# Phylogenetic Trees Made Easy

## A How-To Manual

**FIFTH EDITION**

## Barry G. Hall

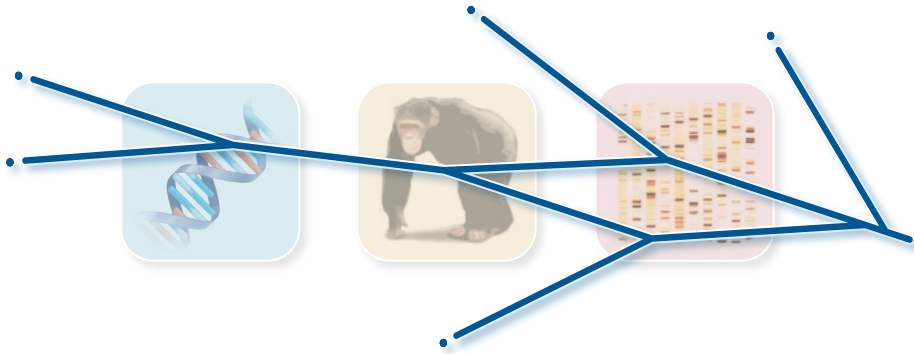# Phylogenetic Trees
# Made Easy

### FIFTH EDITION

# Phylogenetic Trees Made Easy

## A HOW-TO MANUAL

FIFTH EDITION

**Barry G. Hall**

*University of Rochester, Emeritus*
and
*Bellingham Research Institute*

SA **Sinauer Associates**

**Downloadable files to be used with this text are available at**
**oup-arc.com/access/hall-5e**

**Notice of Trademarks**
Throughout this book trademark names have been used and depicted, including but not limited to Macintosh, Mac, Windows, and Adobe. In lieu of appending the trademark symbol to each occurrence, the author and publisher state here that these trademarked product names are used in an editorial fashion, to the benefit of the trademark owners, and with no intent to infringe upon the trademarks.

*To  Miriam Barlow and John Logsdon*

# Acknowledgments

The dedication of this book requires some explanation:

Professor Miriam Barlow is in large part responsible for the existence of the First Edition of *Phylogenetic Trees Made Easy* (*PTME*). In 2000, Miriam was a Ph.D. student in my lab and had just successfully completed a course in Phylogenetic Analysis with Dr. John Huelsenbeck, but she was struggling with the mechanics of actually using software to make a phylogenetic tree. I gave her a fistful of notes I had written to myself when trying to learn phylogenetics a year earlier. She shared them with some fellow students who were similarly struggling, and soon came back to me virtually demanding that I complete them as a book. I did, and *PTME1* was the result.

Professor John Logsdon is similarly responsible for this, the Fifth Edition of *PTME*. John had been an enthusiastic supporter of earlier editions of *PTME*, and he contacted me in the Fall of 2016 asking when a new edition would be available, or more specifically, whether it would be available in time for a course he would teach in 2018. I hadn't thought much about another edition of *PTME*, but after some discussion I concluded that I could add sufficient new material to justify a new edition. This manual is the outcome of those discussions.

So, thanks to Miriam and John for being the bookends that initiated and concluded the *PTME* series. I would not have done it without you.

As always, writing *PTME5* depended on the active cooperation of the authors of some of the software described herein. My thanks to Sudhir Kumar for explaining the implementation of Timetrees in MEGA7, and to Glen Stecher of the MEGA team for responding so quickly to my queries and bug reports. My deep thanks to Andrew Rambaut for his help while learning BEAST and its associated programs. My thanks to Steven Salipante for his collaboration on MSTgold. Steve began as an undergraduate in my lab and is now a Professor at the University of Washington. Working with him has always been a joy.

Author John Scalzi once pointed out that a book is much more than the words written by the author; it is also the result of the efforts by the entire publication team. That is certainly the case for this edition of *PTME*.

Kathaleen Emerson, my Editor, has greatly improved my words and not just by catching typos and spelling errors. She has spotted inconsistencies between text and figures and rearranged the order of paragraphs to improve clarity. Most of all, she has played the role of "generic target reader" and made sure that the text always keeps that reader in mind. She is greatly responsible for making this book the best it can be.

Chris Small and Donna DiCarlo are responsible for the visual appearance of the book, its "look and feel." I am hopeless at anything visual, so I am very grateful to Chris and Donna for making *PTME5* not only informative but aesthetically pleasing.

