# AN IMAGE RETRIVAL ANALYSIS BY USING A PATTERN SIMILARITY SCHEME

## Dr. Madhu Babu Sikha

Malla Reddy Engineering College (Autonomous) Department of Electronics & Communication Engineering

**Abstract—**
In this paper, we propose a novel scheme for efficient content-based medical image retrieval, formalized according to the PAtterns for Next generation DAtabase systems (PANDA) framework for pattern representation and management. The proposed scheme involves block-based low-level feature extraction from images followed by the clustering of the feature space to form higherlevel, semantically meaningful patterns. The clustering of the feature space is realized by an expectation–maximization algorithm that uses an iterative approach to automatically determine the number of clusters. Then, the 2-component property of PANDA is exploited: the similarity between two clusters is estimated as a function of the similarity of both their structures and the measure components. Experiments were performed on a large set of reference radiographic images, using different kinds of features to encode the low-level image content. Through this experimentation, it is shown that the proposed scheme can be efficiently and effectively applied for medical image retrieval from large databases, providing unsupervised semantic interpretation of the results, which can be further extended by knowledge representation methodologies. Forum (CLEF). Since 2004, a medical image retrieval task has been added. Goal is to create databases of a realistic and useful size and also query topics that are based on real{world needs in the medical domain but still correspond to the limited capabilities of purely visual retrieval at the moment. Goal is to direct the research onto real applications and towards real clinical problems to give researchers who are not directly linked to medical facilities a possibility to work on the interesting problem of medical image retrieval based on real data sets and problems. The missing link between computer science research departments and clinical routine is one of the biggest problems that becomes evident when reading much of the current literature on medical image retrieval. Most databases are extremely small, the treated problems often far from clinical reality, and there is no integration of the prototypes into a hospital infrastructure. Only few retrieval articles specifically mention problems related to the DICOM format (Digital Imaging and Communications in Medicine) and the sheer amount of data that needs to be treated in an image archive ($> 30:000$ images per day in the Geneva radiology ).

## 1. INTRODUCTION

ONE of the primary tools used by physicians is the comparison of previous and current medical images associated with pathologic conditions. As the amount of pictorial information stored in both local and public medical databases is growing, efficient image indexing and retrieval becomes a necessity. During the last decade, the advances in information technology allowed the development of content-based image retrieval (CBIR) systems, capable of retrieving images based on their similarity with one or more query images. Indicative examples of such systems are QBIC [1], SIMPLicity [2], and FIRE [3]. It is interesting that more than 50 CBIR systems are surveyed in [4].

The benefits emanating from the application of content-based approaches to medical image retrieval range from clinical decision support to medical education and research [5]. These benefits have motivated researchers either to apply generalpurpose CBIR systems to medical images [3] or to develop dedicated ones explicitly oriented to specific medical domains. Specialized CBIR systems have been developed to support the retrieval of various kinds of medical images, including highresolution computed tomographic (HRCT) images [6], breast cancer biopsy slides [7], positron emission tomographic (PET) functional images [8], ultrasound images pathology images and radiographic images.

Common ground for most of the systems cited earlier is that image retrieval is based on similarity measures estimated directly from low-level image features. This approach is likely to result in the retrieval of images with significant perceived differences from the query image, since low-level features usually lack semantic interpretation. This has motivated

researchers to focus on the utilization of higher-level semantic representations of image contents for content-based medical image retrieval. Recent approaches include semantic mapping via hybrid Bayesian networks semantic error-correcting output codes (SECC) based on individual classifiers combination and a framework that uses machine learning and statistical similarity matching techniques with relevance feedback. However, these approaches involve supervised methodologies that require prior knowledge about the dataset and introduce constraints to the semantics required for the image retrieval task. Content{based visual information retrieval (CBVIR) or content{based image retrieval (CBIR) is an extremely active domain in the multimedia and computer vision _elds [1, 2, 3, 4]. An ever{increasing amount of multimedia data (images, video, music, ...) is produced and made available in digital form. Almost every modern computer user has most of its hard disk _lled with multimedia data (images, video clips, mp3 music, ...) but tools to manage these data well are scarce. Most web pages become increasingly mixed{media documents integrating images, animations, texts, etc. The medical _eld is no exception to this trend. There is an increasing amount and variety of visual data being produced for the diagnostic process and the role of images in the diagnostic process is increasing. Currently, these visual or multimedia data are mainly used for the treatment of a single patient, only. Much of the diagnostic process of medical doctors (MDs) is based on comparing a current case with experience from past cases. To support the memory concerning images, many medical doctors store interesting or typical cases with a textual description and the images on their hard disk or in a teaching _le such as myPACS1 or casimage 2 [5]. Having a larger source of images and descriptions available for all medical doctors can make this stored information and experience available to a larger audience, but the rising number of images requires good tools to not only store the data. Quick search and retrieval tools are needed for these growing databases to _nd relevant information quickly. Then of course, tools are necessary to anonymise the images as the use of images out of the pure diagnostic or treatment planning process is often not permitted, even within a single institution.

