# Bioinformatics

## FOR DUMMIES®

**2nd Edition**

**Updated to cover multiple new genomes and databases**

## A Reference for the Rest of Us!®

FREE eTips at dummies.com®

**Jean-Michel Claverie, PhD**
*Research Director, France's Centre National de la Recherche Scientifique (CNRS)*

**Cedric Notredame, PhD**
*Professor of Bioinformatics, Switzerland's Lausanne University and the CNRS*

# Bioinformatics

## FOR DUMMIES®

**2nd Edition**

Updated to cover multiple new genomes and databases

# A Reference for the Rest of Us!®

**FREE eTips at dummies.com®**

**Jean-Michel Claverie, PhD**
*Research Director, France's Centre National de la Recherche Scientifique (CNRS)*

**Cedric Notredame, PhD**
*Professor of Bioinformatics, Switzerland's Lausanne University and the CNRS*

# Bioinformatics For Dummies®, 2nd Edition

# by Jean-Michel Claverie, PhD and Cedric Notredame, PhD

contained herein may not be suitable for every situation. This work is sold with the understanding that the publisher is not engaged in rendering legal, accounting, or other professional services. If professional assistance is required, the services of a competent professional person should be sought. Neither the publisher nor the author shall be liable for damages arising herefrom. The fact that an organization or Website is referred to in this work as a citation and/or a potential source of further information does not mean that the author or the publisher endorses the information the organization or Website may provide or recommendations it may make. Further, readers should be aware that Internet Websites listed in this work may have changed or disappeared between when this work was written and when it is read.

# About the Authors

**Jean-Michel Claverie** is Professor of Medical Bioinformatics at the School of Medicine of the Université de la Méditerranée, and a consultant in genomics and bioinformatics. He is the founder and current head of the Structural & Genomic Information Laboratory, located in Marseilles, a sunny city on the Mediterranean coast of France. Using science as a pretext to travel, Jean-Michel has held positions in Paris (France), Sherbrooke (PQ, Canada), the Salk Institute (La Jolla, CA), the Pasteur Institute (Paris), Incyte pharmaceutical (Palo Alto, CA); and the National Center for Biotechnology Information (Bethesda, MD). He has used computers in biology since the early days –– his Ph.D. work involved modeling biochemical reactions by programming an 8K Honeywell 516 computer right from the console switches! Although he has no clear recollection of it, he has been credited with introducing the French word "bioinformatique" in the late eighties, before involuntarily coining the catchy "bioinformatics" by mistranslating it while giving a talk in English!

Jean-Michel's current research interests are in microbial and structural genomics, and in the development of bioinformatic methods for the prediction of gene function. He is the author or coauthor of more than 150

scientific publications, and a member of numerous international review panels and scientific councils. In his spare time, he enjoys the relaxed pace of life in Marseilles, with his wife Chantal and their two sons, Nicholas and Raphael.

**Cedric Notredame** is a researcher at the French National Centre for Scientific Research. Cedric has used and abused the facilities offered by science to wander around Europe. After a Ph.D. at EMBL (Heidelberg, Germany) and at the European Bioinformatics Institute (Cambridge, UK) under the supervision of Des Higgins (yes, the ClustalW guy), Cedric did a post-doc at the National Institute of Medical Research (London, UK), in the lab of Willie Taylor and under the supervision of Jaap Heringa. He then did a post-doc in Lausanne (Switzerland) with Phillip Bucher, and remained involved with the Swiss Institute of Bioinformatics for several years. Having had his share of rain, snow, and wind, Cedric has finally settled in Marseilles, where the sun and the sea are simply warmer than any other place he has lived in.

Cedric dedicates most of his research to the multiple sequence alignment problem and its many applications in biology. His friends claim that his entire life (past, present, future) is somehow stuffed into the T-Coffee multiple-sequence alignment package. When he is not busy dismantling T-Coffee and brewing new sequences, Cedric enjoys life in the company of his wife, Marita.

# Dedication

This is for my parents Monique and Jack, for keeping me in school, and for Chantal, for keeping me happy — in and out of the lab. It's also for my daughter Vanessa, and my sons Nicholas and Raphael, for reminding me that not *everything* in life is scientific.

–– J-MC

This is for my wife Marita, my daughter Lina, my mother Marie and in memory of my grandparents, Simone and Louis.

–– CN

# Authors' Acknowledgments

The entire Wiley staff did a great job pulling together to publish this book on tight deadlines. We'd especially like to thank our tireless project editor, Paul Levesque, and Barry Childs-Helton, who did a great job copyediting a text full of obscure biochemical words.

We'd also like to thank Amey Godse, our technical editor. Amey nailed down major and minor inaccuracies alike. His many suggestions did much to improve the book.

