# Short Contents

# Contents

CHAPTER 1

# Introduction to Species Tree Inference

*L. Lacey Knowles and Laura S. Kubatko*

## 1.1 Introduction

Estimation of the evolutionary relationships among a collection of organisms remains a central focus of much of evolutionary and ecological study within the field of biology as these relationships provide the background for testing hypotheses in these fields. For example, support for different hypotheses about early animal evolution, and in particular the evolution of sophisticated cell types such as nerve and muscle cells, was contingent upon the phylogenetic relationships among the earliest diverging animal lineages. Especially important in addressing these questions was the placement of Ctenophora because of their shared complex cell types with bilaterians [642]. As another example, accurate time and rate estimation forms the basis for questions in ecology and evolution [468], with shifts in rates being central to tests about the drivers of diversification (e.g., [143, 596]). Clearly, accurate estimation of phylogenetic relationships that can leverage all available data within a firm inferential framework are crucial to addressing questions such as these.

Within the last 20 years, the field of phylogenetics has grown rapidly, both in the quantity of data available for inference and in the number of methods available for phylogenetic estimation. Our first book, *Estimating Species Trees: Practical and Theoretical Aspects*, published in 2010, gave an overview of the state of phylogenetic practice for analyzing multilocus sequence data at the time, but much has changed since then. Indeed, the rapid pace at which the field has advanced in the intervening time has led to the need for an updated reference. We intend this book both to serve as an update on current practice within the field and to provide a timely look toward the future.

We begin this chapter with a brief recap of the history of species tree estimation, including definitions and basic terminology. We next discuss both opportunities and challenges in the field. This discussion includes a critical look at the limitations currently imposed by data availability and computational power and how these might be expected to change in the future, but it also addresses uncertainty surrounding sampling and data analysis in the wake of the big data wave sweeping phylogenetics. We then consider inference beyond the species tree, highlighting the important problems that a genome-scale phylogeny and underlying data allow us to address in a rigorous inferential framework. We conclude with an overview of the book and its organization.

## 1.2  Background and Terminology

Prior to the routine collection of DNA sequence data, the fields of population genetics and phylogenetics were largely viewed as distinct as they addressed questions at different evolutionary time scales. Much of the mathematical and statistical development of models at the within-population scale was undertaken in the 1980s, through contributions by Kingman [364, 365, 363] and others (e.g., [746, 745]) that resulted in what is now known as *Kingman's coalescent model*, a continuous-time approximation of the Wright–Fisher (and other) population-level models. Kingman's coalescent today forms the theoretical basis for many of the methods used for species tree inference.

Following these developments, several authors noted that when Kingman's coalescent model was applied across species, inferred evolutionary relationships might vary from gene to gene. Important contributions to the development of these ideas, including mathematical details, were provided by [743], [784], [744], and [559], among others. However, much of this work went unnoticed by the phylogenetics community until the mid-1990s, when a seminal paper by Maddison [455] provided clear descriptions of the possible causes of differences in gene-level and species-level phylogenies. This coincided with a decrease in the cost of DNA sequencing, and the subsequent availability of multi-locus sequence data prompted several authors to highlight the need for new inferential frameworks to accommodate these data properly [813, 538, 633, 634].

Importantly, the potential for differences between gene trees and species trees were also recognized to result not only from the coalescent process but also from other evolutionary processes, such as horizontal transfer and gene duplication and loss. By the early 2000s, several papers highlighted the possibility of variation in the evolutionary history across the genome in carefully annotated empirical data sets (e.g., [134, 630, 213]), and the need for methodology that specifically aimed to estimate species-level phylogenetic trees became well accepted by many in the community.

### 1.2.1 DEFINITIONS AND TERMINOLOGY

A species tree or species phylogeny can be defined as a rooted bifurcating phylogenetic tree for which the tips of the tree represent species and the internal nodes represent speciation events. The times associated with internal nodes of the tree represent the times of speciation events, and branch lengths along the species phylogeny represent the amount of time between speciation events. Speciation times are often given in coalescent units, which can be defined as the number of $2N_e$ generations, where $N_e$ is the effective population size. The advantage of using coalescent units to describe speciation times is that a standardized unit can be discussed in such a way that characteristics associated with this unit can be "translated" to any species of interest once the generation time in years and the effective population size are specified. When $N_e$ varies across the tree, it may be more difficult to define an appropriate unit (number of generations is a reasonable choice, see [446]). Mutation units, the unit commonly used for gene tree inference that is given by the number of substitutions per site per unit time, are also sometimes used. Figure 1.1 shows an example species phylogeny for three taxa, labeled $A$, $B$, and $C$ (shaded, thicker tree in each panel).

