

## Estimating Species Phylogeny from Gene-Tree Probabilities Despite Incomplete Lineage Sorting: An Example from *Melanoplus* Grasshoppers

BRYAN C. CARSTENS AND L. LACEY KNOWLES

Department of Ecology and Evolutionary Biology, 1109 Geddes Avenue, Museum of Zoology, University of Michigan, Ann Arbor, MI 48109-1079, USA; E-mail: knowlesl@umich.edu (L.L.K.)

**Abstract.**— Estimating phylogenetic relationships among closely related species can be extremely difficult when there is incongruence among gene trees and between the gene trees and the species tree. Here we show that incorporating a model of the stochastic loss of gene lineages by genetic drift into the phylogenetic estimation procedure can provide a robust estimate of species relationships, despite widespread incomplete sorting of ancestral polymorphism. This approach is applied to a group of montane *Melanoplus* grasshoppers for which genealogical discordance among loci and incomplete lineage sorting obscures any obvious phylogenetic relationships among species. Unlike traditional treatments where gene trees estimated using standard phylogenetic methods are implicitly equated with the species tree, with the coalescent-based approach the species tree is modeled probabilistically from the estimated gene trees. The estimated species phylogeny (the ESP) is calculated for the grasshoppers from multiple gene trees reconstructed for nuclear loci and a mitochondrial gene. This empirical application is coupled with a simulation study to explore the performance of the coalescent-based approach. Specifically, we test the accuracy of the ESP given the data based on analyses of simulated data matching the multilocus data collected in *Melanoplus* (i.e., data were simulated for each locus with the same number of base pairs and locus-specific mutational models). The results of the study show that ESPs can be computed using the coalescent-based approach long before reciprocal monophyly has been achieved, and that these statistical estimates are accurate. This contrasts with analyses of the empirical data collected in *Melanoplus* and simulated data based on concatenation of multiple loci, for which the incomplete lineage sorting of recently diverged species posed significant problems. The strengths and potential challenges associated with incorporating an explicit model of gene-lineage coalescence into the phylogenetic procedure to obtain an ESP, as illustrated by application to *Melanoplus*, versus concatenation and consensus approaches are discussed. This study represents a fundamental shift in how species relationships are estimated—the relationship between the gene trees and the species phylogeny is modeled probabilistically rather than equating gene trees with a species tree. [Coalescent; gene trees; incomplete lineage sorting; species phylogeny.]

Systematic investigations of recently derived species are often complicated by incongruence between the phylogenetic tree estimated from a particular locus (i.e., a gene tree or genealogy) and the actual species phylogeny (e.g., Avise et al., 1990; Hughes and Vogler, 2004; Morando et al., 2004; Crews and Hedin, 2006; Buckley et al., 2006; Ford et al., 2006). Various processes may lead to incongruence between the species tree and a gene tree (Maddison, 1997; Slowinski and Page, 1999), including gene duplication (Fitch, 1970) and horizontal transfer (Cummings, 1994), irrespective of the timing of species divergence. However, at the initial stages of species divergence, incomplete sorting of ancestral polymorphism is a ubiquitous source of discord (Tajima, 1983; Takahata and Nei, 1985; Hudson, 1992). Consequently, for recently formed species, the lack of reciprocal monophyly at any particular locus and incongruous genealogies among loci (Hudson and Coyne, 2002; Hudson and Turelli, 2003) present researchers with a dilemma of how to infer the species phylogeny.

Different approaches have been applied to infer species phylogenies when the species relationships are not strictly apparent from the structure of the gene tree. Species relationships have been inferred from concatenated data when there is genealogical discordance among loci (e.g., Lundrigan et al., 2002; Rokas et al., 2003; Steppan et al., 2005). An implicit assumption in this approach is that, given enough data, the dominant signal of the species phylogeny will emerge and the discord among the different gene trees will, at worst, lower nodal support values. Others have adopted the use of consen-

sus tree approaches (e.g., Jennings and Edwards, 2005), where the topology most frequently observed among the reconstructed gene trees for each locus is chosen as the species tree. Neither approach solves the problem of inferring a species phylogeny when the species are not reciprocally monophyletic—a relatively common condition that is not apparent unless multiple individuals are sampled from each species (e.g., Funk and Omland, 2003). However, it should not be assumed that a species phylogeny cannot be estimated before ancestral polymorphism has fully sorted. A simulation study showed not only that the historical signal of a species phylogeny persists despite the lack of reciprocal monophyly at any particular locus and substantial discord among loci but also that the species tree could be accurately estimated by incorporating the process of gene lineage coalescence into the phylogenetic procedure (even without a full probabilistic framework; Maddison and Knowles, 2006).

What is currently lacking is a demonstration that the empirical trees from these approaches accurately reflect the species tree, as with the concatenation or consensus approaches (e.g., Lundrigan et al., 2002; Rokas et al., 2003; Jennings and Edwards, 2005), or in the case of the coalescent-based approach (Maddison and Knowles, 2006), that the stochastic variance of genetic processes will not overwhelm the phylogenetic signal in an empirical data set. The study presented here couples (a) an approach for generating an estimate of the species phylogeny (an ESP) when there is both incongruence among gene trees estimated from each locus and a lack of reciprocal monophyly at any single locus (arguably

one of the most difficult scenarios for inferring species relationships), with (b) a simulation study to gauge the accuracy of the ESP generated when a coalescent-based approach is incorporated into the phylogenetic procedure.

Because the process of lineage sorting can be incorporated into the phylogeny estimation procedure (Degnan and Salter, 2005; Maddison and Knowles, 2006), an estimate of the species tree is possible even when there is genealogical discord among loci and between the gene tree and species tree. Maddison (1997) suggested such a likelihood framework, and specifically:

$$L(ST) = \prod_{\text{loci}} \sum_{\substack{\text{possible} \\ \text{gene trees}}} [P(\text{sequences}|\text{gene tree}) \\ \times (\text{gene tree}|\text{species tree})].$$

The expression  $[P(\text{sequences}|\text{gene tree})]$  is the familiar statement on the likelihood of the gene sequences given a model of evolution that can be readily calculated using a variety of phylogenetic software packages (e.g., PAML, PAUP\*, PHYLIP). In the second expression  $[P(\text{gene tree}|\text{species tree})]$ , coalescent theory can provide the probability that gene copies would coalesce to yield a particular gene tree (Pamilo and Nei, 1988; Takahata, 1989), given a specific species tree with defined branch lengths and widths (i.e., the number of generations occurring along a branch versus the effective population size of the branch). Assuming that the pattern of incomplete lineage sorting, and hence the discord between a gene tree and species tree, reflects the retention of ancestral polymorphism, the probabilities of different gene trees depends on the pattern of gene lineage coalescence (e.g., the probability of a gene tree decreases depending on its deviation from the theoretical expectation). However, searching for the maximum-likelihood species tree would involve a search not only over species trees (which includes the two parameters length and width specified independently for each branch) but also the entire gene-tree space (both topology and branch lengths) for each locus for every species tree (Maddison, 1997). Evaluating this sum for each locus (i.e., the integral over gene-tree space) is not computationally feasible. Although the search might be possible (see Discussion) using approximate likelihoods (see Felsenstein, 1992) or a Bayesian approach (Liu and Pearl, 2006), here we use a straightforward procedure for calculating the maximum-likelihood species tree—the species tree that confers the highest probability on the observed gene trees (Maddison, 1997). In this approach, the species tree is modeled probabilistically from a set of estimated gene trees. Because the likelihood of particular species relationships (i.e., the topology of a species tree) differs under varying branch lengths (relative to the effective population size along the length of the branch), the probability of obtaining the estimated gene trees is calculated from coalescent theory for a range of branch length parameters of the species trees.

Computing the probabilities of gene trees individu-

ally considers the phylogenetic signal from different loci independently, thereby permitting the evaluation of the likelihood of an ESP when stochastic variation produces discord between the topologies of the gene tree and the species tree. Here we apply this approach to generate an ESP for *Melanoplus* grasshoppers where traditional inferences about species relationships based on the estimates of genealogies (as reviewed in Felsenstein, 2004) are complicated by incomplete lineage sorting and incongruence among gene tree topologies. The recency of speciation, coupled with an apparent radiation coincident with the Pleistocene glacial cycles in the western montane *Melanoplus* species (Knowles, 2001a, b), creates a setting where topological congruence between gene trees and species trees is not likely (Knowles and Carstens, 2007). With a glacial-interglacial periodicity on the order of 100,000 years during the Pleistocene (Gibbard and Van Kolfschoten, 2004), successive speciation events associated with different glacial cycles will occur well before ancestral polymorphism from previous speciation events has fully sorted (Takahata and Nei, 1985; Pamilo and Nei, 1988; Takahata, 1989; Hudson and Coyne, 2002). Phylogenetic procedures that incorporate a model of gene lineage coalescence are expected to be ideal for estimating species phylogeny under these conditions. The ESP obtained from the coalescent-based approach used here is compared to the tree estimated from a concatenation approach, where the data is combined and treated as one locus with a single history of descent. The accuracy of these approaches is also explored with simulations. The strengths and potential problems of incorporating an explicit model of gene-lineage coalescence into the phylogenetic inference procedure (as applied here, and more generally) versus concatenation and consensus approaches are discussed.

