# The National Medical Knowledge Bank

Warren Sterling

NCR Parallel Systems

100 N. Sepulveda Blvd. El Segundo, California 90245

Warren.Sterling@ElSegundoCA.ncr.com

## Abstract

This paper describes an advanced development program to create a medical information system called the National Medical Knowledge Bank (NMKB). This five year program is sponsored in part by a grant from the National Institute of Standards and Technology Advanced Technology Program. The goals of the program, covering computer-assisted diagnosis, medical training, remote consultation, and medical records storage, are defined. The web-based architecture of the medical knowledge bank is presented, including the Teradata Multimedia Services, an object/relational database which serves as the central data repository for medical data stored in multiple data types. Also described are the applications of physician support, including case-based reasoning and image analysis for determining case similarity; virtual medical conferences; and initial/continuing medical education.

## 1 Introduction

The National Institute of Standards and Technology Advanced Technology Program (NIST ATP) sponsors high-risk technology programs with the potential to benefit a broad segment of American society.

Program grants are usually awarded to consortiums of for - profit and not -for -profit entities, with the goal of commercializing successful programs. Consortium members retain rights to all intellectual property developed during the term of the grant. NIST ATP partially funded a five-year program to revolutionize healthcare management through the creation of a medical information system called the National Medical Knowledge Bank (NMKB). The program is being implemented by a consortium including Allegheny Health Education & Research Foundation (AHERF), NCR Parallel Systems, NCR Human Interface Technology Center (HITC), AT&T Business Markets Division, and AT&T Solutions. Carnegie-Mellon University is also supporting the program with research in the areas of video segmenting and indexing, medical content analysis, and case-based reasoning for medical diagnosis. Teradata Multimedia Services (TMS), a massively parallel object-relational database, will serve as the Multimedia Registry, or central data repository, for the knowledge bank. The ambitious goal of this grant is to revolutionize healthcare through the development of advanced tools for: computer-assisted diagnosis utilizing case-based reasoning (CBR), initial/continuing medical education (IME, CME), remote consultation and diagnosis with medical specialists, surgical rehearsal utilizing virtual reality sessions, and storage and management of patient records for illustrative cases. Effectively, the project is applying business process reengineering to healthcare. The key challenge for database people is the need to develop a medical multimedia data warehouse system that will improve productivity as well as patient care.

Initial work on the program has focused on the development of knowledge bank components such as the Multimedia Registry, simple user interfaces for healthcare practitioners, medical diagnosis case-based reasoning (CBR), the formatting of patient medical records to support CBR, automated analysis of MRI brain scans, and virtual conferences. Almost all work to date has been limited to the domain of neurology. This is expected to expand to other medical domains as the knowledge bank components are proven and the techniques for capturing and formatting patient medical records for inclusion in the knowledge bank are verified.

**Proceedings of the 24ᵗʰ VLDB Conference, New York City, 1998.**

The knowledge bank is being developed, tested, and deployed in several stages over the five-year period, with each stage being commercially viable in its own right. Initially, the goal is to develop a complete online end-to-end system where the data is stored in the multimedia database. The information will be accessed using Web-based browser tools, either across a high-speed campus Intranet or across the public Internet.

# 2 Architecture

## 2.1 Three Tiered: Web Client, Application Server, Data Server

Figure 1 shows the three-tiered architecture of the knowledge bank: web client tier, application server tier, and the data server tier. The web client uses standard browser technology augmented by Java Applet and appropriate plug-ins or players for displaying various types of multimedia data.

The application server tier consists of a standard web server and Publisher, a component which extends web servers to handle database references within a web page.

The data server tier consists of the Multimedia Registry and MedWeb, a component which converts Publisher requests into database queries. The Multimedia Registry includes TMS and a media server which can stream video and audio to clients over the Intranet/Internet. TMS must store, retrieve and analyze the following: medical imagery such as MRI scans, x-ray images, and EKG graphs; videos of surgical procedures, medical conference presentations, patient presentations (real and staged for case studies), and synthetic interviews (pre-recorded video segments providing answers to frequently answered questions posed in natural language); and medical case records which include medical images and video. To deliver video to clients across the Intranet/Internet with minimal latency and small client buffer space, TMS utilizes the media server, which is capable of generating concurrent streams of video to clients. The Multimedia Registry, MedWeb and Publisher can communicate on a separate high speed sub-network to reduce latency. This is particularly important in the case of video. Selected video must be retrieved from the database and cached on the media server, from where it is streamed to the user. This eliminates the need to support video Quality of Service from the database directly but adds the latency of the full transfer of the video to the media server; hence the need to minimize this latency with a high speed sub-network.

The query language for TMS is based on SQL3, a nickname for a series of evolving standards supporting

the object-relational model (the addition of objects to the relational model).

