# A Web Database for Computer-Aided Detection and Diagnosis of Medical Images

Dave Tahmoush and Hanan Samet

University of Maryland, College Park, Computer Science Department,
College Park, MD, USA
`tahmoush@cs.umd.edu`

**Abstract.** Building effective Computer-Aided Detection and Diagnosis (CAD) systems involves the combination of running experiments, image markup, security, analysis, evaluation, feature extraction, and feature combination in order to capture and evaluate medical images effectively. This requires the involvement of a large community of experts across several fields. We have created a CAD development system which integrates the community together without requiring disclosure of sensitive details. The system design enables researchers to upload their feature sets and quickly compare the effectiveness of their methods against other stored feature sets. Additionally, research into the techniques used by radiologists is possible through double-blind radiologist comparisons based on their annotations and feature markups. This research archive contains the essential technologies of secure transmission and storage, textual and feature searches, spatial searches, annotation searching, filtering of result sets, feature creation, and bulk loading of features, while creating a repository and testbed for the community.

**Keywords:** CAD, medical image database.

## 1 Introduction

This Breast cancer remains a leading cause of cancer deaths among women in m any parts of the world. In the United States alone, over forty thousand women die of the disease each year [1]. Mammography is currently the most effective method for early detection of breast cancer [7]. For two-thirds of the women whose initial diagnosis of their mammogram is negative but who actually have breast cancer, the cancer is evident upon a second diagnosis of their mammogram [7]. Computer-aided detection (CAD) of mammograms could be used to avoid these missed diagnoses, and has been shown to increase the number of cancers detected by more than nineteen percent [4]. Improving the effectiveness of CAD could improve the detection of breast cancer, and could improve the survival rate by detecting the cancer earlier.

Detecting breast cancer is challenging because cancerous structures have many features in common with normal breast tissue. This means that a high number of false positives or false negatives are possible. Improving CAD can help reduce the number of false positives so that true positives are more obvious. The majority of work on CAD analysis of mammograms has focused on determining the contextual similarity

to cancer, finding abnormalities in a local area of a single image [5,9]. CAD work has used methods ranging from filters to wavelets to learning techniques, but a detailed discussion of various imaging techniques is beyond the scope of this paper. Problems arise in using filter methods [5] because of the range of sizes and morphologies for breast cancer, as well as the difficulty in differentiating cancerous from non-cancerous structures. The size range problem has been addressed by using multi-scale models [9]. Similar issues affect wavelet methods, although their use has led to reported good results [6] with the size range issue being improved through the use of a wavelet pyramid [8]. Learning techniques have included support vector machines [2] and neural networks [6].

There are several different types of research that go into developing an effective CAD system for medical images. The primary research is done by radiologists, who perform the medical scans as well as provide diagnoses. This system can help radiologists by organizing their images, capturing their patient notes and digitizing their image annotations, speed up the analysis of experiments, and enable quick comparisons of different radiologist techniques such as comparing double-reading to single reading of mammograms. Web publication of research into radiologist techniques can be simplified using built in annonimizing and web publishing. The next stage of research is feature extraction and analysis. An example of this is the measurement of spiculation [10] as a feature used to aid in the detection of spiculated lesions. The system can help with research into feature extraction by providing annotated images, comparison features, and comparison results sets, as well as data analysis and feature combination. Once the images have features and diagnoses associated with them, all of the pieces are available for research into CAD techniques.

The rest of this paper is organized according to the steps in the process of creating a CAD system. Section 2 discusses the data collection and analysis by radiologists. Section 3 reviews the unique feature upload, storage, and sharing capabilities. Section 4 draws conclusions and discusses future work.

