

Applications of Image Understanding in Semantics-Oriented Multimedia Information Retrieval

Zhongfei (Mark) Zhang
Computer Science Dept.
SUNY Binghamton
zhongfei@cs.binghamton.edu

Rohini K. Srihari Aibing Rao
Dept. of Comp. Sci. & Eng.
SUNY Buffalo
{rohini, arao}@cedar.buffalo.edu

Abstract

This paper focuses on research in development of semantics-oriented multimedia information retrieval techniques. Semantics-oriented information retrieval addresses the effectiveness of the retrieval. With the goal of significantly improving retrieval precision and query capability, image understanding technologies are applied to develop effective retrieval techniques. Specifically, three techniques are developed by combining IU technologies (i.e., face detection) with other existing or developed techniques (i.e., text indexing, conventional similarity matching, and face identification) for image retrieval. These techniques are evaluated to demonstrate the significance of the semantics-oriented multimedia information retrieval.

Keywords: Multimedia Information Retrieval, Image Retrieval, Image Understanding.

1 Introduction

Multimedia information retrieval is an important and challenging research area in multimedia software engineering. Today it has become a reality that information exists in many modalities. Due to this fact, conventional methods [5] in information retrieval may not always work in all the applications, and in certain applications where information in modalities other than text is concerned, these methods are not appropriate at all. This motivates the research in multimedia information retrieval.

Ideally, a good retrieval technique must be both efficient and effective. The efficiency concerns the response time after a query is posed. The effectiveness addresses the issue of relevancy between the query and the retrieved documents in semantics. However, in the reality, there is no technique that satisfies the both constraints, due to the fact that given the current research status in media understanding (image understanding, text understanding, etc.), semantics-based

retrieval techniques are typically expensive and unreliable. Users normally do not have the patience to wait for hours in order to have a precise retrieval. Instead, they would rather have a retrieval that consists of documents that might not be all relevant to the queries, and that might not be exclusive (i.e., there might be other relevant documents in the database that fail to be retrieved); users are willing to pay the price of going through the retrieved documents for another round of manual retrieval, and the risk of missing some relevant documents. Based on this consideration, the majority of the proposed retrieval techniques focuses on the efficiency issue.

In addition, most of the existing techniques only allow queries in one modality (typically just textual key words). In many applications, queries need to be expressed in more than one modality. For example, in the law enforcement or intelligence community, we may have photos of an unknown individual. If we intend to query all the materials related to this person under certain occasions, we will have to use text to describe the contextual occasions, and to use the photos as the pointer to indicate that this is the person we are looking for. Thus, it is important that queries be expressed in different modalities at the same time.

The goal of this research is to propose a multimedia information retrieval paradigm that focuses on the effectiveness of the retrieval while at the same time the efficiency is also balanced. Specifically, we aim at significantly improving the retrieval precision, and at expanding the query capability by allowing more modalities. This is achieved by applying image understanding (IU) research in multimedia information retrieval to develop semantics-oriented techniques. We show three techniques in this effort: *face detection combined with text indexing*, *face detection combined with conventional similarity matching*, and *face detection combined with face identification*.

While in general image understanding technologies are not considered to be robust given the current sta-

tus of the research in this area, research shows [1] that under certain application domains and for certain clearly defined functionalities, image understanding technologies can deliver robust performance. On the other hand, even if an employed IU technique fails to deliver a 100% robust performance, say the performance has dropped to 50%, at least we would have half of the retrieved images that are relevant, which brings up to 50% precision. That is considered to be a very high precision in real applications.

In this paper, we define our application domains as news photos and consumer photos, and show that in these domains, robust semantics-oriented multimedia information retrieval is feasible. Since all the three techniques shown in this paper require a robust face detection system, we assume that there is such system available. In fact, in the IU community, there are many robust face detection systems that we can directly take off the shelves (e.g., [4]).

2 Face Detection with Text Indexing

In this technique, we assume that there is certain collateral textual information available accompanying the images, such as caption. This is a reasonable assumption as almost all the news photos come with caption, and many consumer photos have annotation text associated with the pictures. Thus, the individual identification is approached through verification at the text level using text indexing as opposed to that at the image level using image indexing. Consequently, no face identification at the image level is required; only face detection is conducted to verify the existence of human beings in the image. Take the example of the query *find images of Hillary Clinton in speech*. Image indexing is only performed for face detection, i.e., verification of the existence of human beings in the images; the verification of the existence of the specific individual, Hillary Clinton, is performed through verification at the text level — looking for the evidence of the presence of the name *Hillary Clinton* or the related equivalents such as *First Lady*, *Mrs. Clinton*, etc. in the collateral text. This approach intentionally avoids the processing of face identification in the image domain, and thus saves significant computation in retrieval. Due to this filtering using face detection, the precision of the retrieval is significantly improved.