Marie Scavotto is responsible for marketing and it is thanks to her efforts on this and earlier editions that *PTME* can be sold at a reasonable price. If she doesn't bring *PTME* to the attention of enough buyers it can't be sold at an affordable price, so I am very grateful to her for her efforts.

Last, as always, Andy Sinauer has been tremendously supportive of every edition of *PTME*. When I first thought of writing *PTME* I asked colleagues whom I should contact about publishing it. Sinauer Associates was the universal suggestion, and the universal comment about Andy was, "You can trust him absolutely." I have found that to be some of the best advice I ever received. Andy is both trustworthy and highly effective. He plans to retire in the near future, and he will be tremendously missed.

Finally, my deepest thanks to my wife of 53 years, Sue Hall, for her patience and encouragement while writing *PTME5*.

# Table of Contents

# Read Me First!

Phylogenetic analysis was once a tool used only by taxonomists, or as they are now called, systematists. Their interests were first in classification, then in elucidation of the historical relationships among organisms. Phylogenetics was considered an arcane and difficult aspect of biology that, fortunately, wasn't widely used by molecular and biochemically oriented biologists. With the advent of DNA sequencing that began to change, and the power of phylogenetics as a tool for understanding biology at all levels became apparent. Today there are few biological journals in which at least some of the published papers do not include phylogenetic trees.

Phylogenetic analysis lies at the core of the fields of genomics and bioinformatics. The emergence of bioinformatics as an important field, however, has had little effect on the perception of phylogenetic analysis as difficult and intimidating. Those who specialize in systematics, phylogenetics, genomics, and bioinformatics have done little to set aside the sense that phylogenetic analysis is a tool best used by specialists, with the result that molecular and cell biologists often seek collaboration with systematists when one of their papers requires a tree. But in most cases, preparing a robust, valid phylogenetic tree for a paper is no more challenging than using word processing software to write the paper or graphics software to prepare the figures. The purpose of this book is to make phylogenetic analysis accessible to all biologists.

*Phylogenetic Trees Made Easy* is a "cookbook" intended to aid beginners in creating phylogenetic trees from protein or nucleic acid sequence data. It assumes basic familiarity with personal computers and with accessing the web using web browsers. I have not attempted to explore all the alternative approaches that might be used, intending only to give the beginner an approach that will work well most of the time and is easy to carry out. I hope the book will also serve the investigator who has a modest familiarity with phylogenetic tree construction but needs to address some aspects and problem areas in more depth.

This book is not intended to be used as the primary text in systematics or phylogenetics courses. It can, however, be used to supplement the primary text and can serve as a tool for making the transition between a theoretical understanding of phylogenetics and a practical application of the methodology.

## New and Improved Software

The six years that have elapsed since I wrote the Fourth Edition have seen significant changes in MEGA, the software for phylogenetic analysis and tree construction that is at the core of this book. MEGA has changed significantly both in terms of increased capabilities and in terms of its interface. Details of using some other programs have also changed significantly. In general, the changes have been positive and especially beneficial to those who are new to phylogenetics.

It is neither possible nor desirable to cover all of the programs available for phylogenetic analysis; that would require a tome too heavy to lift and too boring to read. The objective here is to describe programs that are sufficient to the task and relatively easy to use. Because I focus more on phylogenetic analysis and less on learning the vagaries of different programs, I try to minimize the number of programs described, and, as time goes on, my perception of the most appropriate programs changes.

The Fifth Edition continues the practice of updating software to reflect the interface of the current version and of adding value in the form of additional advanced chapters. MEGA7 adds two important new functions: Time tree estimation and Batch processing.

The Third and Fourth Editions each introduced "advanced" topics that went beyond just estimating phylogenetic trees per se. In keeping with that practice, the Fifth Edition adds a chapter on estimating phylogenetic trees from whole genome sequences without aligning those sequences. That chapter will be particularly important to those who consider phylogenies of bacteria and viruses based on whole genome sequences. It also adds a chapter on Minimum Spanning Trees (Chapter 14), an alternative to phylogenetic trees when too much homoplasy has reduced the phylogenetic signal below an acceptable level. Finally, because so many valuable and important (and mostly free) programs are command-line programs that do not use the familiar GUI interface, I have added an appendix on the command-line interface and an appendix on installing and running command-line programs.