## 2   Data Collection and Analysis

Our system is designed to be able to run and analyze double-blind studies of radiologist techniques. It can be configured so that multiple radiologists can annotate the same image under different conditions without viewing the biopsy-based ground truth or the other annotations. Their input can then be viewed by a user which can only view the doctor information after it has been automatically annonimized to maintain the integrity of the study. The data can be analyzed in several ways. A particular data set can be isolated using the dataset name in the image annotations. This data set can be filtered to show only particular types of cases, such as malignant cancers, normals, or benign cancers. A particular (annonimized) doctor's diagnosis can further refine the results, showing the percentage of accurate diagnoses on different classes of images. The position of the actual cancer can be compared to the biopsied "truth" position using a spatial search, finding all of the cases where the doctor's diagnosis is within a specified distance of the "truth" position. The analysis can then be finished in a matter of hours. The results can then be stored, and the study can be published by allowing guest user access.
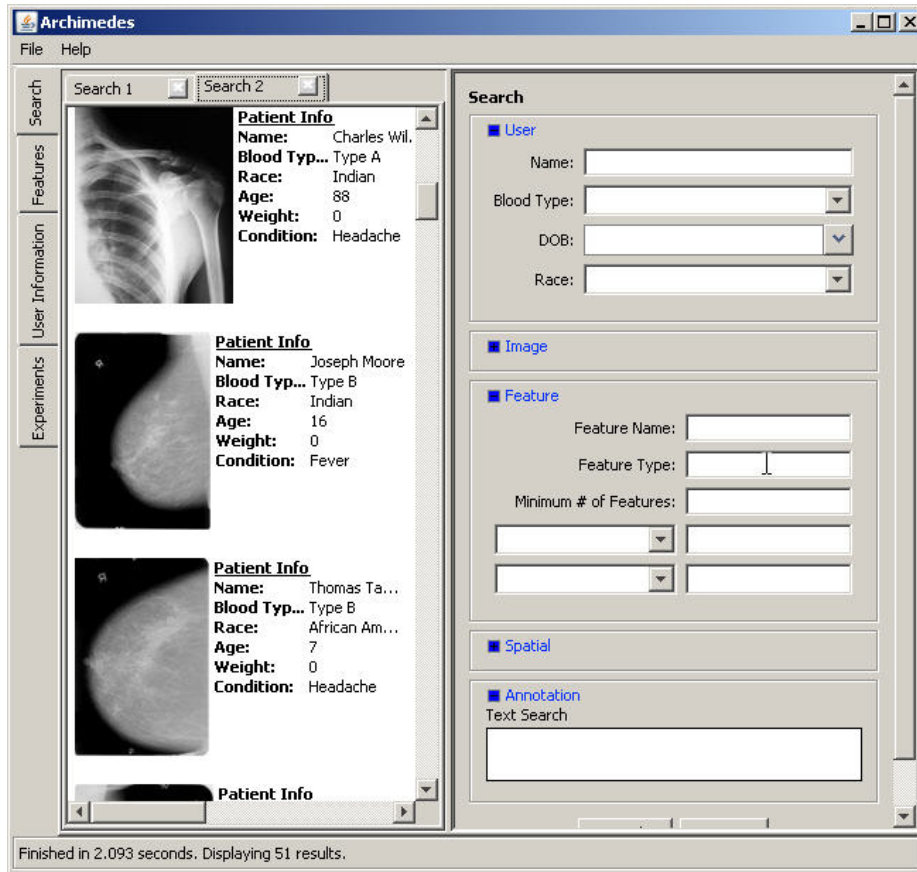
**Fig. 1.** The database GUI. The database searching can be done over the particular user, image characteristics, features or clusters of features, as well as spatial and annotation searches. The images are returned as thumbnails with some info, and then the particular image of interest can be viewed in greater detail.