## MATERIALS AND METHODS

### *Species*

The focus of this study is on four species that diversified across the sky islands of the northern Rocky Mountains during the Pleistocene: *M. montanus*, *M. oregonensis*, *M. marshalli*, and *M. triangularis*. These species are morphologically distinguishable (Knowles, 2000), and exhibit significant differences in the shape of the male genitalia (Marquez and Knowles, in preparation). The species are members of a radiation of *Melanoplus* grasshoppers that took place during the dynamic Pleistocene, which has made it difficult to attain reliable estimates of species relationships (Knowles and Otte, 2000; Knowles, 2000). Computational constraints limit the number of taxa that can be analyzed with the coalescent-based approach described below; with four species all possible species-tree topologies could be examined. Because the four species studied were selected as members of a larger group of taxa, it is possible that these taxa may not represent a monophyletic clade because of the limited taxon sampling. Nevertheless, the results are important for evaluating the feasibility of estimating *Melanoplus* species relationships—which has defied

traditional phylogenetic approaches because of the problems of separating historical signal from the stochastic sorting of gene lineages in these recently diverged taxa. Future analyses will be aimed at extending the taxonomic coverage of this species-rich genus.

These attributes make *Melanoplus* an ideal group for estimating species relationships using the coalescent-based approach presented here. Previous phylogeographic studies (Knowles, 2001a; Knowles and Carstens, 2007; Knowles et al., 2007) have shown that estimates of effective population sizes are large in comparison to the timing of species divergence (Carstens and Knowles, 2007; Knowles et al., 2007)—conditions for which topological equivalency between the gene trees and the species tree are unlikely (Takahata, 1989; Maddison, 1997; Rosenberg, 2002). Estimates of intraspecific migration (Knowles and Carstens, 2007; Knowles et al., 2007), as well as interspecific migration (Knowles et al., 2007; Carstens and Knowles, 2007), also indicate that gene flow is not predominate in this system—which is consistent with geographic isolation of the grasshopper species among the sky islands. Although this suggests the patterns of incomplete lineage sorting observed reflects the sorting of ancestral polymorphism, the possibility of gene flow cannot be eliminated, especially given the distributional shifts associated with the frequent glacial cycles over the Pleistocene. This emphasizes the importance of sampling multiple individuals per species when estimating relationships based on the pattern of gene lineage coalescence (Maddison and Knowles, 2006; Knowles and Carstens, 2007), as discussed below.

#### Sampling and Data Characteristics

Five alleles per species were used in order to have a balanced sampling design, because unequal sampling affects the probability of observing reciprocal monophyly (Rosenberg, 2003) and because the computational demands of calculating gene tree probabilities limits the number of individuals. For each individual, sequences from the mitochondrial gene cytochrome oxidase I (COI) and five anonymous single-copy nuclear polymorphic sequences (SCNPSs) were analyzed, for a total of 5.5 kb and 377 variable sites (for accessions, see GenBank EF217516 to EF218053). Alleles were determined either

by PCR subcloning or with PHASE version 2.0 (Stephens and Donnelly, 2003). Variable sites occur primarily at silent 3rd codon positions in the mitochondrial COI locus and the SCNPS loci are noncoding; furthermore, there is no evidence of directional selection or recombination based on estimates of Tajima's *D* and the four-gamete test, respectively (see also Knowles and Carstens, 2007; Carstens and Knowles, 2007; Knowles et al., 2007).

Sequences from *M. triangularis* were collected using the primers and protocol described in Carstens and Knowles (2006). Alleles exhibiting the largest average intraspecific genetic divergence across loci were selected from *Melanoplus oregonensis*, *M. montanus*, and *M. marshalli*. Because this sampling would introduce a bias for population genetic estimates that rely on the frequency distribution of alleles at polymorphic sites (e.g., Wakeley et al., 2001), all estimates of population genetic parameters referred to in the study are based on analyses of the complete sampling of individuals (Knowles and Carstens, 2007; Carstens and Knowles, 2007; Knowles et al., 2007), not the subset of individuals used here. Considering multiple individuals (as opposed to relying on the most recent interspecific coalescence; Takahata, 1989) minimizes the potential for low levels of gene flow to influence the estimate of the species phylogeny.

#### Phylogenetic Estimation of Gene Trees

Gene trees were estimated using maximum likelihood (ML) in PAUP\* 4.0 (Swofford, 2002). DT-ModSel (Minin et al., 2003) was used to select a model of sequence evolution (Table 1), and PAUP\* was used to conduct heuristic searches under the chosen model using TBR branch swapping, 10 random-addition replicates, and a random starting tree. Nodal support was assessed using a non-parametric bootstrap analysis with 1000 replicates (Felsenstein, 1985). Gene trees were rooted at the midpoint.

#### Estimating a Species Phylogeny Using the Coalescent-Based Approach

The computation and evaluation of the likelihood of different species tree with the coalescent-based approach involved: (1) estimating a gene tree from the nucleotide

TABLE 1. The length (BP) and number of variable sites (VS), estimates of the model of sequence evolution, and lnL score of the gene tree for each locus, as well as the concatenated data.

Data	BP	VS	Estimated model of sequence evolution		lnL
COI	1147	110	HKY+I	$\pi_A = 0.33$ ; $\pi_C = 0.175$ ; $\pi_G = 0.14$ ; $\pi_T = 0.355$ ; Ti/Tv = 4.0029; PINV = 0.851	-2413.89579
SCNPS-2	956	56	HKY	$\pi_A = 0.342$ ; $\pi_C = 0.166$ ; $\pi_G = 0.184$ ; $\pi_T = 0.308$ ; Ti/Tv = 1.0619	-1588.54188
SCNPS-73	853	11	F81+I	$\pi_A = 0.274$ ; $\pi_C = 0.18$ ; $\pi_G = 0.267$ ; $\pi_T = 0.279$ ; PINV = 0.9839	-1282.80789
SCNPS-85	826	41	JC+I	PINV = 0.9289	-1488.48994
SCNPS-89	581	44	HKY	$\pi_A = 0.289$ ; $\pi_C = 0.164$ ; $\pi_G = 0.259$ ; $\pi_T = 0.288$ ; Ti/Tv = 1.969	-1115.51451
SCNPS-211	1188	115	HKY+G	$\pi_A = 0.282$ ; $\pi_C = 0.202$ ; $\pi_G = 0.246$ ; $\pi_T = 0.27$ ; Ti/Tv = 1.136; a = 0.07133	-2397.05639
Concatenated loci	5551	377	HKY+I+G	$\pi_A = 0.305$ ; $\pi_C = 0.19$ ; $\pi_G = 0.213$ ; $\pi_T = 0.292$ ; Ti/Tv = 1.549; PINV = 0.8497; a = 1.2387	-11375.0932

data using ML in PAUP\* 4.0 (Swofford, 2002); (2) computing the probability of the gene tree for each specified species tree using the program COAL (Degnan and Salter, 2005); (3) calculating the likelihood of a species tree from the products of the probabilities of the gene trees given the species tree; and (4) using a likelihood-ratio test (with 1 degree of freedom), with a correction for multiple comparisons (e.g., Anisimova and Gascuel, 2006) to assess whether the species tree with the highest likelihood (highest  $\ln L$  score) was significantly better than the other species trees. This approach will be referred to as the coalescent-based approach, rather than repeating these details for each of the different contexts in which the approach was applied as part of this study.

The likelihood of each of the possible 15 rooted ultrametric species phylogenies (Fig. 1) and an unresolved (star) phylogeny was computed using the coalescent-based approach at total species tree depths ranging between  $1N$  to  $10N$  (where time is measured in coalescent units, or  $t/N_e$ , where  $t$  is the time of species divergence in generations and  $N_e$  is the effective population size of the species; e.g., a total tree depth of  $1N$  corresponds to a  $t$  and  $N_e$  of equal values). The probability of the gene tree topology given the species tree was calculated with the program COAL (Degnan and Salter, 2005). This program uses gene tree topology, but not branch lengths, to evaluate the probability of the gene tree. It is noteworthy that although branch lengths of the gene trees were not used, the simulation study nevertheless confirms the accuracy of the ESPs using this approach. This might reflect the large effect of the pattern of gene lineage coalescence on the probability of gene trees for the very shallow histories considered here; the importance of considering gene-tree branch lengths might become more apparent with older species divergence times that encompass a broader range genetic distances among alleles (see Rannala and Yang, 2003).

