

# Video Searching and Browsing Using ViewFinder: Interactive Search Experiment for TRECVID-2003

**Dan Albertson**

Laboratory for Applied Informatics Research, Indiana University, Bloomington  
1320 E. 10<sup>th</sup> St, LI011, Bloomington, IN 47405. Email: daalbert@indiana.edu

**Javed Mostafa**

Laboratory for Applied Informatics Research, Indiana University, Bloomington  
1320 E. 10<sup>th</sup> St, LI011-D, Bloomington, IN 47405. Email: jm@indiana.edu

**John Fieber**

Laboratory for Applied Informatics Research, Indiana University, Bloomington  
1320 E. 10<sup>th</sup> St. LI 011, Bloomington, IN 47405. Email: jfieber@indiana.edu

**This research project explores the topic of video information retrieval in conjunction with the task definitions and data provided by the Text REtrieval Conference's (TREC) 2003 Video Workshop (TRECVID-2003). Included in this paper, we discuss our processes and various phases in participating with TRECVID-2003. Specific sections discussed include database development, data indexing and retrieval approaches, development of user-interface and client side features, interactive search experiments, results, and conclusions.**

## Introduction

Everyday more and more video is being digitized and made available through various information systems and/or the World Wide Web. As a result, demands for video resources have increased significantly and such querying is becoming more prevalent in everyday information seeking [4]. Spink et al. (2001) observed this by examining Excite query logs over the span of 3 years (1997 to 1999) and found that searches for video content actually doubled. For these reasons, along with other similar findings, there suggests a growing importance in the exploration of problems and questions surrounding video retrieval; thus, there is a need for members of the research community to

collaborate and learn from one another through professional and academic forums (such as TREC).

In order to participate in the 2003 Text REtrieval Conference's (TREC) video workshop (TRECVID-2003), several researchers from the Laboratory for Applied Informatics Research (LAIR) at Indiana University, Bloomington developed a video retrieval system named ViewFinder. This is the second consecutive year that results from ViewFinder have been entered in TRECVID; however, numerous modifications had to be made from last year's system in order to conform to this year's participation requirements. Factors that significantly contributed to this year's system and experimental adjustments include an entirely new video and image dataset, automatic speech recognition (ASR) and closed-caption (CC) outputs (provided by the workshop), and stricter task definitions.

The problem at hand attempts to explore query modeling and user-interfaces (of video retrieval systems) by enabling users to search and browse through the assigned TRECVID data. There were several major components that we fulfilled for the purpose of exploring this problem and finalizing the TRECVID-2003 experiments.

First, we concentrated on indexing the ASR data, and applied an appropriate weight for each keyword. This data would be utilized in a keyword

search feature, in which users can formulate queries which consists of search terms of their choosing.

Next, we gave users the option to browse video data without having to perform a formal keyword search. As a result, the system offers several major headings in which the user can browse the associated keyframes. Moreover, ViewFinder's interface displays the keyframes designated for each video shot, thus allowing the user to browse by visual clues as opposed to text-based (although further textual information is available to the user upon request).

In regards to the user experiments, we participated and submitted results for 1 run which fulfilled the interactive task as defined by the workshop.

Further details of the above information (including system development and search experiments) will be discussed in the following sections. Also included will be a discussion of related research and literature, the results of the search experiment, and conclusions.

## Related Works

As previously mentioned, the problem in which we intend to explore through the TRECVID forum is the design of user-interfaces and query modeling for video retrieval systems. There have been numerous works regarding these (user) factors in video retrieval, and that will be the focus of the following section. *Note that although there is a significant number of research questions and problems that are involved in participation with TRECVID (data indexing and representation, experimental design, etc.), this section will reflect previous work that relates to our primary research interests and problem.*

In regards to previous research covering users, query modeling, and interaction with video data there are several areas in which we hope to benefit from. For example, what types of queries should be supported in video retrieval systems? A number of works discuss using low and high-level query features of video and image data (where low-level features represent content-based information and high-level features represent user's expressions in terms of a keyword search) [6].

