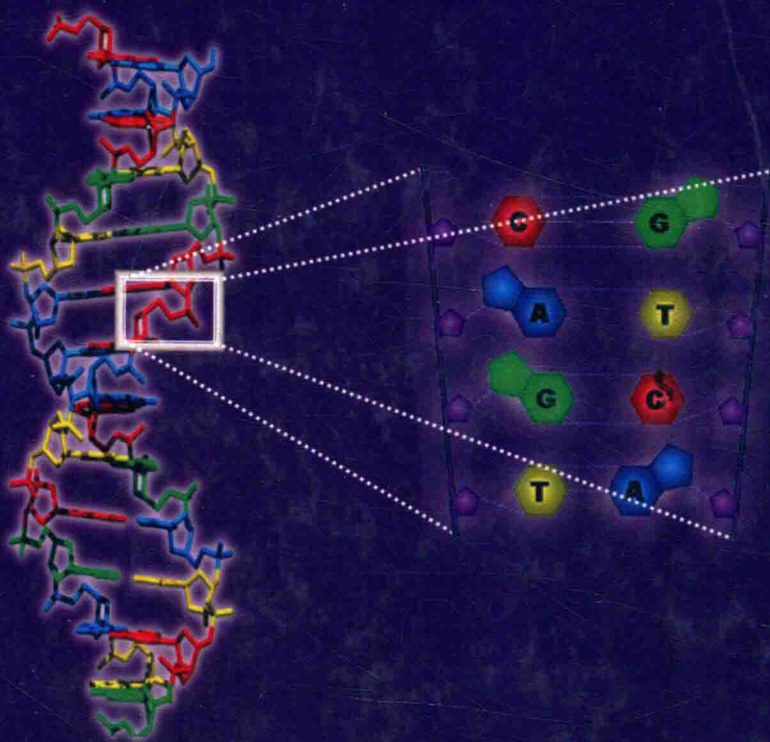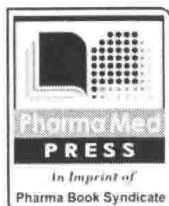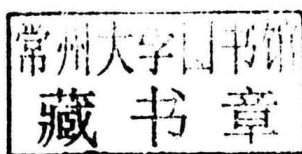# Elementary
# Bioinformatics

**Imtiyaz Alam Khan**

# Elementary Bioinformatics

## Imtiyaz Alam Khan

*M.Sc., PGADB*

Bioinformatics tutor, Microbiologist.
Ex-Lecturer, Dept. of Biotechnology,
Rai University, Hyderabad.

# Elementary Bioinformatics

# Preface

This book presents good coverage of Bioinformatics concepts, concise coverage of multiple sequence alignment practicals, gene identification practicals, domain identification practicals. More illustrations have been included in most of the chapters to emphasize the importance of topics under study. This particular edition is prepared to facilitate students, academicians and research workers. Genetic Engineering which is a part of Bioinformatics, is also explained very clearly. This book is designed keeping in view of all Universitie's Syllabus; basically meant for JNTU B.Tech Students. Anybody interested in learning basic Bioinformatics can use it.

This book will serve well as a supplement to either a formal Bioinfomatics Course or an integrated curriculum. We hope that senior students, graduates, postgraduates, academicians and researchers will use this book as a reference for Bioinformatics information.

## TO THE READER

"Elementary Bioinformatics" is a very handy book, which is helpful in understanding Bioinformatics from its grass root levels. As we all know, bioinformatics is a highly emerging field. It is now in its primordial state and within few years it will reach climax. Understanding Bioinformatics at this level that too in a very conceptual, concise and vivid manner is what the requirement of present scenario. Every topic is explained very clearly, especially sequence alignment, multiple sequence alignments, phylogenetics predictions, gene predictions and so on.. After reading them the reader will get clear concept of those topics. Few practicals are also included in order to make them comfortable with online Bioinformatics tools.

"Elementary Bioinformatics" is designed to provide a logical framework for organizing, learning, reviving and applying the conceptual and factual information required.

# Acknowledgement

# Contents

# 4  Computational Biology
## (Biological Database Management System)

# 5  Computational Biology :
## Information Retrieval From Databases

# 6 Phylogenetic Analysis

# 7 Phylogenetic Prediction

# 8     Genetic Engineering

# 9     DNA Sequencing

# 10 Genome Mapping

# 11 Gene Identification

# 1

# Internet

## 1.1 Scope of Bioinformatics

Bioinformatics has emerged as a distinct discipline that straddles the interface between the traditional biological sciences and the computer sciences and advance computational methodologies. It is rapidly becoming a powerful new approach to understanding life, and it may well reverse the reductionist paradigm that has held sway in molecular biology ever since Erwin Schrodinger turned on a generation of physicists to biology with the publication of "What is life ? more than 50 years ago.

The new discipline of bioinformatics promises to provide the tools needed to attack the complexity of conducting holistic biological research, "because of this complexities, biology will eventually become the most computational science, surpassing physics," said Delici, who predicts that within the next 10 to 15 years bioinformatics will become the integral part of biology.

