# A case study of Image Retrieval on Lung cancer chest X-ray pictures

Gile Narcisse Fanzou T.[a], Wang Ning[a], Nathalie Cindy K.[c], François Siewe[b], Lin Xudong[a], Xu De[a]

[a] School of Computer & Information Technology
Beijing Jiaotong University
100044 Beijing China
fanzounar2002@yahoo.fr

{nwang,dxu,linx}@bjtu.edu.cn

[b] Software Technology Research Laboratory
School of Computing
De Montfort University
Leicester LE1 9BH
UNITED KINGDOM

fsiewe@dmu.ac.uk

[c] Department of Computer Science
University of Yaoundé 1
812 Yaoundé
Cameroon
nathkuicheu@yahoo.fr

**Abstract:** This paper presents a case study of an image retrieval system based on a notion of similarity between images in a multimedia database and where a user request can be an image file or a keyword. The CBIR (Content Based Image Retrieval) system, the current System of Search for Information (SSI) --e.g. PEIR, MIRC, MIR, IRMA, and Pathopic-- and the Current Search Engines (CSE) --e.g. Google, Yahoo and Alta Vista-- make image search possible only when the query is a keyword. This type of search is limited because keywords are not expressive enough to describe all important characteristics of an image. For example, an exact match request cannot be formulated in such systems and in SSI system, users should know natural language (e.g. English, French or German) used. We used XIRS (an XML Image Retrieval System) to set up a similarity distance between images, then to compare the request image with those in a database. An experimentation of XIRS on lung cancer diagnosis is presented. The statistics show that our system is more efficient than leading CBIR systems such as ERIC7, PEIR, PathoPic and CSE.

**Keywords:** XML, Image retrieval, similarity search, diagnosis, web, Medical Information systems.

## I. INTRODUCTION

Users start with *information needs*, which they translate into *query representations*. Similarly, there are *documents*, which are converted into *document representations*. The role of an Information Retrieval (IR) system is to extract from the document representations the information needed by the users and stated in the query representation. The purpose of the search process is to obtain user's needs from a database by comparing the user's requirements with available information. This comparison is carried out by a System of Search for Information (SSI) [3], which is a set of programs with the goal to return to the user the maximum relevant documents available that meet his needs.

The SSI, CBIR (Content Based Image Retrieval) system and the CSE (Current Search Engine) make image search possible only when the query is a keyword. This type of search is limited because these keywords are not expressive enough to describe all important characteristics of an image. To resolve this problem, ERIC7 [6] which is a CBIR system compatible with the MPEG-7 Multimedia standard proposed to the user to search images by features. Hence, in ERIC7 the user can choose between 15 features by navigating within XML files using a tool that generates UML diagrams. However, ERIC7 is limited because the user should be an expert in search for images to recognize these features. He should also be able to read and understand XML files and UML diagrams. We also observe that an exact match request cannot be formulated in such systems.

MPEG-7[8] is a standardization of XML metadata structures called Descriptors (D) and Description Schemes (DS), which are used to describe and annotate multimedia information [11]. The Ds and DSs are defined using the MPEG-7 Description Definition Language (DDL), which is based on the XML Schema Language. Many technologies still need to be developed around the MPEG-7 for extracting, searching and querying multimedia databases, which involves similarity matching including features, content and semantics.

In this work, using XIRS[4], we present a case study of image retrieval in which a request might be an image file or a keyword. We describe an image as an XML document using MPEG-7 standard. We have defined a similarity distance between images which is used to compare the features of the request image to images stored in a multimedia database. The statistics show that our system is more efficient than leading content based image retrieval systems such as ERIC7, PEIR, Pathopic and the CSE. Posting an image for the similarity search in a Database can have an importance in Hospitals to find the diagnosis of the radiographic stereotypes [2], and also used to implement iconic communication systems [7]. As application, an assistant software for lung cancer diagnosis is presented.

This paper is organized in the following way: Section 2 describes the XIRS system; Section 3 is devoted to the case study of image retrieval on lung cancer diagnosis and the discussion.

## II. XIRS (XML IMAGE RETRIEVAL SYSTEM)

This section gives a brief description of XIRS. Readers are referred to [4] for a full presentation of the system. XIRS is a set of 3 components: the XIRS Mediator, the interrogation module, and the XIRS Server. Starting from the feature extraction and annotation process of a multimedia asset, the XML documents are generated and stored in a repository.