A gene tree represents the evolutionary history for an individual gene, where a gene is defined as a stretch of contiguous sequence of any length. The tips of a gene tree represent sequences collected from individuals sampled from a particular species, while the internal nodes represent gene divergence times (looking forward in time) or common

**(a)** Gene tree-species tree match

**(b)** Incomplete lineage sorting

**(c)** Horizontal transfer

**(d)** Hybrid speciation

**(e)** Gene flow after speciation

**(f)** Gene duplication and loss

Figure 1.1. Relationships between gene trees and species trees. In each panel, the species tree is represented by the shaded, thicker tree. Speciation events are indicated with horizontal dotted lines, and the length of time between speciation events is denoted by $t$. Gene divergence, or coalescent, events are indicated in panel (a) by black circles. Each panel shows a possible relationship between the gene tree and the species tree resulting from a specific evolutionary process: (a) The gene tree and species tree share the same topology. (b) The topologies of the gene and species trees are discordant due to incomplete lineage sorting. Tracing the lineages sampled from species $B$ and species $C$ back in time, we see that they fail to coalesce in the immediately ancestral population, and instead the lineage sampled from species $C$ coalesces with that sampled from $A$ in the common ancestral population. (c) Genetic information is transferred horizontally across the phylogeny from species $A$ to species $C$, leading to a gene tree that is discordant with the species tree. (d) A species network in which species $C$ is a hybrid of species $A$ and $B$ is shown. For the particular gene sampled, species $C$ inherited its genetic material from species $A$. Owing to the hybrid speciation event, it is possible for $C$ to inherit genetic information directly from either $B$ or $A$, even in the absence of incomplete lineage sorting. (e) Gene tree discordance due to gene flow from $A$ to $C$ following speciation. (f) A gene duplication event, marked by a star, occurs after the separation of the lineage leading to $A$ from the ancestor of $B$ and $C$; the duplicated lineage is sampled in $A$ and $C$, while the original lineage is sampled in $B$, leading to discordance between the gene tree and species tree. See also figure 7.1.

ancestor events for the sampled sequences (looking backward in time). These are sometimes also called coalescent events. A gene tree may have many more tips than a species tree because multiple individuals may be sampled within each species included in the species phylogeny. A gene tree may differ from the species tree that gives rise to it both in terms of its topology (branching pattern) and in terms of the times associated with its nodes. Differences in topology between gene trees and the species tree can result from many different evolutionary processes. For example, incomplete lineage sorting (i.e., the failure of lineages to coalesce in their immediately ancestral population) can lead to gene trees with topologies that differ from the species tree (see figure 1.1b). This form

of gene tree discordance is typically modeled by applying Kingman's coalescent across the phylogeny (which is then commonly referred to as the *multispecies coalescent*) and is well studied; in particular, the probability distributions of both gene tree topologies [179] and gene genealogies [601] have been derived.

Horizontal transfer (figure 1.1c) is another evolutionary process that is well-known to generate discord between gene trees and the species tree and refers to any process by which genetic information is moved from one species to another by means other than modification with descent. For example, in bacteria, horizontal transfer occurs when distinct bacterial strains recombine to generate unique sequences that include genetic material from both strains. In sexually reproducing organisms, horizontal transfer can occur when a virus or other vector moves a segment of DNA from one species' genome to another. Hybridization (figure 1.1d) and introgression/gene flow (figure 1.1e) can also be thought of as forms of horizontal transfer, in that these processes both involve the exchange of genetic material between distinct, contemporaneous species (i.e., "horizontally" along the phylogeny) rather than through a process of descent with modification within a single species. Regardless of the precise mechanism by which the horizontal transfer occurs, such processes can result in portions of the genome that are inherited differently than others. For example, introgressed loci will show a pattern of inheritance from a species different than that of the majority of the genome if the introgression occurs between non-sister taxa (e.g., figure 1.1c). In the absence of other processes, the extent of discordance due to horizontal inheritance will depend on the extent to which genetic material has been transferred from one species to another throughout the evolutionary history of the set of species under consideration.

The process of gene duplication and loss (figure 1.1f) provides another evolutionary mechanism that results in differences between gene trees and species trees. When a gene is duplicated in a genome, the two versions of the gene subsequently evolve independently of one another, and in descendent species one or both versions of the gene may be present in the genome being sampled. Depending on which copy is sampled, the gene tree for the locus under consideration may differ from the true species-level relationship. Loss of one copy of a duplicated gene may also lead to incongruence between the gene tree and the species tree, or may result in missing data for the locus under consideration, depending on the time that has passed since the duplication and loss events. Gene duplication and loss is prevalent in many species and provides an important mechanism for the generation of new gene function (e.g., a duplicated copy of the gene is under less evolutionary constraint and may evolve to provide a new function in the organism). Thus, consideration of this evolutionary process at the stage of species tree inference is crucial, and many methods have been and continue to be proposed for inference in the presence of duplication and loss.