Zhou and Huang (2002) explored the use of content-based information (low-level features) in

conjunction with keyword-based representation (high-level semantics), and conclude that low-level features are usually not sufficient enough in generating relevant results. However, they also found that when used in combination with high-level semantics, they (low-level features) can be very valuable (i.e. in associating vague keywords and concepts to certain images) [6].

However, the issue of content-based searching is far from being resolved. In fact, some researchers have found that interaction with image/video retrieval systems is directly influenced by (searching) context. For example, the usefulness of content-based searching can be influenced by whether a user has a specific information need or a general need [3]. This has encouraged some researchers to explore other non-traditional means in which to query for visual content, including spatial based querying where the users are actually allowed to formulate their query in the form of a sketch [1]. The results of such a querying model was shown promising, especially in the case where the user had a solid mental image of what their information need was [1]. *There are no immediate plans to incorporate such a feature into ViewFinder, but we consider this interesting grounds for discussion and thought.*

Now that we have discussed the problem of query modeling in video/image retrieval systems, we need to also think about how to present this information to the user (i.e. how should we develop our user-interface (UI)). Moreover, aside from what searching (querying) features to include in the UI, we need to also consider how to best present the visual search results to the users, which has also been explored in previous works.

One interesting application of this research question was to organize images according to similarity (by textual and non-textual clues), which, in a sense, arranges images to form clusters within the UI [2]. This was shown to have both positive and negative side effects. For example, sorting according to content-based information allowed for easier scanning for relevant images; however, some users reported that images had a tendency to "merge" thus possibly resulting in overlooking relevant images [2]. The same research also found that even an interface where images are randomly arranged can be useful, especially when user's don't have a specific information need in mind [2].

Additional research explores this (arrange images by similarity) concept even further, and actually conveys image similarity by the distance between the images in the UI [3]. Moreover, the actual distance between the images is very meaningful in that content-based similarity was calculated and used to form certain lengths (in centimeters) [3].

The research detailed in this section will influence our future work with ViewFinder, and how we attempt to explore these questions will be discussed in the following sections of this paper.

## Methods

Building upon previous research and experiences, we have employed certain methods in which we believe to be suitable for participating in TRECVID-2003. This section will cover specific aspects of our methodology including system development (system and client side) along with the experimental design used to carry out TRECVID-2003 experiments.

### *Data and Keyword Indexing*

Considering an entirely new data set was issued for this year’s TRECVID, the tasks of creating a database schema and indexing keywords had to be performed. The contextual data provided by TRECVID-2003 includes video and image data. More specifically, around 133 hours of video data derived from CNN Headline News, ABC World News Tonight, and CSPAN (only around 13 hours worth of CSPAN), which resulted in approximately 125 thousand keyframes (images) to represent the individual shots. All CNN and ABC video was originally broadcast during the span of January 1998 to June 1998, while CSPAN video ranges from 1998 to 2001.

Accompanying the visual data, TRECVID also provided an assortment of textual information. One such example of this data corresponds to the collection of video files as a whole (i.e. individual information regarding all video files). This data was issued in XML format in which we extracted (using Java’s XML API) and indexed (using JDBC) to form the “Video Table” (see Table 1 for database schema and corresponding attributes of Video Table).

**Table 1: Database Schema of ViewFinder**

Table Name	Attributes
Video Table	video_id, video_url, video_use, video_source, video_date, num_of_shots
Shot Table	video_id, video_filename, video_start_time, video_duration, shot_id, shot_start_time, shot_duration, image_url, time of shot
Keyword Table	video_id, shot_id, keyword, weight, freq_per_shot, freq per video, freq per dataset
Unique Terms Table	video_id, keyword, num_of_shots, idf

Next, we made use of textual data which comprised the common shot boundary directory (also issued by TRECVID). This data contained a separate XML file for each video and includes textual information corresponding to each shot (residing in each video). Considering the format of this data (XML), we parsed and indexed it in a similar fashion to the video collection data. The resulting data from this process can also be found in Table 1, and indexed under the “Shot Table”.