## 2. AXES OF RETRIEVAL EVALUATION

This section explains several of the axes that we regard as important for creating the tasks for ImageCLEF to satisfy various research directions but also to stick to our goal by creating a research environment to prepare medical image retrieval for the use in a real{world setting. Much of the outline and form of the ImageCLEF evaluation is based on the experiences of the TREC workshops and will not be detailed in this article.

### 2.1 User– vs. system–centered evaluation

User{centered evaluation is evaluating how a user judges the results of an information retrieval system. This includes more than only technical aspects as the user judges what he receives as a result interactively, and a large number of factors together influence the user's judgement on the entire retrieval system. Query speed and ease of use and layout of the interface are extremely important (an example interface for visual queries can be seen in Figure 1). On the other hand, the evaluation can be subjective as several users might judge the same result in a different way. Even the same user might judge the same result differently at different times. User{centered evaluation is also relatively \expensive" as it does include the time of real system users and cannot be automated. Each new setting of parameters requires a new interaction circle with the users. System{centered evaluation is less costly as it can be automated and does not necessarily require user interaction. Normally, query topics are formulated in advance, and then system developers can tune their system and submit results that are subsequently evaluated against a ground truth, which is usually created after submission. This means that a large number of system variations can be evaluated with low cost but on the other hand only a part of the system parameters is taken into account, the technical parameters, and important parts such as query speed and the user interface are not analysed at all. Both TREC and CLEF run mainly system{centered tasks.

### 2.2 Visual vs. textual vs. mixed retrieval

One of the  first questions regarding image retrieval is to choose whether a purely textual image retrieval based on available meta data is planned or whether visual data is to be used for the retrieval [1]. Based on the chosen application scenario, only one or the other is really possible. If only very limited meta data is available for retrieval and if many images do not contain any annotation, a keyword search will not be successful but a search with an image example can allow navigation in the database. If good meta data is available text allows to search for semantics and concepts which is usually what a user is looking for. Purely visual retrieval is currently limited to extremely simple concepts and a fairly limited

number of concepts as well. On the other hand, visual content and textual context of the images are most often very complementary. Even if the query is only in one media, the other media can be used in a combined visual/textual approach to improve the final results.

## 2.3 Multilingual vs. monolingual retrieval

Most experience in information retrieval is de nitely available on monolingual and mostly on English retrieval. Still, in elds such as web search a large number of users existwho might want to use a query language other than English but still retrieve English documents. Most image collections are actually understandable without the text, so searching in a multilingual collection for images is also possible, even if the language can not be understood. In multi{lingual environment such as the European Union or Switzerland, multi{lingual information retrieval is indispensable.

## 2.4 Classification vs. Information retrieval

An often discussed topic is whether information retrieval is basically the same thing as classification or not. Often, we can see an information retrieval problem as a two{class problem with the class of relevant and the class of non{ relevant items maybe with a third class of partially relevant items, and without having any learning data. Still, in most cases, when we think about information retrieval, we have very large collections in mind on which we do not have have much information concerning the content, groups of images or documents, etc. Then, we would like to satisfy the information need of a user and nd documents that (s)he is interested in for a particular search. Through the use of frequency{based feature weights some information on the distributions of words or features within the database are extracted in an automated fashion. Judgement of the entire collection for relevance is often impossible due to the large size, so incomplete relevance sets are often based on pooling methods.
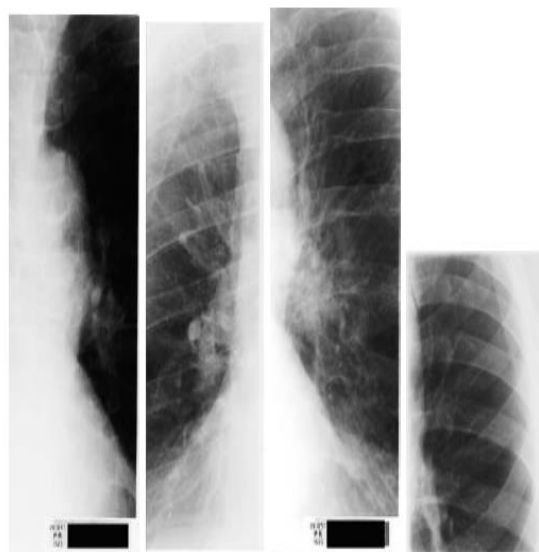