Our system was originally designed to do double-blind research studies, but it also has an image analysis and patient records management tool that can be used by the medical research community. The graphical user interface for the database is shown in Figure 1. Radiologists can store and organize medical images such as x-rays, mammograms, CAT scans, MRIs, and any other image that is stored in a DICOM format. Radiologists can rapidly retrieve images and patient records, and can also find patients with similar images, conditions, or annotations to compare treatment successes. The software archives the addition of markups and notations to images, as well as associating text and patient info with images. The image GUI is shown in Figure 2, and the zooming effect is shown in Figure 3.
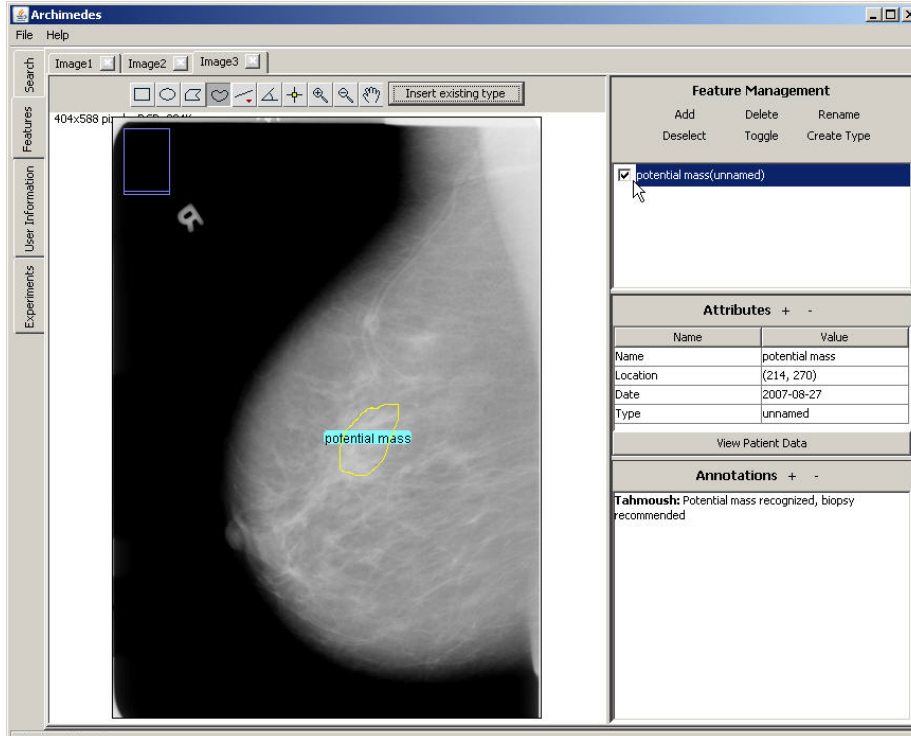
**Fig. 2.** The image GUI. The image set can be viewed and annotated in full-screen mode, and features can be outlined or drawn with a variety of tools. Features types can be created and reused for other images.

Medical professionals can view and manipulate images on the system and tab through set of images. There is a zooming interface in order to focus in on interesting parts of the images. Point features can be inserted and described as overlays to the images. Text annotations can be entered. Multiple overlays can be captured for each image, allowing double reading of images. This type of detailed image information is essential for the design and evaluation of CAD systems. The radiologists' diagnoses are captured and the data can be accessed remotely, thereby allowing tele-medicine applications.

The system helps store, annotate, and retrieve data and images for the radiologists, who are the primary data collection agents in the development of a medical image CAD system. Improving their work and giving them incentives is one of the key steps in creating a collaborative environment to create CAD. AN overview of the system design is shown in Figure 4.
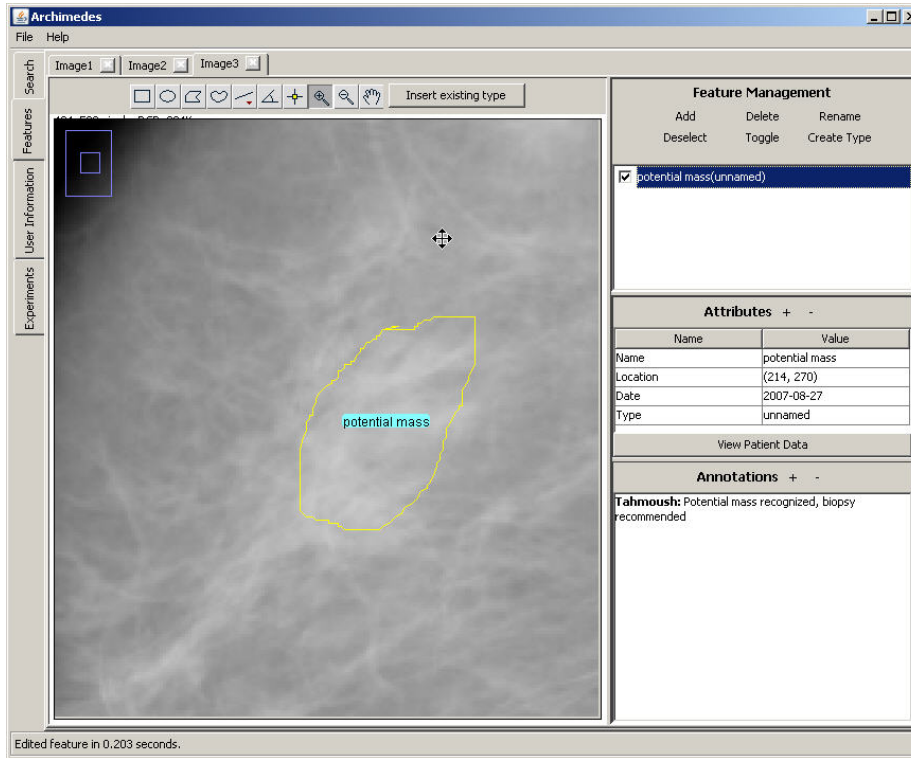
**Fig. 3.** The zoomed image GUI. The image details can be viewed using a zooming tool to examine the details of the image. Note that the other images in the set are accessible through tabs across the top.