To keep things simple, and to make my life easier, I will from now on refer to the Fifth Edition of *Phylogenetic Trees Made Easy* as "*PTME5.*"

## Just What Is a Phylogenetic Tree?

To a mathematician, a phylogenetic tree is an abstract construct embodied in a special type of directed or undirected graph. To systematists and most evolutionary biologists, a tree is a representation of the relationships of different species to their ancestors. To a molecular biologist, a tree is a representation of

the relationships of gene or protein sequences to their ancestral sequences. All agree, however, that a tree consists of nodes connected by branches (except that mathematicians tend to refer to branches as edges).

The tips of a tree, sometimes referred to as **leaves**, are the **external nodes** and biologically they represent existing taxa. Systematists and evolutionary biologists generally think of taxa as being synonymous with species, while molecular biologists think of taxa as being synonymous with sequences. In either case, these leaves/taxa/sequences are the only entities in a tree that we can be sure about. They represent the real data—factual information—from which everything else in the tree is inferred. That factual information can be the states of morphological characters or it can be the states of nucleotide or amino acid characters in macromolecular sequences. The discussion in this book will be limited to molecular sequence characters.

Even with sequence data, however, much is assumed. We may use sequences of the cytochrome *c* gene to make a tree in which each sequence is from a different species, so we may refer to the tips as the "cow sequence," "dog sequence," etc. In fact, the cow sequence is not the cytochrome *c* sequence for all cows; it is the sequence of that gene from one particular cow. We can use that sequence to represent all cows because we assume (1) that any variation in cow sequences occurred after cows diverged from dogs, and (2) that any one randomly selected cow sequence is as representative of cows as the next. Both these assumptions arise from a fundamental assumption of phylogenetics, that of genetic isolation—that there is no genetic exchange between taxa. When this fundamental assumption is violated, different parts of the information used to make the tree (gene sequence, set of sequences, full genome) can have different histories and valid trees cannot be constructed. This issue becomes particularly important when we try to estimate a phylogeny based on whole genome sequences. For microorganisms it is almost always the case that different parts of a genome have different evolutionary histories. Two alternatives to phylogenetic trees, minimum spanning trees and phylogenetic networks, are less disrupted by incomplete genetic isolation.

The interior nodes of the tree represent *hypothetical* ancestors. We don't know the character states of those ancestors and in most cases cannot know those characteristics because the ancestors no longer exist. We can only infer the characteristics of ancestral sequences from information about their descendants (the leaves of the tree). We assume a process of modification with descent in which mutations accumulate and are inherited by their descendants. **Branches** connect nodes and represent that descent, and **branch lengths** represent the amount of change between an ancestor and its descendant.

The most important fact to keep in mind about any phylogenetic tree is that it is not fact, it is an estimate—an inference. Many evolutionary biologists would say that a tree is a hypothesis, but I disagree. A hypothesis is useful only if it can be tested, and we cannot test the validity of evolutionary trees. We can estimate trees, we can estimate how accurate they might be, and we can change our estimates (and the trees) when new data are available, but that's it. An estimate is a fine and useful thing; just keep in mind that it is not absolute truth.

## Estimating Phylogenetic Trees: The Basics

Chapters 2–10 of *PTME5* cover the basics of four distinct and equally import-
ant steps that are involved in making a phylogenetic tree based on molecular
sequence data:

1. Identify and acquire the sequences that are to be included on the tree.
2. Align the sequences.
3. Estimate the tree by one of several methods.
4. Draw the tree and present it to an intended audience.

Chapters 2–10 describe implementing these four steps using the MEGA7 soft-
ware package. The data acquisition, alignment, and tree-drawing functions of
MEGA7 are so elegantly implemented and easy to use that this is the program
of choice for most phylogenetic methods. Understanding these chapters will
permit confident estimation of valid phylogenetics trees by a well-accepted,
reliable method.

   Chapter 2 is a tutorial that takes the reader through each step, using MEGA7
to construct a simple phylogenetic tree. The primary purpose of Chapter 2 is to
familiarize you with both the mechanics of implementing the steps necessary
to make a tree and the basics of the MEGA7 software. Chapter 2 will *not* lead to
the optimal tree based on the data provided, so details of each step are covered
in the subsequent chapters.

   Chapter 3 deals in more detail with Step 1: identifying the sequences that
might be included on a tree, deciding which to include and which to exclude,
and downloading those sequences from international databases operated by
government agencies.