The last set of textual data that was indexed includes the automatic speech recognition (ASR) output (provided by TRECVID along with the video data). This data had a different format than what was previously mentioned (i.e. not XML), so different techniques had to be used to extract the keywords. This procedure included simple string comparison and modification techniques as offered through the Java API.

Embedded in this data (along with the keywords) were timestamps where the length of time and the time in which each keyword was spoken. Moreover, other tags indicated a timestamp for a certain block of time (i.e. for a “statement”) as opposed to a single time of occurrence for each and every keyword. We would use this timestamp (for “statements”) for keyword indexing and shot association purposes, and is discussed in more detail below.

The ASR output was utilized using two different approaches, and indexed accordingly. By extracting all the lines from the ASR output and comparing the timestamps (of the ASR files) with timestamps

within the shot boundary directories, we were capable indexing all the keywords and associating each with a corresponding shot and video ID. This process resulted in the formation of our “Keyword Table” (see Table 1 and Keyword Table). Note that certain timing (compliance) calculations had to be performed in order to make the two timing formats comparable.

As just mentioned, in the “Keyword Table”, all keywords were extracted, indexed, and assigned a video and shot ID. In the case that the same keyword appeared in the same shot (of the same video) the redundant use of the keyword was disregarded, but a (keyword) frequency per shot integer was incremented and indexed accordingly. Moreover, redundant keywords in a video file were still indexed; however, they are distinguished by different shot IDs and weights (which is discussed below).

Next, the ASR data was used to form a table of unique terms (see Table 1 and “Unique Terms” table). Here, each unique term was indexed per video. In this instance, if the same term appears multiple times in the same video, instead of re-indexing it, the number of shots the keyword appeared in was tracked and indexed along with the keyword.

After populating the “Keyword” and “Unique Term” tables (with the data described above) we were capable of applying certain weights to each keyword. First, an *idf* weight was given to each term located in the “Unique Terms” table. The calculation used to formulate the *idf* weight is seen directly below in Equation 1.

$$idf = \log_2(N/n)$$

$N$  = total number of shots in a video file

$n$  = total number of shots in which the term appears

**Equation 1: idf Used in ViewFinder**

Once the *idf* value for each unique term was stored, an overall *tfidf* weight was then calculated and assigned to each keyword (appearing in the “Keyword Table”). This weight consists of the product of the *idf* calculation mentioned above and the term frequency per shot (previously stored in the “Keyword Table”).

*User-interface and Client Side Features*

The graphical features and user-interface of ViewFinder were constructed and operate using Java’s Swing API. The interface itself is made up of two primary panels, which include a results display panel and a searching features (querying) panel (See Appendix A for snapshot of ViewFinder interface).

The results panel takes up approximately the left half of the interface, and has several functions associated with it. First, it is used to display keyframes of individual shots returned after the user has queried the system; thus, allowing the user to visually browse the search results. The results panel can display up to 8 keyframes (results) at a time, with results being ranked from most relevant (upper-left corner) to least relevant (bottom-right). (The displayed keyframes were generated from the images issued by TRECVID and were reduced to approximately ¼ their original size (i.e. thumbnails) for display purposes).

The results panel also offers the user several other features including the option to view further textual information regarding a specific keyframe (shot), and the option to expand upon the results of a previous search. These options are presented to the user in a series of drop down menus located directly below the 8 displayed keyframes (where each menu corresponds to the keyframe located directly above it).

The options included within the menus are “Details” and “Promote”. By selecting “Details” the system will be prompted to retrieve textual data such as video source, video date, video ID, shot ID, and a larger sized image of the keyframe (i.e. the video details) and display the information in a separate window.

