

**BOOKS: MOLECULAR BIOLOGY AND COMPUTING** 

## How to Make Sense of Sequences

### Carol J. Bult

multitude of computational algorithms and electronic resources are available to biologists for the analysis of DNA or protein sequences. But finding the right tools and then interpreting the results can be bewildering and frustrating. Bioinformatics resources for

Bioinformatics A Practical Guide to the Analysis of Genes and Proteins Andreas D. Baxevanis and B. F. Francis Ouellette, Eds.

Wiley-Interscience, 1998. New York, 384 pp. \$149.95, £135. ISBN 0-471-32441-8. Paper, \$59.95, £38.95. ISBN 0-471-19196-5. sequence evaluation are scattered in various places and function on different computing platforms. Some resources can be used directly through World Wide Web interfaces, some must first be downloaded and configured for the user's specific computing environment, and some are only avail-

able in integrated packages of analysis tools. After the software is selected, differences in the underlying logic or methodology of algorithms for aligning sequences, choosing phylogenetic trees, detecting functional patterns, or predicting protein sequences can produce conflicting interpretations of the data.

In Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins, Andreas Baxevanis (from the National Human Genome Research Institute at NIH) and B. F. Francis Ouellette (from the National Center for Biotechnology Information, NCBI, also at the NIH) provide a way for a scientist to chart a course through these confused seas. Despite its title, this useful compendium of resources for sequence analysis is not so much a "how to" guide as it is a companion for sequence analysis. It provides concise overviews of the structure, logic, content, and access methods for many widely used bioinformatics resources-with a strong emphasis on the "tools and databases developed

sis on the "tools and databases developed at the National Center for Biotechnology Information." *Bioinformatics* is composed of 14

chapters contributed by 16 authors, most

of whom are NCBI staff members. Although the book is not organized by themes, the chapters fall into three general categories: bioinformatics infrastructure, databases, and concepts and tools for sequence analysis.

Four chapters are devoted to the infrastructure of bioinformatics resources. Baxevanis introduces network architecture and Internet protocols. Ostell and Kans describe the NCBI data model for sequence-related information. Baxevanis offers an overview of methods for retrieving information from NCBI data re-

sources such as GenBank. Kans and Ouellette review procedures for getting sequences and their annotations into public databases. Although these chapters are not about sequence analysis per se, they provide useful insights into the development and integration of bioinformatics resources. They will help many readers develop strategies for using the online resources available at NCBI more effectively.

Three chapters describe sequence and mapping databases. Ouellette provides a detailed dissection of the format of GenBank flatfiles, the organizational core of this database's massive

assembly of gene sequences. He also relates GenBank to its international counterparts, the European Molecular Biology Laboratory (EMBL) database and the DNA Database of Japan (DDBJ). Chapters on structure databases (by Hogue and Bryant) and physical mapping databases (by Stein) provide examples on how these resources are accessed and used. (Details on using the GenBank database are included in a separate chapter by Schuler.)

The other seven chapters consider specific concepts and tools for sequence anal-

ysis. Two describe well-known sequence analysis packages: the commercially available SeqLab interface to the Wisconsin package of the Genetics Computer Group (GCG), and the freely distributed ACEDB (a Caenorhabditis elegans database). Butler succinctly summarizes the salient features of the SeqLab environment and provides an annotated list of programs in the package organized by topics (such as pairwise comparison, multiple comparison, evolution, and gene prediction). Walsh, Anderson, and Cartinhour describe the data model concepts and the sequence analysis capabilities of ACEDB. Readers unfamiliar with the concepts and jargon of object-oriented data modeling may find the first part of their chapter difficult going. Fortunately the authors recommend several on-line tutorials on implementing customized versions of ACEDB.

BOOKS AND NEW MEDIA



**Protein portraits.** Alternative representations of a portion of the barnase structure 1BN1 (*1*) produced with RasMol software (clockwise, from upper left): wireframe image with element-based coloring; space-filling model; cartoon format showing secondary structure;  $\alpha$ -carbon backbone colored by residue type.

The chapters that describe and compare tools for analyzing sequences are especially useful (and practical). Veterans and beginners alike will benefit from the concise descriptions of methods and the discussions on interpreting results of different algorithms for sequence analysis. Schuler's excellent introduction to sequence alignment and database searching successfully combines reviews of both concepts and methods. Baxevanis's chapter on multiple sequence alignment covers the major algorithms for aligning sequences, methods for

The author is at the National Center for Geographic Information and Analysis, University of Maine, Orono, ME 04469, USA, and in the Mouse Genome Informatics Group, The Jackson Laboratory, 600 Main Street, Bar Harbor, ME 04609, USA. E-mail: cbult@spatial.maine.edu

finding motifs and patterns, and means for informatively presenting multiple alignments (an often overlooked topic). Hershkovitz and Leipe contribute a very readable survey of the concepts and methods used in phylogenetic analyses, but their chapter lacks practical advice on using phylogenetic trees as predictive frameworks in molecular biology. Fickett describes the conceptual and logical framework of current approaches for finding genes and regulatory regions in DNA sequence data. Focusing primarily (but not exclusively) on methodology, Baxevanis and Landsman's chapter on predictive methods for protein sequences provides brief sketches and comparisons of tools for characterizing proteins and predicting structures.