### II.1 XIRS Mediator

An image is represented as a set of descriptors (features) which are structured as XML nodes and stored in a XML document (see Figure 1). The image is stored in a multimedia database and the XML document is then stored in the XML repository. The XML document used by the XIRS Mediator is obtained by combining two parts:
- «Visual Descriptors» extracted from the image by MPEG-7,
- «Metadata descriptors»: the XML document is completed with some information describing the semantic and contents (e.g. keywords, its author, its size…) coming from the database.
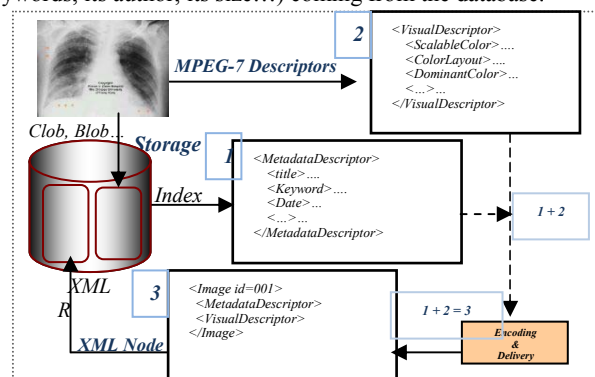


**Figure 1. XIRS Mediator**

A fix DTD is used by XIRS Mediator to construct XML documents. Once XIRS Mediator has described an image in XML node, the node is categorized (to prevent too bulky XML documents) and stored in XML documents of the collection. The role of XIRS Mediator is thus to define an image in XML and vice versa.

## II.2 Interrogation Module

The data model of the XIRS interrogation module is a simplification of XPath data model presented in [1], where a structured document is a tree, composed of simple nodes, sheet nodes and attributes. A node can be a document, an element, a text, a namespace, an instruction or a comment. Two cases of request arise.

### II.2.1 The request is a keyword

A request is a conjunction of sub-requests. We have the following illustration:

Request → sub-request AND sub-request | sub-request OR sub-request | NOT sub-request.

Hence, the similarity distance between an image node $N_I$ and a request node $N_q$ is defined as:

$$\Phi(N_q, N_I) = \begin{cases} \dfrac{|N_q|}{|N_I|} & \text{if } N_q \text{ matches } N_I \\ 0 & \text{otherwise} \end{cases}$$

Where $N_q$ (resp. $N_i$) is the number of sub-nodes and $|N_q|$ (resp. $|N_I|$) is the number of sub-nodes+1 in the query node and image node respectively. $N_q$ matches $N_I$ iff $N_I$ belong to the set described by $N_q$. Note that if $N_q$ matches $N_I$ then $|N_q| = \Phi(N_q, N_I) |N_I|$ and $\Phi(N_q, N_q) = 1$.

### II.2.2 The request is an image

The comparison between an image and a request amounts calculating a score. The image relevance with respect to the request is calculated by a similarity function noted d(q, I), where q is the

request image and I is an image of the Database. It thus leads to calculate a similarity distance between two XML nodes. Lets I =

$(I_1, I_2,…,I_m)$ an image set and $T = (t_1, t_2,…, t_n)$ a keyword set. We describe the image Ij as a vector : $I_j = (w_{1,j}, w_{2,j}, . . . , w_{i,j},…, w_{n,j})$ where $w_{i,j} \in \{0, 1\}$ is the term-weighting. $f_i$ denote the function that returns the associated weight of the term ti : $f_i(\vec{I_j}) = w_{i,j}$

The XML node produced by the XIRS Mediator and corresponding to the request image is regarded as a block of requests (like a system of equation with several unknown factors), in which each sub-node (features) is seen as a request. It is thus a question of reassuring when one has a node coming from a XML document of the Database that both sub-nodes are similar.

If a feature of an image is indexed by tj and if tj < tk then it is also indexed by tk. Therefore, one can extend the vector Ii so that: $w_{k,i} = 1$ if $w_{j,i} = 1$ and tj < tk, otherwise $w_{k,i} = 0$. The usual similarity measure used in XIRS is given in *Formula 1*, where $q_n$ and $S_n$ are XML nodes representing the query image and one image of the database. $S_{sn}$ and $q_{qn}$ are sub-nodes of $S_n$ and $q_n$ respectively. V is the vocabulary of non-structural terms; weight $(S_n,t,S_{sn})$ is the weight of term 't' in XML context $S_{sn}$ in node $S_n$.