Closely related to the concept of a species tree is that of a species network, in which relationships between species are depicted by a sequence of speciation events, as in a species tree, but in which species may arise from more than one immediately ancestral species. This may result from evolutionary processes such as hybrid speciation (figure 1.1d), extensive gene flow between distinct species (figure 1.1e), or other forms of horizontal transfer. Much recent work has focused on carefully defining species networks and developing methods of inferring such networks from phylogenomic data, often within a coalescent framework (see, e.g., [845, 841, 843, 713, 861], as well as chapters 5 and 6 of this volume).
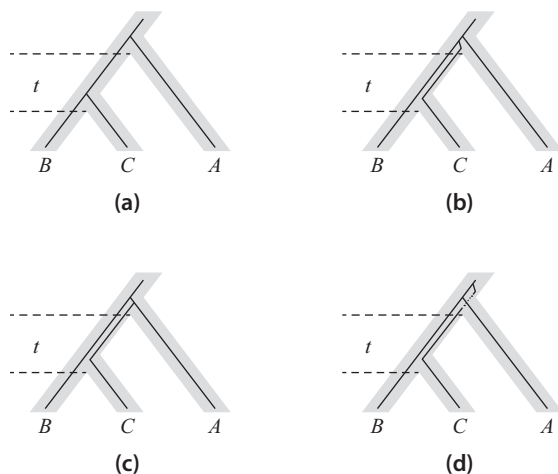
Figure 1.2. Four coalescent histories compatible with a three-taxon species tree. Note that the histories in (a) and (b) share the same topology as the species tree, while those in (c) and (d) do not.

## 1.2.2 AN INTRODUCTION TO THE MULTISPECIES COALESCENT

As mentioned in the previous section, the multispecies coalescent model underlies many of the methods for species tree inference that are commonly applied to multilocus data. Rather than provide a complete mathematical description of this model, we provide here an introduction to the main ideas for three-taxon trees. Readers wishing to see a more full description can consult [383, 770, 289].

Figure 1.2 shows the same three-taxon species trees as shown in figure 1.1. Embedded within the species tree are the four possible *coalescent histories* consistent with this species tree, where coalescent histories refer both to the gene tree topology and the species tree branch lengths along which coalescent events occur. Note that the history in figure 1.2a is the only one in which the first coalescent event occurs within the species branch of length $t$. Under Kingman's coalescent, times to coalescent events follow an exponential distribution with rate given by $\binom{n}{2}$ when $n$ lineages are available to coalesce. Since $n = 2$ lineages are available to coalesce in the interval of length $t$ in figure 1.2a, the probability of observing this history is the probability that an exponential random variable with rate 1 is less than $t$, which is $1 - e^{-t}$.

Since the probability associated with all four histories must sum to 1, this leaves $e^{-t}$ of the probability to be distributed over the other three histories, shown in figure 1.2b–d. Note that these three histories all involve the first coalescent event occurring above the root of the species tree, and all three lineages are available to coalesce within this ancestral population. Under Kingman's coalescent, each pair of lineages is equally likely to be the first to coalesce, and thus each of these histories has probability $\frac{1}{3}e^{-t}$.

Finally, we note that the first two histories (figure 1.2a and b) have the same gene tree topology. Thus to derive the probability distribution of gene tree topologies, we can add these two probabilities. The coalescent model then specifies that for three species, the gene tree topology that matches the species tree occurs with probability $1 - \frac{2}{3}e^{-t}$, while the two nonmatching gene trees each have probability $\frac{1}{3}e^{-t}$. Noting that $1 - \frac{2}{3}e^{-t} \geq \frac{1}{3}e^{-t}$ with equality only when $t = 0$, we can identify a common pattern for which the coalescent model is a good fit: a dominant gene tree topology that occurs with highest frequency (the one matching the species tree) with the two alternative topologies occurring in lower and approximately equal frequencies. Such a pattern has

been observed for empirical data [565, 145], and deviation from this pattern has been
used as evidence for introgression [652].

### 1.2.3 DATA TYPES AND TECHNOLOGIES FOR GENERATING PHYLOGENOMIC DATA

New data collection techniques have driven shifts in not only the quantity of data
but also in the types of data available for phylogenetic inference, with a variety of high-
throughput phylogenetic data collection technologies to choose from (table 1.1). These
range from different types of targeted sequencing technologies (e.g., hybrid enrichment
strategies; [422, 611]) to random genomic sequencing (e.g., reduced representation
restriction site-associated DNA sequencing [RADseq]) or targeted genotyping-by-
sequencing (GBS) (e.g., RAPTURE; see [8, 60]) and whole transcriptome or genome
sequencing.