## 3  Features

One of the main challenges in feature extraction is finding a large enough set of images of the same exact type of cancer in order to focus in on its particular characteristics. But there are a few large databases that do provide these images, for example with lung cancer images [3]. However, in order compare the effectiveness of one feature versus another on the same images, the comparison research has to be replicated. This problem is eliminated by allowing researchers to store their features for comparison along with the or independently from the images from which they were extracted. The input, storage, and sharing of features is one of the design choices that make the system unique. The open sharing of features makes comparisons possible without the need for inaccurate replication of older work, as well as enables research into feature combination both faster and more effective. Features can be combined using spatial search and then fused into a new feature type. A spatial search can be used to quantify the effectiveness of the feature at predicting the position of the cancer by comparing the feature position with the biopsied "truth" cancer position.
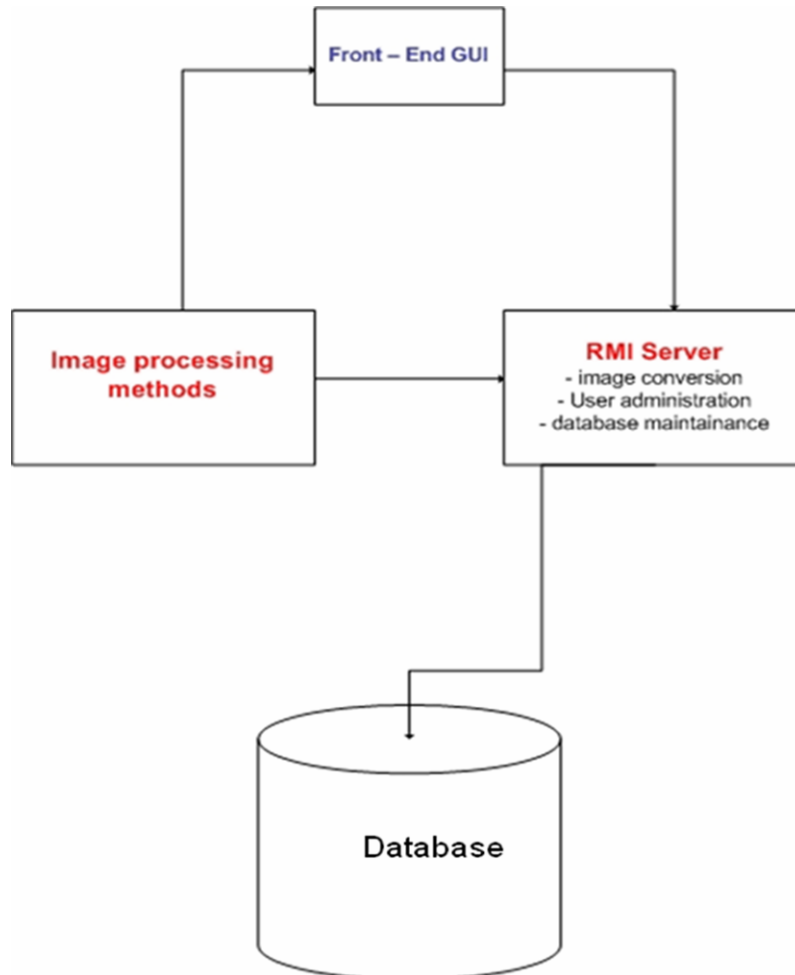
**Fig. 4.** The architecture uses a client-based GUI that connects to a RMI server and database. The image processing methods are separate from the server in order to maintain the flexibility of the design and allow additional methods to be easily incorporated.

Evaluating the effectiveness of features can be done using the system's extensive query capability. Finding cases when a "test" feature is near a "truth" feature can be done using the spatial query capability, as well as those "truth" features that are not near a "test" feature and vice versa. The spatial querying is tunable with a variable input distance for greater flexibility. The image categories can be adjusted with filters to isolate cases that are malignant, or other medically relevant characteristics.

The system can also store and share features that are not associated with a particular position. The flexible design of the feature storage and upload make it capable of handling most features. The system handles feature areas and feature volumes as well as point features.