On the other hand, “Promote” will retrieve the keywords associated with that particular shot (which exceeds a certain *tfidf* weighting threshold) and compare them to all the other shots in the database, then return shots which have matching keywords. Moreover, the system will perform a Boolean ‘OR’ search therefore shots which contain any of the promoted keywords will be returned. In addition, shots which have 2 or more matching keywords have their *tfidf* weights added together resulting in an overall weighting boost for that particular shot. All returned shots are then sorted and returned according to relevance (i.e. by higher shot weighting). Once a “Promote” search has been performed, the

corresponding keyframe (which has been promoted) is transferred to the middle image position (within the results panel) for visual reference for the user.

The searching/querying panel (appearing on the right-hand side of the ViewFinder interface) offers several ways in which users can formulate queries and search/browse the video data. For searching, a text box where terms can be entered and compared with the keywords indexed (in the “Keyword Table”) is available. Similar to the “Promote” search feature mentioned above, if there are 2 or more search terms in which to compare, the system will perform an ‘OR’ search; thus, returning all shots that contain any of the entered keywords. In addition, the same procedure applies when multiple keywords match for an individual shot (i.e. term weights are added together as mentioned above). *Note that considering all ASR keywords contain only capital letters. As a result,, the keyword search feature is not case sensitive as all queried terms had to be transformed to all caps for comparison purposes.*

Aside from the keyword searching function, the system also allows for certain types of video browsing. The browsing options are presented in a drop down menu appearing at the top of searching panel (top right of the ViewFinder interface). By clicking on the menu, the users can choose from video date, video source, and date + source in which to browse. After selecting one of the options, a series of choices are then retrieved and returned to the user and presented in the list box located directly below the drop down menu (*See Appendix A for snapshot of ViewFinder interface*). The user can then select one choice and hit search, which will retrieve the results and display the corresponding keyframes in the results panel.

Other features of the searching/querying panel include the “More” button which becomes available in the case that more than 8 shots are returned after a search; thus, allowing the user to browse all returned shot if necessary, and the “Back” button where the user can re-view previously viewed search results. Also, a feedback field, which will display the last performed query and the number of results returned is located in this panel.

### *Search Experiment Design*

Our experimental designed consisted of performing 1 interactive search run. The interactive run

complied with the mandatory run detailed in the participation requirements, which was to only include experiments regarding the ASR outputs. For classification purposes, ViewFinder was categorized as a ‘C’ system, as it was trained according to the methodology mentioned above, and didn’t meet the criteria of a category ‘A’ or ‘B’ system as described in the requirements.

For this run all 24 search topics (which was designated for the interactive task) was completed in a sequential order. We employed 1 search subject which completed all the topics over 2 testing sessions. The subject was given a maximum of 15 minutes in which to complete the searching topic. The overall average for each topic ended up being 10.4 minutes per topic.

Considering ViewFinder has no contextual searching capabilities, no experiments involving combining ASR data with visual data, or experiments involving exclusively visual data could be performed.

## **Results**

The results discussed in this section involve the submitted run as described above in *Search Experiment Design*, and only discusses those results which were made available from TRECVID (no other result analysis is included). The measurements of mean averaged precision, interpolated recall precision, and precision at  $n$  shots were performed by assessors at the National Institute for Standards and Technology (NIST), and can be further explored in the proceedings of TREC-10 [5].

Out of the 24 search topics designated for the interactive task, there was a total of 2067 relevant shots identified by TRECVID, in which ViewFinder (after completing all 24 search topics) ended up retrieving 282 (13.6%) of them. This came out to an average of 11.75 relevant shots per topic where a range of 58 (max) to 0 (min) was observed.

Our results can also be reflected by the mean averaged precision measured at 0.030 and by the mean precision at the total of relevant shots at 0.051. The mean precision for each search topic had a range of 0.169 (0.169 to 0.000).

Other results issued by TRECVID include the interpolated recall precision and the level of precision at  $n$  shots. The results of these two

measurements are summed up in the following table (Table 2: Summary of Interactive Search Results).