One important factor in deciding among the different technologies is the differ-
ences in their costs, both in terms of the initial time investments and expense but
also associated costs when expanding to large numbers of taxa (or individuals). For
example, amplifying targeted amplicons involves substantial costs for setup, but it is
relatively inexpensive to capture sequences, whereas random genomic sequences from
RADseq technologies are economical and provide a universal approach for collecting
comparative genomic data. As sequencing costs drop, whole transcriptome and genome
sequencing are becoming more widely applied [853, 425]. Alternatively, RADseq can
generate very large numbers of loci (i.e., in the thousands to millions of loci) while
being scalable to large sample sizes [414], including hundreds of thousands of individ-
uals with targeted genotyping-by-sequencing, and, because of the short sequence reads,
they are amenable for applications to museum specimens for which DNA degradation
can preclude large amplicons [837].

Another primary consideration for choosing a technology (besides the cost and ease
of setup) is differences in their utility. For example, the very large numbers of loci
generated by technologies like RADseq become highly desirable for estimation of phy-
logenetic relationships at recent time scales (e.g., [475, 465]). However, their utility
drops as the evolutionary distances between taxa increase (but see [779]) because of
allele dropout (but see [210]), which will result in missing data among more distantly
related taxa (i.e., homologs will not be sequenced in some taxa because of mutations
in the enzyme cutter sites, although new technologies guard against allelic dropout; see
[84]). Decisions about what threshold of missing data to use for analysis of RADseq
data is complicated. Eliminating loci with a lot of missing data can result in a biased
data set with an overrepresentation of loci with low mutation rates [318], which means
the data set may not contain the actual loci that are phylogenetically informative for
resolving relationships among taxa that diversify rapidly—that is, loci with the high-
est rate of evolution. On the other hand, discordant relationships have been shown to
be disproportionately represented among loci with missing data [413], suggesting that
they may be less reliable for phylogenetic inference. Whole-genome or transcriptome
sequencing has the appeal of providing not just a lot of data for phylogenetic inference
but also information to address questions provided by the phylogenetic framework,
including questions about genome evolution [428]. However, in addition to assembly
challenges, such data also pose new challenges because of the potential heterogeneity
of processes contributing to genomic differences among taxa, making model misspec-
ification a more pressing problem compared with the relatively small data sets (e.g.,
hundreds to a few thousand loci). In contrast, targeted amplicon approaches such as

**Table 1.1.** Summary of sequencing technologies.

| Method | Description | Data | Reference |
|---|---|---|---|
| Targeted enrichment | Also known as "hybrid enrichment," "anchored enrichment," and "enrichment"; may also be referred to as "capture." Probes or baits complementary to the sequence of interest are used to hybridize with the target sequence, which is then enriched via PCR and sequenced on a HTS platform. | Typical enrichment/capture data sets for phylogenomics contain tens to hundreds of taxa and hundreds to thousands of loci. | Mertes et al. 2011 (https://doi.org/10.1093/bfgp/elr033); Cronn et al. 2012 (https://doi.org/10.3732/ajb.1100356); Folk et al. 2015 (https://doi.org/10.3732/apps.1500039) |
| Transcriptome sequencing | This method involves sequencing RNA reads (typically via RNA-Seq on a short read HTS platform) and assembling the reads into the transcriptome. When transcriptomes are used for phylogenomic inference, a necessary step is orthology inference to ensure that markers have sequence homology. | This method can generate hundreds to thousands of loci. The number of taxa that can be included is more variable among studies but is typically fewer than 100 due to costs. | Wen et al. 2013 (https://doi.org/10.1371/journal.pone.0074394); Shen et al. 2017 (https://doi.org/10.1093/gigascience/gix116) |
| Whole-genome sequencing | Whole-genome sequencing for phylogenomics has the advantage of obtaining an unbiased representation of genomic data relative to reduced representation methods. As costs for genome sequencing fall, and computational efficiency increases, this method is becoming increasingly accessible to many researchers. Orthology assessment is also a necessary step when using loci from whole-genome sequencing. | Currently, whole-genome uses similar amounts of data as transcriptome sequencing—thousands of loci for tens to hundreds of individuals. | Allen et al. 2017 (https://doi.org/10.1093/sysbio/syw105) |

