# A Framework for the Design of Knowledge-Based Systems in Structural Biology

James F. Brinkley, Jeffrey S. Prothero, John W. Prothero, and Cornelius Rosse

Department of Biological Structure SM-20
University of Washington
Seattle, Washington 98195

## ABSTRACT

Structural biology is a basic medical science dedicated to understanding the function of the body in terms of its three-dimensional organization. Since both normal and abnormal function are fundamentally determined by structure a knowledge base of structural biology should lead to causally-based problem-solving modules in basic science and clinical medicine. Such a knowledge base, and the associated knowledge-based systems, will be large and will need to be developed and coordinated through many research groups. We describe a distributed framework for allowing such sharing to take place, in which individual modules running on widely scattered and diverse machines are made available as knowledge and data servers communicating via remote procedure calls or similar mechanisms. We also describe potential modules for structural biology and present results from our current modules for computing three-dimensional anatomic reconstructions from serial sections.

## Introduction

One of the more important principles of artificial intelligence research is that the effectiveness of an intelligent system, whether human or computer, is proportional to the amount of knowledge to which it has access. With the rapid proliferation of biomedical information this principle is especially relevant to medical care and was recognized in the theme for last year's SCAMC meeting, "Knowledge is Power", and in the increasing numbers of medical informatics papers related to knowledge-based systems [1].

Much of the work in knowledge-based systems in medicine has dealt with clinical diagnosis, utilizing shallow rule-based or probabilistic approaches [2]. However, the brittleness of rule-based systems has led researchers to explore underlying causal models in order to give the computer deeper understanding that it fall back on when simple rules do not suffice [3]. These attempts are an implicit recognition that the basic medical sciences provide the foundation for understanding in medicine, without which a physician is a technician.

Structural biology can be defined as a basic medical science studying the three-dimensional organization of the body at levels ranging from gross to molecular. Since a fundamental hypothesis of biology is that structure determines function it is reasonable to argue that any causal biomedical knowledge base should be built around a framework of biological structure. Once structural knowledge at the level of gross anatomy, histology, and molecular biology is adequately represented, it will be easier to add normal and abnormal functional knowledge (e.g., physiology and pathology). The resulting (large) knowledge base could then be accessed by knowledge-based modules in a scientist's workbench for adding to the store of structural knowledge, in a physician's workbench for using the structural knowledge to diagnose and treat illness, and in a student's workbench for learning about structural biology. The knowledge base would be an example of an emerging trend in artificial intelligence to build large, reusable knowledge bases [4].

The major issue in any knowledge-based system is representation. Structural knowledge can be classified in two ways: symbolic and spatial. Symbolic knowledge is the textual information found in a textbook, such as the names of anatomic objects and their hierarchical relationships. Methods for representing this type of knowledge have been studied in artificial intelligence, but satisfactory solutions have not yet been found [3]. Spatial knowledge is the three-dimensional (3-D) geometric shape and range of variation for anatomic objects, as well as their spatial relationships. Geometric representations have been studied primarily in computer vision for man-made objects, but again satisfactory representations are not yet available for biologic objects [5]. There have been few attempts to couple symbolic and spatial knowledge in a meaningful way.

Although adequate representations do not yet exist for capturing the range of variation for anatomic shapes, many representations have been developed for conveying the shape of particular examples of objects. These representations are in the form of 3-D reconstructions that can be shown on a graphics display and are becoming an important means for conveying structural data to researchers and clinicians. As these reconstructions are accumulated they form the data from which a spatial knowledge base can be generated.

Over the past five years we have developed a facility for generating and displaying three-dimensional reconstructions of anatomy, and we are rapidly building a large database of three-dimensional macroscopic and microscopic reconstructions [6,7,8]. In this paper we describe the current status of our system and propose a design for extending the anatomic database towards knowledge-based systems in

structural biology. The major advantages of this design are its modularity, extendibility, and flexibility, so that we can build incrementally more useful systems as advances are made in the difficult areas of knowledge representation.