This technique allows query in text mode only. Table 1 presents the results of an experiment conducted on 198 images that were from our news photo database. There were 10 queries, each involving the search for pictures of named individual(s); some spec-

ified contexts also, such as *find pictures of Hillary Clinton at the Democratic Convention*. Due to the demands of truthing, we report the results for one query; The 3 columns in Table 1 indicate precision rates using various criteria for content verification: (i) using text indexing alone, (ii) using text indexing and manual visual inspection, and (iii) using text indexing and automatic face detection. As the table indicates, using text alone can be misleading — when inspected, many of these pictures do not contain the specified face. By applying face detection to the result of text indexing, we discard photographs that do not have a high likelihood of containing faces. The last column indicates that this strategy is effective in increasing precision rates. The number in parentheses indicates the number of images in the given quantile that were discarded due to failure to detect faces.

Table 1
Retrieval Results

	Text only	Text + Man Insp	Text + Face Det
At the first 5 docs	1.0	1.0	1.0 (3)
At the first 10 docs	1.0	.70	1.0 (3)
At the first 15 doc	.80	.75	1.0 (2)
At the first 30 docs	.77	.67	NA

3 Face Detection with Conventional Similarity Matching

Similarity based matching is a typical statistical approach to multimedia information retrieval. Algorithms in this approach exploits different statistical features in the image domain for indexing the images, and the retrieval is conducted in the corresponding feature space. While the conventional similarity based matching techniques typically work well for scenery images, they do not work well for images containing human beings. To show this, we have two testing image databases, one of 58 pure scenery images and the other of 66 images containing human beings, both from our consumer photo database. Fig. 1(a) shows the performance comparison of the retrieval technique for the different databases, using color histogram as a representative similarity measure. It is obvious that the performance for the image databases containing human beings is much worse than that for the purely scenery image database.

On the other hand, in certain application domains such as consumer photos and news photos, it is typical that images contain human beings. We propose an effective solution to this problem based on the general theory we proposed earlier [6] of combining object detection with the conventional similarity based

matching techniques for better performance in complex images containing “semantic” objects such as human beings. We directly apply this theory to the special case of background similarity matching, and we only consider the object of human faces. The next section focuses on the opposite scenario, where only foreground similarity is considered in image retrieval.

This technique allows image query mode only. In order to demonstrate and evaluate the proposed approach, this new approach is compared with the conventional approach using color histogram for the same databases, and Fig. 1(b) shows the average precision-recall curves for both approaches. It is noticeable that the performance of this proposed approach is about 20% improvement over that of the conventional approach on the tested databases.

4 Face Detection with Identification

In this section, we focus on the scenario where only foreground similarity is concerned in image retrieval. Again, here we only consider human face object. This functionality is desired in many applications, including law enforcement and intelligence. Using the same face detection technique, the focus of this task is to connect it to a robust face identification technique “automatically”. Note that this technique is complementary to the technique of combining face detection with text indexing presented in Section 2. If the multimedia database is a pure image database with no collateral text, querying images of a specific individual must be done by this technique. On the other hand, if the multimedia database contains information in modalities of both text and images, we can use either of the techniques, or a combination of both techniques (e.g., use face detection with text indexing as the first round of search as it is less expensive, and then use face detection with identification to refine the retrieved documents from the first round of search). In terms of the query mode in this case, queries may be entered in text if we know the names of the specific individuals, or may be entered in images otherwise. Queries may also be entered in both text and images if the queries contain certain description of context (e.g., *find images of this person in hiking* with an image of the person in question).

Although many successful techniques have been proposed in certain specific applications, face identification, in general, still remains as a difficult and challenging problem in the IU community. Almost all the existing methods are based on a fundamental assumption – all the face images (either the images in the database or the query images) must be in the same

dimensions, so that the feature vectors are all in the same feature space; if face images are in different dimensions, which is a common situation in typical image retrieval scenarios, a normalization in scale to a pre-defined space must be conducted. While normalization in scale is a feasible and a common approach, it may lose precious information which otherwise may be used to better achieve identification precision; normalization may also cause extra computation. Based on these considerations, a new approach to face identification, *subspace morphing* (SM) theory [8], is developed. SM theory allows direct comparison between two face images in different dimensions, and thus, completely eliminates the need of normalization in scale. The technique presented in this section is the direct application of the SM theory to the scenario of retrieving images containing the same individuals (i.e., the foreground similarity retrieval).