Figure 1: Images representing one of the smallest classes in the IRMA task of ImageCLEF 2005. Most often for classication, information on class membership of the entire collection is known and well defend, which allows the use of machine learning techniques and system optimisations. An example for images belonging to the same classes is Figure 1 taken from ImageCLEF 2005. To evaluate algorithms there are several methods that are commonly used based on the available training data. Leaving{

one{out means that algorithm training is done on all images but the image under test, making available a maximum of test data. The process is repeated such that all images serve once as tests, and the mean error rate over all experiments can be determined. Classification error rate can be used as performance measure for these completely annotated databases .

## 2.5 Object recognition vs. visual appearance

These two fields are both very active in the domain of computer vision for a variety of application, and both can be very beneficial for image retrieval. Whereas object recognition tries to identify a generally limited number of concepts or objects in an image and label them by techniques such as template matching.

Similarity search by visual appearance in contrast to this takes into account either global features representing the

entire image or features representing the layout of an image such as a smaller representation of the image

itself. Segmentation can also give access to visual appearance search based on regions.

## 3. RESOURCES MADE AVAILABLE

One of the biggest problems when working on medical image analysis is the access to data. As all images are patient data, we need to be careful with them to respect their privacy and everything used for research needs to be anonymised carefully. The advent of the digital radiology and cheap storage capacities have made the exchange and

sharing of images much easier than in the film{based days. Teaching _les are created in many medical institutions and quite a few of these are made available publicly. One of the larger initiatives to publish images on the Internet is the MIRC12 (Medical Image Resource Center) project. In this project, a common access method to teaching fles is created based on the XML standard. Software for clients and servers is made available free of charge and cross{platform in the form of a Java program. Currently, more than 15 databases are accessible in this format to be searched by keywords via the MIRC web page.



Show me all x–ray images showing fractures.
Zeige mir Röntgenbilder mit Brüchen.
Montres–moi des radiographies avec des fractures.

Figure 2: A query that requires more than visual retrieval but visual features can deliver some hints to good results as well.

Still, often images are only stored on local hard disks and much knowledge could be extracted from these images if they were available. One of the databases that is accessible via MIRC is the casimage dataset that contains almost 9.000 images of 2.000 cases and that was used in the ImageCLEFmed 2004 competition [5]. It is also part of the 2005 collection.

Images present in the data set include mostly the radiology department, but also photographs, powerpoint slides and illustrations. Cases are mainly in French, with around 20% being in English. We were also allowed to use the PEIR13 (Pathology Education Instructional Resource) database using annotation from the HEAL14 project (Health Education Assets Library, mainly pathology images). This dataset contains over 33.000 images with English annotation, with the annotation being in XML per image and not per case as casimage. The nuclear medicine database of MIR, the Mallinkrodt Institute of Radiology15 was as well made available to us for ImageCLEF. This dataset contains over 2.000 images mainly from nuclear medicine with annotations per case and in English. Finally, the PathoPic16 collection (Pathology images) was included into our dataset. It contains 9.000 images with an extensive annotation per image in German. Part of the German annotation is translated into English, but it is still incomplete. This means, that a total of more than 50.000 images was made available with annotations in three different languages. Two collections have case{based annotations whereas two collections have image image{based annotations. Only through the access to the data by the copyright holders, we were able to distribute these images to the participating research groups. The automatic annotation task was organised by the IRMA group and based on their datasets. This database is an-notated according to the four{axes IRMA code. To simplify the task in the first year of existence, a subset of 57 classes was chosen that all have at least 5 images in the class. The database contains a total of 10.000 images. 9.000 images representing the 57 classes were given out with class labels as training data. The remaining 1.000 images were given to participants without a class label for classification. The IRMA code in English and German was also made available to the participants.