The XIRS grammar gives a complete description of the request language used. The axiom of the grammar is **Query**, non-terminal symbols are in **bold**, terminal symbols (tokens) are in *italic* and the production rules are described as follow(see table 1)

---

**Query → r1 | r2**
**r1 → ExpressionA ExpressionB**
**ExpressionA →** *keyword* **SuiteExpressionA** | *( keyword )* **SuiteExpressionA**
**SuiteExpressionA → ExpressionA** | ε
**ExpressionB → BooleenOperator r1** | ε
**BooleenOperator →** *OR* | *AND* | *NOT* | ε
**r2 → ExpressionStructure SuiteExpressionStructure**
**ExpressionStructure →** *elementName[* **Condition** *]*
**Condition →** *@attributName = keyword* | **r1** | ε
**SuiteExpressionStructure → BooleenOperator ExpressionStructure** | ε

**Caption:**
*ε* denotes an empty string
*keyword*: terminal symbols representing a keyword
*elementName:* terminal symbols representing a name of tag
*attributName:* terminal symbols representing a name of attribute

**Table 1. XIRS Grammar**

---

$$d(q_n ; s_n) = \sum_{s_{sn} \in S_n} \sum_{q_{qn} \in q_n} \Phi(s_{sn}, q_{qn}) \sum_{t \in V} weight(q_n, t, q_{qn}) \frac{weight(S_n, t, s_{sn})}{\sqrt{\sum_{\substack{s_{sn} \in S_n \\ q_{qn} \in q_n \\ t \in V}} weight(S_n, t, s_{sn}) \times weight(q_n, t, q_{qn})}}$$
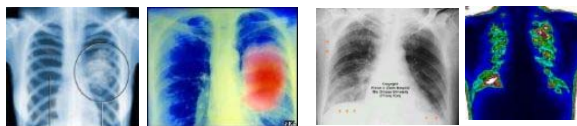
**Formula 1. Similarity distance between two nodes.**



**Figure 2. Lung Cancer**          **Figure 3. Lung**

## III. CASE STUDY AND DISCUSSION

### III.1 Application on lung cancer diagnosis.

### III.1.1 Interface (see Figure 4).

For the experiment, we applied our system on diagnosis search, especially on lung cancer. Lung medicine presents a lot of diseases and each disease has its own chest X-ray and diagnosis.

Our target when setting up this decisional software is to help users (doctors, medical students/researchers or patients) to check if they have a lung disease by analyzing their chest x-ray pictures (see Figure 2 and 3) to produce a diagnosis containing their possible treatments and the way to avoid a lung disease.

For over 100 years, The CLA (Canadian Lung Association) [2] has been dedicated to its mission of promoting and improving lung health. According to the CLA, there are about 39 lung diseases : *Acute bronchitis, Asbestosis, Asthma, Avian flu, Bronchiectasis, Bronchitis, Bronchopulmonary dysplasia (BPD), Chronic cough, Severe acute respiratory syndrome (SARS), Lung cancer, Tuberculosis, etc.* A **Chest x-ray** exams can help the Doctor to confirm if a patient has or not a lung disease. The final diagnosis depends on many tests. - *Medical history, Sputum analysis, Bronchoscopy, Needle Biopsy and Mediastinoscopy…*
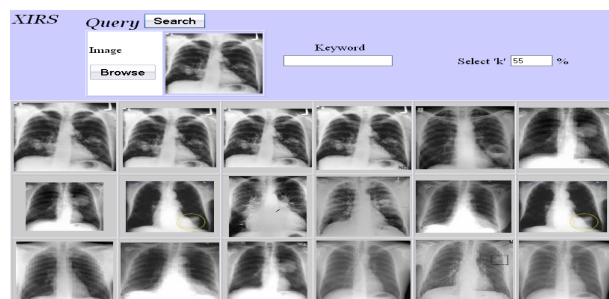


**Figure 4. XIRS Interface, the request is an image file.**

**Figure 5. Specifications when one's double click a picture.**

## III.1.2 Implementation

Our application is quiet simple to use. Users post an image (*a chest x-ray picture for example*) or enter a keyword and specify a degree of similarity k, then our system returns a set of images similar to the post image (or related to the entered keyword) according to k (see Figure 4). One can double click on any picture of the result set to show up the specification of the picture. Specifications are the information about the disease presented by the picture which could be the disease name, signs, symptoms, treatment, and the way to prevent it (see Figure 5).