## System Design

### Design Objectives

The objectives of our research are 1) to extend our methodology for generating and displaying 3-D reconstructions of anatomic objects from serial sections, 2) to create, in collaboration with others, a 3-D spatial database in which the whole human body is represented at gross resolution (about 1 mm), and many parts of the body are represented at microscopic (1 um) and ultrastructural (10 nm) levels of resolution, 3) to develop symbolic and spatial knowledge and databases based on these reconstructions, and 4) to develop, in collaboration with others, modules which utilize these knowledge bases to solve problems such as graphics generation, image segmentation, tutoring, diagnosis, treatment planning, molecular structure determination, national database retrieval, and experiment design.

The reason it is possible to consider so many potentially useful knowledge-based modules is that structural knowledge is fundamental to biology. Thus, the major issue will be how to represent and utilize this knowledge. Because this issue is so difficult we cannot expect to build a system that immediately solves the problem, yet we need to build practical systems before the final answers are in. We also recognize that other groups have worked on knowledge-based systems and that we could never implement all these systems ourselves. Thus, the objectives are to design a framework that 1) is flexible and extendable, allowing incremental improvements in knowledge representation and knowledge-based modules to be incorporated, 2) allows sharing of large databases, knowledge bases, and knowledge-based modules both within our group and among other groups, and 3) recognizes the sociology of research groups, so that groups can maintain independence while collaborating with other groups to build systems that are too large for any one group to build alone.

### General Approach

These design objectives lead to a distributed system in which individual modules running on different machines communicate with each other over the network. As modules are developed, communication can progress from files to remote procedure calls to message passing or other more exotic techniques, such as distributed blackboard systems.

The advantage of this approach is that, once the proper interfaces are developed, the individual modules can be built and changed independently by separate research groups located at widely scattered locations. As long as the interfaces are relatively stable the modules themselves may be changed as techniques improve, and different modules may be run on those machines best suited for them. For example, it can be argued that artificial intelligence modules should be run on specialized Lisp machines, whereas graphics modules should be run on graphics engines.

The distributed approach is much more feasible now than it was just a few years ago. Advances in hardware have led to faster networks and more powerful individual workstations so that data can be communicated quickly among machines. Software advances such as the Networked File System (NFS), Remote Procedure Calls (RPC) or X windows are rapidly improving; processes soon will be able to communicate easily over the national network. What is needed now is to add application-level layers on top of the lower level network communication layers so that knowledge-based communication can be established. These interfaces will allow the development of knowledge servers, database servers, and modules that access them.
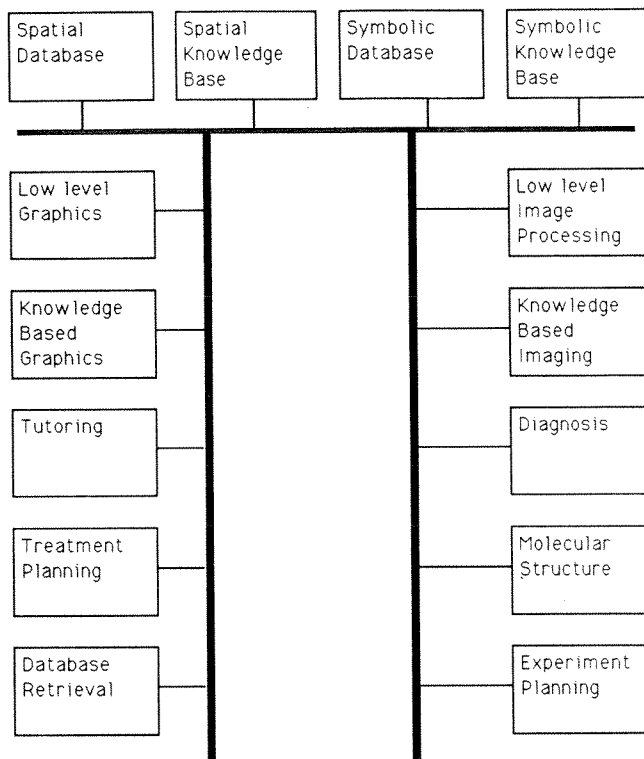


Figure 1: System design.