*(continued)*

**Table 1.1.** (*continued*)

| Method | Description | Data | Reference |
|---|---|---|---|
| RAD-Seq / GBS | This is a reduced representation library approach in which 1–3 restriction enzymes are used to fragment the genome. Adapters with built-in barcodes are ligated to the fragments, then the fragments are pooled and size selected. The pool is typically sequenced on a short-read HTS platform (e.g., Illumina). | The number of loci sequenced at a given depth depends on the genome size and number of individuals multiplexed in a sequencing run. For example, a user could reasonably expect to get 5,000–20,000 loci with 10X coverage when multiplexing 1000 individuals from a species with a 1 Gb genome in a lane of Illumina HiSeq. | Andrews et al. 2016 (https://doi.org/10.1038/nrg.2015.28) |
| Parallel microfluidic PCR + HTS | This approach uses primer pairs designed for PCR and utilizes microfluidic PCR to amplify loci in parallel using HTS. The primers can be designed for organelles or the nuclear genome, but they must have similar annealing temperatures since the PCR will be processed simultaneously in parallel. This method bypasses traditional library preparation. | Typical studies use 48–96 loci for approximately 50–100 taxa. | Uribe-Convers et al. 2016 (https://doi.org/10.1371/journal.pone.0148203); Kates et al. 2017 (https://doi.org/10.1016/j.ympev.2017.03.002) |

*Note:* HTS = high-throughput sequencing, RADseq = restriction site-associated DNA sequencing, and GBS = genotyping by sequencing.

hybrid enrichment approaches avoid the problems of missing data by relying on con-
served sets of priming sites to amplify sequences. They also present less of a challenge for
assembly, modeling, and analysis compared with technologies like RADseq and whole-
genome/transcriptome sequencing. However, they also result in substantially fewer loci,
and because they rely on specific priming sites, they are nonrandom samples of the
genome, which may make them less desirable for some questions.

These different data set properties (e.g., SNP-based information content, or inherent
heterogeneity in underlying evolutionary model with genomic-scale sampling, and/or
differing amounts and distributions of missing loci in data sets) are likewise driving dif-
ferent analytical and theoretical areas in phylogenetic inference. These new areas range
from exciting new approaches for phylogenetic estimation and the evaluation of the
confidence of such relationships (e.g., assessing phylogenetic signal; [423, 771]) to deter-
mination of the different processes contributing to locus-specific patterns of ancestry
(e.g., [88, 371, 771]) and identification of subsets of data for phylogenetic inference from
genome-scale data sets [675, 192, 319]. The analytical methods that might be applied
will also differ depending on the technology used to generate the data. For example,
the short sequence reads of RADseq means that they are not generally amenable to
gene tree estimation but instead are analyzed as SNP data, whereas standard gene tree
estimation methods are applied to sequences generated from technologies like hybrid
enrichment because those technologies target specific genomic regions of longer read
lengths. Likewise, with genome-scale data sets, computational challenges restrict the
types of analyses that might be done [503].

The new technologies and unprecedented abundance of data they generate is chang-
ing phylogenetic inference and no doubt providing better resolved and more reliable
phylogenetic inference in some cases. However, recalcitrant nodes persist (e.g., [798,
590]). Moreover, with phylogenetic estimates differing as a function of analysis, data
set design, or inclusion/exclusion of loci, genome-scale data sets are raising many ques-
tions with no clear answers. For example, how might genome-scale data be analyzed to
provide reliable phylogenetic estimates? If subsets of the data are to be analyzed, how
should such data be identified (both in terms of loci and taxa)? These are some of the
questions that are explored in this book, as researchers contend with the uncertainty
surrounding sampling and data analysis in the big data era. Despite these unknowns,
it is clear that along with these complicated questions come some amazing opportuni-
ties that extend beyond a focus on the species tree itself. As we look to the future, and
in the following chapters, we emphasize this expanded role of genome-scale data—that
is, next-generation inference, which will no doubt become the new focus of researchers
as next-generation sequencing becomes routine (table 1.1).

## 1.3  Overview of Current Methods for Species Tree Inference

Given the processes described above, the precise mechanism by which data arise
must be taken into consideration in the development of methods for inferring species-
level phylogenies. Regardless of the process(es) responsible for gene tree–species tree
discordance, it is usually assumed that gene trees arise from evolutionary processes
occurring along the species tree, and DNA sequence data are subsequently generated
from the gene trees associated with individual loci. Thus, DNA sequences observed
from loci that are freely recombining can be viewed as conditionally independent of one
another, where the conditioning is based on their underlying gene trees arising from a
shared species phylogeny. Inference then proceeds in the "reverse" direction—that is,
given a set of observed DNA sequence data from multiple loci, it is desired to obtain an

estimate of the species tree. Although gene trees are not directly observed, it is clear that
they play an important role in the data-generation mechanism. For this reason, methods
for estimating species trees are commonly categorized according to how they account
for uncertainty in the gene trees in carrying out inference.