The present application is based on tree participants: *Part1- an XML-Enabled data source*, Oracle 8i, in which we store images and information about the disease specification; this data source contains about 1200 Images and 15 XML documents. *Part2- Web server:* built with Oracle 8i, Apache, and PHP 5.0 for web pages management and XIRS interface. *Part3- Software package*: the XIRS mediator, the interrogation module, ConstS and CalWeight.

### III.1.2.1 Example of XML Document

For an example, let us concentrate on the image (**Chest x-ray** of a Lung cancer patient (see Figure 6).
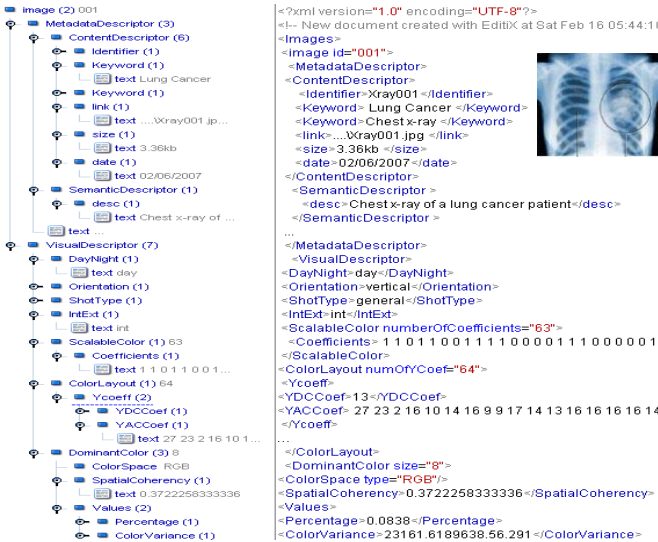


**Figure 6. XML node representing the**

### III.1.2.2 XIRS Principle: Search for images by similarity

The image request is a node; it is a question of returning all the nodes of the XML documents of the collection which are similar to the request node according to a precision ''k''. (See Figure 7)

*To return efficiently results, XIRS uses two algorithms, one to construct the set of results, ConstS* (Sn, W[iq], k), *and another to calculate* $W_{iq}$, *CalWeight* (Ssn, W[j,f]). These algorithms are made using these definitions:

**Definition 1:** Two XML nodes are k-similar if 'k' percent of their sub-nodes (features) are identical.

**Definition 2:** A node belongs to S *(request Set of results),* iff this node is K-similar to the node described by the request image, ie: if $d(q_n, S_n) \geq k.$
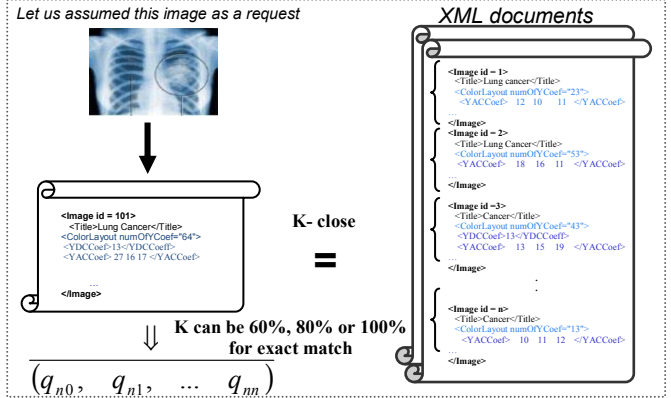


**Figure 7. XIRS principle.**

A similarity distance "d" between two nodes is defined by: *let's use the case of the first node of our XML file as show on the figure 7.*

$$d : N \times N \to D$$

$$\left( \begin{pmatrix} s_{n0} \\ s_{n1} \\ \dots \\ s_{nn} \end{pmatrix} , \begin{pmatrix} q_{n0} \\ q_{n1} \\ \dots \\ q_{nn} \end{pmatrix} \right) \mapsto d\left( \begin{pmatrix} s_{n0} \\ s_{n1} \\ \dots \\ s_{nn} \end{pmatrix} , \begin{pmatrix} q_{n0} \\ q_{n1} \\ \dots \\ q_{nn} \end{pmatrix} \right) \qquad d\left( \begin{pmatrix} LungCancer \\ 05 \\ \dots \\ 12.10.11 \end{pmatrix} \begin{pmatrix} LungCancer \\ 13 \\ \dots \\ 27.16.17 \end{pmatrix} \right)$$

