

Challenging Semantic Concept for Systemic Component of CBIR.

Tohid Aribi*

Department of Electrical Engineering, Miandoab Branch, Islamic Azad University, Miandoab, Iran.

Short Communication

Received: 22/02/2013

Revised: 05/05/2013

Accepted: 12/07/2013

*For Correspondence

Department of Electrical Engineering, Miandoab Branch, Islamic Azad University, Miandoab, Iran.

Keywords: Multimedia, Databases, Medical, Informatics Applications

ABSTRACT

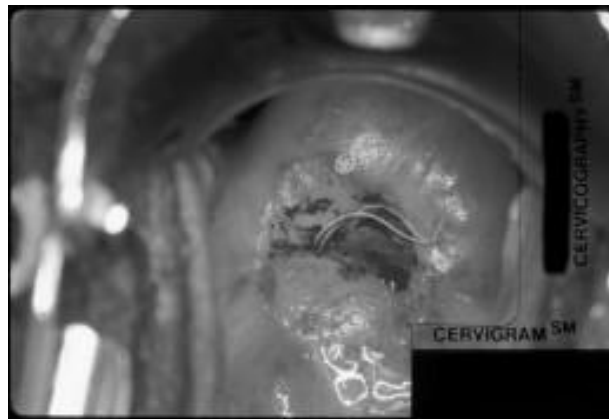
Advances in medical imaging have led to growth in large image collections. We are conducting research on CBIR for biomedical images. We maintain an archive of over 17,000 digitized x-rays of the cervical and lumbar spine from the second National Health and Nutrition Examination Survey (NHANES II). In addition, we are developing an archive of a large number of digitized 35mm color slides of the uterine cervix. Our research focuses on developing techniques for hybrid text/image query-retrieval from the survey text and image data. In this paper we present the challenges in developing CBIR of biomedical images and results from our research efforts.

INTRODUCTION

A common drawback of such systems is that the annotations are imprecise with reference to image feature locations, and text is often insufficient in enabling efficient image retrieval. Even such retrieval is impossible for collections of images that have not been annotated or indexed. Additionally, the retrieval of interesting cases, especially for medical education or building atlases, is a cumbersome task. CBIR methods developed specifically for biomedical images could offer a solution to such problems, thereby augmenting the clinical, research, and educational aspects of biomedicine. For any class of biomedical images, however, it would be necessary to develop suitable feature representation and similarity algorithms that capture the "content" in the image [1,2,3,4,5,6,7].

The Lister Hill National Center for Biomedical Communications, a research and development division of the U.S. National Library of Medicine (NLM), maintains a digital archive of 17,000 cervical and lumbar spine images collected in the second National Health and Nutrition Examination Survey (NHANES II) conducted by the National Center for Health Statistics (NCHS). Classification of the spine x-ray images for the osteoarthritis research community has been a long-standing goal of researchers at the NLM, and collaborators at NCHS and the National Institute of Arthritis and Musculoskeletal and Skin Diseases (NIAMS). Also, the capability to retrieve these images based on geometric characteristics of the vertebral structures is of interest to the vertebral morphometry community. Medical experts have identified visual features of the images specifically related to osteoarthritis, but the images have never been manually indexed for these features which include anterior osteophytes, disc space narrowing for the cervical and lumbar spine, spondylolisthesis for the cervical spine, and spondylolisthesis for the lumbar spine. Another archive of 100,000 digitized 35mm color slides of the uterine cervix is being created in collaboration with the National Cancer Institute (NCI) (figure 1). Researchers at NCI would like to enable use of these images for research and training at sites around the world. The design of a system to achieve these ends relies on research in image compression, database management, and CBIR for image query on the uterine cervix images. Automated or computer-assisted classification, query, and retrieval methods for large medical image archives are highly desirable, since such methods offset the high cost of manual classification and manipulation by medical experts. We are investigating automated or computer-assisted methods that use image features for indexing and retrieval of these images in a manner acceptable to the biomedical community. In addition, we are devoting research efforts into classification of pathology, such as the detection of presence of anterior osteophytes, disc space narrowing, spondylolisthesis in spine images; and squamo-columnar junction boundary, regions with acetowhitening, vasculature, mosaicism and punctation, on the uterine cervix images. As an initial step, we have implemented a modular prototype CBIR system for a subset of the spine x-rays and the associated health survey text data [6]. The system supports retrieval based on shape similarity to a sketch of a complete or partial vertebra, an example vertebral image, as well as conventional text retrieval. In this paper we present the technical

considerations in developing a system for CBIR of medical images, open research problems, and the lessons learned from our research efforts.