One class of methods for species tree inference is referred to as *summary statistics
methods* or *summary methods* because these methods carry out species tree inference
in two distinct steps, the first of which represents a summarization of the data. In this
first step, a gene tree is estimated for each locus in the data set using one of the standard
methods for phylogenetic tree estimation (e.g., maximum likelihood). The gene trees
estimated in this first step are then used as input to the second step of the procedure,
and a species tree estimate is obtained using only the information contained in these
input gene trees. Such methods have the advantage of being computationally efficient.
In the first step, the gene trees for the individual loci can be estimated in parallel, as
each depends only on the sequence alignment for that gene under the conditional inde-
pendence assumption mentioned in the previous paragraph. The second step is typically
carried out by assuming some model for the relationship between the gene trees and the
species tree. The features of the models that are commonly used for inference typically
lead to computationally tractable algorithms for inference. The drawback of summary
methods is that uncertainty in the gene tree estimates is typically unaccounted for in the
second step of the procedure, making estimation of the uncertainty in the species-level
phylogenetic estimate difficult to quantify. Though some suggestions have been made
to remedy this (e.g., using as input bootstrap samples of gene trees rather than only a
single point estimate of the gene trees), these strategies have not definitively been shown
to improve the overall inference.

The second class of methods for species tree inference is referred to as *coestimation
methods*, because they jointly estimate gene trees and the species tree under probabilis-
tic models. The current methods within this class employ a Bayesian framework and
use Markov chain Monte Carlo (MCMC) to carry out inference. Such methods have the
advantages that fairly complex models can be fit and that estimates of all model para-
meters and associated measures of uncertainty are naturally obtained as part of the
MCMC inference procedure. However, these methods can be computationally inten-
sive, particularly as the size of the data and the complexity of the model increase. Most
current methods cannot be feasibly run on genome-scale data for more than 20 or so
species, but as computational power continues to grow, the ranges of data set sizes and
models that can be successfully analyzed in this framework will continue to expand.

A third class of methods are those based on site pattern frequencies, which generally
involve the use of genome-scale data to compare data features expected under various
evolutionary models to those found in observed data sets. Importantly, site-pattern-
based methods are distinct from summary methods as they are applied directly to the
sequence data with the need to estimate gene trees first. Such methods have been shown
to be computationally feasible for large data sets and thus show promise for carrying
out efficient inference for large data sets for which analysis in the Bayesian coestimation
framework is computationally prohibitive.

Across all of these classes of methods, the most common model assumed for the
relationship between gene trees and the species tree is the coalescent process; however,
other processes have been considered. For example, in the class of summary statistics
methods, a model of duplication and loss can be used to carry out inference under
the parsimony criterion by selecting the species tree that minimizes the number of
gene duplication or loss events required to explain the set of estimated gene trees.

Of course, important issues arise, such as whether duplications and losses should be weighted equally and how the space of possible species trees should be searched to find the tree that minimizes the number of duplications and losses. A challenge for the future will be the development of methods that can simultaneously account for many of the evolutionary processes known to generate gene tree–species tree discord in a computationally feasible manner.

## 1.3.1 CONTROVERSIES IN THE ESTIMATION OF SPECIES TREES

Careful testing that employs both simulated and empirical data and that includes comparisons between methods are important components of the development of methods for the problem of species tree inference, and work in this area has predictably led to disagreements among researchers about the appropriateness of various methods. One source of controversy has involved the utility of the multispecies coalescent model in improving the accuracy of estimated species-level phylogenies. Two observations have contributed to this viewpoint. First, when the possibility of discord in the gene trees is ignored and the data from all of the sampled loci are concatenated and analyzed with a method for gene tree inference, such as maximum likelihood, the resulting tree is often an accurate estimate of the species tree topology, particularly when the length of time between speciation events is large. The second, and related, observation is that incomplete lineage sorting is most commonly observed in empirical data for recent speciation events that have occurred in quick succession (e.g., species are often not monophyletic when multiple individuals are sequenced), which has led to the speculation that incomplete lineage sorting, and therefore the need for methods that explicitly model the coalescent process, is not relevant for nodes "deeper" in a tree. However, the coalescent model predicts that the amount of incomplete lineage sorting depends upon the length of time between speciation events, regardless of whether speciation was recent or occurred deep in the past (i.e., regardless of the depth in the species tree). Thus, it is not reasonable to conclude that the possibility of incomplete lineage sorting can be disregarded in the inference stage for deeper histories without knowledge about the time separating speciation events.

In regard to the observation that concatenated methods often infer the species tree topology with high accuracy, a more nuanced view is required to relate this observation to the performance of such methods. First, concatenating multilocus data and carrying out, for example, a maximum likelihood analysis on the resulting concatenated alignment makes the assumption that all loci share a single underlying phylogenetic tree from which data evolve according to the nucleotide substitution model specified during the analysis. Yet several empirical examinations (e.g., [134, 630, 213]) have established that the phylogenetic history varies across the phylogeny in a way that is often consistent with the multispecies coalescent. While one might argue that phylogenetic inference commonly requires simplifying assumptions (for example, methods for estimating gene trees from single-locus alignments typically assume that sites evolve independently), these assumptions are generally necessary because computationally efficient methods that implement more appropriate models are lacking. The variety of computationally tractable methods described within this book demonstrates that this is not a limiting factor in the case of modeling gene tree discordance.