We also have to thank the bioinformatics community for creating the many great Web resources that we describe in this book and for making them available for free over the Internet. We personally know a number of the folks who keep these sites up and running –– and salute all of them for their hard work, enthusiasm, and dedication. Topping this list are the staff members of the Swiss Bioinformatics Institute, who run the ExPASy and the Swiss EMBnet Web server. They always went out of their way to answer any query regarding their site. The NCBI folks have also been very helpful, and we thank them for that.

We also want to pat each other on the back for making the writing of this book great fun!

Finally, we'd like to thank our families and friends, who put up with missed dinners, extra child care, changing

deadlines, late nights, and the many other demands of a project like this. We really appreciate their patience –– and promise that we won't do another one . . . at least not anytime soon!

# Publisher's Acknowledgments

We're proud of this book; please send us your comments through our online registration form located at [www.dummies.com/register/](www.dummies.com/register/).

Some of the people who helped bring this book to market include the following:

**Acquisitions, Editorial, and Media Development**

**Project Editor:** Paul Levesque

**Acquisitions Editor:** Melody Layne

**Senior Copy Editor:** Barry Childs-Helton

**Technical Editor:** Amey Godse

**Editorial Manager:** Leah Cameron

**Media Development Specialists:** Angela Denny, Kate Jenkins, Steven Kudirka, Kit Malone

**Media Development Coordinator:** Laura Atkinson

**Media Project Supervisor:** Laura Moss

**Media Development Manager:** Laura VanWinkle

**Editorial Assistant:** Amanda Foxworth

**Sr. Editorial Assistant:** Cherie Case

Cartoons: Rich Tennant ([www.the5thwave.com](www.the5thwave.com))

**Composition Services**

**Project Coordinator:** Jennifer Theriot

**Layout and Graphics:** Carl Byers, Lavonne Cook, Barbara Moore, Shelley Norris, Barry Offringa, Laura Pence

**Proofreaders:** Susan Moritz, Charles Spencer, Rob Springer, Techbooks

**Indexer:** Techbooks

**Anniversary Logo Design:** Richard Pacifico

# Publishing and Editorial for Technology Dummies

**Richard Swadley,** Vice President and Executive Group Publisher

**Andy Cummings,** Vice President and Publisher

**Mary Bednarek,** Executive Acquisitions Director

**Mary C. Corder,** Editorial Director

**Publishing for Consumer Dummies**

**Diane Graves Steele,** Vice President and Publisher

**Joyce Pepple,** Acquisitions Director

**Composition Services**

**Gerry Fahey,** Vice President of Production Services

**Debbie Stailey,** Director of Composition Services

# Contents

# Part II : A Survival Guide to Bioinformatics

# Part III : Becoming a Pro in Sequence Analysis

# Introduction

Welcome to the second edition of *Bioinformatics For Dummies!*

In the first edition, we presented bioinformatics as a brand new discipline on the rise. How right we were! Since then, it has become so prominent that anybody with an interest in biology, biotechnology, modern medicine, or (for that matter) genetically engineered food or drugs simply cannot afford to remain ignorant about the topic. With this book, you've come to the right place to quickly learn the basics.

But wait — if you expect something complicated, you're in for a (good or bad) surprise: Bioinformatics is nothing but good, sound, regular biology, appropriately dressed so it can fit into a computer.

Bioinformatics is about searching biological databases, comparing sequences, looking at protein structures, and (more generally) asking biological and biomedical questions with a computer. The bioinformatics we show you in this book can save you months of work in the lab at the minute cost of a few hours' work with your computer.

Although you'll find standard biological terms throughout, don't look here for long equations and computer-geek gibberish. The purpose of this book is to

show you quickly and plainly how to use the bioinformatics programs that you need to get your work done. On every page, we give you tricks and treats to get the most out of existing tools. If you didn't know that you can use the most sophisticated programs for free over the Internet — and that you can do this (sometimes) without installing anything on your own computer — then stay tuned: You're in for many more good surprises.

# What This Book Does for You

This book is here to help you get things done. For every standard bioinformatics task you may want to undertake, you'll find detailed steps that you can use to quickly produce the result you need.

To use most of the tools we describe in this book, you don't need to install any program on your computer. Everything we show you here runs over the Internet via your Internet browser.

If you know what you want to do — or at least know the task by name — going through the Table of Contents is the best strategy for finding exactly what you need. If you have an idea of what you want to do but you're not sure how to express it with words, Chapter 2 is here to help you decide which part of the book will suit your needs.