   Chapter 4 deals with the critical problem of aligning sequences (Step 2).
Sequence alignments, whether of nucleotides or proteins, are the data upon
which phylogenetics programs operate to estimate a tree. If you don't get this
part right, nothing you do in later steps will matter; your tree will be worthless.

   Chapter 5 very briefly discusses and compares the major methods for esti-
mating phylogenetic trees. Chapter 6 then deals with Neighbor Joining (NJ), the
most widely used tree estimation method, and discusses the often intimidating
issue of bootstrapping to assess tree reliability.

   In many ways, Chapter 6 is the meat of this book. The major phylogenetic
concepts are covered in this chapter. Subsequent chapters dealing with other
phylogenetic methods are mainly technical and simply illustrate how to accom-
plish the same goals with other methods and software.

   Chapter 7 describes how to draw trees in various ways and how to choose
which style will make it easiest for your audience to correctly interpret the tree
and understand the biological point you are making.

   Chapter 8 discusses using MEGA7 and SeaView to estimate trees by Maxi-
mum Parsimony (MP). MP is one of the first methods to have been applied to
molecular sequence data.

   Chapter 9 describes how to use MEGA to estimate trees by Maximum Like-
lihood (ML). ML is a statistical method that has long been appreciated, but was

not widely used outside of the fields of systematics and phylogenetics because it had been perceived as being intimidating, slow, and not applicable to large data sets. Recent advances in software have solved the speed problem and have made ML applicable to large data sets. I hope that after reading Chapter 9 you will find ML no more intimidating than NJ or MP.

Chapter 10 discusses constructing trees by Bayesian Inference (BI) using BEAST. Like ML, BI is a powerful statistical method and it has proven to be slightly more accurate than the other methods (Hall 2005).

Chapter 11 considers ways to decide which of the major methods to use.

Chapter 12 discusses some issues that arise when working with MEGA7 on its non-native platforms, Mac OS X and Linux.

| TABLE 1.1 |
|---|
| **Some Conventions Used in this Book** |

| CONVENTION | DESCRIPTION |
|---|---|
| Click | Use the mouse to position the cursor over the indicated button on the screen (as in "click the OK button" ), then depress and quickly release the mouse button. |
| Double-click | Click twice rapidly without moving the mouse. |
| Drag | Position the mouse and, while holding down the mouse button, move the mouse to another position |
| Select | Highlight a menu item, section of text, or an object on the screen by dragging the mouse over or clicking (or double-clicking) on the desired operation. In the text these operations are indicated in **Blue, Bold Sans Serif** type. |
| Screen display | Text in **Black, Bold Sans Serif** type indicates information you will see on your screen, but no action is required. |
| Enter | For command line programs such as MSTgold and kSNP3 and when using this book's utility programs, text shown in the `Courier typeface` indicates commands that you will type into an input file (or see on the screen if you have downloaded a utility file). |
| ⬇ Download | Indicates a file in the package that you can download from the PTME5 website |

## Beyond the Basics

Those readers without previous experience with phylogenetics will probably find the information in Chapters 2–12 sufficient to meet their needs for some time to come. Increased experience with phylogenetics, however, often leads to recognition of the utility of more advanced aspects of phylogenetic analysis. Some of those advanced topics are discussed in Chapters 13–19.

Chapter 13 concerns phylogenetic *networks* as opposed to phylogenetic trees. In some circumstances phylogenetic trees are an inadequate means of describing the historical relationships of taxa. For instance, trees based on different genes from the same set of species are often significantly different because those genes have different evolutionary histories as the result of recombination, horizontal transfer, or other "reticulate" events. In such cases phylogenetic

networks more realistically describe evolutionary history. An interior node of a phylogenetic tree may have multiple nodes descending from it, but it can have only one immediate ancestral node. In a phylogenetic network, a node can have both multiple descendants and multiple immediate ancestors.

Chapter 14 discusses another alternative to phylogenetic trees, minimum spanning trees (MSTs). MSTs do not attempt to estimate relationships of taxa to hypothetical common ancestors, only relationships of taxa to each other. They consider identity by state, not identity by descent, and thus really only reveal clustering. In effect, they only ask, "Which sequences are most alike?," ignoring whether they are alike because of descent from a common ancestor, because of convergent evolution, or because of recombination. When phylogenetic signal is weak, MST can still offer some insights into relative similarity of sequences. Chapter 14 discusses the program MSTgold for estimating MSTs.