Second, and more importantly, while the topology may be estimated fairly accurately, other important quantities may not be. One example is the quantification of uncertainty in the estimated tree, which is typically carried out by bootstrapping when a method like

maximum likelihood is used to analyze the concatenated alignment. Two issues arise when using the bootstrap on the concatenated alignment. First, the bootstrap assesses uncertainty in the repeated application of a particular method to data from the population of interest and thus cannot address uncertainty that results from inaccurate assumptions in the model used [222]. Given that concatenation represents an incorrect modeling assumption (i.e., that all loci share a common gene tree), the bootstrap support values are difficult to interpret. Second, bootstrap support from large concatenated alignments will tend to overestimate the actual support for a node. For example, consider a node supported by 55% of the sites in a data set while 45% of the sites favor some other arrangement, and suppose that the data contain one million sites (thus, 550,000 bp favor the node of interest, while 450,000 bp favor the other relationship). Most bootstrap samples will contain a majority of sites favoring the node of interest, and thus the bootstrap support for the node under consideration will be near 100%. Yet, the data show much more even support for the two alternative relationships. For example, significant underlying conflict was masked by high support values of concatenated analyses when conflicting phylogenetic relationships were actually strongly supported (e.g., [675, 774, 426]).

Finally, species divergence times (or speciation times) will be inaccurately estimated using a concatenation method as such a method assumes that the common ancestor for all loci in the concatenated alignment was identical. However, under a model such as the coalescent, it is clear that the gene divergence events (i.e., common ancestor times for the individual loci) must all predate the time of the speciation event in the absence of gene flow or some other form of horizontal transfer (see figure 1.1). Use of a coalescent model allows estimation of the speciation time after accounting for this, while concatenation methods do not. An example of the potential consequences of this is given in [390] (their table 5) for an empirical data set of *Sistrurus* rattlesnakes, for which the speciation times estimated for a concatenated versus coalescent-based analysis differ by as much as 70%. In addition, tree-associated parameters other than the speciation times, such as overall evolutionary rates, effective population sizes, and rates of gene flow, could potentially be affected by model misspecification when loci are concatenated.

## 1.4  A Look to the Future

### 1.4.1 CURRENT LIMITATIONS AND FUTURE PROSPECTS

The primary challenge facing current methods of species tree inference arises from the conflict between the desire to fit increasingly realistic, and therefore complex, models to the data and the computational resources required to fit such models as the size of data sets grows both in terms of taxa and genomic coverage. Overcoming this challenge will require new approaches to these problems, and these approaches will need to be designed specifically for the problem at hand. Here, we highlight several of the important issues to be addressed when the goal is to infer a species tree using phylogenomic data.

First, we note that new methods must be designed by carefully considering the data characteristics specific to phylogenomic data. In particular, properties of methods developed for inference of gene trees from single loci may need to be re-evaluated when applied to genome-scale data. As a first example, we note that when carrying out analyses of single loci, it is common to use a model selection procedure (e.g., ModelTest [584]) to choose a model of nucleotide substitution prior to inference of the gene tree. However, for a phylogenomic data set of several hundred or several thousand genes, selection of specific models for individual genes to be specified in the downstream species tree

inference procedure might yield little increase in the statistical power for tree inference at the expense of computational time in the model-fitting stage.

As a second example, consider the case of *phylogenetic invariants*, which were proposed in the late 1980s by [398], [124], and [125] as a possible method for inferring gene trees for samples of three or four taxa. Although promising from a theoretical standpoint, later work by [333] showed that invariants-based methods lacked power for gene tree inference and were outperformed by other methods in common use at the time. One explanation for the result of Huelsenbeck and Hillis is that phylogenetic invariants are formed from polynomials in the site pattern frequencies, where the site pattern frequencies are estimated from the data. For single loci that may be only a few hundred to a few thousand base pairs in length, accurate estimation of site pattern frequencies may be difficult, and polynomials formed from these estimates may have high variance, making them ineffective at differentiating among trees. For genome-scale data, however, a wealth of data are available, and estimates of site pattern frequencies are expected to be much more accurate, making invariants a reasonable tool for examining species-level phylogenetic relationships. The success of invariants-based methods, such as the ABBA-BABA test [208], in addressing complex problems (i.e., hybridization) indicates that the performance of such methods warrants new examination in light of the very different data structure provided by multiple loci.

