# PHYLOGENOMIC DATA ACQUISITION

# PHYLOGENOMIC DATA ACQUISITION

## PRINCIPLES AND PRACTICE

W. BRYAN JENNINGS

The cover image shows a computer model of a Tn5 synaptic complex, which is comprised of a Tn5 transposase enzyme (blue) bound to the ends of a DNA transposon (red and green). These enzymes play a vital role in some Next Generation Sequencing methods. (image credit: Laguna Design, Science Photo Library/Getty Images.)

# CONTENTS

# PREFACE

Phylogenomics intersects and unites many areas in evolutionary biology including molecular and genomic evolution, systems biology, molecular systematics, phylogeography, conservation genetics, DNA barcoding, and others. Although these disciplines differ from each other in their study questions and methods of data analysis, they all use DNA sequence datasets. Phylogenomics is moving forward at a dizzying pace owing to advances in biotechnology, bioinformatics, and computers. This is reminiscent of what occurred two decades ago when the field of molecular systematics was coming of age. Another major factor that undoubtedly helped spur the growth of molecular systematics was the arrival of *Molecular Systematics*, 2nd edition (Hillis et al. 1996), a book that allowed me (along with countless others) to jump into this exciting field. As phylogenomics has grown substantially since the dawn of the genomics era—in large part due to the advent of next generation sequencing—the time is right for a book that presents the principles and practice of obtaining phylogenomic data.

This book enables beginners to quickly learn the essential concepts and methods of phylogenomic data acquisition so they can confidently and efficiently collect their own datasets. Directed at upper level undergraduate and graduate students, this book also benefits experienced researchers. The inference of gene trees from DNA sequence data represents one of the fundamental aspects of phylogenomic analysis. Accordingly, because

robust gene tree inferences are generally made using longer DNA sequences (e.g., ~200–2,000 base pairs long), this book focuses on methods for obtaining sequences in this length range.

This book is organized as follows. Chapter 1 introduces phylogenomics within a historical context, points out connections between DNA sequence data and gene trees, discusses gene trees versus species trees, and provides an overview of the methods used today to acquire phylogenomic datasets. Chapter 2 describes the landscapes of eukaryotic genomes followed by discussion of molecular processes that govern the evolution of DNA sequences. Chapter 3 continues the discussion about properties of DNA sequence loci by reviewing six common assumptions that pertain to data characteristics before describing the different types of DNA sequence loci used in phylogenomic studies. Chapter 4 covers DNA extraction methods including high-throughput methods. Chapter 5 reviews PCR theory, discusses applications in phylogenomics, and considers high-throughput workflow. Chapter 6 describes Sanger sequencing including high-throughput sequencing. Chapter 7 explains Illumina sequencing technology and how it is used to obtain phylogenomic datasets. Chapter 8 reviews theory and methods for designing novel DNA sequence loci. Finally, Chapter 9 offers a vision of the future in phylogenomic data acquisition.

Most of the information contained in this book can be found elsewhere, but it is worthwhile to

# AUTHOR

**W. Bryan Jennings** is a foreign visiting professor in the Department of Vertebrates and Post-Graduate Program in Zoology at the National Museum of Brazil and Federal University of Rio de Janeiro. He earned his BA in zoology from the University of California at Santa Barbara; MS in biology from the University of Texas at Arlington; and PhD in ecology, evolution, and behavior from the University of Texas at Austin. He was a postdoctoral fellow in the Department of Biology at the University of Washington, and in the Museum of Comparative Zoology and Department of Organismic and Evolutionary Biology at Harvard University. He was then appointed teaching fellow for one year in the Department of Molecular and Cellular Biology at Harvard before becoming an assistant professor of biology at Humboldt State University. At Humboldt, he taught genetics labs, bioinformatics, biogeography, introductory molecular biology, and introductory biology for nonbiology majors. In 2010, he moved to the National Museum of Brazil to accept a CAPES foreign visiting professorship. At the National Museum, he cofounded the Molecular Laboratory of Biodiversity Research, teaches a graduate course in phylogeography, and mentors masters students, doctoral students, and postdocs. Studies in his lab are focused on phylogenomics of vertebrates with an emphasis on phylogeographic and conservation genetics studies.