**Table 2: Summary of Interactive Search Results**

Interpolated Recall Precision		Precision at $n$ Shots	
0.0	0.5835	5	0.2250
0.1	0.0816	10	0.1333
0.2	0.0473	15	.01028
0.3	0.0047	20	0.0896
0.4	0.0006	100	0.0446
0.8	0.0006	500	0.0163
1.0	0.0006	1000	0.0118

In depth result comparison (with other systems) is not yet available for the TREC notebook paper, but will be included in the proceedings paper following the conference.

## Conclusions and Future Improvements

For this year’s TRECVID experiments, ViewFinder only made use of textual data and by analyzing the results we can draw several conclusions regarding this approach. From first glance, our *tfidf* weighting seems to be somewhat pertinent considering the number of relevant shots returned (See Results section above). However, we realize that several adjustments need to be made to our application of the algorithm.

Although we were somewhat pleased with the percentage of relevant shots returned by ViewFinder, our mean average precision obviously suffered. As a result, we are beginning to explore how to better limit the search results (i.e. attempt to only include relevant shots).

One such possibility includes incorporating a stop word list, which wasn’t used for this year’s ViewFinder. This could reduce the number of returned shots by disregarding the use of widely used terms (the, and, on, etc.). Specific characteristics of such a stop word list (one for the purpose of video retrieval) have yet to be discussed.

Next, we would like to incorporate additional Boolean search options, instead of narrowing the search to only include an ‘OR’ search. Here, users would be capable of further limiting their results by

using the ‘AND’ and ‘NOT’ operators along with ‘OR’.

Finally, we would like to make the search and browse functions of ViewFinder cross compatible. Moreover, currently with ViewFinder, each searching feature operated independently from one another (i.e. the keyword search will take precedent over the browse features if search terms have been entered into the keyword field). Also, there is no way to search within a search (i.e. once a browse function has been performed), and this could be very useful in limiting the number of returned shots.

As for contextual based searching, our initial goal for TRECVID-2003 was to also submit a run based solely on graphical features (i.e. image analysis). However, due to time constraints, we were unable to complete the image processing and database population tasks.

We still plan on exploring video retrieval in this fashion, and have done some preliminary experiments using Java’s Advanced Imaging (JAI) API. With JAI, we were capable of extracting color histogram information from keyframes (which was taken from 10 TRECVID-2002 videos) and incorporated a search by “Histogram” function into a prototype of ViewFinder. The preliminary results were somewhat satisfying, but we feel that additional content-based search features need to be incorporated (along with color histogram) to make for a practical search function. Such other content based search features may include an edge detection algorithm. By analyzing this year’s search topics issued by TRECVID, we feel that a content-based search feature is necessary to participate in future TRECVIDs, which our goal is to have such a prototype completed and functioning by TRECVID-2004.

## References

- [1] Jose, J. M., Furner, J., & Harper, D. J. (1998). Spatial querying for image retrieval: a user-oriented evaluation. *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, Melbourne, Australia*, 232 – 240.
- [2] Rodden, K., Basalaj, W., Sinclair, D., & Wood, K. (2001). Does organization by similarity assist image browsing? *Proceedings of the*

*SIGCHI Conference on Human factors in Computing Systems, Seattle, WA, 190 – 197.*

- [3] Santini, S., & Ramesh, J. (2000). Integrated browsing and querying for image databases. *IEEE Multimedia*, 7(3), 26 – 39.
- [4] Spink, A., Goodrum, A., & Hurson, A. R. (2001). Multimedia we queries: Implications for design. *Proceedings of the International Conference of Information Technology: Coding and Computing, Las Vegas, NV*, 589 – 593.
- [5] Vorhees, E. M., & Harman, D. K. (Eds.). Common Evaluation Measures. (2001). *NIST Special Publication 500-250: The Tenth Text Retrieval Conference, Gaithersburg, MD*, A14 – A23.
- [6] Zhou, X. S., & Huang, T. S. (2002). Unifying keywords and visual contents in image retrieval. *IEEE Multimedia*, 9(2), 23 - 33.

Appendix A: Snapshot of ViewFinder user-interface.