## 4. APPLICATION OF THE AXES

### 4.3.1 User vs. system–centered

ImageCLEF has an interactive (user{centered, non{medical) task since 2004, but participation is still fairly low containing s2{5 submissions, mostly due to the high cost of user involvement and the lack of experience in this domain. The task measures how many steps a user needs to find several images by keyword search and using relevance feedback. Still, most of the tasks are clearly system-centered, and all the medical tasks currently are.

### 4.3.2 Textual vs. visual vs. mixed

ImageCLEF covers all three fields but has a main focus on mixed retrieval as this is a field where still a lot of research is needed and much less experience is currently available. To ease such a combination, visual retrieval results were made available and in the next year it is planned to make also textual retrieval results available to all topics for participants mainly working in one of the two fields. In 2004, the medical task had an image as query, only, as shown in Figure 3, whereas the ad hoc query task was a text accompanied by a single image. In 2005, a purely visual medical image annotation task was added (IRMA task). On the other hand, the medical retrieval task contains one or several images plus text in three languages (English, French, German) and has thus a small visual component. Several topics are expected to be solvable with a visual system such as the example in Figure 4, whereas other topics are more semantic and text processing appears to be necessary. This focus towards more semantic queries was based on critics in 2004 with the goal to have more realistic topics that are useful in a clinical setting. The 2005 topics are based on a real user survey among medical professionals.

### 4.3.3 Multilingual vs. monolingual

The medical task in 2005 models the scenario of a collection in several languages, currently English, French and German. This is also a fairly common and realistic scenario as medical doctors often annotate their cases in their mother tongue, whereas they might understand enough in another language as well to use the images of a case. Thus, for the medical retrieval task 2005, query topics were made available in the same three languages as the collection, and queries also contain one or several query images (Figure 4). Techniques for multilingual retrieval include the translation of the queries to a unique language, translations of the documents or the extraction of concepts in multilingual ontologies such as MeSH (Medical Subject Headings).

### 4.3.4 Classification vs. information retrieval

In the context of ImageCLEF, the classification task is actually called automatic annotation task, which is a very similar problem because the classes actually correspond to a text that can be added to the image collection. The IRMA code [45] to which the classes correspond actually exists in several languages, so such a classification and annotation can further{on be used for multilingual retrieval as well. We distribute a learning set of images and then an evaluation set that the evaluation is performed on, so participants have no idea about class memberships of the images

to be categorized but can use the entire training data for system optimisation. The main retrieval task is a typical information retrieval task with 25 query topics that correspond to an information need of a user from a very large data set. The relevance judgements are done on the first $N = 40$ images of all system submissions so results stay comparable even if relevance is not judged on the entire dataset. As training data, only the topics from 2004 were made available that were not really corresponding to the 2005 topics and underline the character of an information retrieval task.
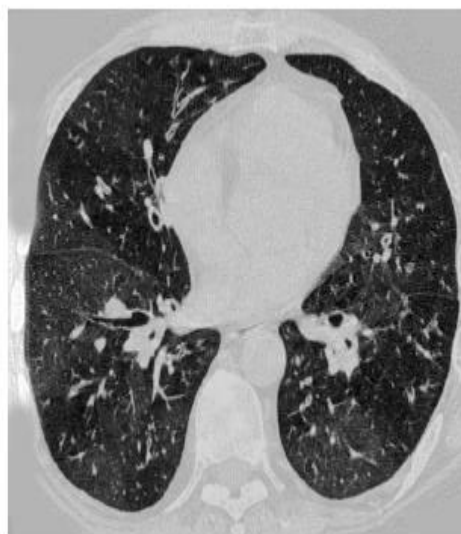


Figure 3: An example query from the 2004 medical task, with the goal to retrieve all images of the same anatomic region, viewing angle and modality. Here, all lung CTs independent of the pathology are expected as result.
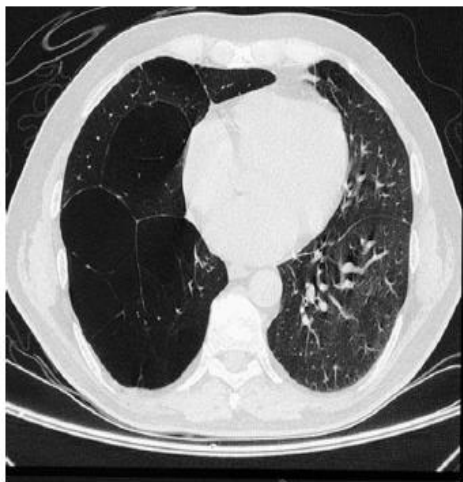
Figure 4: An example of a query that is solvable visually, using image and text as query. Still, the use of the annotation can augment the retrieval quality. The query text is presented in three languages.

