John A. A. Sillince · Maria Sillince

# Molecular Databases for Protein Sequence and Structure Studies

John A. A. Sillince · Maria Sillince

# Molecular Databases for Protein Sequences and Structure Studies

An Introduction

With 27 Figures

Springer-Verlag

Berlin Heidelberg NewYork
London Paris Tokyo
Hong Kong Barcelona Budapest

## Dr. John A. A. Sillince
Lecturer in Management Information Systems
British Sheffield University
Management School
Crookesmoor Building
Conduit Road
Sheffield S10
1FL, UK

## Dr. Maria Sillince
Assistant Subject Librarian
Wolverhampton Polytechnic
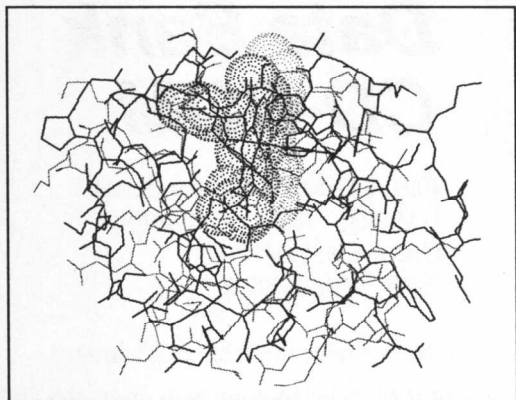Wolverhampton
UK

# Acknowledgements

# Preface

Molecular information is now too vast to be acquired without the use of computer-based systems, which can select according to supplied criteria. Developments in programming have created the ability to extend molecular science in ways that would have been impossible without their help. New databases are being established which enable previously unanswerable questions to be considered. One of these questions is whether or not one can predict the three dimensional structure of a protein from information about its sequence of amino acids.

In order to help the reader to understand this new field, several topics are explained in this volume. The structure and function of proteins and nucleic acids are described, in order to emphasise the way in which three dimensional structure reflects a protein's role in the organism. Also it is important to consider what is involved in molecular data, and how it is represented and registered in software and on the screen. Another aspect to consider is how computer-based research tools are used in molecular science, in particular for manipulating sequence and structure information. Sequence and structure are at the centre of research problems in molecular science, in the identification of a new protein (14000 are known so far) or its three dimensional structure (only 400 are known so far), in patent writing and patent searching, and in modelling proteins.
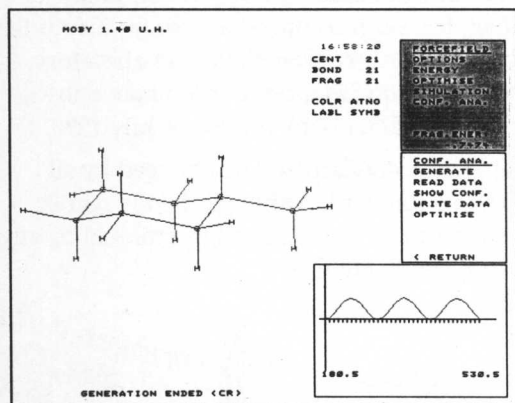
There is also a description of the state of the art in what data-banks exist, both for sequences and for structures, and what types of system are available for using them, for modelling, for searching, and for integrating the operations of the online database and the local system such as a PC in a laboratory. New developments in knowledge-based systems and database technology are described. The case study on protein structure prediction, which includes developments in the specification of an expert system for such a problem, is intended to exemplify the integrated nature of modelling and search (both computer-based) and laboratory experiment in molecular science.

# A <u>New</u> Molecular Modeling Program for the PC

## U. Höweler ___ MOBY ___ Version 1.4



Structure analysis of a hemoglobin subunit of Cytochrome C Van der Waals representation of the hemoglobin structure embedded in the peptide environment. The iron atoms is coordinated by the four nitrogen atoms of the ring system, by a nitrogen of a histidine residue and the sulfur atom of a methionine side chain.



Confirmational analysis of methyl cyclohexane which shows MOBY menu at right, additional information about the molecule at left and rotational profile at bottom right.

Moby is a Molecular Modeling Program for IBM PC and compatible computers.

It provides the following functions:

● 3 D graphic display for up to 2000 centers

● structure and property analysis and comparison

● force field calculations for 150 centers interacting with up to 2000 centers

● geometry optimization and conformation analysis and molecular dynamics simulation

● semiempirical quantum chemical calculation (MNDO, AM1)

● MOBY reads and writes standard structure file formats (e. g. Protein Structure Database format)

● MOBY reads and writes 3 D geometries in any format

● MOBY writes HPGL plot files and generates hardcopy output

### Hardware requirements:
IBM PC or compatible computer, 640 kB RAM, 80x87 arithmetic coprocessor, MS-DOS version 2.x or higher, hard disk, 1.5 MB free disk space, graphics card EGA or VGA, HERCULES supported, mouse optional.

# Contents

# List of Figures