(a)

Figure 1: Example images of uterine cervix.

The CBIR Trail

Content-based image retrieval hinges on the ability of the algorithms to extract pertinent image features and organize them in a way that represents the image content. Additionally, the algorithms should be able to quantify the similarity between the query visual and the database candidate for the image content as perceived by the viewer. Thus, there is a systemic component to CBIR and a more challenging semantic component. As a first step, the images must be indexed, at least, for the pathologies of interest.

INDEXING TRAIL

The indexing trail has been presented from a systemic viewpoint. A graphic user interface (GUI) of the indexing system allows the users to index the text and image data. In indexing images, visual features that correspond to the pathology of interest are segmented (extracted) from the image. Shown as the “Segmentation” block in the figure, this step is synonymous with “Feature Extraction”. The output of the segmentation step is usually in the form of image components such as subimages, edges, boundary contours, color/intensity measurements, texture measurements, etc. Feature extraction is usually done at the local region of interest. In case of the digitized spine x-ray image collection, the only features of interest are the shape of the vertebra and the positional relationship with other vertebrae. At the end of the segmentation step the resulting data is a vector of 2D coordinate points describing the vertebra boundary outline. Segmentation techniques include variants of active contour segmentation [7] and active shape modeling [3]. In the case of the uterine cervix images, the extracted features include, in addition to shape, color and texture measurements on homogeneous regions that are determined using a *k*-Nearest Neighbor classifier. The segmented features, then, need to be represented in a form suitable for indexing and similarity computation. This task is done at the “Feature Representation” stage. “Feature Vector Computation” is often coupled with this stage.

The output of these stages is a feature vector that is often in a form invariant to rotation, translation, scale, and also possesses many other properties such as *stability, uniqueness, etc* [6]. Image similarity is, then, defined as the distance between the feature vectors for two images. Also, each feature representation algorithm may have to use a corresponding similarity measure. A side effect of feature representation is loss of information incurred due to approximation which is done for algorithmic or efficiency reasons, or to avoid the *curse of dimensionality*.

Representations, such as histograms or color averages, approximated boundaries, are often sufficient to enable some form of CBIR and are found in many prototype systems discussed in the literature [4]. The cost in loss of representation of subtle variations in image features, however, can lead to poor retrieval quality. Storing higher dimension feature vectors, while enabling query of subtleties in image content, can cause problems for indexing, creating a Catch-22 situation. We are experimenting with a variety of shape representation techniques for segmented vertebrae that include polygon approximation, Fourier descriptors, shape properties, invariant moments, and Procrustes distance [6, 7].

An outstanding problem in the extraction of feature vectors from the raw boundary data is development of an effective shape representation and similarity method that provides for data reduction while simultaneously preserving the shape characteristics that are essential for reliable indexing and retrieval.

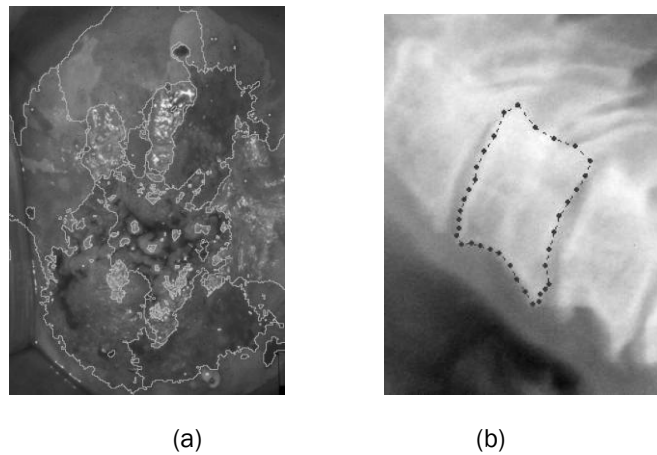


Figure 2: Examples of segmented features: (a) c-spine vertebra boundary; (b) regions on the uterine cervix