*Note that function "d" is calculated according to formula 1*

As we can see, ( $s_{n0}$ , $s_{n1}$ , … , $s_{nn}$ ) =(Lungcancer, 05, …, 12 10 11) is a first node (*representing the features of the image*) coming from XML documents of our database, N is a set of Nodes and D is a set of distances. The image request (see Figure 7) being an XML node, ( $q_{n0}$ , $q_{n1}$ , … , $q_{nn}$ ) =(Lungcancer, 13, …, 27 16 17) are fixed and are query sub-nodes; ( $w_0$ , $w_1$ , … , $w_n$ )=(100%, 38.46%, …, 70.7%) are weight (*similarity distance between features*) associated to the sub-node $s_{nl}$ compared to the request $q_{nl}$ with $l \in [0, n]$, '*l*' is the number of sub-nodes of a given node. *In fact, the sub-query here is "Lung cancer" and $d_{Lungcancer}(Lungcancer)=1$.*

Example, if k=100%, $\quad w_0 = \dfrac{s_{n0}}{q_{n0}} \quad = \dfrac{Lungcancer}{Lungcancer} = 100\%$

### III.2 Discussion

When working on medical image (patient data) analysis, the access to data is really a problem. From the Internet, there are many institutions which publish images [5]:

- The PEIR (Pathology Education Instructional Resource), their database use annotation from the HEAL project. This dataset contains over 33.000 pathology images with English annotation, the annotation being in XML per image. *http://peir.path.uab.edu/*

- The MIR (Mallinkrodt Institute of Radiology). This dataset contains over 2.000 images mainly from nuclear medicine with annotations per case and in English. *http://gamma.wustl.edu/home.html*

- The PathoPic collection (Pathology images). It contains 9.809 images with an extensive annotation per image in German. *http://alf3.urz.unibas.ch/pathopic/intro.htm*

- The MIRC (Medical Image Resource Center) project. Cross-platform is available. Currently, more than 15 databases are accessible to be searched by keywords via the MIRC web page. One of the databases is the "casimage" dataset that contains almost 9.000 images in French. *http://mirc.rsna.org/*

- The IRMA. This database of 10.000 images is annotated in English and German, it is organized in subset of class.

These databases available from Internet are all using only keywords to query images. In this case, users should initially know the language. Because in English, the French word "Cancer

du poumon" is "Lung Cancer". Hence, the keyword is limited for those systems. Using XIRS, users could query a database using a posting image or a keyword, if the query is a posting image, XIRS will returns a set of images as result of his query. So the user does not have to know the database language.

### III.3 Experiments

The evaluation of XIRS was conducted on a computer with Intel Pentium 4 clocked at 3.00 GHz(2CPUs), 100 Gb of hard disk and 1520 Mb of main memory. The O/S was Windows XP SP2. We used XIRS, ERIC7, Google Search, Yahoo Search, PEIR and PathoPic. We measured **i)** *the precision of retrieval (percentage of similarity between the query and the result (PR))*. **ii)** *Time of results*. We believe that a better and more accurate measure could be achieved by using these metrics. We considered the types of queries used in basic processing operations of search.

#### i)The precision of retrieval (PR)

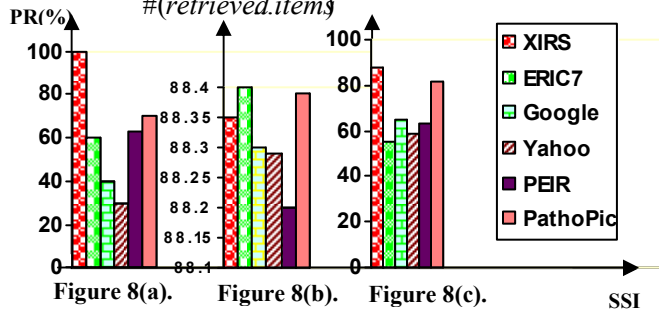$$PR = \frac{\#(relevant.images.retrieved)}{\#(retrieved.items)} = PR(relevant | .retrieved)$$



**Figure 8(a).    Figure 8(b).    Figure 8(c).**

- ✓ **Exact match search** (see Figure 8a): when the value of k is equal to 100%, XIRS returns only the XML nodes identical to the XML node of the request image and thus the returning images are the one identical to the image request. In 100 images returned by ERIC7, 40 are totally different to the request image depending of the features given by the user. In CSE, 70% of returned images are not similar. PathoPic returns 70.11% exact images when PEIR returns 62.92 %.
- ✓ **Full text search** (See Figure 8b): the PR of ERIC7 and PathoPic is closed to 88.4% while that of XIRS is 88.35%, due to the database clustering done by ERIC7 and the classification in PathoPic.
- ✓ **Semantic search** (See Figure 8c): XIRS is about 35 % more efficient than ERIC7, due of the semantic descriptors insert in the XML Nodes by XIRS Mediator. PathoPic is closed to XIRS because of the multiple annotations of PathoPic databases.