The probabilities of the gene trees were also evaluated with randomized data sets to confirm that there was phylogenetic signal in the data. Five individuals were drawn at random for each species from the randomized nucleotide data, based on reshuffling individuals across species in MESQUITE (Maddison and Maddison, 2004). For each locus, the gene trees were estimated from 100 reshuffled data sets, and the probability of these genealogies was calculated for the single species tree with the highest likelihood from the coalescent-based analysis above.

#### Evaluating ESP Accuracy with Simulations

A simulation study was used to evaluate the ability of the coalescent-based approach to recover an accurate ESP given the data, and specifically given the level of genetic divergence, the number of loci, and amount of sequence data collected in the grasshopper species. Using MESQUITE, 100 replicate nucleotide data sets were simulated under a model that matched the model selected using DT-ModSel from the empirical data, both in terms of the model of sequence evolution and the total number

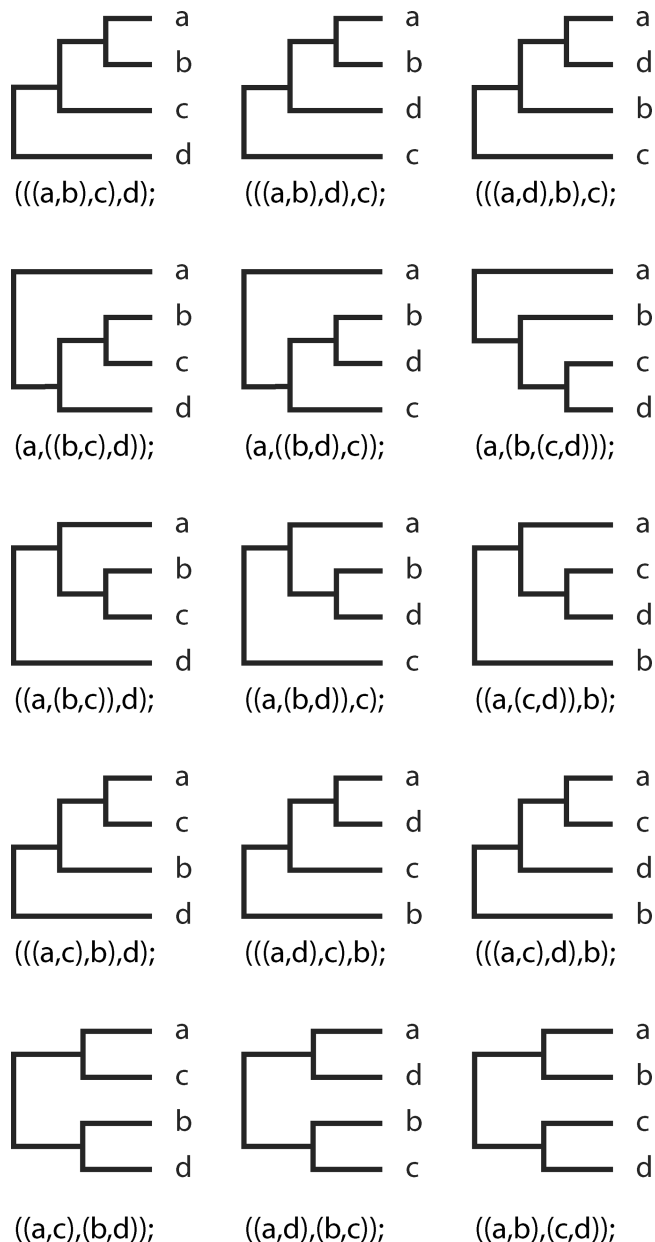


FIGURE 1. The 15 rooted topologies for four species ( $a = M. montanus$ ;  $b = M. oregonensis$ ;  $c = M. triangularis$ ; and  $d = M. marshalli$ ); the likelihood of each topology was evaluated based on the probabilities of the gene trees for each the possible ultrametric species trees over a range of tree depths (see Materials and Methods).

of variable sites for each of six loci, on the gene genealogies simulated by a neutral coalescence process. The gene genealogies were simulated with an  $N_e$  of 140,000 (the average  $N_e$  previously estimated for the separate species; Carstens and Knowles, 2007; Knowles et al., 2007) under the species tree with the highest likelihood based on the coalescent-based analysis of the empirical data. Gene trees for the simulated nucleotide data sets and the likelihood of the known species tree (i.e., the species tree under which the data were simulated—the parameters

for which matched the empirical data, specifically with regards to  $N_e$  and timing of species divergence) were estimated using the coalescent-based approach. Accuracy of the coalescent-based approach was evaluated by recording the proportion of data sets in which the true species tree had the highest  $\ln L$  score compared to the alternative species trees.

The ability to accurately estimate a species phylogeny was also evaluated across a range of total species-tree depths with different relative branch lengths across the tree (in contrast to the ultrametric species trees described above; see Fig. 1). Data were simulated using an unrooted species tree of ((a, b), (c, d)) at five total tree depths (1N, 2N, 3N, 4N, and 8N) with an internal and tip branch ratio of 1:1, 2:1, and 1:2, which results in differing amounts of topological discordance between the species tree and gene trees. For each of these 15 species trees, replicate nucleotide data sets were simulated on 100 genealogies simulated by a neutral coalescence process using models of sequence evolution that matched those selected using DT-ModSel from each locus of the empirical data. Gene trees for the simulated nucleotide data sets were estimated and their probabilities calculated for each of the possible unrooted species trees using the coalescent-based approach. The accuracy of the species tree estimate was evaluated by recording the proportion of data sets in which the gene trees probabilities were highest for the species tree under which the data were simulated compared to the alternative species trees.

#### Comparing Different Approaches for Inferring a Species Tree

A ML tree was estimated from the concatenated empirical data, and bootstrap values were calculated using 1000 replicates. A parametric bootstrap analysis (Huelsenbeck et al., 1996) was used to evaluate whether the tree with the highest  $\ln L$  score estimated from the concatenated data differed significantly from the ESP with the highest  $\ln L$  score from the coalescent-based approach. Using MESQUITE (Maddison and Maddison, 2004), 1000 nucleotide data sets were simulated under the ESP from the coalescent-based method, with model parameters and branch lengths estimated using the ML model selected from the concatenated data set using DT-ModSel (Minin et al., 2003). The likelihood of the reconstructed gene tree for each simulated data set (from a ML search with PAUP\*; Swofford, 2002) was calculated for an unconstrained and constrained search, where the ESP from the coalescent-based analysis served as the constraint tree. A null distribution from the difference in log-likelihood scores ( $\ln L_{\text{unconstrained}} - \ln L_{\text{constrained}}$ ) for the replicate data sets was used to assess the difference in likelihood scores between the trees estimated from the concatenated empirical data versus the coalescent-based approach.

To test whether concatenation provides an accurate inference about the species tree for the shallow depths of species divergence considered in this study, the simulated data sets used to evaluate the accuracy of the coalescent-based approach were also analyzed using the

concatenation approach. The accuracy of the concatenation approach was evaluated by recording the proportion of data sets in which the ML tree estimated from the concatenated data matched the species tree under which the data were simulated.

#### RESULTS

A single ML tree was identified for all loci (Fig. 2), except SCNPS-85 and SCNPS-89, where 10 and 2 trees were found, respectively. Species were not reciprocally monophyletic at any locus, and the degree of species monophyly varied among species and among loci. For example, *M. montanus* and *M. triangularis* were monophyletic at three of six loci, but not the same subset of loci, *M. marshalli* was monophyletic at one locus, and *M. oregonensis* was not monophyletic at any locus. Consequently, the species phylogeny is not obvious based on consideration of the topologies of the gene trees alone.

#### The Likelihood of the Species Phylogeny Based on the Probabilities of Gene Trees

The most likely species tree given the gene trees ( $\ln L = -200.887$ ) had a topology of (*M. triangularis*, (*M. montanus*, (*M. marshalli*, *M. oregonensis*))) (Fig. 3). The best  $\ln L$  scores from each of the 16 topologies occurred at different total tree depths (Table 2). Based on the likelihood ratio test, the most likely topology (Fig. 3) was significantly better than all but two alternate topologies. In one alternate topology *M. montanus* was basal to the clade containing *M. triangularis*, *M. oregonensis*, and *M. marshalli*, and in the other *M. montanus* and *M. triangularis* were sister (Table 3).

The probability of the gene trees given the species tree differed considerably among loci, and on average the loci with more variable sites (Table 1) had higher probabilities

TABLE 2. The highest  $\ln L$  score for each of the possible rooted topologies, as well as a star phylogeny, is shown. Significant differences in the likelihood between the most likely tree and the suboptimal species trees (shown in bold) was assessed using likelihood-ratio tests, after correcting for multiple comparisons with a Bonferroni correction. Species names are abbreviated as follows: *Melanoplus montanus* (a); *M. oregonensis* (b); *M. triangularis* (c); and *M. marshalli* (d).