Modern day biotech and drug discovery industry has witnessed the development of high through put automated equipment, which enables amassing of data faster than it can be analyzed and the utilized. Pharmaceutical research will clearly be the one major benefactor of development in bioinformatics. 100 of new drug targets have been identified by using computational techniques, which involves searching for genes similar to known protein encoding genes. In the future, virtual toxicology screening may be the first setup in predicting the effects of new chemical or complex metabolic pathways. In addition, bioinformatics will likely provide the methodology, to make highly accurate predictions about protein tertiary structure based on amino acid sequences and a viable means to design drugs based on computer simulation of the docking of small molecules to the predicted protein architecture.

According to figures from Framingham mass based market analysis firm IDC Bio-IT market will reach $38 billion in 2006. The market includes pharmaceutical, healthcare research and biotechnology companies, as well as Govt. linked institutions. A number of recent work force studies have shown that there is a high current and unmet demand for people trained to various levels of expertise in bioinformatics, serve the upcoming biotech and biopharmaceutical industry which has observed significant growth in genomic era.

## 1.2  Introduction (Elementary Commands and Protocols)

The Internet has become an important tool for biological and biomedical research scientists. Using the Internet, it is possible to perform a number of kinds of analysis on research data and to search for and obtain information. Over the last several years, the number of tools and the amount of information relevant to biologists available on the Internet has grown and the ease of use of these tools has grown as well. As a result of both of these trends, the value of Internet resources for biologists now significantly outweighs the costs in time and money of using it. The overall goal of this chapter is to help biologists use the Internet effectively and to illustrate to computer scientists, how biologists are currently using the Internet.

***This chapter has three specific goals :***

1. To provide background information which will help demystify computer network usage.

2. To provide an introduction to the resources available to biologists over the Internet in sufficient detail to allow the students in this course to explore and learn how to use these resources on their own.

3. To provide practical instruction to these students on using the specific network resources needed during the remainder of the course.

It is assumed that the students can use a Web Browser (e.g. Internet Explorer or Netscape Navigator) to access the contents of the course.

## 1.3   How Different Internet Services are Used

As noted above, different Internet services are characterized by client software, server software, a set of capabilities they agree upon (e.g. text, pictures in the GIF format, etc.), a protocol by which they communicate (i.e. how the data is encoded in a stream of bytes), and a port on which the client contacts the server to begin communication. The distinction between different services can be blurred because different services can perform similar functions, because different services can share capabilities, and because of the existence of multifunction clients (and servers). Specifically, many modern Web clients are highly multifunctional, being gopher, ftp, e-mail and net news clients in addition to being Web clients. Finally, note that although this chapter deals with services delivered via the Internet, some of these same services can be delivered via other, very different kinds of networks, Internets, or dedicated connections.

What is presented here is a very superficial overview of a few of the available Internet services.

## 1.4   Telnet

Telnet is one of the oldest of the network services and perhaps the easiest to understand. Telnet allows one computer to "log on" to another computer as if it were a terminal. Once logged on, you frequently will have all the privileges of a local user; you can run programs, create and delete files. This is probably the most common way that users with accounts will use a computer.

Although "full service logins" as is described above are perhaps the most common use of the telnet protocol, in fact as much control as the host's system administrator desires may be imposed on a telnet connection. Thus, a telnet service may be advertised with a public login name and password. Login with this name, however, is likely to be restricted to a limited number of commands. The National Institutes of Health in the United States used, at one point, such a telnet login to disseminate information as to the membership of study sections. Such specialized telnet services have become much less common since the rise in popularity of the Web.

A telnet session can negotiate a range of different protocols, but this almost always includes ASCII text. Because many protocols for other services (e.g. SMTP, HTTP) are encoded as ASCII text, a telnet client can sometimes be used to connect to a server for these other protocols. Most people will use a telnet client the first time connecting to a MOO, and some people will continue to use telnet as their client, although most of us find dedicated clients to be significantly more convenient. Similarly, it is possible to connect to a Web server with a telnet client if you understand the syntax of HTTP. This is almost never done to *use* a Web server, but is occasionally done when debugging.

From a practical point of view, every telnet host will be different, and thus you will need to learn about each one as you have occasion to use it.

## 1.5  Ftp

Telnet is useful for interactive computer access, but is much less useful for transferring files. Ftp is an older service designed specifically for file transfer. Originally like telnet, it was intended for account owners. However, as it became apparent that it was useful to make files available to the world at large without giving all those wanting the files an account, the variant of "anonymous ftp" developed. In this variant, logging in with a "magic" user name (most commonly "anonymous" or "ftp") eliminates the requirement for a password.

In 1996, "To a large extent, use of the World Wide Web has rendered (direct) ftp access obsolete." Although there was and is some truth to that statement (especially given that files on an ftp server can be retrieved by a Web client), the need for ftp clients has not vanished. Some users will choose to avoid them, preferring the simplicity of dealing with a single piece of software, but within their domain ftp clients are more versatile than Web browsers, in some cases one has more control with an Ftp client, and for simple file transfer they are quicker and more convenient.