Techniques that combine features can also be stored as features as well. Currently the spatial search can be used to combine features, but more complex approaches have to be done offline and uploaded as features. One planned upgrade is a learning package built into the system that would simplify the development of classifications and analysis of medical images. This would allow radiologists to use the system to explore relationships and use the learning package to optimize the approach, allowing the research to be done by just one expert instead of requiring expertise in multiple different fields.

The combination of advanced querying and feature sharing enables rapid analysis of features and combinations of features for CAD and the comparison of computed features to "truth" features defined by a radiologist. By providing a platform for the analysis and comparison of features, we encourage the collaboration between researchers designing features as well as researchers building CAD systems.

## 4   Conclusions and Future Work

We have created a secure web-enabled database for the storage, retrieval, manipulation, and annotation of medical images and medical records for the development and evaluation of CAD methods. The most unique quality is the ability to input, store, and share multiple feature sets and result sets for each image, thereby allowing greater flexibility for CAD and allowing web collaboration in the development of CAD. Each expert needed for the development of CAD gains advantages in their individual work by collaborating, while improving the project overall. The advanced querying and feature storage capabilities provide rapid analysis and comparisons radiologist techniques, medical image features and CAD techniques.

Future work on this project on the hospital applications includes changes to easily incorporate the OCR converted scanned forms for older patient data as well as associating a scanned permission forms with images that are to be made public or semi-private. A test site has been deployed and we are awaiting deployment at a larger hospital. On the other hand, future work on the research applications side includes support spatial search on features with certain value ranges, and the incorporation of an arbitrary number of feature values.

The incorporation of XML as well as DICOM capabilities improved the flexibility of the system and allows compatibility between systems that are DICOM based and newer systems that might incorporate XML. The inclusion of XML was simple because of the use of Java which contains much of the required capability.

The incorporation of filtering of search results was useful in analyzing trends in the medical data. For example, a particular radiologist's diagnoses can be checked by searching for their diagnosis feature, then filtering with malignant to see how often the diagnosis was correct. Missed diagnoses can be checked by comparing all of a radiologist's patients without a diagnosis feature that were eventually malignant.

## References

1. American Cancer Society. Breast Cancer Facts and Figures 2005-2006. American Cancer Society, Inc., Atlanta, GA (2006)
2. Campanini, R., Bazzani, A., Bevilacqua, A., Bollini, D., Dongiovanni, D., Iampieri, E., Lanconelli, N., Riccardi, A., Roffilli, M., Tazzoli, R.: A novel approach to mass detection in digital mammography based on support vector machines. In: Proceedings of the 6th International Workshop on Digital Mammography (2002)
3. Clarke, L.P., Croft, B.Y., Staab, E., Baker, H., Sullivan, D.C.: National Cancer Institute initiative: Lung image database resource for imaging research. Academic Radiology 8(5), 447–450 (2001)
4. Freer, T.W., Ulissey, M.J.: Screening mammography with computer aided detection. Radiology 220, 781–786 (2001)
5. Heath, M.D., Bowyer, K.W.: Mass detection by relative image intensity. In: Proceedings of the 5th International Conference on Digital Mammography, Medical Physics Publishing, Madison (2000)
6. Kalman, B.L., Kwasny, S.C., Reinus, W.R.: Diagnostic screening of digital mammograms using wavelets and neural networks to extract structure, Technical Report 98-20, Washington University (1998)
7. Linda, J., Burhenne, W., Wood, S.A., D'Orsi, C.J., Feig, S.A., Kopans, D.B., O'Shaughnessy, K.F., Sickles, E.A., Tabar, L., Vyborny, C.J., Castellino, R.A.: Potential contribution of computer-aided detection to the sensitivity of screening mammography. Radiology 215, 554–562 (2000)
8. Lui, S., Babbs, C.F., Delp, E.J.: Multiresolution detection of spiculated lesions in digital mammograms. IEEE Transactions on Image Processing 6, 874–884 (2001)
9. Sajda, P., Spense, C., Parra, L.: Capturing contextual dependencies in medical imagery using hierarchical multi-scale models. In: Proceedings of the IEEE International Symposium on Biomedical Imaging, pp. 165–168 (2002)
10. Sampat, M.P., Bovik, A.C.: Detection of spiculated lesions in mammograms. In: Proceedings of the 25th Annual International IEEE Conference on Engineering in Medicine and Biology Society, vol. (1), pp. 810–813 (2003)