Topology	Depth	$\ln L$	LRT
<b>((a,(b,d)),c)</b>	4N	<b>-200.887</b>	<b>0</b>
<b>(a,(b,d),c)</b>	4N	<b>-202.181</b>	<b>2.588</b>
<b>((a,c),(b,d))</b>	3N	<b>-203.969</b>	<b>6.164</b>
Star	5N	-206.391	11.008
(a,(b,c),d))	3N	-207.667	13.56
((a,d),b),c)	3N	-207.863	13.952
((a,b),d),c)	3N	-207.908	14.042
((a,(b,c)),d)	6N	-208.122	14.47
((a,b),c),d)	7N	-208.817	15.86
((a,c),b),d)	6N	-209.477	17.18
(a,(b,c,d)))	9N	-210.417	19.06
((a,d),(b,c))	5N	-210.713	19.652
((a,b),(c,d))	8N	-211.349	20.924
((a,(c,d)),b)	4N	-212.4	23.026
((a,d),c),b)	3N	-213.265	24.756
((a,c),d),b)	6N	-214.213	26.652

TABLE 3. The probabilities of each gene tree for each of the alternative species trees that were statistically indistinguishable from the species tree with the highest likelihood that is in bold (i.e., probability (GT|ST)); the likelihood (see last column) of each of these species trees ( $\ln \text{LST|GT}$ ) was calculated as the product of the probabilities of the individual gene trees is also shown; a = *M. montanus*; b = *M. oregonensis*; c = *M. triangularis*; and d = *M. marshalli*.

Species tree	Depth	SCNPS-2	SCNPS-73	SCNPS-85	SCNPS-89	SCNPS-211	SCNPS-COI	$\ln \text{LST GT}$
<b>((a,(b,d)),c)</b>	<b>4N</b>	<b><math>7.70 \times 10^{-11}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>7.40 \times 10^{-14}</math></b>	<b>-200.887</b>
	5N	$1.03 \times 10^{-10}$	$<1.0 \times 10^{-15}$	$<1.0 \times 10^{-15}$	$<1.0 \times 10^{-15}$	$<1.0 \times 10^{-15}$	$4.60 \times 10^{-14}$	-201.072
<b>(a((b,d),c))</b>	<b>4N</b>	<b><math>1.26 \times 10^{-11}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>1.24 \times 10^{-13}</math></b>	<b>-202.181</b>
	5N	$6.09 \times 10^{-11}$	$<1.0 \times 10^{-15}$	$<1.0 \times 10^{-15}$	$<1.0 \times 10^{-15}$	$<1.0 \times 10^{-15}$	$2.43 \times 10^{-13}$	-202.236
<b>((a,c),(b,d))</b>	<b>3N</b>	<b><math>2.81 \times 10^{-12}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>&lt;1.0 \times 10^{-15}</math></b>	<b><math>9.30 \times 10^{-14}</math></b>	<b>-203.969</b>

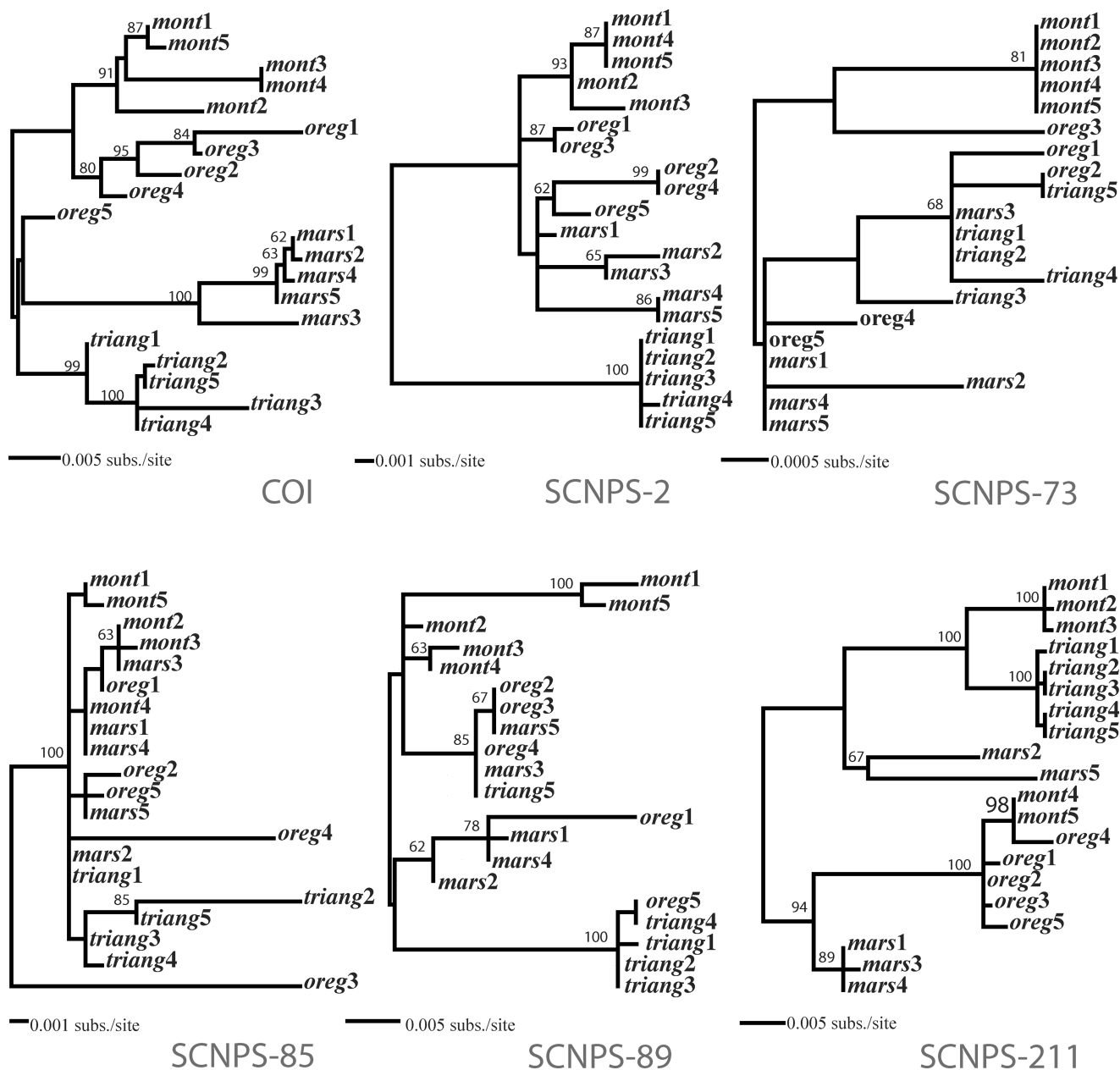


FIGURE 2. Maximum likelihood estimates of the gene trees (rooted at the midpoint) for each of the six loci; bootstrap values  $>50$  are shown above the nodes. Species names are abbreviated: *M. montanus* = mont; *M. oregonensis* = oreg; *M. triangularis* = triang; and *M. marshalli* = mars.

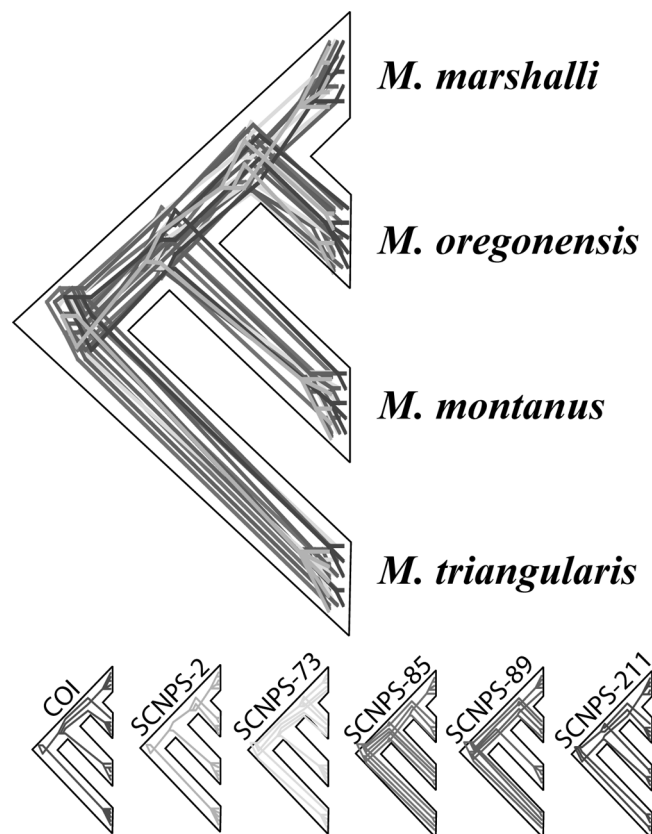


FIGURE 3. Estimate of the species phylogeny (ESP) with the highest likelihood score ( $-\ln L = -10,487.1938$ ) calculated with the coalescent-based approach from the probabilities of the gene trees (see Fig. 2), with each gene tree shown within the species tree.

