction	2-4-5-9-4	8-2-3-6-5

Table 3: Results of the post-topic questionnaires: counts for each step on the scale of 1 to 5.

cated that they were satisfied with their results for these topics, had sufficient time, and found it easy to find relevant shots (by rating these statements with 4 or 5 on a scale of 1=poor to 5=good).

These cases are taken to suggest that searchers' criteria for accepting shots were not strict; they seemed to accept the shots that best matched the topic, even when that match was only partial. The reason for this is not clear. A possible explanation is that the system's MAP score was not too high, which may have made finding relevant shots relatively difficult for users. They therefore may have been willing to make compromises by accepting shots that partially matched the topics. Moreover, it may mean that users need to be instructed more stringently to follow the topic exactly. Finally, shots that were somewhat unclear – which made it more difficult to correctly judge their contents – might have benefited from the option to play the video for disambiguation.

Zooming in on the queries formulated by users, we found that of a total of 1053 queries to 48 topics (24 per system variant) 44 typos and 8 language errors were made. A 'language error' is defined as either an error due to incorrect transfer of a Dutch word or a misspelled word that was not corrected by the searcher (which a recognized typo would be). Moreover, 128 queries were repeated literally within search tasks, and 8 times participants used operators (AND, NOT, -)<sup>7</sup>. Mainly the high number of repeated queries makes search inefficient, which was already remarked by the participants: they suggested to include the functionality of showing search histories to prevent this.

As for the post-topic questionnaires (i) the ease to find results, (ii) the sufficiency of time to complete search, and (iii) overall satisfaction with the result were rated (on a scale of 1=poor to 5=good). The results are shown in table 3. The median ratings are consistently higher for ASR-based than for concept-based search for each of the subquestions, but nonparametric statistical tests revealed no differences. With respect to the individual topics, users found topics 0197, 0202, 0203, 0208, 0210, and 0211 especially difficult, rating the ease to find relevant shots at 1 or 2. On the other hand, topics 0199, 0204, 0212, and 0213 seemed relatively easy.

The post-test questionnaire addressed the UI's usability by asking about learnability, satisfaction, ease of use, and interface design on a scale of 1 (=poor) to 5 (=good). The median for ease of use was high, i.e. 5, but overall satisfaction was below average at 2.5. We think that these findings are due to the combination of a basic design adhering to known functionalities that however did not lead to satisfying numbers of relevant shots. Furthermore, the system was judged relatively easy to learn (3.5) and also its design was rated positively (4).

Finally, users were asked for any remarks on the system or suggestions for its improvement. According to three out of six participants improvement was mainly needed

<sup>7</sup>The system was not designed to use operators, but participants were not informed of this beforehand.

in the match between the shot and its associated concepts. Users might have found the idea of associated concepts useful for their searches, but perceived the accuracy of the concept assignments as too low. None of them made remarks about the relatively poor quality of the ASR text in the detailed view, but it has been shown that low-quality ASR does not help users, e.g., [10], and may even hinder them [9]. Another suggestion for improvement made by the users, again three out of six, was to include more information on their current location in the result set.

## 7 Conclusion

We conclude that we achieved in the official runs around the median of the other systems. Later we found that stemming and query stop words improved the results significantly. The usage of English or Dutch ASR (or machine translated ASR) did not yield a significant difference. In comparison to combination methods the performance was worse. To incorporate them as an artificial detectors score into the combination lowered MAP. Finally, we found that our method strongly depends on the kind of detector source.

Given the evaluation of the interactive search system presented in this report there is room for improvement of the UI, but also of the search engine itself. Reasons for why we performed comparatively bad in the Interactive Search task are:

- The current functionality of the interface is very basic;
- Our users were novice users of the system, who were not information professionals, and also second language users;
- The low MAP scores obtained on the data make the task difficult to begin with.

Therefore, to improve overall performance, both the UI and the search engine should be developed further.

As for the UI, relations between shots should be exploited. As we mentioned in section 5.1, we did not have the time to include such relations. For instance, once a user finds a relevant shot, it is likely that neighboring shots are relevant as well, but in our system this temporal dependency was not exploited. In addition to temporal relations, semantic and visual similarity could also be used, either by grouping shots on screen according to these criteria, or by allowing users to use a ‘give me more of this’ functionality. A second improvement of the UI would be to show search histories per search task, so that queries are not formulated twice, because the searcher did not remember he/she had already tried that particular combination of search terms. Thirdly, user support for orientation within the result pages should be improved, for instance by including a page index.

The search system showed median performance in the automatic run. This, however, did not seem to result in a precision that was high enough from the viewpoint of the user. For instance, Web users at most look at 2 to 3 pages of results [8], which is only few dozens of results. With a MAP in the range of a few percents, users will not be satisfied with those first few dozens, and therefore may perceive their task as

effortful and difficult if the UI does not support efficient exploration of the collection. This calls for improvements in baseline search performance.