At the end of most chapters you'll find a convenient "Doing It for Free over the Internet" section, where we list a few carefully chosen Web sites that are similar to those we describe in the rest of the chapter. Treat this information as a spare wheel! If the main site is down, this section probably lists a convenient replacement.

# Foolish Assumptions

Putting a project's assumptions right up front is just good policy. While writing this book, we have assumed that

✔ You have a PC running Microsoft Windows.

✔ You have an Internet connection (a fast one if possible, but not necessarily).

✔ You likely have a background in molecular biology. If you don't — or if you need to brush up on your molecular biology — Chapter 1 gives you a brief overview of the basics.

✔ You know how to use an Internet browser but not much more about computers.

✔ You don't want to become a bioinformatics guru; you simply want to use the right tools for your problem and not spend days finding out about things you don't need!

✔ Most private biotech companies consider it unsafe to send data over the Internet. We assume here that the data you want to analyze over the Internet is *not* very confidential. Also, some of the "public" databases and services listed in this book require commercial users to enter into a license agreement.

# How This Book Is Organized

Bioinformatics is a broad field, with many nooks and crannies, hills and dales, and other charming features. Rather than present the whole vast discipline in one fell swoop, we've divided our discussion into five (more manageable) parts.

# Part I: Getting Started in Bioinformatics

If you have less than an hour to find out what bioinformatics can do for you, Part I is the right place for you! It tells you everything you need to know in order to actually *do* something with bioinformatics. In Part I, we also remind you of just those bits of molecular biology that you'll need to know when you do sequence analysis. We show you here how to run the main bioinformatics tools so that you know what's in store for you.

# Part II: A Survival Guide to Bioinformatics

If you want to find out everything that's ever been published on your sequence, this part is for you. It shows you how you can deal with the bioinformaticist's bread and butter: *DNA or protein sequences and their databases.* Here we tell you where you can find all the available sequences, and how to find the one you really need among zillions of irrelevant others. We also show you how to gather everything that's known in the universe about this special sequence that interests you so much (at least all of it that's available online).

# Part III: Becoming a Pro in Sequence Analysis

If you want to compare sequences, this is the part for you. Here we show you how to search databases for sequences that are similar to yours, as well as show you how to compare two or more sequences. This part also tells you how to gather hints about the function of a gene, through sequence comparisons. Finally, we give you pointers on how to produce, edit, and beautify your multiple sequence alignments so you can show them in presentations and publications.

# Part IV: Becoming a Specialist: Advanced Bioinformatics Techniques

To take full advantage of this part, you should have a pretty good idea of what you're looking for. Heavy stuff is going on here: how to predict a protein structure, how to predict an RNA structure, and how to do phylogenetic analysis. These are complicated subjects; it's simply amazing what you can do with a simple PC, thanks to the Internet resources we describe in this part.

# Part V: The Part of Tens

Welcome to our bazaar! If you haven't found what you were looking for in the other parts, you're now in the right place. The wealth of online resources that exist in bioinformatics is extraordinary — and almost overwhelming. With every student and his or her cousins putting semester reports online, finding exactly what you need with a simple keyword search can be a daunting task. In the Part of Tens, we give you a list of central resources that you can use as a starting point. Chances are that the program or server you're looking for is only one or two clicks away. In this part, we also give you ten important pieces of advice to make sure that your lab work can safely depend on your Internet work.

# Icons Used in This Book

Always eager to please, we've decided to use a series of icons in the margins of this book as a way to help you key in on important information. We came up with four, which seemed like a nice, round number.

Some particularly technoid information is coming up. You can skip it and nothing terrible will happen. Yet, if you want to be in full control of what you're doing, reading this may help! Your call. . . .

This icon shows you something simple, or smart, or a cute shortcut. In any case, it's something that can save you time and trouble.

There are many booby traps around when you use Internet servers. This icon warns you when some ambiguity surrounds what the server you're using is up to — or when disaster is only one (wrong) mouse click away. Treat the Warning icon with respect — especially in a steps list!

This icon indicates something you should remember. It can be one of the few important principles that you need to know, or it can be a very special tip — the kind that can save you three days of work (or drive you nuts if you forget it). You may assume that the head of your institute/company got to the top by discovering and applying one or more of pearls of wisdom in these very special tips!

# Where to Go from Here

If you know nothing about bioinformatics, this book is here to reassure you. Bioinformatics is a much simpler subject than you ever thought possible. For most people new to this field, the main difficulty is finding out the kind of questions they can ask with these new tools. If you're a biologist, don't let the computer scare you; bioinformatics is nothing more than good, sound, regular biology hidden inside a computer.

The magic thing about bioinformatics is that, with a simple Internet connection, you can browse databases that contain the sum of our entire human biological knowledge — and you can do this with the most