A second challenge in developing methodology for species tree inference is that the availability of large quantities of sequence data provides the opportunity to use subsets of the data selectively to address specific questions of interest. Using only portions of the data can clearly result in increased precision in the resulting inferences (e.g., excluding loci with errors in alignment, assembly, or orthology detection; see [95]), but it also risks the introduction of bias when the data are not appropriately sampled. As described above, sequencing technologies in current use often result in large quantities of data, but the quality of the data for phylogenetic inference may differ. For example, some loci might be characterized by large amounts of missing data, or some loci may be involved either in a disproportionate amount of discord throughout the tree or discord in parts of the tree may arise by processes other than those captured by the phylogenetic model used for inference (e.g., [700, 95]). Depending on the method used for species tree inference, the investigator may be required to make a decision about whether to exclude taxa and/or loci because of the pattern of missing data (see, e.g., [638]). However, few methods currently incorporate a model of missing data explicitly in the inference procedure. Similarly, although data filtering is also becoming more common during the preprocessing of data in preparation for species tree inference, methods that explicitly incorporate the filtering step into formal inference model are lacking. The filtering may involve inclusion/exclusion of genes due to level of variation, the possibility of horizontal transfer or gene duplication, or numerous other reasons. However, when the process of sampling data is not correctly incorporated into the underlying models used for inference, the inferred phylogeny may be biased. It is unknown how much bias may be introduced by the nonrandom inclusion of data, and there is as yet no consensus on how filtering decisions should be made [371].

## 1.4.2 BEYOND THE SPECIES TREE

Even though a major challenge to species tree resolution is the gap between the data we collect for phylogenetic analyses (i.e., large-scale transcriptomic and genomic data) and the methods that accommodate the inherent complexity of big data (i.e., processes in addition to incomplete lineage sorting that contribute to discord among loci), such

data is also an unprecedented opportunity to better illuminate the processes that shape the tree of life. That is, big data and all its complexities when studied in a phylogenetic framework open new opportunities to address questions beyond the primary goal of resolving the species tree.

The characterization of patterns of discordance and the contribution of different processes to this discordance is itself of interest for generating hypotheses about the role of lateral gene transfer, gene duplication, and incomplete lineage sorting during the divergence of different taxa [627, 697]. Such hypotheses include those about the distribution of duplication across the tree of life and its potential association with the shifts in diversification rates or its concentration at the origin of major clades (as opposed to being dispersed across taxa). Likewise, by focusing on the processes that lead to gene tree discord, we can test whether lateral gene transfer is commonly associated with hypothesized ecological transitions versus evidence of convergent molecular evolution (e.g., [249]).

In the future, and as comparative genomics expands [425], the species tree and all its ancillary applications will only grow (e.g., dissecting evolution and disease; [491]). Such novel biological questions depend upon seeking evolutionary explanations for the distribution of discord in the gene trees, which means that the heterogeneity of comparative phylogenetic data sets should be embraced (as it is arguably the key to accurate resolution of recalcitrant phylogenetic relationships).

## 1.5  Organization of This Book

In the book we highlight, by example, not only how species tree estimation differs from a phylogeny estimated from concatenated multilocus data but also the issues that arise more generally from a mismatch between phylogenomic data, with all its complexities, and the models used for inference. This includes both conceptual and practical issues related to improving species tree estimates, as well as inferences that are nonbifurcating (i.e., networks). The book devotes five chapters to methodological developments, whereas the latter portion of the book, a total of eight chapters, focuses on empirical applications, including those that consider questions beyond the species tree. Some of the contributors were participants in a workshop offered at the 2018 Society of Systematic Biology Standalone meeting, whereas others were invited to cover topics identified by workshop attendees from questionaires they completed to assure that the book reflects the experiences, interests, and concerns of the diverse community that is engaged in species tree inference.

Through the set of chapters (authors representing their own perspective on aspects of species tree estimation relevant to their individual research programs), a diversity of perspectives and backgrounds are represented. This diversity means that the book speaks to people with varying levels of familiarity with the topic of species tree estimation, but it does not (nor is it intended to) provide a comprehensive overview of the subject. The combination of theoretical and empirical work is meant to provide readers with a level of knowledge of both the advances and limitations of species tree inference that can guide researchers in applying the methods while also inspiring future advances among those researchers with an interest in methodological development. Such cross talk (between empiricists and theoreticians/mathematicians) is vital to the growth of phylogenomics as it refocuses attention on the biological history of diversification (i.e., the timing and pattern of species divergence), the processes generating the observed patterns of genetic variation (e.g., sorting of ancestral polymorphism and gene flow, in addition to mutation models of nucleotide evolution), and the vast opportunities of study that includes and goes beyond a focus on the species tree itself.

# Index