Finally, the fact that our users were novice, non-native searchers may have affected their search capabilities somewhat. Though the number of observable language errors was low, the participants' ability to successfully rephrase queries was probably less than in natives. The use of student participants instead of information professionals has the disadvantage that the former group has probably developed less strategies to work around quirks of search systems.

In future work on UIs for searching audiovisual archives we will incorporate the recommendations made by our users, and we will further investigate the question of how to improve textual representations of the spoken content.

## References

- [1] R. B. N. Aly, D. Hiemstra, and R. J. F. Ordelman. Building detectors to support searches on combined semantic concepts. In *Proceedings of the Multimedia Information Retrieval Workshop, Amsterdam, The Netherlands*, pages 40–45, Amsterdam, August 2007. Yahoo! Research.
- [2] M.G. Christel. Evaluation and user studies with respect to video summarization and browsing. In *Proceedings of IS&T/SPIE Symposium on Electronic Imaging*, 2006. San Jose, CA.
- [3] Christiane Fellbaum. *Wordnet: An Electronic Lexical Database*. The MIT Press, 1998.
- [4] C. Hauff, R. B. N. Aly, and D. Hiemstra. The effectiveness of concept based search for video retrieval. In *Workshop Information Retrieval (FGIR 2007), Halle, Germany*, volume 2007 of *LWA 2007 Lernen - Wissen Adaption*, pages 205–212, Halle-Wittenberg, 2007. Gesellschaft fuer Informatik.
- [5] A.G. Hauptmann, R. Baron, M. Christel, R. Conescu, J. Gao, Q. Jin, W.-H. Lin, J.-Y. Pan, S. M. Stevens, R. Yan, J. Yang, and Y. Zhang. Cmu informedias TRECVID 2005 skirmishes. In *Proceedings of the 3rd TRECVID Workshop*, 2006.
- [6] Djoerd Hiemstra, Henning Rode, R. van Os, and J. Flokstra. Pftijah: text search in an xml database system. In *Proceedings of the 2nd International Workshop on Open Source Information Retrieval (OSIR), Seattle, WA, USA*, pages 12–17. Ecole Nationale Supérieure des Mines de Saint-Etienne, 2006.
- [7] Marijn Huijbregts, Roeland Ordelman, and Franciska de Jong. Annotation of heterogeneous multimedia content using automatic speech recognition. In *Proceedings of the second international conference on Semantics And digital Media Technologies (SAMT)*, Lecture Notes in Computer Science, Berlin, December 2007. Springer Verlag.

- [8] B.J. Jansen, A. Spink, and T. Saracevic. Real life, real users, and real needs: a study and analysis of user queries on the web. *Information Processing & Management*, 36:207–227, 2000.
- [9] C. Munteanu, R. Baecker, G. Penn, E. Toms, and D. James. The effect of speech recognition accuracy rates on the usefulness and usability of webcast archives. In *Proceedings of CHI 2006*, pages 493–502, Montreal, Quebec, Canada, April 22-27 2006.
- [10] A. Ranjan, R. Balakishnan, and M. Chignell. Searching in audio: the utility of transcripts, dichotic presentation and time-compression. In *Proceedings of CHI 2006*, 2006.
- [11] Nicu Sebe. The state of the art in image and video retrieval. In *Image and Video Retrieval*, volume Volume 2728/2003, pages 1–8. Springer Berlin / Heidelberg, 2003.
- [12] Alan F. Smeaton, Cathal Gurrin, Hyowon Lee, Kieran Mc Donald, Noel Murphy, Noel O’Connor, D O’Sullivan, B Smyth, and D Wilson. The Físchlár-News-Stories system: Personalised access to an archive of TV news. In *RIAO 2004 - Coupling Approaches, Coupling Media and Coupling Languages for Information Retrieval*, 2004.
- [13] Cees G. M. Snoek, Marcel Worryng, Jan C. van Gemert, Jan-Mark Geusebroek, and Arnold W. M. Smeulders. The challenge problem for automated detection of 101 semantic concepts in multimedia. In *MULTIMEDIA '06: Proceedings of the 14th annual ACM international conference on Multimedia*, pages 421–430, New York, NY, USA, 2006. ACM Press.
- [14] Thijs Westerveld, Tzvetanka Ianeva, Liudmila Boldareva, Arjen de Vries, and Djoerd Hiemstra. Combining information sources for video retrieval: The lowlands team at TRECVID 2003. In *In Proceedings of the 12th Text Retrieval Conference (TREC-12) Video Evaluation Workshop*.
- [15] Rong Yan. *Probabilistic Models for Combining Diverse Knowledge Sources in Multimedia Retrieval*. PhD thesis, Canegie Mellon University, 2006.