Image content represented by feature vectors is then indexed. Unlike in a traditional database, however, it is difficult to develop unique keys from the feature vectors. One approach is to use hierarchical cluster trees (figure 2). This approach links images to leaf nodes in a tree. It then clusters “similar” images together and assigns their cluster centroids as parent nodes. This process is repeated on these centroids until only one remains. Such a hierarchical organization strategy can be very efficient and significantly reduces image search times. In addition, it supports both *target queries*, where one matching image is sought, as well as *range queries*, where images that match a certain feature measurement range are sought. There are, however, some shortcomings with this approach. First, it requires that the similarity measure be a metric and most effective similarity measures are relative. Second, the index tree is optimized for a single type of query, e.g., in spine x-ray images, the tree might be optimized for queries on anterior osteophytes. For other query types a new index tree would be necessary. This limits the types of queries possible on a dataset and is not directly helpful to the long term goals of CBIR. As an initial step, however, we have adopted this approach [4] for organizing indexing trees and optimizing the node structure with the spine x-ray images shapes. In general, organization of image features for CBIR is an open research problem. The survey text data that accompanies the images is indexed in a traditional RDBMS. Our current implementations do not link images indexed in the hierarchical cluster tree to the RDBMS text data, though such an approach is conceivable. Currently, we link the image to the text data by the image name.

Open Problems

While we have developed reasonably functional solutions for segmentation and representation of vertebrae, these remain open problems for the uterine cervix images. Organization of feature vectors comprising multiple features also remains an open problem. Specifically, it is necessary to decrease the dependence of hierarchical cluster trees on similarity distance metric. Finally, an important problem little discussed in the literature but of much importance, is that of validation of retrieval results. For example, how can we justify calling one set of shape retrieval results better than another? How can we compare results among different shape representations and similarity measures? The validation of the query results in either a quantitative sense or with a non-quantitative approach that will justify confidence in the results using a particular method remains a critical issue for this work. Beyond the important issue of what we have called “engineering validation” of results, there remains the further issue of biomedical validation, for the biomedical community is the system end user. There is a critical need for more extensive expert data to enable development of better algorithms. A key requirement of these data sets is that they should be collected by multiple expert observers; only then can the performance of computerized methods relative to human performance be evaluated. With the NCI uterine cervix image set, an open problem is the technique for combining multiple image features. The images are rich in color, exhibit texture in pathology and also have boundaries. Techniques for effective and efficient combination of features need to be developed.

CONCLUSIONS

This paper presents our experiences, techniques adopted, and in continuing research towards enabling CBIR for large medical image archives. The paper presents the CBIR trail and presents some results from our work in each task on the trail and open research problems.

REFERENCES

1. Antani S, Kasturi R, and Jain R. A survey on the use of pattern recognition methods for abstraction, indexing and retrieval of images and video. *Pattern Recognition*. 2002;35(4): 945-65.

2. Tagare HD, Jaffe CC, Duncan J. Medical image databases: a content-based retrieval approach. *J Am Med Inform Assoc.* 1997; 4(3):184-98.
3. U Sinha, A Ton, A Yaghmai, RK Taira H Kangarloo. Image Content Extraction: Application to MR Images of the Brain. *Radio Graphics.* 2001; 21(2): 535-547
4. ME Mattie, L Staib, E Stratmann, HD Tagare, J Duncan, PL Miller. Path Master: Content-based Cell Image Retrieval Using Automated Feature Extraction. *J Am Med Inform Assoc.* 2000; 7(4): 404-415.
5. Long LR, Antani S, Lee DJ, Krainak DM, Thoma GR. Biomedical information from a national collection of spine x-rays: film to content-based retrieval. *Proc SPIE Med Imaging: PACS and Integrated Med Info Sys: Design and Evaluation, 15-20 Feb 2003, San Diego, CA, SPIE Vol. 5033: 70-84.*
6. Antani S, Long LR, Thoma GR. A biomedical information system for combined content-based retrieval of spine x-ray images and associated text information. *Proc. 3rd Indian Conf on Computer Vision, Graphics, and Image Proc., Ahmedabad, India, 16-18 Dec 2002, 242-47.*
7. Kass M, Witkin A, Terzopoulos D. Snakes: Active contour models. *Int J Computer Vision.* 1987;1:321-331.