In news and consumer photos, it is common that human faces are “buried” under a background of other objects or scene. Since SM allows direct comparison between face images in different dimensions, we can first apply face detection to an image to automatically crop out all the face candidates, and then feed them directly to SM for identification without necessarily normalizing the cropped face candidates in scale. This approach is called *streamlined approach*. Fig. 2 shows an example of this approach. This new technique is evaluated using a model database of 112 samples of 25 individuals, and a testing database of 226 samples of the same 25 individuals. All the face samples are collected in different scales and in different poses to show the varying appearances. The evaluation is done against the four classic identification methods: nearest neighbor based eigenface (NE) [7], linear combination based eigenface (LE) [3], nearest neighbor based SVD (NS) [2], and linear combination based SVD (LS) by combining [2] and [3]. Since all the classic techniques require normalization in scale, the 338 face images are scaled to a common dimension space of 120×80 to generate a separate version of the two databases. Table 2 shows the performances between SM and the four classic face identification techniques. Each image in the testing database is taken as a query to the model database, and the first retrieved image from the model database is taken as the identified image to report the identification rate for each of the evaluated techniques. Given the complicated variations of the appearance in the testing images w.r.t. the model images, SM is obviously better than the four conventional identification techniques.

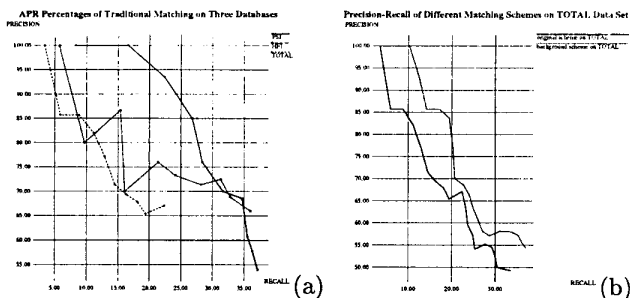


Figure 1: (a) Average Precision-Recall for different databases where the number of retrieved images in each query is 10. (b) Average Precision-Recall for different matching schemes on a combined database consisting of pure scenery images and consumer images where the number of retrieved images for each query is 20. PSI: pure scenery image database, HFI: image database containing human beings, TOTAL: combined database.

Table 2

Face Identification Results

NS: Number of testing images succeeded
 NF: Number of testing images failed
 IR: Identification rate

Methods	SM	NE	LE	NS	LS
NS	208	92	104	102	56
NF	18	134	122	124	170
IR	92%	41%	46%	45%	25%

5 Conclusion

This work is aimed at developing techniques for semantics-oriented multimedia information retrieval. Specifically, image understanding (IU) techniques are used in specially defined application domains to deliver semantics-oriented multimedia information retrieval. Three techniques are demonstrated using IU technologies combined with other indexing tools for effective image retrieval. Experimental results show that these proposed techniques significantly improves the retrieval precision at the reasonable cost of efficiency, and significantly expands the query capability by allowing more modalities.

References

[1] D.A. Forsyth. Computer vision tools for finding images and video sequences. *Library Trends, Special Issue on Progress in Visual Information Access and Retrieval*, 48(2):326–355, 1999.

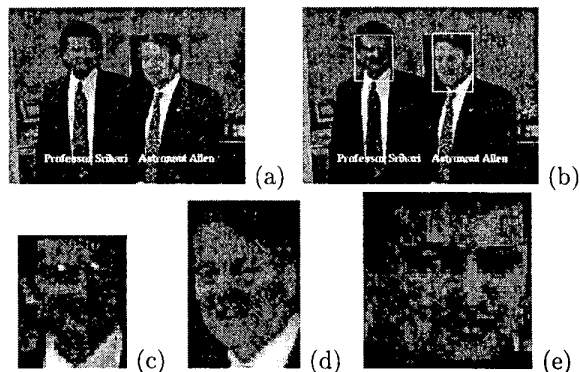


Figure 2: An example of the result of the streamlined approach to face identification. (a) an original color image. (b) face detection result. (c) the first face candidate returned for face identification. (d) the second face candidate returned for face identification. (e) another face image of the same individual as in (c).

[2] Z. Hong. Algebraic feature extraction of image for recognition. *Pattern Recognition*, 24, 1991.

[3] S.Z. Li. Face recognition based on nearest linear combinations. In *Proc. CVPR*. IEEE, 1998.

[4] H.A. Rowley, S. Baluja, and T. Kanade. Neural network-based face detection. *Trans. on PAMI*, 20(1), 1998.

[5] G. Salton. *Automatic text processing: the transformation, analysis, and retrieval of information by computer*. Addison Wesley, 1989.

[6] R.K. Srihari, Z. Zhang, and A. Rao. Image background search: Combining object detection techniques into content-based similarity image retrieval(cbsir) systems. In *Proc. of IEEE International Workshop on Content-Based Access of Image and Video Libraries*. IEEE Press, June 1999.

[7] M. Turk and A. Pentland. Eigenfaces for recognition. *Journal of Cognitive Neuroscience*, 3(1), 1991.

[8] Z. Zhang and R.K. Srihari. Subspace morphing theory for appearance based object identification. In *Proc. of International Conference on Advances in Infrastructure for Electronic Business, Science, and Education on the Internet*, L’Aquila, Italy, 2000. Scuola Superiore G. Reiss Romoli.