## 2.2 Information Flow

A user initiates a request for information using the client. The request is sent to the web server using standard HTTP protocol. In the case of a request for a static page, the web server will process it directly. Requests involving dynamic web pages are transferred to Publisher. These requests include "template" pages. Publisher determines if data is required from the Multimedia Registry in order to build the dynamic page. If so, it formats the query and accesses the Multimedia Registry through MedWeb. MedWeb queries the Multimedia Registry and returns results to Publisher. Publisher uses the results to replace data references in "template" pages with actual data. It then returns the page to the client via the web server. For streaming media requests, MedWeb retrieves a media object from the Multimedia Registry and stages it on the media server, where it will be accessed by the client and displayed at the client, using a browser plug-in if necessary. Note that objects can remain on the media server after staging, or can reside on the media server permanently in order to assure minimum latency for a client accessing large objects such as videos.
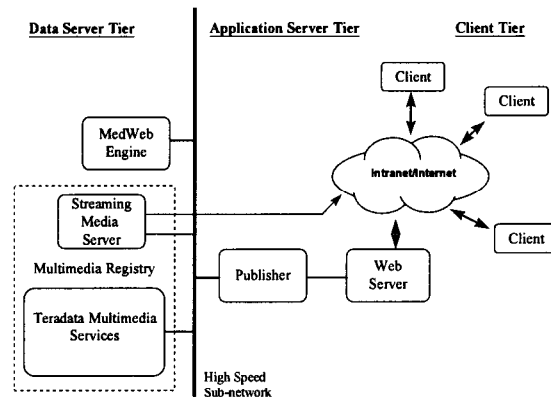


**Figure 1.** Architecture Block Diagram for the NMKB

## 2.3 User Defined Functions

A key feature of the TMS is the ability to execute User Defined Functions (UDF) which can manipulate and analyze objects stored in the database. An example of a UDF is an image processing algorithm that analyzes MRI brain scans to detect the presence of a tumor. UDFs are written in C or C++, loaded into the database, bound to particular object types, and then executed as a function call within a query. For the knowledge bank, UDFs will be developed for MRI brain scan analysis, metadata extraction from DICOM-3 images (DICOM-3 is a standard format for medical imagery), case-based reasoning application to

638

medical case records, and video indexing and segmentation.

## 2.4 Scalability and Parallelism

Another important feature of TMS is scalability - the ability to grow in size incrementally with increases in database size or concurrent usage, and maintain performance. TMS runs on the NCR WorldMark massively parallel processing system. Scalability is achieved by adding nodes to the system as required. Further information on TMS and the NCR WorldMark computers can be found in [CS98].

Massive parallelism in the object/relational database, which serves as the Multimedia Registry, is a key technical objective. Fully deployed, a medical knowledge bank is expected to host large numbers entities containing multimedia data: electronic medical case records, medical journal articles, conference presentations, healthcare training courses, surgery videos. On-line storage requirements can easily reach multiple terabytes. Realistically, hundreds (potential thousands) of users could be accessing the Multimedia Registry concurrently. Queries launched against this massive database will require searches against entire tables, and will include content analysis of multimedia data. The database storage requirements and query processing power requirements will increase over time and the Multimedia Registry will need to grow, scaling linearly to accommodate these storage and processing power increases. Massive parallelism is the best solution to this problem [CSK95] [OIS96]. As storage requirements grow, nodes can be added to the system to increase storage. Because database tables are distributed across all nodes, queries against full tables can be decomposed into parallel steps, each running on a separate node. Increasing the number of nodes increases the level of parallelism which can be applied against a query. Increased parallel execution of query steps allows more queries to execute concurrently. Even for relatively simple transactional queries, increasing the number of nodes allows the system to handle an increased number of user sessions. Furthermore, for high volumes of transactional queries involving large object retrieval (such as video objects), the number of independent transmission paths can be increased as the number of nodes increases, minimizing transmission bottlenecks.

## 3 Multimedia Warehouse

The Multimedia Registry will operate in a typical data warehouse environment for this application. For the prototype implementation at AHERF, NMKB will serve multiple campuses across Pennsylvania via a high bandwidth ATM network connecting the campuses. It will also serve the world outside AHERF via the Internet. Tools are being created to extract patient data from hospital legacy systems and format it for the registry. Data stored in the knowledge bank will be available to all healthcare practitioners and other interested individuals to support a wide variety of medical applications.