### Structural Biology Modules

Figure 1 shows a design for knowledge-based systems in structural biology. All modules are shown connected to other modules via a network. The modules will often run on different machines and will communicate primarily via remote procedure calls. In theory, each module can be both a client and a server, so routines in the module can be accessed by other modules and each module can in turn access other modules. The communication protocol will primarily be Lisp functions parsed by the server to call individual routines in the module. The modules themselves can be written in whatever language is most convenient (currently Commonlisp and C). The user will interact through the user interface, which will be hypertext-based and which will call other modules via remote procedure calls.

This server protocol is modeled after a 3-D protein structure system developed at Stanford called PROTEAN, in which a knowledge-based module running on Xerox D-machines accessed a computational geometry server running on a Silicon Graphics Iris [9].

The currently defined modules are:

1. Spatial database: Three-dimensional coordinate tables describing reconstructions of specific biologic objects, later to be augmented with images as well as higher level geometric entities such as surfaces and volumes.

2. Symbolic database: An accounting database for keeping track of the 3-D reconstructions. A relational database will be used initially, followed by an object-oriented visual database coupling the symbolic and spatial databases.

3. Symbolic knowledge-base: Symbolic information about anatomic objects and their relationships. Since this is a difficult problem, we will begin by developing anatomical taxonomies for use as indices into the database. Several expert systems shells, such as KEE or BB1 [10], will be tried. The knowledge that we don't yet know how to put into the symbolic knowledge base will be stored in hypercards associated with the objects in the knowledge base. As we come to understand the nature of the knowledge it will be transferred to the knowledge base.

4. Spatial knowledge-base: Generic spatial models of anatomic objects using information gained from the individual models in the spatial database. These models will contain the essential shape, as well as the range of variation, for an anatomic object referenced in the symbolic knowledge base.

5. Low level graphics: a graphics server for drawing the 3-D models, based on our current display program.

6. Low level image processing: an image processing server, adapting one of the many packages already available.

7. Knowledge-based graphics: A module that, by utilizing the symbolic knowledge base, responds to high level scripts for displaying anatomic reconstructions. This module will, in turn, call the low level graphics server.

8. Knowledge-based image processing: The spatial knowledge base will be utilized in a model-based vision system to aid the segmentation of object borders on images.

9. Other modules will be developed as the knowledge bases are improved.

## Current Status

We currently have in place rudimentary modules for the spatial database (files of 3-D coordinates), low level image processing (manual contouring), and low level graphics (3-D display program). The main activity in our group has been the generation and display of 3-D reconstructions from serial sections [6,7,8]. We routinely carry out reconstructions at the macroscopic and microscopic levels of resolution. Cadaver or tissue segments are thinly sliced and photographed; then their contours are manually traced. Slices are also obtained from clinical CT or MRI images. External fiducials (or markers) are used to align adjacent slices and to aid the computation of 3-D coordinates of tissue contours. The contours are then edited and displayed with an interactive 3-D display program.

The display program, called SKANDHA, was written by one of the authors (Jeffrey Prothero). Given the 3-D coordinates of various contours the program calculates a triangular tiling between adjacent contours, computes the surface normals, and displays the object using Z-buffering and Gouraud shading. An extensive set of interactive menus allows the user to edit contours or tilings, to store and retrieve large sets of contours grouped in separate files, to move objects via arbitrary transforms, to change the color and lighting of objects or subobjects, and to control a video tape recorder for automatically generating animated sequences. SKANDHA is written in C and runs on both the Silicon Graphics Iris and enhanced AT's.