### 4.3.5 Object recognition vs. visual appearance

In ImageCLEF 2005, both of these techniques have very useful applications and can well improve retrieval quality. A typical example for an object recognition topic can be seen in Figure 7, where all images showing faces are wanted as a response. For several other queries, object recognition can be useful through very specific detectors but in general the variability of medical images in our database and the variability of query topics is extremely large and constructing one detector per topic is tedious. Thus, for most of the topics, query by visual appearance can deliver overall acceptable results in addition and as complement to the textual queries, although query by visual appearance is much less specific. Many of the queries are very hard for object recognition as well as for search by visual appearance, which makes the use of text important to complement the two. Whereas object recognition can be important if almost no annotation is available to extract semantics, the visual appearance is important where textual information is available. This can for example be used to rank images within a group of semantically related images, such as ranking all images with a text containing the word emphysema based on the similarity with a lung CT.

### CONCLUSION

We presented a novel scheme for efficient content-based medical image retrieval. This scheme utilizes rich-in-semantics *pattern* representations of medical images, defined in the context of PANDA, a framework for representing and handling data mining results. The theoretical contributions of this paper are validated by comprehensive experimentation on the IRMA reference collection of radiographic images. The results advocate both its efficiency and effectiveness in comparison with state of the art.

Future perspectives of this paper include: 1) systematic evaluation of the proposed scheme for the retrieval of various kinds of medical images, such as endoscopic [29] and ultrasound images [43] according to their pathology; 2) the enhancement of the retrieval performance by using image indexing techniques based on specialized data structures; and 3) the integration of the proposed scheme with ontology-based information extraction and data mining techniques for the retrieval of medical images using heterogeneous data sources. By storing the semantically rich patterns along with low-level features in a unified way, according to the PANDA framework, will enable the extension of the CBIR methodologies with knowledge representation techniques for semantic processing and analysis.

### REFERENCES

[1] Arnold W. M. Smeulders, Marcel Worring, Simone Santini, Armarnath Gupta, and Ramesh Jain. Content{based image retrieval at the end of the early years. IEEE Transactions on Pattern Analysis and Machine Intelligence, 22 No 12:1349{1380, 2000.

[2] Yong Rui, Thomas S. Huang, Michael Ortega, and Sharad Mehrotra. Relevance feedback: A power tool for interactive content{based image retrieval. IEEE Transactions on Circuits and Systems for Video Technology, 8(5):644{655, September 1998. (Special Issue on Segmentation, Description, and Retrieval of Video Content).

[3] Serge Belongie, Chad Carson, Hayit Greenspan, and Jitendra Malik. Color{ and texture{based image segmentation using EM and its application to content{based image retrieval. In Proceedings of the International Conference on Computer Vision (ICCV'98), pages 675{682, Bombay, India, 1998.

[4] Amarnath Gupta and Ramesh Jain. Visual information retrieval. Communications of the ACM, 40(5):70{79, May 1997.

[5] Antoine Rosset, Henning M•uller, Martina Martins, Natalia Dfouni, Jean-Paul Vall_ee, and Osman Ratib. Casimage project { a digital teaching _les authoring environment. Journal of Thoracic Imaging, 19(2):1{6, 2004.

[6] Henry J. Lowe, Ilya Antipov, William Hersh, and Catherine Arnott Smith. Towards knowledge{based retrieval of medical images. The role of semantic

indexing, image content representation and knowledge{based retrieval. In Proceedings of the Annual Symposium of the American Society for Medical Informatics (AMIA), pages 882{886, Nashville, TN, USA, October 1998.

[7] Stelios C. Orphanoudakis, Catherine E. Chronaki, and Despina Vamvaka. I2Cnet: Content{based similarity search in geographically distributed repositories of medical images. Computerized Medical Imaging and Graphics, 20(4):193{207, 1996.

[8] Hemant D. Tagare, C. Ja_e, and James Duncan. Medical image databases: A content{based retrieval approach. Journal of the American Medical Informatics Association, 4(3):184{198, 1997.