Chapter 15 discusses time tree estimation. Time tree estimation means estimation of the divergence times of interior nodes based on outside information about the times at which a few of the interior nodes diverged. Time trees are becoming increasingly frequent and important in the literature.

Chapter 16 deals with the reconstruction of the amino acid sequences of ancient ancestral proteins and nucleotide sequences of ancestral genes. This subject, which some call *paleobiology*, permits the estimation of ancestral sequences as a necessary prelude to actual synthesis of ancestral proteins. Biochemists and others interested in understanding protein function often have an interest in synthesizing, then studying, a protein they believe to be the common ancestor of a group of proteins with distinct, but mechanistically related, functions. The first step in the process is to estimate a reliable phylogenetic tree, and the second step is to estimate the most likely sequence of that common ancestor.

Chapter 17 deals with another advanced topic, detecting *adaptive evolution* (i.e., purifying or diversifying selection) along phylogenetic trees.

Chapter 18 discusses the problems associated with estimating phylogenetic trees from whole genome sequences (WGS). Because microbial genomes of the same species are replete with rearrangements, major deletions, and major insertions it is impractical to align more than about 35–40 WGS of microorganisms. There are now tens of thousands of microbial WGS in the databases, and there is a need for practical ways to estimate phylogenies of hundreds of such genomes. This chapter discusses an alternative to genome alignment and its implementation in the program kSNP3 (Gardner and Hall 2013; Gardner et al. 2015).

Chapter 19 discusses why it is important for you to learn at least the rudiments of programming and offers some suggestions about how to approach that learning.

Appendix I discusses file formats and their interconversion.

Three new appendices discuss the command-line environment. The command-line is a very old interface that predated the current GUI with its windows, menus, mouse, etc. Most free academic programs, such as MSTgold, and kSNP3, use the command-line environment and it is important to become comfortable in that environment. Appendix II discusses text editors that are used to make input files for command-line programs, and to read their output files.

Appendix III discusses navigating the command-line environment. Appendix IV describes how to install and run command-line programs.

Appendix V briefly discusses additional programs that are often useful.

Appendix VI answers some frequently asked questions—or at least some questions that should be frequently asked.

## Learn More about the Principles

Just as it is possible to implement molecular methods without understanding them by following the protocols in commercial "kits," it is also possible to implement phylogenetic methods without understanding them by following the protocols in this book. Most of us insist that our students understand the principles underlying the methods implemented by these kits because we know that without such an understanding it is impossible to spot and troubleshoot many problems.

It is in this spirit that the reader will find *Learn More* boxes scattered throughout the text. These boxes present somewhat more detailed background on the biological and/or mathematical principles underlying the various methods and suggest further reading. It is not necessary to read the boxes to be able to estimate a valid phylogenetic tree, but understanding these principles can help troubleshoot the phylogenetic problems that arise when estimating trees from molecular alignment data.

Readers who want to go beyond the *Learn More* boxes will find Dan Graur and Wen-Hsiung Li's *Fundamentals of Molecular Evolution* (2000) and Masatoshi Nei and Sudhir Kumar's *Molecular Evolution and Phylogenetics* (2000) very helpful and enjoyable. Li's *Molecular Evolution* (1997) and Chapters 11 (Swofford et al. 1996) and 12 (Hillis et al. 1996) of *Molecular Systematics* (David Hill, Craig Moritz, and Barbara Mable, eds.) provide more detailed insights into these topics. Joe Felsenstein's outstanding book, *Inferring Phylogenies* (2004), is very technical, very mathematical, and not for the faint of heart. At the same time, it is delightfully written and well worthwhile for anyone who really wants to understand phylogenetic theory. *Phylogenetic Networks*, by Daniel Huson, Regula Rupp, and Celine Scornavacca (2011) provides an extremely detailed explanation of the concepts, algorithms, and applications of this relatively new approach (described briefly in Chapter 13 of this book). Like Felsenstein's book, *Phylogenetic Networks* is highly technical and mathematical, but it is also a very worthwhile presentation of this recent alternative to phylogenetic trees for understanding relationships among taxa.

## About Appendix VI: F.A.Q.

Well, not really "Frequently Asked Questions," but "Sometimes Asked Questions," or more likely, "Questions That Should Be Asked but Probably Aren't."