Figure 2:     Three-dimensional reconstruction of a cadaver skull from 1.5 mm CT scans.

Figure 2 shows a cadaver skull, reconstructed from CT scans 1.5 mm apart, and is an example of the output of the 3-D graphics program. In addition, we have a fairly extensive collection of 3-D reconstructions and images of details of the brain, the hand, the pelvis, the eye and orbit, and the human embryo. At present, the images are used primarily for teaching via animated videotape sequences. Three anatomists devote full time to this project, and numerous others from our department and elsewhere are also involved.

Our current hardware configuration consists of workstations connected via Ethernet to each other and to the rest of the world. Several "anatomist's workstations" are IBM AT's fitted with a digitizing tablet for manual contouring. Four Silicon Graphics Personal Irises (4D/20) and a Silicon Graphics 3020 support the graphics program SKANDHA. A Silicon Graphics 4D/70 is used for graphics production and for generating video tapes. The "symbolic" workstations consist of three Apple MacIntosh IIX's, one of which has a Texas Instruments MicroExplorer board running LISP. Our fileservers are a Sun 3/160 with a 1.2 gigabytes disk, a Microvax II with a 2.2 gigabyte optical WORM drive, and a NeXT.

## Implementation Plans

Our current hardware is adequate for initial development of many of the modules. We plan to proceed incrementally, gradually adding capabilities as we develop more robust knowledge bases. Since our primary activity is currently graphics production, the first knowledge-based modules will be devoted to speeding up that activity. The SKANDHA program will be made available as a server (the low level graphics module in figure 1), so that it can be called by other modules. A commercial relational database system will be embedded in a server to form the symbolic database. This database will contain "accounting" information such as the source of the specimen, tissue codes, the date of entry, etc. One of the fields in the database will be the name of a SKANDHA file in the spatial database of 3-D coordinates, thereby providing a link between the symbolic and spatial databases. Another field will be the name of the anatomic object, which will provide a link between the symbolic database and generic anatomic objects in the symbolic knowledge base. The database servers will run on the Sun.

The symbolic knowledge base will initially only contain an object-oriented glossary of anatomic terms and synonyms, as well as PART-OF relationships from several points of view. It will, thus, be used at first as an index into the database, but as the anatomic hierarchies are developed, additional symbolic information will be added. The symbolic knowledge-base server will run on the MicroExplorers.

Many of the objects in the knowledge base will contain pointers to cards in HyperCard stacks on the MacIntosh. These cards will contain knowledge that we do not yet know how to put into the computer-readable symbolic knowledge base. The user interface will also be written in HyperCard.

These initial modules will give us the rudiments of a knowledge-based information retrieval system for accessing our 3-D reconstructions via anatomically-meaningful terms in the symbolic knowledge base. We will then develop knowledge-based systems for high level graphics and image processing. These modules will speed up the production of graphics images by partially automating two of the main bottlenecks we currently face, namely, generation of graphics displays and extracting meaningful contours from the original images.

The graphics module will be an expert system that uses knowledge of anatomy from the symbolic knowledge base, as well as heuristics from the anatomists, to display images automatically from high level scripts. This module will access the database server to find the reconstructions, and will call the low level graphics server, SKANDHA, to draw the images. The high level graphics module will initially be accessed via the user interface, but it may later be accessed by a knowledge-based tutor or other modules requiring graphics expertise.

The knowledge-based image-processing module will utilize both symbolic and spatial knowledge to aid the process of segmenting contours on actual images. This module will call on routines in the low level image processing server, which can be built from one of the many existing image processing packages developed in other departments here or elsewhere. Some of these packages are very large and involve many man-years of effort. By making that work available as a server on the machine it is best suited for, we are able to take advantage of the substantial previous work in this area.