#### ii) T*ime results*

The *Response time(R)* is discussed in this section. We used the keyword **"lung cancer "** as a query**.**
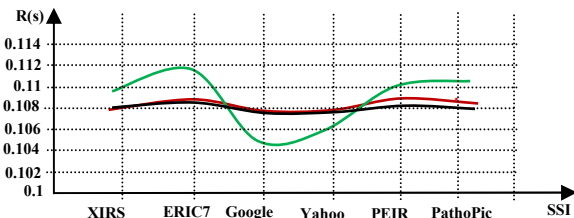


**Figure 9. Time results of the query "lung cancer".**

- ✓ **Exact match search**(see black chart): According to the chart, Google is a little bit fast(**0.1070** seconds) with 15 results pages of about 320**,**000 images where only **212 are exacts** for the query "lung cancer", when XIRS gives back 2 result pages of 96/1200(*96 images over 1200 images*) **exact images** in 0.1078 seconds, ERIC7 returns 5 result pages of 850/1200 images with only 25 exact images in 0.1080 seconds.

- ✓ **Full text search** (See blue chart ):As we can see from the chart, Google stills the faster (0.105s) with 21 result pages of 531,000

images, when XIRS return 4 result pages of 165/1200 images in 0.110 seconds and ERIC7 gives 7 result pages of 1000/1200 images in 0.112 seconds.

- ✓ **Semantic search** (See red Chart):The chart shows that our six search systems are running at the same response time; Google returns 19 result pages of 431,000 images at 0.1080s, when XIRS return 2 result pages of 102/1200 images in 0.1080s seconds and ERIC7 gives 3 result pages of 600/1200 images in 0.1082 seconds.

In general, XIRS is running a bit fast than PEIR, PathoPic and ERIC7. Google returns results a little bit fast but those results (images) are not all close to the query, the similarity distance between returning images and the query is very high.

### IV. CONCLUSION

In this paper, we have presented a case study of image retrieval when the request is an image file or a keyword. The user has the possibility to formulate his requirements in information using a given precision K.  We used XIRS to define a similarity distance between two images by defining the similarity between two XML nodes representative the two images.  An evaluation of XIRS shows the effectiveness of this system towards the CBIR systems, the Current Search Engines (e.g. Google and Yahoo), PEIR and PathoPic as for the search for images.

As future works, this system should let users give another query on the set of results; hence, they could have a result close as possible to their needs. The consideration of more than one image in a request (e.g. iconic sentences) and the consideration of heterogeneous sources of images (data every where and pay as you go system) are also very important for this type of system.

### V. ACKNOWLEDGMENTS

### VI. REFERENCES

[1] Boag S.  and al. (Eds), XQuery 1.0 : An XML Query Language, *W3C Working draft*, 2003.

[2] Canadian Lung association *http://www.poumon.ca/diseases-maladies/a-z_f.php*, 1750 Courtwood Crescent, 300 Ottawa.

[3] Fanzou T.G.N., XIRL : XML Information Retrieval , *Mémoire de DEA, University of Yaounde 1*, Cameroon, 2006.

[4] Fanzou T. Gile N., Xu De, Wang N., Francois Siewe, XIRS: an XML-Based Image retrieval system, *WSEAS International Conference on MULTIMEDIA, INTERNET & VIDEO TECHNOLOGIES (MIV '07)*, Beijing-China September 17-18, 2007, pp. 233-238.

[5] Henning M., Paul C. Evaluation Axes for Medical Image Retrieval Systems — The ImageCLEF Experience; *MM'05-ACM,* November 6–11, 2005, Singapore.

[6] L. Gagnon, S. Foucher, V. Gouaillier, ERIC7: An Experimental Tool for Content-Based Image Encoding and Retrieval under the MPEG-7 Standard, R&D Department, *Computer Research Institute of Montreal*, 2004.

[7] N. C. Kuicheu, P. L. Fotso, F. Siewe. Iconic Communication System by XML Language (SCILX). *Proceedings of the 2007 ACM International Cross-Disciplinary Conference on Web Accessibility, Banff, Canada,* 7-8 May 2007.

[8] Kosch H., MPEG-7 and Multimedia Database Systems, *SIGMOD Record, Vol. 31, No. 2*, June 2002.