(Table 3). For example, the average probability of the COI gene tree given the species trees was  $2.59 \times 10^{-12}$  compared to less than  $1.0 \times 10^{-17}$  for SCNPS-73. However, the probabilities of the gene trees from the randomized data sets were never larger than those for the empirical gene trees, indicating that there is a significant signal of phylogenetic relationships—otherwise, the probabilities of the gene trees from the randomized and empirical data would have been of similar magnitude.

#### Evaluating the Accuracy of the ESP with Simulations

Comparison of the likelihood of the species phylogeny for the simulated data sets under the species tree used in the simulations to the alternative species trees demonstrates that the approach is able to accurately recover an ESP under conditions corresponding to the empirical data. In 94% of the replicate data sets, the  $\ln L$  score of the true species tree was greater than the  $\ln L$  scores of the alternate species trees.

A similar pattern was observed across the range of total tree depths considered, where increasing levels of topological incongruence between the species tree and the gene trees are observed for the shallower tree depths. As the total depth of the species phylogeny increases (i.e., with the older timing of species divergence), the accu-

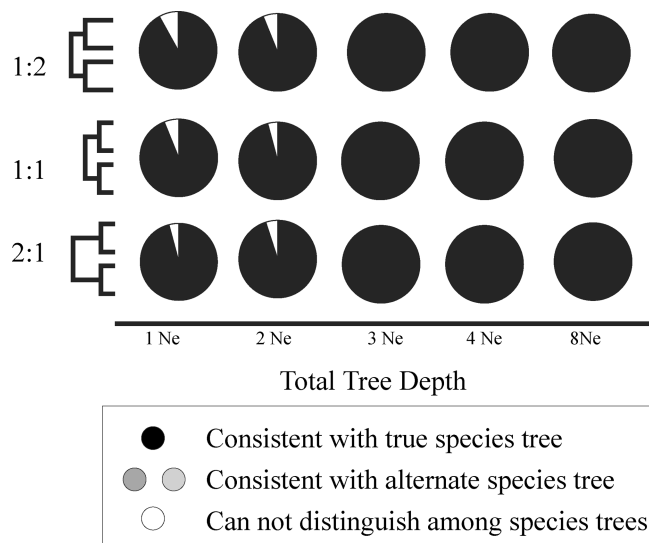


FIGURE 4. Accuracy of the species trees estimated from the simulated data using the coalescent-based approach. Accuracy of the species trees at each tree depth is represented by a pie chart, showing the proportion of simulated data sets in which (a) the true species tree was correctly identified as most likely (in black); (b) an alternate species tree was identified as most likely (in grey); and (c) it was not possible to distinguish among the alternative species trees (in white). Note that, in contrast to concatenation-based approaches (see Fig. 6), there were no cases where an alternate species tree was identified as most likely (i.e., no grey regions).

racy of the phylogenetic estimate improved, regardless of the relative branch lengths of the underlying species tree (Fig. 4). When the total tree depth reaches  $3N$ , the species tree with the highest probability matches the simulation tree under all tree shapes and for virtually every replicate. At the lowest divergence levels (e.g.,  $1N$  and  $2N$ ), estimates of the species phylogeny are still remarkably accurate using the coalescent-based approach, even when the relative length of internal to tip branch differs (Fig. 4). With just five individuals per species and five loci, the true species tree was recovered in more than 90% of the simulated data sets under the most difficult condition of a  $1N$  total tree depth. In the few replicates where the true species tree was not identified, the method never identified an incorrect topology as the most likely species tree but instead indicated that there was insufficient data to distinguish among the alternative topologies (i.e., the  $\ln L$  scores of the alternate species trees did not differ significantly).

#### Comparing Different Approaches to Inferring a Species Phylogeny

Analysis of the concatenated data produced a tree in which *M. marshalli* and *M. triangularis* were monophyletic, but *M. oregonensis* and *M. montanus* were not (Fig. 5). Moreover, *M. montanus* is split into two well-supported but unrelated clades, to a large degree mirroring the gene tree estimated from SCNPS-211 (Fig. 2). The SCNPS-211 locus contains about 25% of the total variable sites across the multilocus data set (Table 1), which

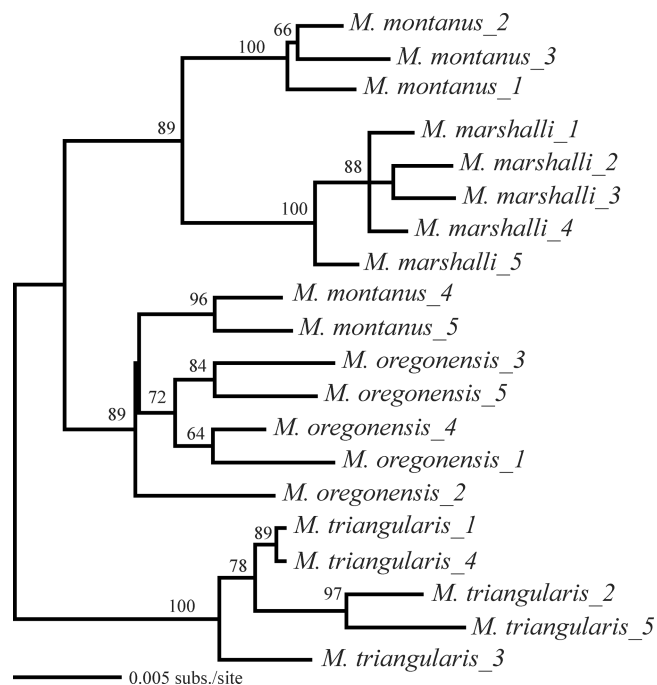


FIGURE 5. Maximum likelihood estimate of the tree from the concatenated data; bootstrap values >50 are shown above the nodes.

may underlie the apparent impact of this locus on the tree estimated from the concatenated data.

The basic topological structure in the tree reconstructed with the concatenation approach might be interpreted as (*M. triangularis*, (*M. oregonensis*, (*M. montanus*, *M. marshalli*))) if the two *M. montanus* samples that form a clade with *M. oregonensis* are ignored (although this non-quantitative decision highlights the problems with inferring species relationships from concatenated data sets when the species are not monophyletic). This tree (Fig. 5) does not match the topology of the species tree with the highest likelihood (Fig. 3) modeled probabilistically from the estimated gene trees or the other equally likely species trees (Table 3). The likelihood score of the concatenated empirical data constrained to fit the ESP from the coalescent-based approach (Fig. 3) decreased from  $-11,375.093$  to  $-11,454.085$ , or  $78.992 \ln L$  units. This decrease is significant ( $P < 0.001$ ), based on the parametric bootstrap, indicating the topology of the tree estimated from concatenating the data differs significantly from the ESPs computed from the gene tree probabilities.

Analysis of the simulated data sets also showed the problem of obtaining a quantitative estimate of the shallow species trees because of a lack of species monophyly, which predominated at total tree depths of  $1N_e$  and  $2N_e$  (Fig. 6). The concatenation approach performed well at the deepest tree depth of  $8N_e$  (which would correspond to a divergence of 800,000 years ago assuming one generation per year and an  $N_e$  of 100,000) and moderately well at the  $3N_e$  and  $4N_e$  depths. However, there were a few very disconcerting (albeit rare) cases in which the species were reciprocally monophyletic, and the tree inferred from the concatenated data actually differed from

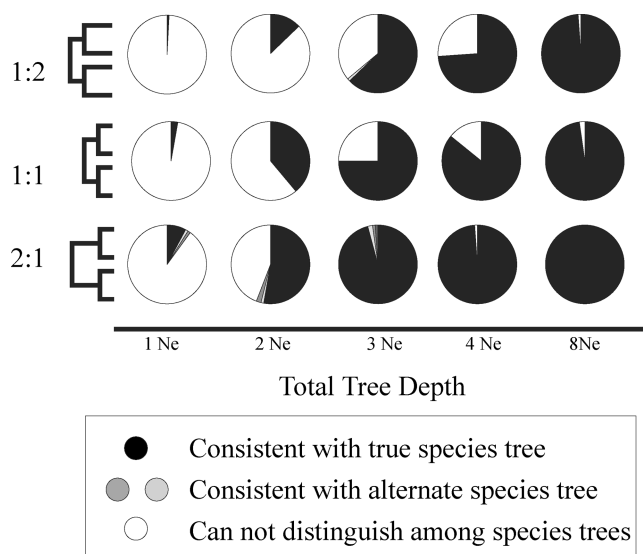


FIGURE 6. Accuracy of the species tree estimated from the simulated data using the concatenation approach. As in Figure 4, accuracy of the species trees at each tree depth is represented by a pie chart, showing the proportion of simulated data sets in which (a) the true species tree was correctly identified (in black), and (b) an alternate species tree was identified (in grey). But the case (c) reflects the proportion of simulated data sets in which a quantitative estimate about the species tree was not possible because one or more species were not monophyletic (in white).

the phylogeny used to simulate the data (see the grey shaded areas of Fig. 6).