The knowledge-based image processing module will require spatial knowledge in the form of expected shape and range of variation for anatomic objects. This knowledge will be in the form of generic models generated from the individual reconstructions in the spatial database and indexed via the symbolic knowledge base and symbolic database. One potential representation for the generic models is geometric constraint networks [5]. However, since the problem of model-based vision is far from solved, we expect the knowledge-based imaging system will be interactive, communicating with the user via the user interface.

As these modules are developed, they will lead to better representations for the symbolic and spatial knowledge bases, allowing, in turn, the development of some of the other modules in figure 1. Given the complexity of each of these modules it is likely that some of them will be developed elsewhere, but the distributed framework will allow us to take advantage of that work, while others take advantage of our work.

## Discussion

In this paper we have outlined a design for knowledge-based systems in structural biology based on a large spatial database of 3-D anatomic reconstructions. These systems will also rely on a large knowledge base of structural biology that will need to be developed and shared over many institutions. The development of specific problem solving modules will lead to incremental improvements in the representation of both spatial and symbolic knowledge of structure. These better representations will, in turn, lead to more powerful problem-solving modules. Given the fundamental role of structure in biology these modules can impact many areas of medicine and biology.

The distributed nature of the design will allow incremental development of progressively more powerful modules and will not require solving all the theoretical problems first. It will also allow sharing of expertise among research groups, since the results of a group's expertise can be made available as a knowledge server over the network.

Our approach is based on the hypothesis that intelligent systems arise out of the coherent interaction of small subparts. This hypothesis has been articulated many places, including artificial intelligence, and evidence for it can be seen at levels ranging from protein molecules to neural networks to rule-based systems to distributed problem-solving systems. The design of our framework follows the distributed problem-solving approach of modules communicating over a network, but by making the human developer and user an integral part of the process we allow the creation of practical yet evolutionary systems.

## REFERENCES

1. Greenes, R.A., "Knowledge is Power", in Proc. Twelfth Annual Symposium on Computer Applications in Medical Care, Washington, D.C., November 6-9, pp.2-3, 1988.

2. Clancey, W.J. and Shortliffe, E.H., Readings in Medical Artificial Intelligence: The First Decade, Reading, Mass., Addison-Wesley, 1984.

3. Haimowitz, I.J., Patil, R.S., and Szolovits, P., "Representing Medical Knowledge in a Terminological Language is Difficult", in Proc. Twelfth Annual Symposium on Computer Applications in Medical Care, Washington, D.C., November 6-9, pp. 101-105, 1988.

4. Lenat, D., Prakash, M., and Shepard, M., "CYC: Using Common Sense Knowledge to Overcome Brittleness and Knowledge Acquisition Bottlenecks", AI Magazine, 6(4):65-85, 1986.

5. Brinkley, J.F., "Knowledge-driven Ultrasonic Three-dimensional Organ Modelling", IEEE Trans. Pattern Analysis and Machine Intelligence, PAMI-7(4):431-441, 1985.

6. McLean, M., Ross, M.A., and Prothero, J.W., "A Three-dimensional Reconstruction of the Myofiber Pattern in the Fetal and Neonatal Mouse Heart", Anat. Rec., In Press.

7. Prothero, J.S. and Prothero, J.W., "Three-dimensional Reconstruction from Serial Sections IV. The Reassembly Problem", Computers and Biomedical Research, 19:361-373, 1986.

8. Stimac, G.K., Sundsten, J.W., Prothero, J.S., Prothero, J.W., Gerlach, R. and Sorbonne, R., "Three-dimensional Contour Surfacing of the Skull, Face, and Brain from CT and MR Images and from Anatomic Sections", AJR 151:807-810, 1988.

9. Brinkley, J.F., Altman, R.B., Duncan, B.S, Buchanan, B.G., Jardetzky, O., "Heuristic Refinement Method for the Derivation of Protein Solution Structures: Validation on Cytochrome b562", J. Chem. Inf. Comput. Sci., 28(4):194-210, 1988.

10. Hayes-Roth, B., "A Blackboard Architecture for Control, Artificial Intelligence, 26(3):251-321, 1985.