One particularly useful feature of *Bioin*formatics is that most chapters include multiple tables and figures highlighting the concepts and resources discussed by the authors. Appendices provide a limited glossary and examples of commonly used formats for sequence files. Each chapter but one includes a compilation of Internet addresses for topics referenced within it. These lists of uniform resource locators (URLs) will help guide readers to further resources, but given the pace of change in bioinformatics some of the addresses will likely become "stale" over time.

The book suffers from the lack of a consistent approach. The level of detail about specific methods and tools varies among chapters. Some contributions (for example, Fickett; Hershkovitz and Leipe) focus primarily on conceptual aspects of sequence analysis while others (including Baxevanis; Baxevanis and Landsman; Butler) emphasize descriptions of methodologies and resources. Additionally, some practical issues facing users of bioinformatics resources are not addressed adequately. For example, most authors did not distinguish analysis and data management strategies useful with a handful of sequences from those required for working with hundreds to thousands of sequences. Readers should also be aware that although the book covers an extensive range of topics, many relevant bioinformatics resources are not included.

Despite these limitations, *Bioinformatics* offers researchers a good general overview of bioinformatics concepts, analyses, and resources. The contributors have created a reference that should be in the personal library of any biologist who uses the Internet for the analysis of DNA and protein sequence data.

#### References

#### SCIENCE'S COMPASS

NEW MEDIA: SOFTWARE

# A Handler for Big Data

### **Tony Cass**

iQ is a technical graphics, analysis, and presentation package from the same company that pro-

duces the LabView data acquisition package. Although HiQ can act as the analysis "front end" for LabView, it also functions nicely as stand-alone software for the analysis and presentation of large data sets, as discussed here.

The basic metaphor used by HiQ is that of a notebook the user creates by assembling individual elements, called objects, on a page. Objects can be data tables, graphs (two- or three-dimensional), HiQ functions, or scripts. HiQ uses Microsoft's ActiveX to work with the objects, which allows the inclusion of non-HiQ objects in the notebook (for example, Excel spreadsheets, Word files, HTML files, Adobe Acrobat documents, Origin graphs, or Adobe Photoshop images). This capability allows the user to assemble notebook components and presenobject and dropped. The resulting graph can then be viewed from different perspectives and the lighting and projections changed.

Although the user may wish to create stand-alone notebooks dedicated to performing a single function (for example, a specific statistical analysis), there are often occasions

> where the restrictions of standalone notebooks are undesirable. In these cases, users will opt for the more standard, interactive notebooks, which allow access to the full range of the program's analytical functions. These are referred to in HiQ as "problem solvers." The program comes

with several problem solvers, including a data fitter for plotting and fitting data (to predefined or user-defined functions), an ordinary differential problem solver, an integration problem solver, and several others. Problem solvers rely heavily on the HiQ scripting language, which is extensively documented both in a printed reference manual and in Adobe Acrobat format. Scripts can be accessed or created by typing in the command window or from a script object. Scripts are automatically highlighted in blue in the command window for easy viewing. Further



HiQ

National Instruments

Corporation

\$495 (students, \$45).

Phone: 512-794-0100

www.natinst.com

tations from many different sources. The latest release also includes tools to allow HiQ users to access matrix data in other formats. For example, HiQ users can directly read files created by the popular matrix manipulation program, MatLab.

Objects for data visualization and analysis would typically include a data matrix, one or more graphs, and regression or other analysis routines. HiQ imports data with the now-ubiquitous "wizard." Data can be brought into the program in binary, Excel, or text formats. Numerous custom data formats are also recognized and converted via custom import scripts. Two- and three-dimensional graphs are available in several standard formats (bar, pie, x-y, and so forth), and the user has numerous options for controlling their appearance. Creation of three-dimensional graphs is particularly easy, as the data can simply be dragged over the graph

The author is in the Biochemistry Department, Imperial College of Science, Technology, and Medicine, South Kensington, London SW7 2AY, UK. E-mail: t.cass@ic.ac.uk information on built-in functions can be quickly obtained via context-sensitive online help. Once a problem solver notebook has been created, it can be worked with interactively and easily modified to take account of changing circumstances.

The power of HiQ comes at a price—a fairly steep learning curve. Users unfamiliar with an object-oriented approach, or with scripts, will need to invest time to learn them. HiQ is not a good solution for people who work with relatively small data sets and carry out a limited number of common data-fitting routines. But for visualization and analysis of large data sets, drawn from diverse sources where results need to be presented in a dynamic and visually appealing fashion (see the figure), HiQ is an excellent choice.

HiQ, along with its associated help files and examples, installs from a CD-ROM. Although the minimum configuration for the Windows version is a 486 system with math co-processor, anything less than a 200-MHz Pentium is likely to prove frustratingly slow, even for simple tasks like scrolling. HiQ is also available in a Macintosh format.

23 OCTOBER 1998 VOL 282 SCIENCE www.sciencemag.org

A. M. Buckle, K. Henrick, A. R. Fersht. J. Mol. Biol. 23, 847 (1993).