## DISCUSSION

The lack of species monophyly and incongruence among the gene trees reconstructed from the multilocus data from *Melanoplus* highlights the difficulties with inferring species relationships for recently derived taxa. For recently diverged species, under traditional phylogenetic treatments where gene trees and species trees are synonymized, the historical signal used to infer species relationships is overwhelmed by the stochasticity of genetic processes because the most recent common ancestor of multiple individuals sampled within a species is not likely to occur within the species (Takahata, 1989; Hudson and Coyne, 2002; Hudson and Turelli, 2003). In contrast, the simulation study demonstrates that an accurate ESP is possible with the consideration of the genetic processes underlying the observed genealogical discord (Maddison, 1997; Degnan and Salter, 2005; Maddison and Knowles, 2006). Despite the messy gene trees (Fig. 2), an estimate of the species phylogeny was obtained for the four *Melanoplus* species (Fig. 3). Given the data, the simulation study also confirms that an accurate species tree is expected using the coalescent-based approach, with the correct species tree estimated in 94% of simulated data sets. Additional simulations verified the robustness of the coalescent-based approach, demonstrating that the species trees were correctly estimated over a range of different species divergence times (Fig. 4). Below we discuss some of the strengths, shortcomings, and extensions of

methods that consider the stochastic sorting of gene lineages during the phylogenetic estimation procedure, as illustrated by application to *Melanoplus*, and comparison with other phylogenetic approaches.

#### *Extracting Phylogenetic Signal When There Is Incomplete Lineage Sorting*

Within a traditional phylogenetic framework, because species relationships are inferred directly from the estimated gene tree (i.e., the topologies of gene tree and species tree are assumed to be the same), it is not clear how to infer the species tree when there is incomplete sorting and incongruence among gene trees (Slowinski and Page, 1999; Parker and Kornfield, 2000; Stepan et al., 2004; Buckley et al., 2006). Various subsampling approaches might be used so that there is a direct correspondence between a species tree and a gene tree in the sense that individuals (or more precisely, gene copies) are interchangeable with species (i.e., the OTUs of the species and gene tree are equivalent). For example, the most recent interspecific coalescence might be targeted as the basis for inferring species relationships (e.g., Takahata, 1989). This approach might perform well when the discordance between the species and gene trees results from incomplete lineage sorting (see Maddison and Knowles, 2006). However, relying on a single individual (or exemplar) per species would lead to inaccurate estimates of species phylogeny if genetic similarities were due to gene flow rather than common ancestry (see Knowles and Carstens, 2007). This problem would also apply to consensus approaches that rely on the most common gene tree across loci estimated from single exemplars from each species (Jennings and Edwards, 2005). In this case, gene flow is expected to influence the entire genome (except for those loci for which selection prevents their movement across species boundaries; Wu, 2001; Rieseberg 2001), thereby producing a positively misleading impression that the most common gene tree topology reflects its higher probability under a specific evolutionary branching pattern when it is caused by the demographics of gene flow. Moreover, relying on single exemplars per species disregards information relevant to inferring species relationships—namely, the pattern of gene lineage coalescence across multiple individuals sampled from each species (Maddison and Knowles, 2006).

The coalescent-based approach has two important and desirable qualities. First, because the likelihood of a species tree is computed from the probabilities of the individual gene trees, there is no assumption of topological equivalency among gene trees or between the species tree and gene trees. For recently derived species (see Fig. 2), such an assumption would clearly be violated (Takahata, 1989; Hudson, 1992; Rosenberg, 2002). Because the distribution of possible gene trees (Degnan and Salter, 2005) and their respective probabilities are derived from coalescent theory (Kingman, 1982), the species tree can be modeled probabilistically from the observed gene trees long before there is reciprocal monophyly at the sampled

loci. For example, accurate ESPs were generated for total tree depths of just 1N generations (or divergence as recent as 0.3 between a pair of taxa) with just five individuals sampled at five loci in the simulated data (Fig. 4), whereas it would take about 10N generations for reciprocal monophyly to be observed between a pair of species if five individuals at a single locus were sampled (Hudson and Coyne, 2002). Second, the analyses produce a comparison of the likelihood of a species tree (Table 2)—there is no nonquantitative guess of what the species tree might be, as when incomplete lineage is not considered. For example, phylogenetic analysis of the concatenated data set for the *Melanoplus* species produced a single tree (Fig. 5). Yet, the species tree still is not clear. Ignoring the two *M. montanus* samples that form a clade with *M. oregonensis* might be interpreted as supporting the species tree (*M. triangularis*, (*M. oregonensis*, (*M. montanus*, *M. marshalli*))), whereas ignoring the three *M. montanus* samples that form a clade with *M. marshalli* might be interpreted as supporting the species tree (*M. triangularis*, (*M. marshalli*, (*M. oregonensis*, *M. montanus*)))—this also involves the decision to overlook the paraphyly of *M. oregonensis*. Neither of these trees matches the species trees estimated using the coalescent-based approach (Table 3) that takes into account the observed gene tree incongruence and incomplete lineage sorting (Fig. 2). For recently derived species, it is extremely unlikely that a consensus approach will resolve this problem since the most probable and common gene tree under such a species history will certainly not be one in which the species are reciprocally monophyletic when more than one individual is sequenced per species (Hudson and Coyne, 2002; Rosenberg, 2003).

This study provides an empirical confirmation of the suggestion by Maddison and Knowles (2006) that the signal of species relationships persists in the pattern of gene lineage coalescence. Computation of the gene tree probabilities that formed the basis for evaluating the likelihood of alternative *Melanoplus* species trees (Table 3) is based on the pattern of gene lineage coalescence throughout a species tree (Degnan and Salter, 2005). As such, this study also reaffirms the importance of sampling design for extracting historical information inherent in the pattern of coalescence for phylogenetic inference. Data from multiple loci are important for providing independent information for inferring a species tree across all species divergence times (Pamilo and Nei, 1988; Takahata, 1989; Maddison, 1997). However, if most gene lineages coalesce deeper than the species divergence (as in recently diverged taxa), by sequencing multiple individuals per species, each gene lineage that coalesces with a gene lineage from a sister species provides information regarding species relationships (i.e., the pattern of deep coalescence retains a signal of the history of species divergence; Takahata, 1989; Maddison and Knowles, 2006).

#### *Implications for Estimating Species Phylogeny*

Given the data, and based on the likelihood-ratio test, we were not able to statistically distinguish between

three possible species tree topologies (out of the 15 possible topologies). This might be viewed as a weakness of the approach. Nevertheless, being able to determine that these topologies are equally likely is arguably preferable to a method that produces a single tree, but one for which there is no way to evaluate its accuracy, such as with the concatenation and consensus methods. For example, when concordance among gene trees is not theoretically expected, or when topological equivalency of the species and gene trees is not a valid assumption, it is not generally known whether concatenation of multilocus data will produce an accurate estimate of the species tree. For example, in the few cases that species were monophyletic in the simulated data sets for the recent divergences (i.e., total tree depths of 1N and 2N; Fig. 6), surprisingly, the species tree inferred from the concatenated data actually did not always match the phylogeny used to simulate the data (Fig. 6; see also Kubatko and Degnan, 2007). The reason for these positively misleading species trees is not clear; in no instances was an incorrect species tree identified using the coalescent-based approach that models the species tree probabilistically from the independent gene trees (Fig. 4). For the data simulated under a model of recent species divergence, the empirical data suggest a possible explanation for the incorrect species tree resulting from concatenation across loci. This explanation relates to the distribution of phylogenetically informative sites across loci. For example, the split of *M. montanus* into two unrelated clades in the tree estimated from the concatenated data set (Fig. 5) mirrors to some degree the topology observed for the SCNPS-211 locus—which contains about 25% of the total variable sites across the six loci—when none of the other gene trees show *M. montanus* split into distant clades (Fig. 2). This would seem to indicate that the independent realizations of history that multiple loci are supposed to provide are not being treated as such under the forced assumption of topological equivalence among loci of methods that concatenate data. There currently (to our knowledge) has not been any systematic investigation of whether the distribution of sites across loci might bias the trees estimated with the concatenation of data, though such a study would no doubt be very useful in the context of its application to recently diverged taxa. Irrespective of the outcome of such investigations, the most frequent problem for inferring a phylogeny of recently diverged species using concatenated data was the ambiguity of species relationships posed by the lack of reciprocal monophyly in both the empirical (Fig. 5) and simulated data (Fig. 6). This contrasts with the simulation results for the older species divergence for which the concatenation approach performed well (Fig. 6).