## 4 Applications

### 4.1 Physician Support

Computer-assisted medical diagnosis involves the retrieval from the knowledge bank of medical cases similar to the case being analyzed by a physician. Case-based reasoning [Sim92] [Bar91] is the base technology being used to determine case similarity. A set of proto-typical medical cases is chosen to cover a particular medical domain. Individual case attributes are assigned weights indicating their correlation to a particular diagnosis, based on expert knowledge of the cases and analysis of the case base (where the diagnosis of each case is known). Of particular interest is the handling of multimedia imagery which can make up a significant part of the medical cases. Initial work is focusing on the automated analysis of MRI brain scans to detect and characterize brain tumors, aneurysms and bleeds [BR95]. Work at Carnegie-Mellon University and AHERF is currently underway to develop these automated analysis techniques [LD98], and then to understand how to utilize these case attributes for determining case similarity. NCR HITC is developing the actual case-based reasoning engine and the user interface.

### 4.2 Virtual Medical Conferences

Virtual conferences deliver live or persistently stored presentations (video and/or audio, illustrations, text versions of the presentations) via the Intranet/Internet. A web site with two such virtual conferences has been prepared and is currently being evaluated by a selected audience of healthcare practitioners. The first is a conference on "brain attack", or stroke, held in May 1997 in Scottsdale, AZ. Twelve prominent researchers and practitioners in the field of brain attack were brought together for the purpose of generating content for this virtual conference. Participants gave two presentations, one targeted to healthcare practitioners and one targeted to the general public. The second virtual conference is based on recordings made at the Conference of Neurosurgeons held in October 1997 in New Orleans, LA. Different post-production methods were used to generate each virtual conference - one labor-intensive, the other relatively automatic. However, access methods for each conference are essentially the same. Users access a combination of text papers and streaming video synchronized with slides

and other images. Currently, the consortium is evaluating feedback on the two virtual conferences in order to determine the minimum post-production effort required, consistent with acceptable presentation quality.

### 4.3 Initial/Continuing Medical Education

Healthcare professionals go through initial medical training, of course, but they must also acquire a minimum number of additional education credits each year. Typically, these are obtained through attendance at conferences or short focused courses. In the case of continuing nursing education, one way to obtain these credits is through nursing periodicals. Training articles are published in the periodicals. Nurses read the articles, then complete a multiple-choice test which is mailed in for grading and credit. The knowledge bank will be used to deliver multimedia-enabled courses for both IME and CME, using a web-based interface. AHERF has developed several multimedia-enabled synthetic interviews and case studies to be used in their medical school classes (IME). Virtual conferences hosted on the knowledge bank can be used for CME credits. Current work is focusing on building a framework for creating and administering CME courses, including testing, in a web-based environment. This framework is being built on an intelligent agent architecture. The student can select a course, which consists of a set of lessons typically augmented with multimedia data. Lessons are indexed by the concepts they teach. A modeling agent creates a student model, indicating such things as background, knowledge state and goals. The student profile (courses completed, lessons mastered, areas requiring remedial review, concepts the student has not yet studied) is maintained by a profiling agent. The student model and profile are used by a lesson planning agent to provide the student with a lesson plan. All course material and course administration data will be maintained in the knowledge bank.

## 5 Conclusion

This paper briefly describes a project to create the National Medical Knowledge Bank operating in an Intranet/Internet environment, using web technologies. This technology is predicated on the existence of a massively parallel and scalable object/relational database, Teradata Multimedia Services, to store medical data in diverse forms. UDFs are used to perform content-based access to the stored data in support of a number of medical applications. Three applications were described here: physician support, virtual conferences, and initial/continuing medical education; however, the architecture supports a number of other applications,

such as synthetic interviews, electronic patient record storage, and remote consultation.

### References

[Bar91]    R. Barletta, An Introduction to Case-Based Reasoning, *AI Expert*, pp. 43-49, August 1991.

[BR95]    P. Black and E. Rossitch, Jr., *Neurosurgery, An Introductory Text*, Oxford University Press, New York, 1995.

[CS98]    F. Cariño and W. Sterling, Parallel Strategies and Concepts for a Petabyte Multimedia Database Computer, *IEEE Parallel Database Techniques*, 1998.

[CSK95]    F. Cariño, W. Sterling, P. Kostamaa, Industrial Database Supercomputer Exegesis - The DBC/1012, the NCR 3700, the Ynet, and the Bynet, *Emerging Trends in Knowledge and Database Systems*, IEEE Computer Society Press, Los Alamitos, California, 1995.

[LD98]    Y. Liu and F. Dellaert, A Classification-Based Euclidean Similarity Metric for 3D Image Retrieval, *Computer Vision and Pattern Recognition Conference*, January 1998.

[OIS96]    W. O'Connell, I. T. Ieong, D. Schrader, *et.al.*, Prospector: A Content-Based Multimedia Object Server for Massively Parallel Architectures, *Proceedings of the ACM SIGMOD Conference on Management of Data (SIGMOD 96)*, pp. 68-78, 1996.

[Sim92]    E. Simoudis, Using Case-Based Retrieval for Customer Technical Support, *IEEE Expert*, pp. 7-11, October 1992.