This book is organized and presented as though you will use MEGA whenever possible for every step in the process of estimating and presenting a tree. What if you want to use MEGA to download sequences, but you want to use

another program for alignment? What if you have estimated a tree using other programs, but you want to use MEGA to draw and present that tree? The more you use phylogenetics the more likely you are to want to explore and use alternative programs. Such issues can interrupt the flow of an individual chapter and be distracting, so I have collected them in Appendix VI. Look there first for tips and hints. I hope that these "What if?" questions will both help and encourage you to expand your horizons beyond the methods that are presented in this book.

## Computer Programs and Where to Obtain Them

Programs are continuously being updated. Readers should always download the latest versions from the sources listed below. Readers should also be aware that new versions might have screens that differ from the examples in this book, or they may accept commands in slightly different ways. Usually the modifications are minor and will present few problems to the careful user.

### MEGA7

MEGA7 (Kumar et al. 2016) is a modern, integrated phylogenetics package that elegantly handles the following: downloading sequences through its own specialized web browser; aligning sequences through its implementations of ClustalW and MUSCLE; estimating phylogenetic trees by a variety of methods including Neighbor Joining (NJ), Maximum Parsimony (MP), and Maximum Likelihood (ML); and drawing those trees in a variety of ways. MEGA7 is free and can be downloaded from www.megasoftware.net.

### BEAST

BEAST(Bayesian Evolutionary Analysis by Sampling Trees)  is used for Bayesian Inference (BI) of phylogenetic trees (Drummond et al. 2012). The BEAST home page beast.bio.ed.ac.uk/ provides helpful information and links to the download site github.com/beast-dev/beast-mcmc/releases/tag/v1.8.4. BEAST is free and is available for Macintosh, Linux, and Windows computers.

The current version is 1.84. The BEAST package includes several programs and some documentation. The BEAST Manual and Practical Beast in the Docs folder are helpful, but they have not been revised in some time so many of the screen shots are not consistent with the current version.

### FigTree

FigTree is an excellent tree-drawing program. It reads tree files in both Nexus and Newick formats. It also has an extended Nexus format that includes such features as fonts, branch colors, etc. FigTree is the default program for drawing the consensus tree written by BEAST. Like MEGA's tree-drawing program, FigTree permits re-rooting, rotating clades around a branch, showing or hiding branch lengths, and other features. The current version is 1.4.3. Download from tree.bio.ed.ac.uk/software/figtree/.

### codeml

Codeml is part of the free PAML (Phylogenetic Analysis by Maximum Likelihood) package. You can download PAML for Macintosh, Windows, and Unix/Linux from **abacus.gene.ucl.ac.uk/software/paml.html**. The current versions of PAML are 4.8a for Mac OS X and Linux, and 4.9.d for Windows. This book discusses the use of PAMLX, a graphical user interface for PAML (Xu and Yang 2013).

### SplitsTree and Dendroscope

SplitsTree and Dendroscope are free programs for estimating and drawing phylogenetic networks. They are available from **www-ab.informatik.uni-tuebingen.de/software/splitstree4/welcome.html** and **www-ab.informatik.uni-tuebingen.de/software/dendroscope/welcome.html**, respectively.

### Graphviz

Graphviz is a program that draws minimum spanning tree graphs from .dot files. It is required to display and print the minimum spanning trees that are discussed in Chapter 14. Graphviz is freely available from **www.Graphviz.org/Download.php**.

### Utility Programs

I have provided a couple of utility programs that are very helpful for changing file formats and for reconstructing the sequences of ancestral proteins and nucleic acids. Those programs, along with various examples, templates, and so forth can be downloaded from the *Phylogenetic Trees Made Easy* website (**oup-arc.com/access/hall5e**). Appendix IV gives details on installing and using these programs.

### Text Editors

Users who intend to use codeml, kSNP3, MSTgold, or the utility programs will need a text editor program to prepare the necessary input files. Word processor files such as those written by Microsoft Word, WordPerfect, etc. will not work. See Appendix II for recommended text editors and instructions on how to use them.

## Acknowledging Computer Programs

Phylogenetic analysis is completely dependent upon computer programs; without those programs, molecular phylogenetics simply would not be possible. The programs described in this book are free—at least they are free to us, the users. The authors of these programs have paid dearly for them in terms of time, effort, and creativity. The only reward the authors get for their effort (and it is an enormous effort) is having that software cited in publications. Please be sure to cite the software you use just as rigorously as you cite published papers. You can find the appropriate citation in the program itself, in the documentation, or at the program's website. If you can't find a paper, simply cite