#### *Considerations for Modeling Species Trees Probabilistically from Gene Trees*

The probabilities of the gene trees for the six loci in this study were calculated for each possible species tree topology (Fig. 1) over a range of branch lengths from 1N to 10N, for a total of 160 different species trees. Al-

though the subset of tree space explored does cover the range of parameter space indicated by the separate population genetic analyses in terms of the total species tree depth (Carstens and Knowles, 2007; Knowles et al., 2007; Knowles and Carstens, 2007), this is clearly a small subset of the possible species trees—the tree depths explored were in 1N increments across ultrametric species trees. Consequently, we were able to determine the maximum-likelihood species tree (Fig. 3) for the branch length/width parameters considered, but a large part of the tree space was not explored, namely varying the relative branch lengths. Although short internal branches might make it more difficult to obtain an ESP, the simulation study demonstrates that such estimates are possible (Fig. 4). The ESPs were accurately identified in more than 95% of the simulated data sets, even when the internal branch was only 0.3N generations, and twice as short as the species branches. Nevertheless, searching species-tree space for the maximum-likelihood species tree is tedious, and it becomes increasingly difficult as the number of taxa increases. With addition taxa the tree-space grows dramatically in terms of possible topologies (Felsenstein 1978; Hillis et al., 1994, 1996), but the set of possible combinations of branch lengths also increases at even a higher rate because each topology will have the full configuration of possible branch lengths.

Another assumption that would influence the evaluation of the likelihoods of the species trees from the gene trees is that the gene trees were estimated without error. Hence, the estimates of the species trees are only reliable to the extent that gene trees are correctly reconstructed—which also applies to the phylogenetic tradition of interpreting the gene tree as the species tree.

It is not clear what effect these assumptions will have on the accuracy of the ESP from the coalescent approach. Their influence is likely to differ among data sets (e.g., depending on the sequence length and number of variable sites) and with specific details of the species tree (e.g., depending on the shape of the underlying gene tree distribution; Degnan and Salter, 2005). Avoiding these simplifying assumptions might (or might not) improve the accuracy of ESPs. For example, some sort of approximation, such as importance sampling (Robert and Casella, 1999), might be used to search across tree-space. Perhaps with such approximations, the coalescent-based approach could also incorporate other biological realities, such as allowing the effective population size to vary by branch (Degnan and Salter, 2005). However, increased model complexity may not necessarily lead to a more accurate historical inference (Swofford et al., 1996; Sullivan and Swofford, 2001; Knowles and Maddison, 2002). Alternatively, rather than considering the species tree as a fixed parameter, a Bayesian framework might be used to estimate a species tree (Liu and Pearl, 2006), where the tree with the highest posterior probability could be interpreted as the best species tree. A prior distribution that includes information regarding both the species tree topology and branch lengths would need to be specified for this approach (Maddison, 1997; Degnan and Salter, 2005). However, establishing a prior for the species tree would

also require assumptions, and again, what effect these assumptions will have on the accuracy of the ESP is not known. Moreover, whether traditional sampling of most phylogenetic studies (i.e., single individuals sampled per species and relatively few loci) would contain sufficient information for estimating the necessary parameters is an important question that requires investigation. Just as the effects of the model of sequence evolution on the accuracy of gene trees has been thoroughly investigated (Swofford et al., 1996; Bruno and Halpern, 1999; Sullivan and Swofford, 2001; Anderson and Swofford, 2004), the modeling of the species tree is an exciting area that will also require critical study.

### CONCLUSIONS

Species relationships estimated from the coalescent-based approach are not only accurate, but the method also provides a direct statistical evaluation of the estimated species phylogeny (the ESP), as opposed to inferring it from the topology of a gene tree under the assumption of topological equivalency. The study also shows that estimated trees will not accurately reflect the species tree if topological congruence is forced—as with the concatenation approach—on recently derived species for which widespread discordance among gene trees and incomplete lineage sorting predominates. Because the species tree is modeled probabilistically from the individually estimated gene trees, the coalescent-based approach does not assume concordance among the gene trees of different loci and between the species tree and gene trees (Maddison, 1997). Moreover, by incorporating a model of gene coalescence into the phylogenetic procedure (Rannala and Yang, 2003; Degnan and Salter, 2005), species relationships can be estimated long before reciprocal monophyly of the species is observed (see also Maddison and Knowles, 2006). The coalescent-based approach used here bridges what has been viewed conceptually as the separate disciplines of systematics and population genetics. What emerges from this synergy is a shift in how species relationships are estimated. Rather than synonymizing gene trees with a species tree for inferring phylogenetic relationships, an estimate of the species phylogeny (an ESP) is chosen that maximizes the probability of the gene trees. This marks a very interesting and significant development in phylogenetics, and an area of largely unexplored potential.

### ACKNOWLEDGMENTS

We thank members of the Knowles lab for their input. The research was funded by the following awards to LLK: a National Science Foundation grant (DEB-04-47224), the Elizabeth Caroline Crosby Fund, National Science Foundation ADVANCE Project, University of Michigan, and a grant from the University of Michigan (Office of the Vice President for Research). We thank the Ohio State University—Mathematical Biosciences Institute for inviting us to a workshop where this and related approaches were discussed. We thank James Degnan for his assistance during the preparation of the manuscript; Jack Sullivan, Noah Rosenberg, and Scott Edwards for discussions related to this topic; and Tim Collins, Robb Brumfield, and one anonymous reviewer for providing suggestions that improved the manuscript.

### REFERENCES

- Anderson, F. E., and D. L. Swofford. 2004. Should we be worried about long-branch attraction in real data sets? Investigations using meta-zoan 18S rDNA. *Mol. Phylogenet. Evol.* 33:440–451.
- Anisimova, M., and O. Gascuel. 2006. Approximate likelihood-ratio test for branches: A fast, accurate, and powerful alternative. *Syst. Biol.* 55:539–552.
- Arvestad, L., A.-C. Berglund, B. Sennblad, and J. Lagergren. 2004. Gene tree reconstruction and orthology analysis based on an integrated model for duplications and sequence evolution. *RECOMB* 04:326–335.
- Avise, J. C., C. D. Ankney, and W. S. Nelson. 1990. Mitochondrial gene trees and evolutionary relationships of mallard and black ducks. *Evolution* 44:1109–1119.
- Bruno, W. J., and A. L. Halpern. 1999. Topological bias and inconsistency of maximum likelihood using wrong models. *Mol. Biol. Evol.* 16:564–566.
- Buckley, T., M. Cordeiro, D. Marshall, and C. Simon. 2006. Differentiating between hypotheses of lineage sorting and introgression in New Zealand alpine cicadas (*Maoricicada* Dugdale). *Syst. Biol.* 55:411–425.
- Carstens, B. C., and L. L. Knowles. 2007. Shifting distributions and speciation: genomic resolution of species divergence during rapid climate change. *Mol. Ecol.* 16:619–627.
- Carstens, B. C., and L. L. Knowles. 2006. Variable nuclear markers for *Melanoplus oregonensis* identified from the screening of a genomic library. *Mol. Ecol. Notes* 6:683–685.
- Crews, S. C., and M. Hedin. 2006. Studies of morphological and molecular phylogenetic divergence in spiders (Araneae: Homalonychus) from the American southwest, including divergence along the Baja California peninsula. *Mol. Phylogenet. Evol.* 38:470–487.
- Cummings, M. P. 1994. Transmission patterns of eukaryotic transposable elements: Arguments for and against horizontal transfer. *Trends Ecol. Evol.* 9:141–145.
- Degnan, J. H., and L. M. Salter. 2005. Gene tree distributions under the coalescent process. *Evolution* 9:24–37.
- Felsenstein, J. 1978. Cases in which parsimony or compatibility methods will be positively misleading. *Syst. Zool.* 27:401–410.
- Felsenstein, J. 1985. Confidence limits on phylogenies: An approach using the bootstrap. *Evolution* 39:783–791.
- Felsenstein, J. 1992. Estimating effective population size from samples of sequences: a bootstrap Monte Carlo integration approach. *Genet. Res.* 60:209–220.
- Felsenstein, J. 2004. *Inferring phylogenies*. Sinauer Associates, Sunderland, Massachusetts.
- Fitch, W. M. 1970. Distinguishing homologous from analogous proteins. *Syst. Zool.* 19:99–113.
- Ford, V. S., J. Lee, B. G. Baldwin, and L. D. Gottlieb. 2006. Species divergence and relationships in *Stephanomeria* (Compositae): PgiC phyleny compared to prior biosystematic studies. *Am. J. Bot.* 93:480–490.
- Funk, D. J., and K. E. Omland. 2003. Species-level paraphyly and polyphyly: Frequency, causes, and consequences, with insights from animal mitochondrial DNA. *Ann. Rev. Ecol. Syst.* 34:397–423.
- Gibbard, P., and T. Van Kolfschoten. 2004. The Pleistocene and Holocene Epochs. Pages 441–452 in *A geologic time scale* (F. M. Gradstein, J. G. Ogg, and A. G. Smith, eds.). Cambridge University Press, London.
- Hillis, D. M., J. P. Huelsenbeck, and C. W. Cunningham. 1994. Application and accuracy of molecular phylogenies. *Science* 264:671–677.
- Hillis, D. M., B. K. Mable, and C. Moritz. 1996. Applications of molecular systematics: The state of the field and a look into the future. Pages 515–543 in *Molecular systematics*, 2nd edition (D. M. Hillis, C. Moritz, and B. K. Mable, eds.). Sinauer Associates, Sunderland, Massachusetts.
- Hudson, R. R. 1992. Gene trees, species trees and the segregation of ancestral alleles. *Genetics* 131:509–512.
- Hudson, R. R., and J. A. Coyne. 2002. Mathematical consequences of the genealogical species concept. *Evolution* 56:1557–1565.
- Hudson, R. R., and M. Turelli. 2003. Stochasticity overrules the “three-times” rule: Genetic drift, genetic draft, and coalescence times for nuclear loci versus mitochondrial DNA. *Evolution* 57:182–190.

- Huelsenbeck, J. P., D. M. Hillis, and R. Jones. 1996. Parametric bootstrapping in molecular phylogenetics: Applications and performance. Pages 19–45 in *Molecular zoology: Advances, strategies, and protocols* (J. D. Ferraris and S. R. Palumbi, eds.). Wiley-Liss, New York.
- Hughes, J., and A. P. Vogler. 2004. The phylogeny of acorn weevils (genus *Curculio*) from mitochondrial and nuclear DNA sequences: The problem of incomplete data. *Mol. Phylogenet. Evol.* 32:601–615.
- Jennings, W. B., and S. V. Edwards. 2005. Speciation history of Australian grass finches (*Poephila*) inferred from thirty gene trees. *Evolution* 59:2033–2047.
- Kingman, J. F. C. 1982. The coalescent. *Stochastic Process Appl.* 13:235–248.
- Knowles, L. L. 2000. Tests of Pleistocene speciation in montane grasshoppers (genus *Melanoplus*) from the sky islands of western North America. *Evolution* 54:1337–1348.
- Knowles, L. L. 2001a. Did the Pleistocene glaciations promote divergence? Tests of explicit refugial models in montane grasshoppers. *Mol. Ecol.* 10:691–701.
- Knowles, L. L. 2001b. Genealogical portraits of speciation in montane grasshoppers (genus *Melanoplus*) from the sky islands of the Rocky Mountains. *Proc. Roy. Soc. Lond. B* 268:319–324.
- Knowles, L. L., and B. C. Carstens. 2007. Inferring a population-divergence model for statistical phylogeographic tests in montane grasshoppers. *Evolution* 61:477–493.
- Knowles, L. L., B. C. Carstens, and M. L. Keat. 2007. Coupling genetic and ecological-niche models to examine how past population distributions contribute to divergence. *Current Biol.* 17:1–7.
- Knowles, L. L., and W. P. Maddison. 2002. Statistical phylogeography. *Mol. Ecol.* 11:2623–2635.
- Knowles, L. L., and D. Otte. 2000. Phylogenetic analysis of montane grasshoppers from western North America (genus *Melanoplus*, Acrididae, Melanoplinae). *Ann. Ent. Soc. Am.* 93:421–431.
- Kubatko, L. S., and J. H. Degnan. 2007. Inconsistency of phylogenetic estimates from concatenated data under coalescence. *Syst. Biol.* 56:17–24.
- Liu, L., and D. K. Pearl. 2006. Species trees from gene trees: Reconstructing Bayesian posterior distributions of a species phylogeny using estimated gene tree distributions. Technical report #53, Ohio State University.
- Lundrigan, B. L., S. A. Jansa, and P. K. Tucker. 2002. Phylogenetic relationships in the genus *Mus*, based on paternally, maternally, and biparentally inherited characters. *Syst. Biol.* 51: 410–431.
- Maddison, W. P. 1997. Gene trees in species trees. *Syst. Biol.* 46:523–536.
- Maddison, W. P., and L. L. Knowles. 2006. Inferring phylogeny despite incomplete lineage sorting. *Syst. Biol.* 55:21–30.
- Maddison, W. P., and D. R. Maddison. 2004. Mesquite: A modular system for evolutionary analysis. Version 1.01. Available at <http://mesquiteproject.org>
- Minin, V., Z. Abdo, P. Joyce, and J. Sullivan. 2003. Performance-based selection of likelihood models for phylogeny estimation. *Syst. Biol.* 52:674–683.
- Morando, M. L. J. Avila, J. Baker, and J. W. Sites. 2004. Phylogeny and phylogeography of the *Liolaemus darwini* complex (Squamata: Liolaemidae): Evidence for introgression and incomplete lineage sorting. *Evolution* 58:842–861.
- Pamilo, P., and M. Nei. 1988. Relationships between gene trees and species trees. *Mol. Biol. Evol.* 5:568–583.
- Parker, A., and I. Kornfield. 2000. Evolution of the mitochondrial DNA control region in the mbuna (*Cichlidae*) species flock of Lake Malawi, East Africa. *J. Mol. Evol.* 45:70–83.
- Rannala, B., and Z. Yang. 2003. Bayes estimation of species divergence times and ancestral population sizes using DNA sequences from multiple loci. *Genetics* 164:1645–1656.
- Rieseberg, L. H. 2001. Chromosomal rearrangements and speciation. *TREE* 16:351–358.
- Robert, C., and G. Casella. 1999. Monte Carlo statistical methods. Springer, New York.
- Rokas, A., D. Kruger, and S. B. Carroll. 2005. Animal evolution and the molecular signature of radiations compressed in time. *Science* 310:1933–1938.
- Rokas, A., B. L. Williams, N. King, and S. B. Carroll. 2003. Genome-scale approaches to resolving incongruence in molecular phylogenies. *Nature* 425:798–804.
- Rosenberg, N. A. 2002. The probability of topological concordance of gene trees and species trees. *Theor. Pop. Biol.* 61:225–247.
- Rosenberg, N. A. 2003. The shapes of neutral gene genealogies in two species: Probabilities of monophyly, paraphyly, and polyphyly in a coalescent model. *Evolution* 57:1465–1477.
- Slowinski, J. B., and R. D. M. Page. 1999. How should species phylogenies be inferred from sequence data? *Syst. Biol.* 48:814–825.
- Stephan, S. J., R. M. Adkins, and J. Anderson. 2004. Phylogeny and divergence-date estimates of rapid radiations in murid rodents based on multiple nuclear genes. *Syst. Biol.* 53:533–553.
- Stephan, S. J., R. M. Adkins, P. Q. Spinks, and C. Hale. 2005. Multi-gene phylogeny of the Old World mice, Murinae, reveals distinct geographic lineages and the declining utility of mitochondrial genes compared to nuclear genes. *Mol. Phyl. Evol.* 37:370–388.
- Stephens, M., and P. Donnelly. 2003. A comparison of bayesian methods for haplotype reconstruction. *Am. J. Human. Genet.* 73:1162–1169.
- Sullivan, J., and P. Joyce. 2005. Model selection in phylogenetics. *Ann. Rev. Ecol. Evol. Syst.* 36:445–466.
- Sullivan, J., and D. L. Swofford. 1997. Are guinea pigs rodents? The importance of adequate models in molecular phylogenetics. *J. Mamm. Evol.* 4:77–86.
- Sullivan, J., and D. L. Swofford. 2001. Should we use model-based methods for phylogenetic inference when we know that the assumptions about among-site rate variation and nucleotide substitution pattern are violated? *Syst. Biol.* 50:723–729.
- Swofford, D. L. 2002. Paup\*: Phylogenetic analysis using parsimony (\*and other methods). Version 4. Sinauer Associates, Sunderland, Massachusetts.
- Swofford, D. L., G. P. Olsen, P. J. Waddell, and D. M. Hillis. 1996. Phylogenetic inference. Pages 407–514 in *Molecular systematics*, 2nd edition (D. M. Hillis, C. Moritz, and B. K. Mable, eds.). Sinauer Associates, Sunderland, Massachusetts.
- Tajima, F. 1983. Evolutionary relationships of DNA sequences in finite populations. *Genetics* 105:437–460.
- Takahata, N. 1989. Gene genealogy in three related populations: consistency probability between gene and population trees. *Genetics* 122:957–966.
- Takahata, N., and M. Nei. 1985. Gene genealogy and variance of inter-population nucleotide differences. *Genetics* 110:325–344.
- Wakeley, J., R. Nielsen, S. N. Liu-Cordero, and K. Ardlie. 2001. The discovery of single-nucleotide polymorphisms and inferences about human demographic history. *Am. J. Hum. Genet.* 69:1332–1347.

First submitted 18 July 2006; reviews returned 29 September 2006; final acceptance 8 January 2007  
Associate Editor: Tim Collins